

# BRIDGING THE GAP

## WITH DATA

A Product Analytics Framework for  
Accountable Business Performance



KANJA SAHA

# **BRIDGING THE GAP**

**WITH DATA**

**A Product Analytics Framework for  
Accountable Business Performance**

**KANJA SAHA**

Copyright ©2026 by Kanja Saha

All rights reserved. No part of this publication may be reproduced, distributed, or transmitted in any form or by any means without prior written permission of the author, except in the case of brief quotations in reviews or articles.

First edition

[www.kanjasaha.net](http://www.kanjasaha.net)

## **Dedication**

*To my parents and my teachers,  
who ignited my lifelong love for numbers and logic.*

*To my mentors,  
who shaped my early career and taught me to build what matters.*

*To my coworkers and business leaders,  
who challenged me to solve real-world business problems with data.*

*To my family and friends,  
who have continually inspired me to push my boundaries.*

## **Author Bio**

I build AI-ready data foundations — Medallion pipelines, dbt-first models, and semantic layers that turn product telemetry into trusted metrics. At AWS, Google, and high-growth SaaS companies, I've automated reporting, accelerated decision cycles, and driven measurable improvements in revenue and retention. I partner with Product, Engineering, Finance, and GTM to define success metrics and deliver scalable self-serve dashboards and analytics platforms. I enjoy blogging and teaching practical techniques for working with data. This book introduces the O2R (Observability-to-Revenue) framework, a way to connect observability signals like error rate and latency directly to customer behavior and revenue impact.

## **Author's Note**

In the age of large language models, a reasonable question is: why read a book at all? LLMs have ingested much of what has been written — frameworks, methodologies, and best practices.

Analytics engineering, however, is a young profession and has not had decades to document practitioner wisdom. Much of what determines the success or failure of an analytics initiative lives in the memory of people who ran it: the negotiations required to align five teams on a single metric, the 2 AM pipeline failure that left morning revenue numbers wrong.

Experience itself cannot be compressed, but this book aims to do just that: distill years of practice into a repeatable, transferable framework grounded in real-world work across diverse teams and roles. It shares guidance for building cross-functional teams, navigating the modern data stack, and key

considerations at each stage of an analytics project turning data into business value. At its core, it is a product analytics framework that links business goals to metrics, metrics to drivers, and drivers to actions. It does this by defining success metrics at the very beginning of the product lifecycle and by treating observability metrics — from error rates and latency to UI behavior, onboarding completion, and escalation patterns — as first-class inputs to the product and revenue story, not just as technical health checks. Throughout the book, I refer to this as the O2R (Observability-to-Revenue) Framework.

## **Intended Audience**

This book is for product and data leaders who collaborate to drive measurable impact and visibility into product and initiative performance across teams. Product leaders set direction and Analytics leaders activate the framework, but analytics is collaborative: responsibility for linking efforts to outcomes is shared across Product, Finance, GTM, Customer Success, Engineering, and Observability/SRE teams. The O2R (Observability-to-Revenue) Framework is written for product leaders, analytics engineers, and operators who need to translate product and operational signals into revenue conversations their executive teams recognize; this book provides a common vocabulary and a shared system so every product team across the organization can turn data into decisions.

The framework is built for scale, designed to be repeatable, automated, and accurate. It serves as a shared resource across business and technical teams, a translation layer between technical data work and business outcomes, providing a common vocabulary and set of expectations around data projects, deliverables, and performance.

## **Acknowledgements**

This book grew from real-world problems, feedback, and collaboration with data, engineering, and business professionals I have worked with. I also wish to acknowledge the numerous authors whose experience and wisdom have shaped my thinking over the years. To accelerate the editing process, I used AI tools, while recognizing that nothing fully replaces professional editors for in-depth technical review. All database schemas and datasets used here were synthetically generated using AI tools to

protect the privacy of professional work.

## **Preface**

We had dashboards and pipelines that were often disconnected and had gaps, built quickly to address ad-hoc requests. At the end of the year, senior leadership asked a simple question: Where is the impact? While data was visible, its impact was not.

The problem wasn't lack of effort, but lack of a common understanding of how data initiatives influenced business decisions. We lacked a unified framework that connected business goals to metrics, metrics to drivers, and drivers to actions. As a result, gaps in data, metrics, and ownership persisted, preventing us from establishing the connections necessary to convert data into decisions. Analytics is a collaborative effort, and without a mutual understanding of how data relates to decisions, teams end up moving in different directions.

This experience is not unique. This same gap between data investment and business value exists across the industry. Four forces shape this challenge:

*The Value Gap:* A 2025 [Gartner](#) survey of over 850 chief data and analytics officers found that only 52% are focused on delivering measurable value, while competing priorities and siloed operating models remain the top barriers. The missing ingredient is not just better data infrastructure — it is a clear way to connect data to decisions.

*Decision Velocity:* [McKinsey](#) research shows that data-driven organizations are 23 times more likely to acquire customers and 19 times more likely to be profitable. Leaders who act confidently on data signals are the ones with a repeatable system behind their decisions.

*The Data-Accountable Leader:* [Salesforce](#)'s State of Data and Analytics report — based on nearly 8,000 leaders globally — finds that 63% of data and analytics leaders say their companies struggle to drive business priorities with data, even as 76% face growing pressure to prove data's value.

*Framework Over Tools:* [Forrester](#)'s 2026 Technology Predictions find that enterprises will delay 25% of planned AI spend into 2027, with fewer than one-third of AI decision-makers able to tie initiatives to P&L changes. A framework-first approach focused on human systems, cross-functional accountability, and repeatable patterns are what turn investment into impact.

A missing piece in many organizations is a clear line between observability signals and business outcomes. Reliability and performance metrics — error rate, latency, uptime, throttle rates, time-to-first-value — are often treated as concerns for SRE and platform teams, while product, GTM, and finance look only at adoption, retention, and revenue. The result is two disconnected conversations about the same system. This book proposes the O2R (Observability-to-Revenue) framework as a bridge: a way to treat observability metrics as Product Performance signals that sit upstream of Customer and Revenue metrics, so reliability work can be prioritized on business impact rather than gut feel.

Any framework has three elements: a clear input, a defined process, and a measurable output. For the product analytics framework proposed in this book, the input is business metrics, specifically success metrics — the definition of what good looks like for a product or initiative. The process is the eight stages of the analytics lifecycle — Define, Collect, Ingest,

Transform, Implement, Visualize, Predict, and Iterate. The output is four types of analytics: descriptive, diagnostic, predictive, and prescriptive — delivered as dashboards that different audiences can act on.

What follows is a practical, repeatable product analytics framework that connects business goals to metrics, metrics to drivers, and drivers to actions — built on a scalable, trustworthy data foundation.

### **Overview of the Next Edition**

This edition provides a comprehensive framework covering the strategy hierarchy, metric layers, analytics lifecycle, and the organizational model that links them all. It introduces the O2R (Observability-to-Revenue) Framework as the connective tissue between observability data, customer behavior, and financial outcomes — turning error rates, latency, UI behavior, onboarding completion, escalation patterns, and time-to-value into signals the business can act on. The upcoming edition will delve into two key areas, each building on the foundation established in this edition:

*Data Engineering: The Extract and Load Layer:* Source connectors, ingestion patterns, data quality contracts, and dbt-expectations.

*AI Agents for Data: The Autonomous Operations Layer:* Agent patterns that monitor pipelines, diagnose failures, trigger remediation workflows, and surface insights, turning analytics engineering practices into partially self-healing systems.

*Data Science and Agents:* With intermediate data science knowledge on Prediction, Classification, Forecasting, and Clustering, these agents act as your co-pilot for prediction, guiding you through feature engineering,

model selection, and deploying forecasts into your analytics workflows.

# Table of Contents

SECTION I: Foundation	1
1. The Hallmarks of a Successful Analytics Project	2
1.1. What Success Actually Looks Like	2
1.2. The Three Lines of Defense, and Drivers	3
1.3. In Practice	5
2. The Power of Data-Driven Decisions	7
2.1. When Data Misleads	7
2.2. In Practice	9
3. The Case for Automation	11
3.1. The Triple Benefit of Automation	12
3.2. Accuracy Starts With Testing	13
4. The lifecycle of Analytics Projects	18
4.1. The eight stages	19
5. The Rise of the Modern Data Stack	24
5.1. The Full-Stack Days	24
5.2. Data Specialization Phase	25
5.3. The Analytics Shift	25
5.4. Emergence of Analytics Engineers	26
5.5. The Modern Data Stack	27
5.6. The Four Data Roles	30
SECTION II: THE ANALYTICS FRAMEWORK	34
6. Analytics Strategy	35
6.1. The Three-Layer Strategy Hierarchy	35
6.2. Product Strategy	36

6.3. Analytics Strategy: Following Product Strategy . . . . .	37
6.4. Data Strategy: Following Analytics Strategy . . . . .	41
6.5. The Importance of the Right Sequence . . . . .	42
6.6. When the Order Was Wrong: A Recognizable Pattern . . . . .	43
6.7. What This Means for Leadership . . . . .	44
6.8. Connecting Strategy to Execution Readiness . . . . .	46
7. The Three Metric Layers – From Observability to Revenue . . . . .	47
7.1. The O2R (Observability-to-Revenue) Framework . . . . .	47
7.2. The Layers Defined . . . . .	49
7.3. Revenue Metrics . . . . .	51
7.4. Customer Metrics . . . . .	56
7.5. Product Performance Metrics . . . . .	59
7.6. Additional Product Performance Metrics . . . . .	64
7.7. What Changes When You Measure All Three . . . . .	67
8. Designing the Analytics Framework . . . . .	69
8.1. The Space Between the Dashboards . . . . .	69
8.2. Closing the Gap: A Shared Framework . . . . .	70
8.3. Product Analytics Framework . . . . .	71
8.4. Success Metrics, Factors, and Drivers . . . . .	72
8.5. Actionable Insights for Success Metric Factors . . . . .	74
8.6. What Changes When Teams Share the Framework . . . . .	78
8.7. Framework at a Glance . . . . .	79
8.8. Adapting the Framework Beyond Enterprise SaaS . . . . .	82
8.9. A Capacity Utilization Platform: From Gut-Feel to . . . . .	84
9. The Deliverables of an Analytics Project . . . . .	87
9.1. What Each Building Block Requires . . . . .	87

9.2. Building Block Readiness: A Pre-Project Reference .....	92
9.3. The Four Types of Analytics .....	93
10. Statistical Inference and Predictive Analysis .....	95
10.1. When Descriptive Analytics Reaches Its Limit .....	95
10.2. Growing a Plant — The Multivariate Problem .....	96
10.3. Navigating Complexity Without Fear .....	97
10.4. Statistical Inference — Is What You Observed Real? .....	98
10.5. Predictive Analysis — What Will Happen Next? .....	100
10.6. The Common Thread .....	103
10.7. Anomaly Detection — Catching What You Were Not Looking For .....	104
10.8. What to Commission First: A Practical Guide .....	105
SECTION III: FRAMEWORK TO ACTION .....	109
11. Building the Data Foundation .....	110
11.1. The Bronze Layer .....	111
11.2. The Silver Layer .....	112
11.3. The Gold Layer .....	113
11.4. The Semantic Layer .....	115
12. Business Metrics in Action .....	117
12.1. The Dataset .....	117
12.2. Financial Metrics .....	122
12.3. Customer Metrics .....	127
12.4. Product Performance Metrics .....	132
12.5. Reading the Metrics Together .....	136
13. The Complete Reporting Framework .....	138
13.1. The Four Types of Reports .....	138

14. Analysis and Insights .....	140
14.1. From Metrics to Questions .....	140
14.2. Exploratory Analysis .....	141
14.3. Hypothesis-Driven Investigation.....	143
14.4. Packaging the Insight.....	146
15. From Insight to Business Impact .....	149
15.1. The Insight-to-Action Gap .....	149
15.2. Operationalizing an Outcome .....	151
15.3. Measuring Downstream Business Impact.....	152
15.4. The Analytics Review Cycle .....	155
SECTION IV: TECHNICAL REFERENCE.....	157
16. SQL and Python Foundations .....	158
16.1. The Core Concept: Tables Are Sets .....	158
16.2. The Three Types of Joins .....	159
16.3. Grouping and Aggregating: Turning Data into Insights.....	160
16.4. Filtering: Getting to the Data That Matters .....	160
16.5. SQL Reference Functions .....	161
16.6. SQL Best Practices.....	162
16.7. Python Foundations for Data Pipelines .....	163
17. Advanced SQL and Python Reference .....	184
17.1. Date and Time Functions.....	184
17.2. Date Dimension Calendar .....	187
17.3. String Functions .....	188
17.4. Window Functions .....	192
18. dbt Command Reference .....	201
18.1. The dbt Project Structure .....	201

18.2. Project Setup Commands .....	203
18.3. Development Commands .....	204
18.4. Testing .....	208
18.5. Data and History .....	211
18.6. Production .....	212
18.7. Inspection and Utilities .....	215
18.8. Jinja and Macros .....	217
18.9. Scheduling with Airflow .....	220
18.10. Troubleshooting .....	222
SECTION V: CASE STUDIES .....	224
19. Building Your Analytics Environment .....	225
19.1. The Local Stack Setup .....	225
19.2. The AWS Stack — Next Edition .....	247
20. The Framework in Practice .....	248
20.1. Case Study 1: Resource Utilization .....	248
20.2. Case Study 2: Customer Journey Milestones .....	253
20.3. Case Study 3: Marketing Attribution and Growth .....	254

# SECTION I: Foundation

This section (Chapters 1–5) presents the core principles and organizational practices necessary for turning data initiatives into measurable business outcomes. It highlights the key factors that consistently drive analytics success: clear problem definition, stakeholder alignment, reliable data systems, and shared ownership across teams.

These chapters show how to convert business questions into actionable outcomes. They explain how to build the processes, systems, and teams that turn insights into action at scale.

The section closes by framing these practices within the modern data ecosystem, clarifying roles, responsibilities, and the trade-offs required to scale effectively.

# Chapter 1. The Hallmarks of a Successful Analytics Project

Analytics is a team sport — inherently collaborative. Its success rests on two foundations: process and people.

The process foundation defines how analytics projects get done reliably and repeatably across five milestones: It starts with a **Frame** — a clear business problem and measurable questions; **Source** — the right data instrumented and captured; **Refine** — data ingested, transformed, and validated into metrics and drivers; **Translate** — metric movements converted into actionable recommendations; and **Drive** — execute decisions, launch initiatives, measure outcomes, and close the loop.

The people foundation comprises three groups with distinct responsibilities: **Product Managers** — define problems, set success metrics, and own outcomes; **Software Engineers** — instrument systems, and ensure reliable data; and **Data Professionals** — ingest, load, and transform data, architect data models, identify drivers, and translate insights into decision-ready recommendations.

These two foundations together form the product analytics framework introduced later, bridging people and process to deliver accountable business performance.

## ***1.1. What Success Actually Looks Like***

Success in analytics projects has two dimensions: technical excellence and business impact. Technical excellence without business impact is expensive engineering effort. Business impact without technical excellence is

unsustainable; it works until it doesn't. Both paths end with a system nobody trusts.

**Technical excellence** goes beyond pipeline failure rates and accuracy. It involves defining the right success metrics and activating processes that turn recommendations into cross-team actions.

**Business impact** is first a measurable shift in metrics, and secondarily the downstream effects: faster self-service, broad trust in data, and a clear, positive return on investment (ROI).

## ***1.2. The Three Lines of Defense, and Drivers***

Analytics projects fail in predictable ways. In product analytics, there are three lines of defense to prevent those failures — each group owning a distinct quality gate, and together covering the pipeline from problem definition to trusted insight. These same three groups also serve as the drivers of impact — each owning a distinct lever that moves metrics and outcomes. The key to achieving business impact lies in alignment and accountability among these three groups.

### **Product Managers — the first line of defense**

Product Managers state the problem, set the goal, and own the outcomes. They define success metrics and success criteria. Their quality gate is bringing alignment on success metrics, working in partnership with GTM, Customer Success, Account Management, and Finance. The success metrics defined must have the potential to drive a positive business outcome. Without a business-grounded metric at the start, projects drift toward technically elegant but irrelevant solutions.

They drive actions and accountability. They take recommendations and convert them into prioritized initiatives and assign accountability among cross-functional teams. In short, Product Managers turn insight into accountable action.

### **Software Engineers — the second line of defense**

Their quality gate is collection integrity. For analytics projects the core risk is signal loss — missing, malformed, or mistimed events. Engineers partner with Product Managers to capture the right signals for success metrics. They make telemetry reliable — preserving signals and increasing downstream trust. Engineering and business teams owning telemetry sign-off in the release checklist ensures issues are caught at source; If data engineers are to find signal loss, there is already data loss.

They drive signal quality and velocity. With business context, engineers don't just implement specs — they challenge ambiguous signals, propose richer events. As owners of operational excellence, they monitor and iterate on instrumentation, automate validation and drift alerts.

### **Data Professionals — the third line of defense**

Data Professionals own data trust across the full pipeline — from raw signal to decision-ready metric. Their quality gate is trust at every handoff: ingestion validation, transformation tests, reconciliation, and business logic checks that confirm the numbers mean what the business thinks they mean. Data and Analytics Engineers build and maintain ingestion and transformation pipelines, implement automated checks, and publish versioned metrics with lineage. Analysts and Data Scientists define edge cases and expected ranges, prototype analyses, and validate metric

behavior against real business scenarios. Without this gate, organizations accumulate data that is technically correct but operationally untrustworthy.

They drive insight, causality, and learning. Analysts and Data Scientists surface causal drivers, prioritize hypotheses, and design experiments that link initiatives to measurable outcomes. Data and Analytics Engineers automate recurring reports and alerts, codify validated business logic into the semantic layer, and ensure findings are reproducible. Together they close the loop — turning raw signals into tested metrics, tested metrics into decisions, and decisions into learnings that make the next cycle faster and more reliable.

Together, the three lines treat every handoff as a verification point — ingestion checks, transformation tests, and reconciliation. When each group both defends its domain and actively drives work, projects move from brittle handoffs to a self-reinforcing cycle of action, measurement, and learning.

### ***1.3. In Practice***

In one role supporting an enterprise SaaS product, the team was asked a simple question: why are customers churning? Revenue was growing but retention was deteriorating quietly underneath it.

The project followed all five milestones. Frame: the problem was defined precisely — three product milestones were identified as success metrics, representing the moments in the customer journey where meaningful usage was established. Source: CloudWatch telemetry and user journey data were combined to map how customers actually moved through the

product toward each milestone.

Refine: the data was analyzed to identify the factors helping customers reach each milestone and the factors blocking them. Some blockers were repeated platform failures — customers hitting the same error in the same workflow, quietly disengaging rather than raising a ticket. Others were subtler — an ambiguous step in the user interface that a portion of customers could not get past, never visible until the journey data made it impossible to ignore.

Translate: each blocking factor was mapped to the team that owned it — Engineering for reliability failures, Product for interface ambiguity, Customer Success for accounts already past the disengagement threshold. The findings did not land as a report. They landed as an ownership map with a clear signal for each team.

Drive: Engineering fixed the recurring failures. Product redesigned the ambiguous step. Customer Success engaged the at-risk accounts. Churn dropped. The milestone completion rate became the shared metric that every team tracked going forward.

The outcome was not just reduced churn. It was a shared framework — product, engineering, GTM, and analytics working from the same milestone view, with the same blocking factors, toward the same outcome. The churn question that started the project had a precise answer. More importantly, the team had a repeatable process for asking the next one.

## Chapter 2. The Power of Data-Driven Decisions

Humans are natural pattern recognizers. From childhood, every decision is implicitly data-driven — we observe, remember, and adjust based on past experiences. Over time those observations solidify into instinct. We make the call first, then look for data to confirm it. Sometimes the instinct is right. Often it is shaped by incomplete information or the most recent anecdote.

What has changed is that we no longer have to rely on instinct only. The scale of data available today — telemetry from millions of events, hundreds of dimensions, thousands of accounts — means the signal exists before the next decision is made. Leaders who build the discipline to read that signal before acting gain something instinct alone cannot provide: confidence grounded in evidence.

That confidence changes how organizations operate. It gives leaders the ability to prioritize initiatives, act decisively, pivot when needed, and measurably improve outcomes and ROI. Decisions grounded in trusted data reduce risk and bias while sharpening accountability — outcomes are tracked against clear metrics and drivers. Over time a data-driven culture compounds: better decisions improve performance, which produces richer data to inform increasingly effective decisions.

### ***2.1. When Data Misleads***

Data is powerful but not infallible. Three failure modes recur in practice.

The first is vanity metrics — metrics that look good but do not reflect the underlying reality — produce data that confirms what teams want to see

rather than what is actually happening. Getting the success metric right is essential. This is why success metrics must be defined before the work begins — not after. A metric defined after the fact is shaped by the outcome, not the objective. When success metrics are set upfront, in partnership with the teams that own the drivers, they are harder to misread and more likely to reflect what the business actually needs to move.

The second is noise and volume — a challenge born from the big data era, and a good problem to have. Rare but high-impact events get drowned out by the volume of ordinary signals. The human eye can spot patterns, but not at that scale — we process one dimension at a time and cannot reliably rank their importance or separate cause from effect. A utilization spike might reflect seasonality, a customer segment, a product change, or some combination. The fix is to break the data down — by milestone, segment, or time window — rather than looking at aggregates. A signal invisible in the total becomes clear in the right slice.

Mistaking correlation for causation is an expensive analytical error a team can make. A metric that rises every Q4 did not respond to the product change shipped in October — it responded to seasonality. A conversion lift after a campaign launch may reflect an improving economy, not your messaging. When decisions stall, go deeper: Advanced Analytics isolates the true drivers of change, and when structured data alone can't answer the question, Data Science earns its place — moving the team from describing what happened to predicting what happens next.

The mitigations across all three require human judgment — to question the assumption, define the right metric, and know which signal to act on.

Ultimately, data is only as powerful as the people who act on it — it surfaces signals, but human judgment and accountability turn those signals into outcomes.

## ***2.2. In Practice***

### **When the data was used: idle capacity**

An organization realized it had idle infrastructure capacity it could not fully account for. The initial instinct was to look at the aggregate — how much capacity was idle and what revenue was at risk. That view confirmed a problem but gave no direction.

The insight came when the question shifted: not just how much, but how long, who owned it, and why. Idle duration became a tracked dimension with a defined threshold — instances beyond an acceptable idle period were flagged automatically. The data stopped being a report and became a decision tool. The product team proposed a success metric. Analytics translated it into a measurable framework. Engineering and finance acted on it. The result was a shared success metric that gave every team the same view of the same problem — and a clear signal for when to act.

This is what data-driven decisions look like when the framework is working. The reports that were not actionable became impossible to ignore once the right metric gets defined.

### **When the data was ignored: the growth rate signal**

In another role supporting a product with worldwide presence, the data showed something easy to dismiss. A revenue growth rate had doubled

year over year — but the absolute revenue was small enough that it never made it onto a leadership agenda. No alert fired. No initiative launched. The signal sat on the dashboard, unread.

What eventually forced the conversation was not the data — it was business professionals in the field reporting unusual customer interest. The opportunity had been visible in the metrics for a long time. The gap was not in the data. It was in how the data was being read.

Data-driven thinking changes what you notice. The discipline is in knowing which signals matter before they become obvious, and the rate of change is one of the signals most commonly overlooked when absolute values are small.

### **The common thread**

Both examples involve signals that were present in the data before anyone acted on them. In the first case, the right question unlocked the signal. In the second, the signal was visible but deprioritized. Again, the difference was human judgment — a data-driven mindset

## Chapter 3. The Case for Automation

In 2014, nearly halfway through my professional life, I moved from a data engineering role into the analytics domain. I had always used engineering as a means to understand the business through data, so the opportunity to focus directly on analysis, especially the predictive analytics element, was exciting.

A month into the role, as I performed analysis on different datasets, I realized how ad-hoc my projects were. I used SQL for data summary, spreadsheets for pivots, dashboards for time series analysis, and Data Science notebooks for correlations and anomaly detection. I was repeating some or all of the above steps for each dataset. In one of my projects, I modified and transformed the dataset more than ten times to get the analysis right. I was spending most of my time preparing and reshaping data instead of interpreting it. I found myself missing the automation I had in engineering.

Most modern data teams have automated data collection, pipelines, and dashboards, but the analytics itself is still largely manual. That work is rarely versioned or refreshed systematically. As a result, the organization can answer "What happened?" once, but it cannot easily track "Is this still true?" or "Did our decision actually change the trajectory?" over time.

Making a case for automation means treating analytics as an engineered product, not a one-time presentation. The transformations, assumptions, and metric definitions used to answer a question live in code, in a single environment, as reusable building blocks that can be rerun daily or weekly against fresh data. This shift — from manual, ad hoc analysis to automated,

repeatable analytics — is what allows organizations to move from decisions made on instinct alone to decisions supported with continuously measured, compared, and improved data.

That is the shift this book is built around — an automated analytics framework that evolves with the business.

### ***3.1. The Triple Benefit of Automation***

Automation delivers three benefits that compound with each other: speed, scale, and accuracy.

Speed is the most visible benefit. Automation eliminates manual effort, reduces redundant work across teams, and removes the waiting that comes with human-dependent processes. Reports refresh when the data refreshes, not when an analyst is available to run the query.

Scale is the benefit that compounds over time. A well-built automation framework is reusable across products, projects, and new additions. When a feature or data source is introduced, it integrates into existing reports quickly with incremental change. Manual processes do not scale this way — every new product or data source requires proportional effort. Automation breaks that linear relationship, allowing analytics to grow with the business rather than catching up with it.

Accuracy is the benefit that is easiest to overlook. A human eye cannot catch every error at scale — and it is unreasonable to expect it to. When code changes, every downstream transformation, metric, and report is potentially affected. No team can manually test every scenario every time a change is made. A well-designed automated test environment runs every

scenario every time — consistently, without fatigue, without the bottleneck of human review.

Manual processes force tradeoffs between them — speed comes at the cost of thoroughness, scale requires proportional effort, and accuracy depends on individual diligence. Automation removes those tradeoff entirely.

### ***3.2. Accuracy Starts With Testing***

Automation without testing is scheduled risk. The same discipline that software engineering applies to code — unit tests, integration tests, regression tests — must be applied to every layer of the analytics stack: pipelines, transformations, metric definitions, and reports. Each layer is a potential point of failure, and a failure at any layer compounds downstream. The goal is — accuracy at speed.

Manual reports get spot-checked by humans who develop intuition about what numbers should look like. Automated reports need systematic validation that does not rely on human oversight.

**Unit tests** validate individual transformations against typical cases, edge cases, and historical scenarios.

**Integration tests** verify end-to-end flows from source systems through transformations to final outputs.

**Shadow runs** compares new reports against existing ones, catching discrepancies within the safety net of the old system.

**Reconciliation checks** automates the comparison between upstream sources and downstream aggregates, flagging mismatches.+

The goal is catching problems before your stakeholders do. A typo in a

document is unprofessional but the meaning survives. A wrong number in a report is a different kind of failure — business decisions get made on it. It is better to have no number than a wrong one. Trust erodes quickly when automated reports silently produce wrong answers.

A filter written as `(status='locked' OR status='deactivated') AND country='USA'` returns only locked or deactivated accounts in the US. Remove the outer parentheses: `status='locked' OR status='deactivated' AND country='USA'` and the result changes entirely. SQL follows operator precedence — AND evaluates before OR — so this reads as all locked accounts globally, plus deactivated accounts in the US only. One missing parenthesis. A completely different answer. A business decision made on the second version would be built on the wrong population. Automated testing catches this before it reaches a report.

### **Guardrails That Fail Fast**

Catching problems before they reach stakeholders is what guardrails are designed to do. Schema validation stops a pipeline when a source system changes unexpectedly — preventing inaccurate data from flowing downstream. Data quality rules surface problems the moment they cross acceptable boundaries. Data contracts make assumptions between data producers and consumers explicit and enforceable — covering format, freshness, and quality expectations. Upstream systems change — and without guardrails in place, those changes flow silently into reports. Guardrails turn silent failures into alerts that get fixed before anyone makes a decision on wrong data.

### **Designing for Production Realities**

Production systems face conditions that manual processes never

encounter. Data arrives late. Sources send duplicates. Bugs get fixed and historical periods need to be reprocessed. Automation designed for these realities handles them gracefully. Pipelines that can safely reprocess historical data mean that fixing a bug is a routine correction. Late-arriving data lands where it belongs. Fixes are validated in isolation before reaching the data that teams depend on.

### **Active Monitoring**

Automation that runs unmonitored accumulates risk invisibly. Tracking pipeline health, data freshness, quality metrics, and usage patterns surfaces degradation early — before it reaches stakeholders.

Runbooks with remediation steps, assigned owners, and time-to-fix SLAs ensure that when something breaks, the path to resolution is clear. Feedback loops that capture user-reported errors and incorporate fixes systematically turn every bug into an improvement to the quality checks. Backfilling and validating historical data after any logic or schema change ensures consistency across time. Active monitoring is what makes automation more reliable over time.

### **Anomaly Detection**

Testing covers what you anticipate. Anomaly detection covers what you do not. Unit tests, schema validation, and reconciliation checks are designed to detect known failure modes — a field going null, a row count dropping to zero, or a value falling outside an expected range. But markets shift, source systems behave unexpectedly, and data can pass every test you have written and still be wrong. Valid, clean data passes all tests — and something is still statistically wrong.

A field that is normally 99% populated suddenly drops to 60% — every row that arrives is valid, no schema violation, no duplicate, but the pattern itself is anomalous. A revenue metric that is normally distributed evenly across five regions suddenly shows 80% concentrated in one region — total revenue is unchanged, every row is clean, but the distribution has shifted dramatically, signaling either a genuine business event or a routing issue in the source system. These are not data quality failures. They are statistical signals that structured tests cannot catch.

### **The Next Layer of Automation**

Active monitoring surfaces failures early — but it still depends on a human to diagnose and resolve them. The next phase of automation closes that gap. AI agents can monitor pipelines continuously, diagnose routine failures, apply fixes, and alert the team — all before anyone notices something is wrong. Human attention is reserved for what genuinely requires judgment. Routine failures resolve themselves.

Two scenarios illustrate what this looks like in practice. When a source system renames a column, an agent detects and remaps the alias to the new column name and the pipeline continues without failure. When a source system sends duplicate records due to a reprocessing bug, an agent detects and triggers deduplication in the staging layer and notifies the source team.

### **Human Accountability**

Automation transfers execution from people to systems — accountability stays with people. Data owners who validate metric changes before they go live are the last line of defense between a pipeline error and a business decision made on wrong numbers. Automated reports carry an authority

that manual reports do not — they look correct because they ran on schedule. Caveats, limitations, and known issues do not surface themselves. Downstream teams depend on the same data — a schema change that takes minutes to deploy can take days to unwind across ten other workflows. The speed of automation makes human judgment more important. Everything the system delivers, with or without AI, is accounted for by a person.

Automation is a tool. The organizations that get the most from it are the ones that pair it with human judgment — clear ownership, active monitoring, and accountability at every layer. Together, they build something more valuable than efficiency — a system that earns trust.

## Chapter 4. The lifecycle of Analytics Projects

A CTO of an early-stage company once reached out to me for an analytics lead role. I was pleasantly surprised — analytics roles at that stage are rare. I told him so, and asked what prompted it. He smiled and said: "I built a product once without it. I learnt this the hard way."

An analytics project lifecycle starts long before the first data request, or even the first data point is collected. It walks in parallel with product inception. Before a team decides what to build, two questions must be answered: what problem will this product solve, and how will we measure success? These questions are not analytics questions — they are product questions. When analytics is part of the conversation from inception, teams build with a shared vision of the definition of success. They know what data they need, how to instrument it, and how to use it to make decisions. When analytics is an afterthought, it struggles to catch up leading to misalignment, expensive rework, and missed opportunities.

For products already in the market, visibility into the product roadmap is the fastest path to bridge the gap. When analytics team knows the upcoming features and initiatives, measurement can be built in before the next release. A seat at the strategy table and a view of the roadmap transforms analytics from a reactive function into a proactive one, with shortened project timelines as an added win.

Every organization uses analytics, and analytics offers two choices: reactive and foundational. Using it reactively — like a walking stick, consulted as needed, when a metric drops, or a decision cannot be reached without it. Or making it foundational — incorporating it into the business,

from product inception to every strategic decision. The foundational path requires deliberate attention and investment upfront, but it is what allows organizations to pick up speed over time, making data and analytics foundation the backbone and central nervous system respectively. The difference between the two approaches reveals itself under stress. That reactive approach though useful short-term, leaves critical blind spots. Without full picture across systems, cross-signal patterns is easy to miss and slower to detect emerging risks or opportunities. By contrast, treating analytics as foundational — investing in data quality, instrumentation, and integrated workflows — gives continuous, contextual visibility across the business. It turns raw volume into coherent signals, enabling earlier detection, better tradeoffs, and faster, more confident decisions.

Analytics is a progression of four questions to guide a business: what happened (descriptive), how and why it happened (diagnostic), what will happen (predictive), and what can we do about it (prescriptive). The eight stages — Define, Collect, Ingest, Transform, Implement, Visualize, Predict, and Act — provide an analytics lifecycle for answering those four questions based on what the business wants.

Knowing where a project stands within that lifecycle at any given time keeps teams aligned and outcomes on track. When priorities shift, the lifecycle provides a clear starting point for realignment.

### ***4.1. The eight stages***

**Define Success Metrics:** This stage begins when a product roadmap is built from strategic planning. The key question is: what does success look like? The goal is to define concrete success metrics, align on goals, identify

the signals that need to be instrumented for data collection, determine the dimensions that enable deeper analysis, and align on the campaigns and initiatives that will drive those metrics. This is also the stage where timeline and resource allocation are decided. Adopting an agile two-week sprint cadence keeps the project modular, makes decoupling between stages easier, and allows teams to adapt as requirements evolve. Investing time here establishes clarity that guides every stage that follows.

**Collect Usage Data:** Based on the success metrics defined in Stage 1, this stage ensures the right signals are instrumented before the product ships. In most organizations, engineering already collects signals for platform reliability and SLA adherence, many of which can be reused for product analytics. Metering, on the other hand, is the tracking of billable consumption — capturing how much a customer uses the product. Signal design is intentional — API names, event names, and error codes can be structured to reflect specific milestones, drop-off points, or customer journey stages rather than generic logging. A signal that captures time-to-value, for example, directly answers whether the product journey unfolded as anticipated. The signals not tracked in operational monitoring are the ones most likely to be missed — yet they are often the most crucial for understanding customer behavior and milestone mapping. Identifying and instrumenting these signals before the product ships is critical at this stage.

**Ingest data:** With the right signals in place, the next step is to reliably move data from source systems into the data platform. Usage data is ingested directly by data engineers, raw and as-is. Modern products generate petabytes of streaming telemetry — the focus here is on building infrastructure that reliably handles volume, ensures data arrives on

schedule, and provides a stable foundation for every stage that follows. Metering data generally passes through a centralized billing data warehouse — raw consumption signals are combined with pricing plans, discounts, and free trial adjustments to produce usage hours, gross and net revenue. This data, along with telemetry, is then consumed by business functions for analysis and reporting.

**Transform data:** Raw data can answer ad-hoc questions, but automated reporting at scale is a different challenge — a single table capturing three months of streaming telemetry can grow to 100B+ records, making transformation essential. Tools like dbt transform raw data across multiple layers and deliver business-ready models. Transformations range from technical — converting timestamps, standardizing string values, deduplicating records, and handling nulls — to structural, such as standardizing customer segments and classifying transaction types. The medallion architecture and star schema provide the structural frameworks that make this possible at scale — both covered in detail later in the book. Well-designed models eliminate redundant work, establish a single source of truth, and are validated with tests and guardrails that ensure every team is working from data they can trust.

**Implement Metrics:** This is where success metrics from Stage 1 become operational — translated into precise, reusable calculations that power every report, dashboard, and analysis across the organization. The semantic layer is where these definitions live centrally — so when business units all pull the same metric, they are aligned. Some tools, like Looker and dbt Cloud, provide built-in semantic layers — making metric definitions a core capability of the analytics stack. In other scenarios, metric definitions

can be built and managed through datasets within the BI tools, scaled across dimensions and granularity to serve every team consistently.

**Visualize Metrics and Deep Dive:** This is the dashboard layer — where metrics and models are delivered in forms that different audiences can quickly understand and act on. Three types of dashboards drive business decisions. Strategic dashboards — commonly known as State of Business reports — give executives a view of business health across financial, product, and customer performance. Tactical dashboards — or Weekly Business Reviews — show performance against weekly or monthly targets: what is on track, what is falling short, and what can be changed. Analytics dashboards go deeper — investigating the factors behind the numbers, how they impact the success metrics, and surfacing root causes when a metric moves unexpectedly.

**Predict and Prescribe:** Once the analytics dashboard is in place, the natural next question is: what can we proactively do to improve success metrics? Prediction does not have to be complex to be valuable. At its simplest, tracking the rate of change of a metric and extrapolating it forward using SQL can provide directional insights into where the business is heading over the next few months — which customers are at risk of churning, which products are likely to grow, where capacity constraints will emerge. This level of prediction is transparent, explainable, and accessible to any analyst without advanced modeling skills. With that knowledge, specific actions can be prescribed that directly impact success metrics. As the data foundation matures and business questions grow more complex, the sophistication of the prediction layer can grow with it — from statistical models like regression and cohort analysis that add precision, to

machine learning models that handle complexity at scale. The progression is deliberate: start simple, stay explainable, and add sophistication only when the business question demands it.

**Act and Iterate:** Act and Iterate: The value proposition for the analytics framework is right here. This is where recommendations become decisions — strengthening what is working well and improving areas with high potential impact — adjusting product features, reallocating resources, modifying campaigns. Predictive model outputs, such as churn scores or onboarding friction scores, are fed back into the data warehouse and surface directly in dashboards — enabling teams to act proactively. A customer with a high churn score triggers an outreach from Customer Success. A customer struggling at a specific point in the onboarding journey gets targeted intervention before they disengage. Since every stage before is fully automated, incremental changes to each stage will quickly get ready for the next iteration.

While eight stages may seem like a long process, upfront planning allows several of them to run in parallel. Once the schema is defined in Stage 2 (Collect Usage Data), real data collection is still being instrumented, Stage 3, 4, 5, 6 (Ingest, Transform, Measure and Visualize) can start simultaneously with synthetic data for every stage. Stages 2 through 6 can largely be prepared in parallel, significantly reducing overall timelines without compromising quality.

## Chapter 5. The Rise of the Modern Data Stack

In the age of AI, data ecosystems have gained widespread attention — reliable data infrastructure is the foundation for real business impact, because even the best models are only as good as the data they run on. At the same time, leaders must measure the return on large-scale data investments: track business outcomes, adoption, and ROI so you fund what drives value and stop what doesn't. What was once a single full-stack responsibility has evolved into a specialized ecosystem of three complementary roles — Data Engineers who build reliable pipelines and platforms, Analytics Engineers who transform raw data into business-ready models, and Analysts and Data Scientists who turn those models into decisions and value. Hiring the right mix at the right time, and shifting focus from infrastructure upkeep to business logic and outcomes, are among the highest-leverage moves a leader can make.

### ***5.1. The Full-Stack Days***

Most of the roles we now consider standard in the data ecosystem didn't exist in the late 2000s when I first moved into a data-focused role.

Our backend engineering team did whatever was necessary to ensure reports were in the hands of stakeholders. We designed the reports, built the application layer with web services, and wrote the SQL queries. We were involuntary database administrators. Backups, schema changes, and performance issues were ours to solve. We modeled the data itself, structuring databases to support reporting while working around the constraints of transactional systems, which were never designed for analytics. There was no separation of concerns.

## ***5.2. Data Specialization Phase***

As organizations grew, the "do everything" model became unsustainable. Those who have worn the database administrator hat know that it is stressful. One wrong command, one missed backup check, one poorly written query, and it can bring down an entire business. And when problems occurred, midday or midnight, they required immediate action to get systems back online.

Databases became more complex, data volumes grew exponentially, and business demands for analytics accelerated. The role started to split. Data Engineering emerged to handle infrastructure, while others focused on insights and analytics products. Discovering insights from data was fascinating — as analytics product engineers, we built products for Enterprise that analyzed internal documents and organizational data to surface connections, helping teams discover each other and collaborate. We even had guardrails in place, much like today's large language models (LLMs), but on a smaller scale in a contained environment.

## ***5.3. The Analytics Shift***

From 2015 onward, the pioneering use of customer data scaled rapidly. Two critical roles emerged to make sense of this data: data analysts delivering business insights to stakeholders, and data scientists building predictive models to forecast and optimize business outcomes.

Data scientists became essential as organizations realized competitive advantage came from predicting customer behavior, optimizing operations, and automating decisions at scale. They combined statistical expertise with programming skills to build predictive models, run

experiments, and extract patterns from massive datasets.

Data analysts were equally critical, translating complex data into actionable insights for business leaders. Advanced analytical techniques were particularly valuable when there wasn't enough data to build predictive models.

Yet both roles faced the same frustration: spending more time preparing data than doing their actual work. Data scientists cleaned data for their models. Data analysts built transformations for their dashboards. Both wrestled with data quality issues and rewrote the same complex joins repeatedly because no one owned the transformation layer.

A simple SELECT statement would time out because the dimensions table has 800 columns. Partitioning the table became mandatory to retrieve data from a table with 50 billion rows. As data volumes grow, transformation complexity increases exponentially. What worked at the gigabyte scale broke at the terabyte scale.

This is when dbt (data build tool) began gaining traction, bringing software engineering principles like modularity, testing, and documentation directly into the SQL transformation layer. As data professionals, we started wearing our engineering hats to build production-ready SQL with macros, data pipelines, orchestration, and maintainable transformation logic as version-controlled code with built-in testing and dependency management.

#### ***5.4. Emergence of Analytics Engineers***

Meanwhile, data volumes kept growing, business demands accelerated, and the gap between raw data and business-ready insights became a

critical bottleneck. The industry was already shifting to address this. At large organizations, SQL Developers and Report Developers evolved from non-technical reporting roles into engineering-focused positions requiring advanced SQL, ETL, and Python. These roles became hybrids between Data Analyst and Data Engineer, building scalable, automated data infrastructure while understanding business context.

This technological shift created space for a new role across the industry: the Analytics Engineer. Analytics Engineers owned the transformation layer — applying software engineering rigor to the SQL transformations that turned raw data into analysis-ready models. They built the foundational layer of clean, tested, business-ready data that both analysts and data scientists could build upon.

It is worth being direct about what made this possible. The tooling problem has largely been solved. dbt, cloud data warehouses, and the modern data stack gave data teams the engineering rigor they needed — version control, testing, documentation, and modular design patterns that scale. This book does not reinvent that foundation. It builds on it. What remains unsolved is the business alignment problem — connecting that technical foundation to business goals, metrics, drivers, and accountable outcomes. That is the layer this framework addresses, and it assumes a solid data foundation is in place, or being built.

The challenge is no longer building data systems — it is ensuring they consistently inform and shape business decisions.

## ***5.5. The Modern Data Stack***

**Early 2000s — On-Premise Infrastructure:** Infrastructure sat on-premise

in company data centers with expensive storage and hardware scaling that took weeks. Databases ran on dedicated servers with manual backups. ETL happened through proprietary tools or brittle scripts. Version control was uncommon — SQL lived in shared drives, collaboration meant emailing files, testing was manual, documentation lived in people’s heads, and deployment required scheduled maintenance windows.

**Late 2000s — Early Cloud:** Cloud computing emerged in 2006, but enterprise adoption was gradual. Early adopters could provision virtual infrastructure in seconds and pay only for usage, but still managed operating systems, databases, patches, and backups themselves. Many organizations used managed hosting with dedicated servers — faster than on-premise but still requiring physical hardware.

**Today — Cloud-Native and Intelligence:** Cloud data warehouses separate storage from compute, handling petabytes with performance in seconds. Transformation tools bring software engineering to SQL — version-controlled, tested, documented. Everything lives in Git with automated testing, deployment in minutes, and AI-assisted development. Analytics Engineers spend 95% of time on business logic rather than infrastructure. Specialized services exist for streaming, ML, governance, and real-time analytics.

This transformation created Analytics Engineering as a discipline bridging engineering rigor with business insight. The modern stack didn’t just make data work easier — it made Analytics Engineering possible.

### *5.5.1. Common Data Stacks*

The modern data stack is not a single prescribed architecture — it varies by

organization size, cloud preference, budget, and existing infrastructure. Dozens of tools exist across every layer, from ingestion to visualization, and the right combination depends on the specific needs of the business. The two stacks described here are ones I have worked with directly — an open source stack built around PostgreSQL and dbt, and an AWS-native stack built around Athena and Glue. Both implement the same foundational principles: reliable ingestion, tested transformations, a semantic layer, and self-serve visualization.

### Open Source Stack

Layer	Tool	Purpose
Ingestion	Airbyte	Moves data from source systems into the warehouse
Orchestration	Airflow + Cosmos	Schedules and monitors pipeline runs
Transformation	dbt	Builds tested, business-ready models across medallion layers
Warehouse	PostgreSQL	Stores and queries structured analytical data
Visualization	Metabase	Self-serve dashboards and reporting
Notebooks	Hex	Collaborative analysis and data science workflows

### AWS Native Stack

Layer	Tool	Purpose
Ingestion	AWS Glue	Serverless ETL for moving and transforming data at scale
Streaming	AWS Kinesis	Real-time data collection and streaming at scale

Layer	Tool	Purpose
Orchestration	MWAA	Managed Airflow — pipeline scheduling without infrastructure overhead
Transformation	dbt + Glue	Business-ready models with Glue as the compute layer
Warehouse	Athena	Serverless query engine over S3-based storage
Visualization	QuickSight	Self-serve dashboards and executive reporting
Notebooks	SageMaker Studio	Machine learning and advanced analytics workloads

## 5.6. The Four Data Roles

**Data Engineering** emerged from the overlap between software engineering and data work — roughly 40–60% of the skills overlap. Data Engineers are essentially software engineers who specialize in data infrastructure. They bring software engineering principles — version control, testing, CI/CD, scalable architecture — and apply them to data problems. Core skills include Python, SQL, cloud platforms, pipeline orchestration tools like Airflow, and distributed systems. Data Engineers typically sit within the Engineering or Platform organization, working closely with Software Engineering and Analytics Engineering.

Organizations need Data Engineers when three patterns emerge: application teams spend more time on data infrastructure than product features, data volumes exceed what scripts can handle resulting in slow systems and unreliable delivery, and multiple data sources require professional management beyond manual solutions.

**Data Analytics Engineering** emerged from the overlap between Data Engineering and Data Analysis — approximately 60–70% of the skills are shared. They bring engineering rigor from Data Engineers (version control, testing, automation) and business acumen from Data Analysts (understanding metrics, stakeholder needs, business logic). They own the transformation layer — turning raw data into business-ready models and semantic models. Core skills include advanced SQL, dbt, data modeling, semantic layer tools, and a strong understanding of business metrics. Analytics Engineers typically sit within the Data or Analytics organization, working closely with both Engineering and business-facing teams.

Organizations need Analytics Engineers when transformation processes break down as data grows, analysts spend more time maintaining SQL pipelines than delivering insights, data engineers lack business context to prioritize transformations, and teams cycle through repeated requests and revisions before getting data right.

**Data Analysis** emerged from the need to translate data into decisions — approximately 50–60% of skills overlap with Analytics Engineering, with the emphasis shifting from building models to interpreting them. Data Analysts sit closest to the business, combining SQL and visualization skills with genuine curiosity about how the business works. They own Stage 6 — Visualization, Reporting and Deep Dive — turning the semantic layer into insights that stakeholders can act on. Core skills include SQL, data visualization tools like Tableau or Looker, storytelling with data, and a strong understanding of business context. Data Analysts typically sit within the business function they support — Product, GTM, Finance, or Customer Success — or within a centralized Analytics organization.

Organizations need Data Analysts when business teams are making decisions without data, reports are being built manually and inconsistently across teams, there is no clear translation between what the data shows and what the business should do, and stakeholders are asking the same questions repeatedly without a reliable, shared answer.

**Data Science** sits at the intersection of statistics, engineering, and business problem-solving — approximately 40–50% of skills overlap with Data Analysis, with the emphasis shifting from describing what happened to predicting what will happen. Data Scientists own Stage 7 — Predict and Prescribe — using historical patterns to forecast outcomes and recommend actions. Strong candidates come from quantitative backgrounds but the differentiator is the ability to frame a business question as a testable hypothesis, build and validate models, and communicate uncertainty honestly to non-technical stakeholders. Core skills include Python, statistical modeling, machine learning, experimentation design, and the ability to communicate probabilistic thinking to business audiences. Data Scientists typically sit within a centralized Data Science or AI organization, working closely with Product and Engineering to move models from experimentation into production.

Organizations need Data Scientists when descriptive and diagnostic analytics are no longer sufficient, business decisions require forecasting and prediction, experimentation and A/B testing need statistical rigor, and the volume and complexity of data exceeds what rule-based analysis can reliably interpret.

Regardless of where data professionals sit in the organization and the depth of their individual skills, collaboration, alignment, and a general

understanding of the data team's complementary capabilities are what build a strong foundation for supporting business functions. A Data Engineer who understands what an Analyst needs, an Analytics Engineer who communicates clearly with a Data Scientist, and a Data Analyst who appreciates the constraints of the pipeline — these relationships are what turn individual expertise into collective impact. The org chart defines the roles. The collaboration defines the outcomes.

# **SECTION II: THE ANALYTICS FRAMEWORK**

This section (Chapters 6-10) defines a repeatable framework that translates business strategy into measurable analytics work. The chapters establish how analytics priorities connect directly to strategic objectives, how to define robust business metrics, and how to apply statistical inference and predictive methods appropriately. The overall emphasis is on creating a coherent framework where every metric traces back to a business goal and every decision traces back to a metric — giving every team a shared language for measuring and improving business performance.

## Chapter 6. Analytics Strategy

Over the last few years, I've often been invited to product strategy sessions, both in startups and startup-style teams. When that happens, the team usually benefits in two main ways.

First, there's better visibility and planning: when analytics leaders understand where the business is trying to go, they can anticipate the instrumentation, infrastructure, and data foundations needed and start preparing before those needs become urgent.

Second, there's a less obvious advantage: cross-product insight and pattern spotting. Because analytics leaders work with data from across all teams and products, they hold a cross-functional, horizontal view that no single product or business leader has. They can see adoption and usage patterns across related and complementary products at the same time, which is hard for individual product owners focused on their own areas. Those cross-cutting signals often highlight broader shifts in user behavior and point to where the market is heading.

### ***6.1. The Three-Layer Strategy Hierarchy***

Every data-driven organization runs three strategies at once—product, analytics, and data—each dependent on the one above it.

**Product strategy** is the high-level plan for how a product will create value and win: what you'll build, for which customers, in which problems or jobs, and how you'll differentiate from alternatives. It turns a broad company mission into a concrete set of choices about markets, segments, and capabilities. It is the top of the hierarchy. Everything else guides and

serves it.

**Analytics strategy** is the shadow of product strategy. It defines which signals confirm that the product is delivering its intended outcomes, which leading indicators give early warning before outcomes are affected, and which questions the organization must be able to answer at any point in time. A product strategy is a vision with a direction. Analytics turns it into evidence.

**Data strategy** is the shadow of analytics strategy. It defines the infrastructure, collection mechanisms, pipelines, governance, and observability required to enable the automated analytics framework.

A data strategy built without reference to the analytics it must support sometimes produces gaps. At scale, crucially, the signals were missed during instrumentation, or the existing data carries quality issues that compound quietly over time. When the data strategy follows the analytics strategy, every collection decision has a clear rationale, and every quality standard has a clear owner. The cascade matters. Product strategy first. Analytics strategy second. Data strategy third. Building in any other order produces a technically sophisticated data platform that cannot answer the questions leadership actually has.

## ***6.2. Product Strategy***

Product strategy is the organization's chosen approach for how the product will create value and win. Different organizations take different paths: some start from customer problems and "jobs to be done," some organize around specific personas or verticals, others focus on building platforms and ecosystems, and many think in terms of positioning and differentiation

(for example, best-of-breed vs. all-in-one) or growth motions like product-led growth and land-and-expand. Most real strategies are a mix of several of these ideas.

Alongside these strategy approaches, there is a mature landscape of metrics and analytics frameworks designed to measure whether a product strategy is working. Teams often use a North Star Metric to capture the core value delivered to customers, OKRs to translate strategy into measurable objectives and key results, and lifecycle frameworks like AARRR/AARM to structure acquisition, activation, retention, referral, and revenue metrics.

Whatever the combination, these choices sit above analytics. The role of analytics strategy is to provide the measurement layer that shows whether the chosen product strategy is working, and why. Beneath that, data strategy is almost universal—the infrastructure, pipelines, and governance that make analytics possible tend to look similar regardless of what is being measured. The next section focuses on an analytics strategy that can sit under any product strategy. Analytics strategy itself requires minimal adaptation to fit the specific product strategy it serves.

### ***6.3. Analytics Strategy: Following Product Strategy***

The analytics strategy and framework proposed in this book rests on three components. The first is the human element — analytics is a team sport. Success metrics are owned across Product, GTM, Finance, and Engineering. Without cross-functional alignment and shared ownership, the stages and the framework cannot deliver. The second is the eight-stage analytics lifecycle — Define, Collect, Ingest, Transform, Implement, Visualize, Predict, and Iterate — which takes success metrics from definition to delivery. The

third is the framework itself — three diagnostic layers of Revenue, Customer, and Product Performance, connected by a causal chain.

The active subset of success metrics in these three layers shifts as the product evolves. At any stage, the signals that matter span all three layers simultaneously — tracking one in isolation produces blind spots.

At the same time, tracking too many metrics at once clouds judgment and scatters focus. In this book, I use a simple constraint: at any point in the analytics lifecycle, teams align on a small set of top-level success metrics (usually no more than five) and treat everything else as supporting detail. This is a divide-and-conquer approach. Each top-level metric has its own supporting metrics behind it — grouped into small, coherent sets — and those in turn may have their own drivers. The work is recursive: start from a focused set of success metrics, and when something moves, drill down into the supporting metrics and drivers that explain why it moved.

Think of growth like a tree: sunlight, water, and soil must all be present — missing one halts growth even if the others are optimal. Revenue, Customer, and Product Performance metrics work the same way — a coordinated set exposes cross-metric trade-offs and root causes that no single metric can surface.

In the early stages, the questions follow a natural sequence across the lifecycle. At Awareness (pre-acquisition), the focus is on who is even finding you — which audiences, channels, and campaigns are creating qualified interest. At Acquisition, the focus shifts to conversion: Weekly Sign-up and Active Accounts begin to move as visitors become identifiable leads and customers. At Adoption (activation), Product Performance

Metrics take center stage — time-to-first-value, onboarding completion, early error rates, and UI friction determine whether new customers ever experience core product value or stall out after sign-up.

As the product matures, the center of gravity moves to Retention and then Churn. Retention is driven by Customer Metrics — usage retention, depth of adoption, Customer Satisfaction, and account health — all of which show stress long before a contract is lost. Churn is where Revenue Metrics finally record the outcome: ARR churn, contraction, and lost CLTV. By the time revenue moves, the causes have usually been visible for weeks or months in Product Performance and Customer Metrics — in rising failure rates, increasing support ticket volume, slowing adoption, or stagnant usage depth. Usage expansion and net revenue retention reveal whether customers are growing with the product or quietly preparing to leave. Pricing comes into focus last: if net revenue and expansion trend together, pricing is working; if they diverge, the model needs revisiting.

Every team contributes to adoption. Every team contributes to retention. Every team contributes to revenue. That is precisely why a shared analytics framework matters — it gives every team visibility into how their contribution is affecting the outcome. The right framing is not "who owns this metric" but "which teams need to see this metric and act on it."

At scale, the platform must support growth without degrading — latency, uptime, and throttle rates become the operational lens. GTM efficiency comes into focus — are the right customers being reached, and is the pipeline converting? The full set of success metrics remains visible throughout. What changes at each stage is which subset drives decisions.

Amazon's product development process starts with a PR/FAQ — a press release and frequently asked questions written before a single line of code is committed. The discipline is to work backward from the customer experience: what does success look like from the customer's perspective at launch, and what would justify the investment?

The same exercise can be applied as an analytics scoping tool. The exercise, done analytically, produces three outputs before any engineering work begins:

**1. The metric list:** Every measurable claim becomes a defined metric with a formula, a grain, and an expected value at launch.

**2. The instrumentation requirements:** Every metric becomes a list of events, properties, and data points that must be collected for the metric to exist.

**The data contracts:** Every instrumentation requirement becomes a specification for engineering: what to log, when to log it, what fields to include, and at what granularity.

If the product strategy defines "customers will activate within their first session," the success metric becomes Time to Value. If the product strategy defines "the platform delivers sub-200ms response times," the success metric becomes latency — tracked at p50, p90, and p99. If it says "customers in the Enterprise tier generate three times the revenue of Growth tier accounts," that is a commitment to track revenue, account tier, and the data model that joins them.

With business-focused metrics defined, the data strategy ensures these signals are reliably captured, modeled, and delivered.

When these three outputs are marked as deliverable in sprint one, the analytics team is not a downstream consumer of what engineering happens to emit. It becomes an upstream contributor to what engineering is required to produce.

#### ***6.4. Data Strategy: Following Analytics Strategy***

Data strategy exists to support two types of decisions. The first is product observability — is the product running correctly? Observability is a built-in engineering requirement. Every software product needs it. Latency, failure rates, error rates, and uptime are the signals that answer this question, and the tools to monitor them are widely available. Many engineering organizations maintain 99.99% uptime SLAs — these signals are already instrumented and operational.

The second is business decisions — is the business healthy and is the customer healthy? These decisions draw from three types of data: revenue, customer, and product performance. Revenue signals show whether the business is growing in the right direction. Customer signals show whether customers are adopting and staying. Product performance signals show whether the product is delivering value.

Some signals overlap. Revenue data is derived from metering — which is part of operational data. But customer journey and usage trends are not operational concerns. A customer dropping off because the UI is confusing is not an observability problem — it is a product design problem. No operational dashboard surfaces that signal.

This is why data strategy must connect back to analytics strategy. When success metrics are defined, the first question is: are the signals needed to

measure them already available in operational data? For the ones that are not, they must be explicitly added to the instrumentation plan before product implementation — so that when the product ships, those signals are collected, transformed, and ready for reporting.

### ***6.5. The Importance of the Right Sequence***

Once, after taking ownership of a product that had launched a year earlier, a product leader asked me to build a Lifetime Value model. The model was straightforward to build — but there were only three weeks of reliable instrumented data. LTV built on three weeks of data is not LTV — it is a guess. A year of customer behavior, the most critical period immediately after launch, was gone. The window to collect it had closed by the time the product shipped, without instrumentation in place.

Retroactively adding event tracking to an existing system is inconvenient but feasible. The real cost is strategic: the period between launch and instrumentation is a blind spot. Decisions made during that window are made without data, and the patterns established without data — the narratives, the assumptions, the instincts about what is working — are the hardest to displace even after measurement begins.

The business that built Time to Value instrumentation from day one builds its onboarding process around data showing exactly where the drop-offs occur. Twelve months after launch, it has a feedback loop that shapes every product decision. The business that skips instrumentation builds on assumptions — and assumptions compound.

Reports can be delayed. Transformations can be built later. Dashboards can wait. But instrumentation cannot. It is a one-way door — once the product

ships without the right signals in place, that window of customer behavior is permanently closed. No retroactive work recovers what was never collected. The chance to understand the full customer journey, forecast lifetime value, or diagnose early churn is lost the moment the product launches without measurement in place.

## ***6.6. When the Order Was Wrong: A Recognizable Pattern***

The pattern described in this chapter is not theoretical. It plays out with enough regularity across enterprise software companies that it has a recognizable shape. It follows the shape of a funnel — the amount of data and analytics instrumentation should grow incrementally after launch, building on the foundation that was in place from day one. When that foundation is missing, the funnel has no base.

A team ships a product. Instrumentation is deferred — there is always a more urgent sprint. The product gains early traction. Leadership forms a view of what is working based on the signals available: sign-up volume, a few customer conversations, and revenue. Those signals are real, but incomplete.

Six months later, retention starts to soften. The analytics team is asked to find out why. They discover that the data needed to answer the question — feature-level usage events, time-to-first-value timestamps, session depth — was never logged—the earliest cohort with reliable behavioral data signed up four months ago. The customers who churned were from the first two months.

Their behavior cannot be reconstructed — though teams that inherit this situation can partially close the gap. Proxy metrics that do exist — login

frequency, support ticket volume, billing events — can reconstruct a partial picture. Segmenting existing customers by signup cohort and comparing behavior patterns can help infer what early churners may have looked like. These are substitutes, not solutions. The data that was never collected cannot be recovered.

The result is not just a delayed answer — it is a permanently unanswerable question. The business will never know what drove early churn, because the window to collect that data closed at launch. The cost is not the engineering hours to retroactively add tracking. It is the absence of the feedback loop that would have shaped the product differently from month two onward. Twelve months of product decisions were made on incomplete information — and that gap compounds. Colin Bryar and Bill Carr captured it precisely in *Working Backward*: instrumentation before innovation. It is not a process improvement. It is the difference between a product that learns and one that guesses.

### ***6.7. What This Means for Leadership***

For business leaders, the implication is a simple question to ask at every product review before a launch commitment is made: What will we be able to measure on day one?

The answer requires three things to be in place: the metric list, the instrumentation requirements, and the data contracts — all defined before sprint one. If the answer is unclear — if metrics are not defined, instrumentation is not planned, or data contracts are not in the sprint — the product is not ready to launch. Not because the features are incomplete, but because the feedback loop that will validate those features

does not yet exist.

Sequencing matters, and it does not necessarily require more budget: the same engineering effort produces far more value when measurement is built in up front. Done before launch, instrumentation creates a day-one feedback loop. Done after launch, it leaves a permanent blind spot during the most consequential weeks when customer patterns form.

Making data instrumentation part of the deployment checklist — with reports already consuming and reflecting test data before launch — is what turns this from a principle into a practice. Two sample checklists are shared in the next section.

#### *6.7.1. Analytics Instrumentation Checklist — Pre-Sprint*

This checklist is used by the analytics and engineering team before a sprint begins.

1. Metric definition: name, owner, formula, grain, aggregation window.
2. Event/field list: field names, types, null rules, and sample values.
3. Data contract: schema, freshness SLA, retention, and error handling.
4. Event frequency: expected event rates and ingestion cadence.
5. Schema versioning: version ID and migration plan.
6. Transformation spec: field mapping, joins, and derived metrics.
7. Tests and validations: unit tests, schema checks, and range assertions.
8. Shadow run and reconciliation: run plan and variance thresholds.
9. Idempotency and deduplication: backfill approach and dedupe keys.
10. Rollback and backfill: steps to revert or reprocess bad data.

11. Monitoring and alerts: freshness, quality, and volume with owner.
12. Security and PII: masking, access controls, and residency notes.
13. Documentation: README, event catalog entry, and sample queries.
14. Sign-off: Team approvals.

### *6.7.2. Deployment Checklist — Release Gate*

The product leader uses this checklist at the release gate.

1. Metric definition: name, owner, formula, grain, aggregation window.
2. Event/field list: field names, types, null rules, and sample values.
3. Data contract: schema, freshness SLA, retention, and error handling.

## **6.8. Connecting Strategy to Execution Readiness**

Strategy defines direction. The cascade from product → analytics → data only works if each stage of the analytics lifecycle is executed reliably. Alongside a well-crafted strategy, the most critical prerequisite for successful execution is the human element.

Every team — Product, Engineering, GTM, Finance, and Data — must align on success metrics, priorities, and deliverables, and commit dedicated resources.

Meetings matter as much as capacity. In a cross-functional initiative, every invitee is part of the chain — and every absence creates a gap that is rarely filled by someone else.

## Chapter 7. The Three Metric Layers – From Observability to Revenue

Metrics in product management are the foundation of every strategic decision — without the right metrics, product leaders are reacting to outcomes they could have anticipated.

The three layers and their relationships are shown in Figure X, with Product Performance (observability), Customer behavior, and Revenue outcomes connected by the O2R (Observability-to-Revenue) Framework.

The Product Analytics Framework introduced in this book is built on three metric layers: Revenue, Customer, and Product Performance. Understanding each layer independently is the necessary first step before seeing how they connect into the broader O2R framework.

### ***7.1. The O2R (Observability-to-Revenue) Framework***

Most product metrics frameworks stop at the customer layer — adoption, retention, satisfaction. O2R goes one layer deeper. It treats observability signals — error rate, latency, throttle rate, uptime, time-to-first-value — as first-class Product Performance metrics that sit upstream of everything else.

I use "observability" broadly to include not only infrastructure metrics, but any product performance signal — from UI behavior to onboarding completion and escalation data — that helps explain why customer and revenue metrics move.

The causal chain runs in one direction:

Product Performance (observability signals from API logs, error logs, UI interaction telemetry, and business process flows) → Customer (adoption, retention, usage depth, satisfaction, escalation rate) → Revenue (ARR, ARPA, CLTV, growth rate, churn)

A spike in error rate today will show up in customer satisfaction next week. Slow time-to-first-value this month will suppress adoption next quarter. A sustained latency degradation will appear in churn next quarter. O2R makes that chain explicit — so engineering reliability, product experience, and revenue growth can be managed as one system instead of three separate dashboards.

O2R also connects observability to the Define stage of the analytics lifecycle. Before any instrumentation is built, the metrics in all three layers — Performance, Customer, and Revenue — are defined together, with data contracts specifying exactly what to log, when, and at what granularity. This ensures that observability data is not a byproduct of engineering decisions, but a deliberate input to the product and business story.

Revenue metrics reflect business outcomes, Customer metrics explain why those outcomes were impacted, and Product Performance metrics reveal why customer metrics moved — each layer telling a distinct part of the same story. Building fluency across all three layers enables the analytics team to shift from explaining what went wrong to anticipating what comes next.

It is not uncommon in a business meeting where revenue came in below target for everyone to scramble to explain why. The data team pulls reports. The sales team points to the pipeline. The product team points to

churn. The team could have seen this coming. The data signals were there. They just were not on the dashboard. And even if the signals were on dashboards, they were not tied to outcomes. Revenue tells you what happened. Not what is happening. Not what could happen.

Revenue is an outcome metric. It reflects decisions customers made weeks or months ago — whether to adopt the product, expand their usage, or renew. By the time those decisions appear on the revenue line, the window to influence them has already closed.

This dynamic is especially true in enterprise SaaS. Enterprise customers do not churn suddenly. They disengage slowly. Usage plateaus. Support tickets increase. The internal champion goes quiet. The renewal conversation gets harder. Each of these is a signal — and each appears in the data long before the revenue impact lands. The question is not whether the signals exist. They always do. The real question is whether you are instrumenting and monitoring those upstream signals so you can act before revenue moves.

## ***7.2. The Layers Defined***

There are three layers of metrics in an enterprise SaaS product: Revenue, Customer, and Product Performance. Understanding how they relate to each other changes how you run the business.

Time is a reporting dimension that runs across all three layers. The examples throughout this framework default to a weekly cadence — Weekly Sign-up, Weekly Revenue, weekly error rates — because weeks offer the right balance of signal frequency and operational response time for most enterprise SaaS teams. But teams should apply the cadence that matches how they operate. GTM may track pipeline weekly, CS may review

account health monthly, and leadership may review NRR and CLTV quarterly. The framework stays consistent — what changes is the window and the granularity each team chooses to surface the signal. Aggregate metrics reveal trends. Dimensional cuts — by region, product, or segment — reveal patterns. Individual account analysis reveals root causes. The analytics deep dive is where aggregate signals become actionable decisions.

Additionally, pricing, a critical component of the product strategy, is not a metric — it is a product decision, a business driver. A pricing model that misaligns with how customers derive value creates friction at the moment of adoption and erodes retention over time, often before any other signal surfaces. Pricing tier is a dimension you apply to existing metrics rather than a separate layer — it adds resolution to signals you are already tracking without requiring a new category.

**Revenue Metrics** are the outcomes. Revenue, growth rate, ARPA, CLTV. These are the numbers that matter most to the board and to investors — and they should. But they are backward-looking by nature. They confirm what already happened. They are the scores at the end of the game.

**Customer Metrics** are the leading indicators. Active accounts, adoption rate, retention rate, usage retention, customer satisfaction. These move before revenue metrics do. When the adoption rate drops this quarter, the retention rate will follow next quarter. When customer satisfaction falls this month, revenue will feel the impact in 3 to 6 months. If you want to know where your revenue is going, look at your customer metrics today.

**Product Performance Metrics** are the diagnostic layer : the earliest warning system available. Error rate, latency, feature abandonment

rate, time-to-first-value, onboarding completion rate, and escalation rate explain why customer metrics move — they are the leading indicators that shift long before revenue does. They are drawn from the instrumentation running beneath your product: API logs, UI interaction logs, error logs, design telemetry, and business process flows. A spike in error rate today will show up in customer satisfaction next week. Slow time-to-first-value this month will suppress adoption next quarter. These signals cut across engineering operations, UI design, and business operations, which is why product, engineering, and customer success must build and interpret them together — the data already exists, but it only becomes a business signal when someone decides to surface and connect it.

Product Performance metrics vary significantly across products. The metrics listed here are generic and apply broadly across enterprise SaaS platforms. Every product has unique operational characteristics that require additional or different metrics. Teams should work closely with Product and Engineering to identify the specific metrics that best reflect their platform's performance and customer experience. The framework is a starting point — not a final answer.

## **7.3. Revenue Metrics**

### **R.1 Year-to-Date Revenue**

**Definition:** The cumulative revenue generated from the first day of the current fiscal year to the current date.

YTD Revenue answers the most fundamental question a business leader asks: how much have we made this year so far? It is the starting point for every revenue conversation because it provides the baseline against which

all targets and growth rates are measured.

**Example:** A SaaS platform closes January with \$1.2M and February with \$1.4M. YTD Revenue is \$2.6M. Against a \$20M annual target, the business has generated 13% of its annual goal in the first two months — giving leadership an early read on whether the year is on track.

## R.2 Weekly Revenue

**Definition:** The total revenue generated within a single calendar week.

Weekly Revenue is the heartbeat metric of the business. It reveals short-term momentum, seasonal patterns, and the immediate impact of product changes, pricing adjustments, or marketing campaigns. Tracking revenue weekly rather than monthly catches problems earlier — a drop in weekly revenue is visible within days, not weeks.

**Example:** A platform generates \$280K in week one, \$310K in week two, and \$265K in week three. The dip in week three prompts an investigation into whether a technical issue, a pricing change, or a drop in sign-ups is to blame. Without weekly tracking, this signal would not surface until the monthly report.

## R.3 Weekly Growth Rate

**Definition:** The percentage change in revenue from one week to the previous week.

Weekly Growth Rate converts raw revenue numbers into a directional signal. It tells you not just how much revenue you made but whether you are accelerating or decelerating. A business generating \$500K per week

with a positive weekly growth rate is in a very different position from one generating the same amount with a negative rate — even though the absolute number looks identical.

**Example:** Revenue grows from \$280K to \$310K — a growth rate of approximately 10.7%. The following week, revenue drops to \$265K — a decline of approximately 14.5%. The weekly growth rate catches this reversal immediately.

#### **R.4 Cumulative Weekly Growth Rate (CWGR)**

**Definition:** The average weekly growth rate sustained over a defined period, typically since the start of the fiscal year.

Where individual weekly growth rates are volatile and noisy, CWGR smooths the signal. It answers the question: at what consistent rate are we growing week over week on average? It is the most honest measure of sustained growth momentum because it removes the noise of individual weeks and reveals the underlying trend.

**Example:** A platform starts the year with \$250K in weekly revenue and reaches \$400K by week 20. The CWGR of approximately 2.4% per week becomes the benchmark against which future weeks are measured.

#### **R.5 Catch-up CWGR**

**Definition:** The weekly growth rate required from the current position to hit the annual revenue target by the end of the fiscal year.

Catch-up CWGR connects current performance to future targets. Given where we are today, how fast do we need to grow each week to hit our

goal? When Catch-up CWGR is significantly higher than the current CWGR, the business needs to act — accelerate growth, adjust the target, or both.

**Example:** Annual target is \$20M. By week 20, the business has generated \$6M with 32 weeks remaining. Catch-up CWGR is approximately 3.8% per week. If the current CWGR is 2.4%, the gap requires a concrete plan, not just optimism.

## R.6 Annual Recurring Revenue (ARR)

**Definition:** The annualized value of recurring revenue generated from active accounts at a given point in time.

ARR is the single most important number in a SaaS business. It tells you the scale and predictability of your revenue engine — not what you earned last month, but what you are on track to earn over the next twelve months if nothing changes. A growing ARR means the business is compounding. A flat or declining ARR means expansion is not outpacing churn. ARR is the north star metric that every other revenue metric orbits around.

*Formula:* Monthly Recurring Revenue  $\times$  12

**Example:** A platform closes February with \$420,000 in monthly recurring revenue across 145 active accounts, giving an ARR of \$5.04M. New business added \$38,000 in MRR during the month, while churn and contraction removed \$12,000 — a net new MRR of \$26,000. Annualized, that growth rate adds \$312,000 to ARR. Tracking ARR movement broken down into new, expansion, contraction, and churn components — known as the ARR waterfall — gives leadership a precise view of where growth is coming from and where it is leaking.

## R.7 Average Revenue Per Account (ARPA)

**Definition:** The average revenue generated per active account within a defined period.

ARPA measures the quality and efficiency of each customer relationship. A rising ARPA means customers are spending more — through tier upgrades, expanded usage, or adoption of additional capabilities. A falling ARPA means the business is acquiring lower-value accounts or existing accounts are contracting. ARPA is one of the earliest indicators of account health — it moves before churn does.

**Example:** A platform generates \$1.4M across 145 accounts, giving an ARPA of \$9,655 per month. New accounts average \$6,200, and existing accounts average \$11,400 — revealing that new accounts need to expand usage to reach the portfolio average.

## R.8 Customer Lifetime Value (CLTV)

**Definition:** The total revenue a business can expect to generate from a single customer account over the entire duration of the relationship.

CLTV is a forward-looking revenue metric. While all other metrics look backward at what has already happened, CLTV projects the value of the current customer base over time. When retention rates improve, CLTV goes up. As ARPA grows, CLTV increases. Every investment in customer success, product quality, and platform reliability ultimately shows up in CLTV.

**Example:** A platform retains accounts for 24 months with an ARPA of \$8K per month, giving a CLTV of \$192K per account. If the cost of acquiring a new account is \$15K, the business is generating nearly 13 times its

investment. If retention drops to 18 months, CLTV falls to \$144K — a 25% reduction from a 6-month reduction in retention.

## **7.4. Customer Metrics**

### **C.1 Active Accounts**

**Definition:** The total number of active paying accounts at any given point in time.

Active Accounts is the foundation metric of the customer base. Every per-account metric uses it as the denominator. Tracked over time, it reveals whether the business is growing, stable, or shrinking at the account level. It is the first place to look when revenue moves unexpectedly — did revenue change because of account growth or because of changes in how much each account spends?

**Example:** A platform ends Q1 with 120 active accounts and Q2 with 145 — a 20.8% increase. Revenue grew only 10% in the same period. The growth in accounts outpaced revenue growth, signaling that new accounts are spending less than the existing base.

### **C.2 Weekly Sign-up**

**Definition:** The number of new accounts that sign up within a single calendar week.

Weekly Sign-up is the primary acquisition metric. It measures the effectiveness of the entire go-to-market motion — marketing campaigns, sales pipeline, and organic referrals — in converting interest into new paying accounts. A spike following a campaign confirms it worked. A

sustained decline signals the top of the funnel needs attention.

**Example:** A platform averages 12 new sign-ups per week through Q1. A developer conference in week 10 drives 28 sign-ups — more than double the average. The analytics team tracks whether conference-sourced accounts show stronger or weaker adoption rates than the baseline.

### C.3 Adoption Rate

**Definition:** The percentage of new accounts that reach a defined level of active usage within a specified time period — typically 13 weeks from sign-up.

Adoption Rate is the most critical leading indicator of long-term retention. An account that reaches meaningful usage within 13 weeks is significantly more likely to renew, expand, and refer others. An account that does not reach that threshold rarely recovers — the window for intervention is narrow. The 13-week mark represents one business quarter, the typical period within which enterprise customers decide whether a product has earned a place in their workflows.

**Example:** Of 48 accounts that signed up in Q1, 34 reached the usage threshold by week 13 — an Adoption Rate of 70.8%. The 14 accounts that did not adopt are flagged for CS intervention. Historical data shows that accounts that receive a health check before week 10 have a 60% recovery rate.

### C.4 Retention Rate

**Definition:** The percentage of accounts that remain active and paying from one defined period to the next.

Retention Rate is the single most important metric for an enterprise SaaS business. A high retention rate means the business is building on a stable foundation. A low retention rate means the business is running to stand still — constantly replacing lost accounts rather than compounding growth.

**Example:** A platform starts Q2 with 145 active accounts and ends with 138 — a retention rate of 95.2%. The analytics team segments the 7 churned accounts by cohort, industry, and usage pattern to shape Q3 product and CS priorities.

### **C.5 Usage Retention Rate**

**Definition:** The percentage change in total usage across existing accounts from one period to the next.

Where Retention Rate measures whether accounts stay, Usage Retention Rate measures whether they grow. An account that stays but reduces usage is a contraction risk. An account that grows usage is an opportunity for expansion. Usage Retention Rate above 100% means growth in existing account usage more than offsets any contraction or churn — the gold standard for enterprise SaaS growth.

**Example:** Total consumption grows from 4.2B tokens to 4.8B tokens across 138 existing accounts — a Usage Retention Rate of 114%. The existing base grew usage by 14% even after accounting for churned accounts.

### **C.6 Customer Satisfaction Score**

**Definition:** A composite score measuring the overall satisfaction of the customer base with the product, platform, and support experience — tracked on a rolling 30-day basis.

Customer Satisfaction Score is the leading indicator of retention. It measures how customers feel before that feeling shows up in churn data. A declining score is an early warning — accounts are becoming less satisfied, and without intervention, they will eventually leave.

**Example:** The score drops from 78 to 71 following two platform incidents. The analytics team identifies 23 accounts below the 65-point churn risk threshold and flags them for immediate CS outreach — preventing potential churn before it materializes.

## ***7.5. Product Performance Metrics***

Within that set, a meaningful distinction exists between observability metrics and operational metrics — and understanding it helps teams instrument the right signals. Observability metrics expose why a system behaves a certain way, surfacing signals inside the engineering stack. Operational metrics measure whether the product is performing as intended, from the perspective of the team running it. In practice, most operational metrics are drawn from the same instrumentation that powers observability — what differs is the question being asked and the audience acting on the answer.

Latency is the clearest entry point. Just as a customer judges a courier service by how fast a package arrives — not by what happened inside the warehouse — users judge a product by how fast it responds. A checkout page that takes four seconds to load does not feel like a technology problem to the user; it feels like a broken product. Latency is measured through API logs, which record the time elapsed between a request entering the system and a response leaving it. When latency spikes, engineering uses those

same logs to trace exactly which internal service call is responsible.

Error rate is where the business impact is most immediate. When one in twenty users cannot complete a core workflow — submitting a form, generating a report, saving a record — the product is functionally unavailable for those users, even if the system shows green on a status page. Error logs capture every failure: what broke, when, and for whom. For engineering, a rising error rate is an observability signal that triggers investigation. For the product and operations team, that same rate tracked week over week is a stability KPI — a measure of whether the platform is holding up under real usage.

Abandonment rate connects product design to business outcomes in a way that surprises most engineering teams. A user who clicks through three screens and then stops is not always encountering a bug — they may be encountering confusion. UI interaction logs record every click, scroll, and navigation path, making it possible to identify exactly where users disengage. When the abandonment rate on an onboarding flow is high, it is rarely a technology failure. It is often a design failure — and no amount of infrastructure investment will fix it. This is the signal that belongs equally in a product design review and a customer success dashboard.

Time-to-value is the metric that connects engineering performance to revenue most directly. In enterprise SaaS, time-to-value measures how long it takes a new customer to reach their first meaningful outcome — their first successful report, their first completed integration, their first team member onboarded. Think of it as the gap between signing a lease and moving into a finished apartment. The longer that gap, the higher the risk of early churn. Time-to-value is calculated from business process event logs

— the sequence of steps a customer completes from contract signature to productive use — and it crosses boundaries that no single engineering system can observe alone.

Adoption depth rounds out the picture. A product can have strong sign-up numbers and still be failing if users activate only one feature and never go deeper. Adoption depth — tracked through UI design instrumentation such as feature usage heatmaps and interaction flows — reveals whether the product’s design is guiding users toward its full value or leaving capabilities undiscovered. For engineering, this data informs where to invest in performance optimization. For product, it reveals where the design is working and where it is not.

### **P.1 Time to Value**

**Definition:** The time elapsed from a new account’s first sign-up to their first meaningful use of the platform — measured at p50, p90, and p100.

Time to Value is the most customer-facing product performance. A short Time to Value is one of the strongest predictors of long-term adoption and retention. Measuring at three percentiles tells three stories: p50 is the typical account, p90 is the account that struggles, and p100 is the worst-case experience.

**Example:** p50 is 3 days, p90 is 11 days, and p100 is 34 days. The gap between p50 and p100 signals that a minority of accounts are struggling, pointing to onboarding or technical-readiness gaps in complex enterprise environments.

### **P.2 Error Rate**

**Definition:** The percentage of platform requests that return an error response within a defined period.

Error Rate is the most direct measure of platform reliability. Every error is a moment the platform failed to deliver. High error rates erode customer trust, inflate support ticket volume, and directly depress Customer Satisfaction Score. Common error types include E429 (rate limiting) and E503 (service unavailable).

**Example:** Error rate spikes from 0.8% to 4.2% following a deployment. The analytics team correlates the spike with a drop in Customer Satisfaction Score across 31 accounts — confirming a customer-visible reliability event.

### **P.3 Latency**

**Definition:** The time taken for the platform to process and return a response — measured at p50, p90, and p99.

For an ML inference platform, latency is not just a technical metric — it is a product quality metric. Customers build workflows on the platform, and their performance depends on how quickly it responds. High latency breaks real-time applications and drives customers to evaluate alternatives.

**Example:** A deployment increases p99 latency from 1,200ms to 3,400ms. While p50 remains stable, three enterprise accounts running latency-sensitive applications immediately raise support tickets — caught before broader impact.

### **P.4 Uptime and Availability**

**Definition:** The percentage of time the platform is fully operational and accessible within a defined period.

Uptime is the foundational reliability commitment of any SaaS platform. For enterprise customers, it is a contractual commitment under service-level agreements. Breaching uptime SLAs has direct revenue consequences and damages the customer relationship in ways that take months to repair.

**Example:** Two incidents in March combined to cause 6.2 hours of downtime, reducing monthly uptime to 99.2% and breaching the SLA for 18 accounts. CS proactively communicates before customers raise formal complaints.

### **P.5 Mean Time to Recovery (MTTR)**

**Definition:** The average time taken to restore full platform functionality following an incident or outage.

MTTR measures how quickly the team resolves platform failures. A low MTTR means incidents are short and customer impact is minimized. MTTR is as important as incident frequency — frequent short incidents may cause less customer damage than rare, prolonged outages. It is also a measure of organizational capability.

**Example:** Four Q2 incidents have recovery times of 45 minutes, 2.5 hours, 20 minutes, and 1.5 hours — an MTTR of 71 minutes. Updated incident runbooks reduce Q3 MTTR to 38 minutes.

### **P.6 Support Ticket Volume**

**Definition:** The total number of support tickets raised by customers within

a defined period.

Support Ticket Volume is a lagging indicator of product and platform health. When ticket volume rises, something has gone wrong — a platform incident, a confusing feature, a documentation gap, or an onboarding failure. Tracked alongside Error Rate and Latency, it confirms whether operational issues are translating into customer pain.

**Example:** Average weekly tickets jump from 42 to 118 in the same week, and the error rate spikes. Ticket categorization confirms 73% are directly related to the deployment incident — providing clear evidence of customer impact.

### **P.7 Support Ticket Resolution Time**

**Definition:** The average time taken to resolve a support ticket from the moment it is raised to the moment the customer confirms resolution.

Fast resolution builds customer trust. Slow resolution compounds the original problem. Resolution Time is a leading indicator of Customer Satisfaction Score — accounts with consistently slow resolution times show lower satisfaction scores and higher churn rates.

**Example:** Average resolution time rises from 18 to 31 hours when engineering focus shifts to an infrastructure project. Customer Satisfaction Score drops 6 points — confirming support responsiveness is a material driver of satisfaction.

## **7.6. Additional Product Performance Metrics**

Product Performance Metrics vary significantly across products. The

examples in this framework are intentionally generic and apply broadly across enterprise SaaS platforms, but every product has unique operational, design, and business process characteristics that require additional or different signals. A workflow-heavy product may lean more on time-to-first-value and onboarding completion; an infrastructure product may lean more on error rate and latency; a self-serve collaboration tool may care most about feature abandonment and depth of adoption. The job of the analytics team is not to adopt this list wholesale, but to work with Product, Engineering, and Customer Success to identify the specific metrics — drawn from API and error logs, UI interaction and design telemetry, and business process flows — that best reflect their platform’s performance and customer experience. The framework is a starting point, not a final answer.

The Product Performance metrics defined above cover the signals common across most enterprise SaaS platforms. The following two metrics appear in the full framework defined in [\[designing-the-analytics-framework\]](#) and are called out here for completeness.

#### **P.4a Throttle Rate**

**Definition:** The percentage of requests that are rate-limited by the platform within a defined period — typically tracked as E429 responses as a share of total requests.

Throttle Rate is a direct indicator of capacity alignment between the platform and its customers. A rising throttle rate means customers are hitting limits that prevent them from using the platform as intended — which erodes perceived value, inflates support ticket volume, and

accelerates churn risk among high-usage accounts. For ML inference platforms, throttle rate is especially consequential: the accounts most likely to be throttled are the accounts growing fastest, making it a leading indicator of both expansion opportunity and retention risk.

**Example:** Throttle rate climbs from 0.3% to 2.1% for accounts in the top usage decile following a pricing tier cap change. The analytics team flags 14 high-usage accounts approaching tier limits and routes them to the commercial team for upgrade conversations before throttling becomes a churn signal.

#### **P.4b Resource Utilization**

**Definition:** The percentage of available platform capacity — compute, memory, or API quota — consumed across the account base within a defined period, typically tracked at p50, p90, and peak.

Resource Utilization is the operational counterpart to revenue forecasting. It answers the question engineering and finance both need answered: are we approaching capacity constraints that will affect service quality or require capital investment? A platform operating at sustained high utilization has less headroom for traffic spikes, increasing the risk of latency degradation and throttling events. Conversely, consistently low utilization signals that infrastructure investment is outpacing actual demand.

**Example:** Peak resource utilization reaches 87% during a product launch week, causing p99 latency to degrade by 40%. The analytics team correlates the event with a spike in support tickets and flags the utilization threshold — 80% — as the trigger for capacity review in future sprint planning.

## ***7.7. What Changes When You Measure All Three***

The metrics were always there, sitting in systems across the organization, waiting to be connected.

What changes when you build a three-layer framework is not the data you have. It is the way you see it. Revenue metrics stop being the whole story and start being the conclusion of a story that began upstream weeks or months earlier. Customer metrics stop being a Customer Success concern and become an early warning system for the entire business. Product performance metrics stop being an engineering dashboard and start being the first chapter of every revenue conversation.

When you track all three layers together, the conversation in a below-target revenue meeting changes completely. Instead of scrambling to explain what happened, you already know — because you watched it unfold in the upstream metrics weeks earlier. When you catch the upstream signals early enough, you can act before the revenue misses the target — sometimes preventing the miss entirely.

Proactive analytics in practice looks like this: upstream signals that predict revenue outcomes weeks in advance. Not a more sophisticated dashboard. Not a bigger data team. Just a clearer understanding of which metrics tell you where you are and which ones tell you where you are going.

The businesses that grow consistently are not the ones with the most sophisticated tools. They are the ones who learned to read the signals before the revenue line moved — and built their decisions around what they saw.

No single team owns these metrics — every team contributes to adoption, retention, and revenue. That’s exactly why a shared analytics framework matters: it gives every team visibility into how their work affects outcomes. When teams align on metrics and commit to reliable instrumentation, analytics stops being a rear-view mirror and becomes a proactive tool: early warning, prioritized action, and measurable impact. The payoff is faster decisions, fewer surprises, and sustained revenue growth driven by informed product choices.

## Chapter 8. Designing the Analytics Framework

The most common discovery in enterprise SaaS is that every team is doing its job — and yet customer outcomes fall short, because no single team has visibility into how their work connects to everyone else's. The Product Analytics Framework makes those connections explicit by linking three layers — Success Metrics that every team shares, Factors and Drivers that identify who owns each lever, and Actionable Metrics that tell each team whether their drivers are working. When that picture is visible to everyone, behavior changes naturally: GTM qualifies leads differently, Engineering prioritizes fixes differently, and Customer Success engages earlier because the data revealed the connection. The framework does not create new data — it creates new visibility that turns well-intentioned teams into an aligned, accountable organization.

There is a moment in almost every enterprise SaaS company where something goes wrong with a customer and every team has a reasonable explanation for why it was not their fault. GTM brought in the account. Product built the features. Engineering kept the platform running. Customer Success was managing the relationship. And yet the customer churned.

This is not a people problem. It is a measurement problem. When each team measures only what they own, they optimize for their piece of the puzzle without seeing how the pieces connect. The gaps between teams are invisible — until a customer falls through one.

### ***8.1. The Space Between the Dashboards***

Every team has a dashboard. GTM tracks pipeline and sign-ups. Engineering tracks uptime and error rates. Customer Success tracks health scores and renewal dates. Product tracks feature adoption and roadmap delivery.

What nobody tracks is how those metrics connect to each other — and to the customer outcome that matters most.

An account that signed up six months ago is showing early warning signs. Adoption Rate is below the 13-week threshold. Support tickets are rising. Usage has plateaued. Customer Satisfaction Score is declining. Each of those signals lives in a different team's dashboard. No single team sees the full picture. And so nobody acts.

This is the gap. Not between the teams themselves — between the metrics they each watch in isolation.

## ***8.2. Closing the Gap: A Shared Framework***

The most effective analytics teams share one thing in common. They do not just measure what they own. They measure how what they own connects to what everyone else owns.

This requires a framework that makes those connections explicit — one that defines not just what to measure but who owns each driver and how each driver connects to the customer outcomes that drive revenue.

The framework has three layers:

**Success Metrics** are the outcomes everyone is working toward — Active Accounts, Adoption Rate, Retention Rate, Usage Retention, Customer

Satisfaction Score. These are shared across all teams. No single team owns them. Every team influences them.

**Factors and Drivers** are the specific levers each team pulls to move the success metrics. Marketing runs campaigns that drive Weekly Sign-up. Engineering maintains platform reliability that drives Adoption Rate. Customer Success runs QBRs that drive Retention Rate. Each driver has a clear owner — and a clear connection to a success metric.

**Actionable Metrics** are the granular measurements that tell each team whether their drivers are working. Is the email campaign conversion rate improving? Is the error rate across weeks one to thirteen within target? Is the QBR completion rate where it needs to be? These are the metrics each team tracks daily — and they connect directly to the success metrics everyone shares.

### ***8.3. Product Analytics Framework***

#### *8.3.1. For an Enterprise SaaS Business Model*

Product Analytics is a team sport. Although the primary owner of all success metrics is the Product team (PRD), the factors driving these success metrics are owned by Go-to-Market (GTM), Customer Success (CS), and Engineering (ENG), in collaboration with Product.

#### 1. Financial Metrics

R.1 Year-to-Date Revenue

R.2 Weekly Revenue

R.3 Weekly Growth Rate

R.4 Cumulative Weekly Growth Rate (CWGR)

R.5 Catch-up CWGR — rate required to hit annual target from current position

R.6 Monthly Revenue

R.7 Average Revenue Per Account (ARPA)

R.8 Customer Lifetime Value (CLTV)

## 2. Customer Metrics

C.1 Active Accounts

C.2 Weekly Sign-up

C.3 Adoption Rate

C.4 Retention Rate

C.5 Usage Retention Rate

C.6 Customer Satisfaction Score

## 3. Product Performance Metrics

P.1 Time to Value

P.2 Error Rate

P.3 Latency

P.4 Throttle Rate

P.5 Uptime and Availability

P.6 Mean Time to Recovery (MTTR)

P.7 Resource Utilization

P.8 Support Ticket Volume

P.9 Support Ticket Resolution Time

## **8.4. Success Metrics, Factors, and Drivers**

Each driver has multiple metrics that can be correlated with the success metrics above. The drivers are owned by cross-functional teams and explain what causes each success metric to move up or down.

#### *8.4.1. C.1. Active Accounts*

**C.1.1. New Account Acquisition** — campaigns, pipeline, organic referrals.

Driver: GTM

**C.1.2. Account Retention** — product value, CS engagement, renewal management. Driver: PRD + CS

**C.1.3. Churn Prevention** — health checks, escalation handling, platform reliability. Driver: CS + ENG

#### *8.4.2. C.2. Weekly Sign-up*

**C.2.1. Marketing** — campaigns, blog posts, events. Driver: GTM

**C.2.2. Sales Pipeline** — lead qualification, outreach. Driver: GTM

**C.2.3. Organic** — word of mouth, customer referrals. Driver: GTM + CS

#### *8.4.3. C.3. Adoption Rate*

**C.3.1. Product Experience** — features, customer journey. Driver: PRD

**C.3.2. Platform Reliability** — error rate and uptime across weeks 1–13.

Driver: ENG

**C.3.3. Support Experience** — tickets, resolution. Driver: CS

#### *8.4.4. C.4. Retention Rate*

**C.4.1. Product Roadmap** — depth of integration, usage. Driver: PRD + ENG

**C.4.2. Product Pricing** — switching cost, multi-model adoption. Driver: PRD

**C.4.3. Support Experience** — tickets, resolution. Driver: CS

*8.4.5. C.5. Usage Retention Rate*

**C.5.1. Use Case Expansion** — new workflows and applications built on the platform. Driver: PRD + CS

**C.5.2. Team Adoption** — number of internal teams within the enterprise using the platform. Driver: CS

**C.5.3. Model Upgrade Path** — migration to higher capability and higher margin models. Driver: PRD + ENG

*8.4.6. C.6. Customer Satisfaction Score*

**C.6.1. Platform Reliability** — error rates, latency consistency, uptime SLA. Driver: ENG

**C.6.2. Customer Success Engagement** — QBRs, health checks, escalation handling. Driver: CS

**C.6.3. Perceived ROI** — measurable business outcomes from the platform. Driver: CS + PRD

**8.5. Actionable Insights for Success Metric Factors**

*8.5.1. C.1. Active Accounts*

**C.1.1. New Account Acquisition** (Driver: GTM)

- New Accounts Added per Week
- Account Activation Rate
- Time from Sign-up to First Payment

**C.1.2. Account Retention** (Driver: PRD + CS)

- Account Health Score
- QBR Completion Rate
- Renewal Rate

### **C.1.3. Churn Prevention** (Driver: CS + ENG)

- At-Risk Account Count
- Churn Rate
- Escalation Resolution Rate

## *8.5.2. C.2. Weekly Sign-up*

### **C.2.1. Marketing** (Driver: GTM)

- Email Campaign Conversion Rate
- Blog Post Traffic to Sign-up Rate
- Event Lead Conversion Rate
- Cost Per Lead

### **C.2.2. Sales Pipeline** (Driver: GTM)

- Lead Qualification Rate
- Outreach Response Rate
- Pipeline to Close Rate
- Average Sales Cycle Length

### **C.2.3. Organic** (Driver: GTM + CS)

- Referral Rate
- NPS Score
- Organic Traffic to Sign-up Rate

## *8.5.3. C.3. Adoption Rate*

### **C.3.1. Product Experience** (Driver: PRD)

- Feature Adoption Rate
- User Journey Completion Rate
- Product Engagement Score
- Time Spent in Product per Week

### **C.3.2. Platform Reliability** (Driver: ENG)

- Error Rate weeks 1–13
- Uptime Percentage weeks 1–13
- Latency p50 and p90 weeks 1–13

### **C.3.3. Support Experience** (Driver: CS)

- Ticket Resolution Time
- First Contact Resolution Rate
- Support Satisfaction Score

## *8.5.4. C.4. Retention Rate*

### **C.4.1. Product Roadmap** (Driver: PRD + ENG)

- Feature Release Velocity
- Feature Adoption Rate post-release
- Integration Depth Score
- API Call Growth per Account

### **C.4.2. Product Pricing** (Driver: PRD)

- Pricing Tier Adoption Rate
- Tier Upgrade Rate
- Price Sensitivity Score
- Multi-model Adoption Rate

### **C.4.3. Support Experience** (Driver: CS)

- Ticket Volume per Account
- Resolution Time
- Escalation Rate

#### *8.5.5. C.5. Usage Retention Rate*

##### **C.5.1. Use Case Expansion** (Driver: PRD + CS)

- New Use Cases per Account per Quarter
- Workflow Expansion Rate
- Application Build Rate on Platform

##### **C.5.2. Team Adoption** (Driver: CS)

- Number of Teams per Account Using Platform
- Internal User Growth Rate per Account
- Cross-team Adoption Rate

##### **C.5.3. Model Upgrade Path** (Driver: PRD + ENG)

- Model Migration Rate
- Higher Margin Model Adoption Rate
- Token Volume Growth per Account

#### *8.5.6. C.6. Customer Satisfaction Score*

##### **C.6.1. Platform Reliability** (Driver: ENG)

- Error Rate
- Latency p50, p90, p99
- Uptime SLA Adherence Rate
- Incident Frequency

##### **C.6.2. Customer Success Engagement** (Driver: CS)

- QBR Completion Rate

- Health Check Frequency
- Proactive Outreach Rate
- Escalation Response Time

### **C.6.3. Perceived ROI (Driver: CS + PRD)**

- Customer Reported ROI Score
- Business Outcome Achievement Rate
- Cost Savings Attributed to Platform
- Revenue Growth Attributed to Platform

## ***8.6. What Changes When Teams Share the Framework***

When GTM sees that their conference-sourced accounts have a lower Adoption Rate than their campaign-sourced accounts, they change how they qualify leads — not because Product told them to, but because the framework makes the connection visible.

When Engineering sees that a p99 latency spike in week eight is correlating with a drop in Adoption Rate, they prioritize the fix differently — not because CS escalated, but because they can see the downstream impact in the shared framework.

When Customer Success sees that accounts with low Feature Adoption Rate in weeks one to four have a significantly higher churn rate at week thirteen, they shift their onboarding engagement earlier — not because Product instructed them to, but because the data told them where the gap was.

This is what closing the gap looks like in practice. Not more meetings between teams. Not more escalations. Just a shared framework that makes

the connections between teams visible — so each team can see how their work connects to the outcomes that matter.

Every team in a well-run enterprise SaaS business is doing its job. The gap does not exist because teams are not working. It exists because each team is working without seeing how their work connects to everyone else's. The adoption signal that CS needed was sitting in the product dashboard. The churn risk that the product needed to prioritize was in the CS health score. The reliability issue that had been suppressing adoption was in the engineering incident log.

The framework does not create new data. It creates new visibility. And when every team can see the same picture — when the connections between drivers and outcomes are visible to everyone — the gap closes not because someone mandated it but because the right people finally had the information they needed to act.

That is what product analytics looks like when it works.

## ***8.7. Framework at a Glance***

The following table summarizes the full three-layer structure of the Product Analytics Framework — Success Metrics, their Factors and Drivers, and the teams that own them. Use it as a single-page reference when introducing the framework to a new team or onboarding a new stakeholder.

<b>Success Metric</b>	<b>Factor</b>	<b>Representative Actionable Metrics</b>	<b>Driver</b>
C.1 Active Accounts	New Account Acquisition	New Accounts per Week, Activation Rate, Time to First Payment	GTM
C.1 Active Accounts	Account Retention	Account Health Score, QBR Completion Rate, Renewal Rate	PRD + CS
C.1 Active Accounts	Churn Prevention	At-Risk Account Count, Churn Rate, Escalation Resolution Rate	CS + ENG
C.2 Weekly Sign-up	Marketing	Email Conversion Rate, Event Lead Conversion Rate, Cost Per Lead	GTM
C.2 Weekly Sign-up	Sales Pipeline	Lead Qualification Rate, Pipeline to Close Rate, Avg Sales Cycle	GTM
C.2 Weekly Sign-up	Organic	Referral Rate, NPS Score, Organic Traffic to Sign-up Rate	GTM + CS
C.3 Adoption Rate	Product Experience	Feature Adoption Rate, User Journey Completion Rate, Engagement Score	PRD
C.3 Adoption Rate	Platform Reliability	Error Rate wks 1–13, Uptime wks 1–13, Latency p50/p90 wks 1–13	ENG

<b>Success Metric</b>	<b>Factor</b>	<b>Representative Actionable Metrics</b>	<b>Driver</b>
C.3 Adoption Rate	Support Experience	Ticket Resolution Time, First Contact Resolution Rate	CS
C.4 Retention Rate	Product Roadmap	Feature Release Velocity, Integration Depth Score, API Call Growth	PRD + ENG
C.4 Retention Rate	Product Pricing	Tier Upgrade Rate, Multi-model Adoption Rate, Price Sensitivity Score	PRD
C.4 Retention Rate	Support Experience	Ticket Volume per Account, Resolution Time, Escalation Rate	CS
C.5 Usage Retention Rate	Use Case Expansion	New Use Cases per Account per Quarter, Workflow Expansion Rate	PRD + CS
C.5 Usage Retention Rate	Team Adoption	Teams per Account, Internal User Growth Rate, Cross-team Adoption Rate	CS
C.5 Usage Retention Rate	Model Upgrade Path	Model Migration Rate, Token Volume Growth per Account	PRD + ENG

Success Metric	Factor	Representative Actionable Metrics	Driver
C.6 Customer Satisfaction Score	Platform Reliability	Error Rate, Latency p50/p90/p99, Uptime SLA Adherence, Incident Frequency	ENG
C.6 Customer Satisfaction Score	CS Engagement	QBR Completion Rate, Health Check Frequency, Escalation Response Time	CS
C.6 Customer Satisfaction Score	Perceived ROI	Customer Reported ROI Score, Business Outcome Achievement Rate	CS + PRD

## ***8.8. Adapting the Framework Beyond Enterprise SaaS***

The framework defined in this chapter is built for an enterprise SaaS business model: named accounts, subscription revenue, a 13-week adoption window, and a cross-functional team structure of GTM, CS, Product, and Engineering. That is the model it was designed to serve — and the one for which the metric definitions, driver ownership, and actionable metrics are calibrated.

Teams in adjacent business models can apply the same three-layer structure — Success Metrics, Factors and Drivers, Actionable Metrics — with three adaptations.

**Product-led growth (PLG).** In a PLG model, adoption is driven by self-serve users rather than CS-managed accounts. The 13-week adoption

window compresses — the critical threshold moves to days or weeks, not a quarter. Weekly Sign-up is replaced by activation rate at the individual user level. CS-owned drivers shift to product-owned drivers: in-app onboarding completion, time-to-first-value at the user level, and feature discovery rate become the primary levers. The Success Metrics remain the same. The driver ownership shifts heavily toward Product and Engineering.

**Usage-based pricing (UBP).** In a UBP model, revenue is not fixed at the account level — it is a direct function of consumption. This changes which metrics matter most. Usage Retention Rate becomes more important than account Retention Rate, because a retained account that reduces consumption generates less revenue even though it does not churn. Throttle Rate and Resource Utilization take on commercial significance — both constrain the consumption growth that drives revenue. ARPA becomes a consumption-based metric rather than a contract-based one, and the ARR waterfall needs to track consumption growth and contraction alongside account-level new and churn.

**Marketplace and multi-sided platforms.** In a marketplace model, there are two customer bases — supply and demand — each with their own adoption, retention, and satisfaction signals. The framework applies to each side independently, with the additional requirement that liquidity metrics (match rate, fill rate, time-to-match) connect supply-side and demand-side performance to the platform's core value proposition. Success Metrics for the platform as a whole typically focus on gross merchandise volume, take rate, and net revenue — but the upstream drivers require a separate factor and driver map for each side of the market.

In all three cases, the logic of the framework is unchanged: define the outcomes everyone is working toward, identify the drivers each team owns, and measure whether those drivers are moving. What changes is which metrics go in each layer and who owns which driver. Teams adapting the framework should start by mapping their own Success Metrics — the three to five outcomes the whole business is aligned on — and work from there.

### ***8.9. A Capacity Utilization Platform: From Gut-Feel to***

#### Data-Driven Efficiency

Consider what actually changes when an organization makes the shift. A capacity utilization platform that once ran on gut-feel operations transformed when trusted data replaced instinct as the basis for decisions. Here is what changed:

**Greater confidence:** When average equipment utilization is visible at 78% with predictable seasonal patterns, leadership can confidently commit to new contracts knowing exactly when capacity opens up. When utilization spikes to 95% with no relief in the pipeline, the signal is clear: invest in additional capacity or risk service quality deterioration. Executives approve major capital expenditures they would never greenlight on instinct alone, simply because utilization trends, bottleneck analysis, and demand forecasts make the case irrefutable. The data does not make the decision — it makes the risk calculable.

**Clear alignment:** Shared definitions eliminate endless debates. When operations, sales, and finance all agree that "available capacity" means equipment that is operational, staffed, and not already committed,

conversations transform. Sales stops promising delivery timelines operations cannot meet. Finance stops questioning why utilization looks different in every department's report. Teams shift from arguing about whose numbers are right to solving the real problem: how to smooth demand peaks and eliminate the bottlenecks.

**Faster learning:** Inefficiencies get identified quickly enough to fix. Scheduling experiments that used to take quarters to evaluate now show results in weeks. A new shift pattern either improves throughput by 12% or it does not — the answer comes within two weeks, not six months. Failed approaches get abandoned before they become entrenched practices. Winning strategies — like dynamic scheduling that reduced idle time by 18% — get rolled out facility-wide while the data is still fresh and the team remembers what changed.

**Better efficiency:** Resources flow to the actual constraints. When data shows that Machine Group A runs at 95% utilization while Group B sits at 55%, the bottleneck is visible. When labor utilization data reveals that afternoon shifts consistently underperform, the investigation points to scheduling, training, or tooling — not because someone complained, but because the numbers pointed there. Operations teams have eliminated millions in unnecessary overtime simply by rebalancing workload based on actual capacity data instead of historical habit.

**Real accountability:** Outcomes replace activity. "We optimized scheduling" is vague. "We increased equipment utilization from 72% to 81% while reducing overtime 15%" is an outcome that flows straight to the P&L. When teams are measured on utilization rates, schedule adherence, and throughput — not just "keeping things running" — behavior changes.

Projects that look good in presentations but do not improve metrics get killed faster. Teams that genuinely move utilization, reduce downtime, or eliminate bottlenecks get recognized for driving real business impact.

**Faster decisions:** Before trusted metrics, operations managers spent days gathering utilization reports from different facilities, manually reconciling conflicting numbers, and debating whether capacity existed for a new client. With trusted metrics, that question gets answered in minutes. Real-time dashboards show current utilization across all assets — machines, labor, facilities — and leadership can immediately see whether the business is at 65% capacity with room to grow or at 92% and approaching constraints. Speed matters when a prospect needs a capacity commitment by end of day.

The payoff is not prettier dashboards or more reports — it is an operations team that makes capacity decisions in hours instead of weeks, commits to growth opportunities with confidence, and optimizes resources based on evidence. That is millions in efficiency gains that would stay hidden without the data.

## Chapter 9. The Deliverables of an Analytics Project

Strategy defines the destination—but most analytics projects fail long before the journey begins. Not because the analysis is wrong, but because the foundational pieces required to connect data work to business outcomes were never put in place.

Before any analytics project begins, six prerequisites must be in place:

A clear business problem worth solving  
Success metrics that define what done looks like  
The right data sources to measure progress  
Stakeholder alignment on priorities  
A cross-functional team with the skills and ownership to deliver  
A shared understanding of how the project connects to the broader strategy

When any one of these is missing, projects drift—solving the wrong problem, measuring the wrong things, or producing outputs that never translate into decisions. These building blocks represent the organizational readiness required for an analytics investment to succeed.

### ***9.1. What Each Building Block Requires***

The six prerequisites are not a checklist to file and forget. Each one has a specific meaning, a clear standard for what good looks like, and a recognizable failure mode when it is absent. Understanding all six — before any work begins — is what separates analytics projects that produce decisions from ones that produce reports.

#### *9.1.1. 1. A Clear Business Problem Worth Solving*

A business problem is not a topic. "Understand churn" is a topic. "Identify which behavioral signals, collected in weeks one through four, predict account churn at week thirteen" is a problem. The difference is precision: a well-formed problem names the metric being affected, defines the time window, and implies a decision the answer will enable.

The test is simple: if the analysis came back tomorrow with a clean answer, would someone change something? If the answer is yes, the problem is worth solving. If the answer is "it depends on what we find," the problem is not yet specific enough.

**What good looks like:** The problem fits in one sentence. It names the affected metric. It has a clear owner who is waiting on the answer to make a decision.

**Failure mode:** The team spends four weeks building an analysis for a question the business had already answered by intuition — or for a question no one is waiting on. The work is technically correct and organizationally irrelevant.

### *9.1.2. 2. Success Metrics That Define What Done Looks Like*

Before a query is written, the team and the stakeholder must agree on what a successful outcome looks like. Not "a dashboard" or "an analysis" — a specific number, threshold, or decision that the work enables. The analysis either answers the question or it does not. That binary is only possible if done was defined in advance.

This prerequisite also protects the team. When done is defined, scope creep has a boundary. When it is not, every stakeholder question becomes a new

requirement.

**What good looks like:** "Done is when we can say, with statistical confidence, whether accounts that complete onboarding step X within seven days have a higher retention rate at week thirteen than accounts that do not." That is a measurable outcome with a clear finish line.

**Failure mode:** The team delivers technically accurate analysis that the stakeholder receives as inconclusive. The stakeholder expected a recommendation; the analyst delivered a table. The gap was not analytical — it was definitional. Done was never agreed.

### *9.1.3. 3. The Right Data Sources to Measure Progress*

Every metric in the success definition requires data. That data must exist, be accessible, and be reliable enough to support the analysis. Discovering a data gap after the project has started is not a data engineering problem — it is a prerequisite failure that should have been caught before sprint one.

Data quality is not assumed — it is confirmed. Row counts, null rates, schema consistency, and join integrity checks are not engineering tasks. They are project prerequisites. An analysis built on unvalidated data produces findings with an asterisk, and the asterisk becomes the headline.

**What good looks like:** Before the project begins, each success metric is mapped to a data source. Access is confirmed. A basic quality check is completed. Any gaps are surfaced and resolved — or the success metrics are adjusted to reflect what the data can actually support.

**Failure mode:** The project delivers findings with a caveat: "Based on available data, which excludes the first six months of accounts due to a

schema migration." The caveat invalidates the finding for the cohort that matters most.

#### *9.1.4. 4. Stakeholder Alignment on Priorities*

Analytics projects compete for analyst time. Without explicit alignment before work begins, the project drifts toward whoever asks the most questions mid-sprint. Alignment means one named stakeholder owns the problem, has committed to acting on the findings, and has a real decision scheduled that the analysis is meant to inform.

The decision deadline matters. An analysis tied to a product review on a specific date has urgency that is connected to an outcome. An analysis without a decision deadline tends to be refined indefinitely — because "done" has no external forcing function.

**What good looks like:** One stakeholder is named. They have committed to act on the findings. A decision — a product review, a QBR, a board meeting — is scheduled. The analysis has a deadline because the decision has a deadline.

**Failure mode:** The analysis lands in an email thread. It generates follow-up questions but no decision. Six weeks later, someone asks for a refresh because "the data might have changed." The work was done; the outcome was not.

#### *9.1.5. 5. A Cross-Functional Team with the Skills and Ownership to Deliver*

Most analytics projects require more than one skill set. A churn prediction model requires data engineering to build the feature pipeline, a data scientist to build the model, and a business analyst to translate the output

into a playbook the CS team can act on. Missing any one of those roles means the work stalls at handoff, gets completed incorrectly, or produces an output no downstream team can use.

Ownership is as important as skill. Knowing who builds the feature pipeline is not enough — that person must have committed time for this project, not just nominal responsibility.

**What good looks like:** Before the project starts, roles are named and time is committed. The data engineer knows which tables they are building. The analyst knows what format the output needs to be in for the stakeholder to act. Every handoff point has an owner on each side.

**Failure mode:** A technically excellent model is built and handed to the CS team without explanation. The model scores every account. It changes no behavior — because no one explained what a score of 0.7 means in practice, or what the CS team is supposed to do with it.

#### *9.1.6. 6. A Shared Understanding of How the Project Connects to the Broader Strategy*

Analytics projects exist in context. A churn analysis is not just a churn analysis — it is either evidence that the product needs to change, that CS needs to engage earlier, or that the acquisition motion is bringing in the wrong accounts. Without knowing which of those it is meant to inform, the team cannot frame the findings appropriately, and the stakeholder cannot act on them correctly.

This connection to strategy is not a formality. It is the frame that determines how findings are presented and to whom. An analysis framed

as a product problem lands differently than the same analysis framed as a CS playbook — even when the data is identical.

**What good looks like:** The project brief includes one sentence connecting the work to a team, a timeline, and a decision: "This analysis will inform the Q3 CS engagement strategy by identifying which behavioral signals, observed in weeks one through four, predict churn at week thirteen." That sentence is the frame for every output the project produces.

**Failure mode:** The findings are accurate but framed incorrectly — presented as a product problem when the audience needed a CS playbook, or scoped to one team when the finding has implications for three. The work is done. It does not land.

## ***9.2. Building Block Readiness: A Pre-Project Reference***

The following table summarizes the standard for each building block and the failure mode when it is absent. Use it as a pre-project readiness check before committing team resources to any analytics initiative.

<b>Building Block</b>	<b>What Good Looks Like</b>	<b>Failure Mode When Missing</b>
Clear business problem	One sentence. Named metric. Owner waiting on the answer.	Vague scope. No decision waiting. Work is technically correct and organizationally irrelevant.
Success metrics defined	Specific threshold or decision the work enables. Binary finish line.	"Inconclusive" findings. Stakeholder expected a recommendation; received a table.

Building Block	What Good Looks Like	Failure Mode When Missing
Right data sources	Each metric mapped to a source. Access confirmed. Quality checked before work begins.	Analysis delivered with a caveat that invalidates the finding for the cohort that matters.
Stakeholder alignment	One named owner. Committed to act. Decision deadline is real and scheduled.	Work lands in an email thread. No decision. Refresh requested six weeks later.
Cross-functional team	Roles named. Time committed. Every handoff has an owner on both sides.	Technically excellent output that changes no behavior because no one knows how to use it.
Strategic connection	One sentence connecting work to team, timeline, and decision.	Accurate findings framed for the wrong audience. Work done. Does not land.

When all six building blocks are in place before the project starts, the analytics team is positioned to deliver an output that connects directly to a decision. When any one is missing, the gap will surface — in scope creep, in inconclusive findings, in a deliverable that no one acts on. The time to resolve that gap is before the first query runs, not after the final presentation.

### ***9.3. The Four Types of Analytics***

Analytics is not a single activity — it is a progression of four distinct types, each building on the previous.

**Descriptive Analytics** answers what happened by summarizing historical data into patterns and trends — your weekly revenue report, your customer growth rate, your adoption metrics.

**Diagnostic Analytics** goes deeper and answers why it happened — drilling

into the drivers behind the numbers to explain what caused a metric to move up or down.

**Predictive Analytics** uses historical data, statistical models, and machine learning to answer what will happen — forecasting future revenue, predicting churn, or estimating customer lifetime value.

**Prescriptive Analytics** answers what should we do — translating predictions into specific actions, pricing decisions, or investment priorities.

A successful analytics project does not stop at description. It moves through all four stages, from understanding what happened to recommending what to do next. This book is structured around that progression — Section II builds your descriptive and diagnostic capabilities, Section III takes you into predictive and prescriptive territory, and the case studies in Section IV show all four types working together in real business contexts.

## Chapter 10. Statistical Inference and Predictive Analysis

Descriptive and diagnostic analytics tell you what happened and why — but a crucial component of any analytics framework is the ability to look forward, and that requires a different set of tools: statistical inference to test whether what you observed is real or random, and predictive analysis to forecast outcomes and predict behavior before it happens. Anomaly detection complements both — automatically flagging unexpected patterns in data before they surface as incidents, missed targets, or customer complaints. Together these three capabilities complete the analytical toolkit, bridging the gap between understanding the past, catching the unexpected, and shaping the future.

### ***10.1. When Descriptive Analytics Reaches Its Limit***

For small datasets and early-stage products, descriptive analytics — charts, tables, weekly metric reviews — is often enough. You can look at the numbers, see the trend, and draw a reasonable conclusion. When your product has a few hundred accounts and a handful of metrics, a sharp analyst and a well-designed dashboard can carry the organization a long way.

That changes as scale increases. At tens of thousands of accounts, millions of events, and hundreds of dimensions — usage patterns, pricing tiers, regions, cohort dates, feature adoption sequences, support history — no one can look at the data and understand it anymore. Not because they lack skill or experience, but because the human eye is not built to detect patterns across hundreds of variables simultaneously. The relationships

that matter most are often invisible to visual inspection precisely because they involve the interaction of multiple variables at once.

This is not a failure of the analyst. It is a fundamental limitation of descriptive methods when applied beyond their intended scale. The moment a business reaches meaningful volume, traditional approaches start producing a familiar and dangerous outcome: confident conclusions drawn from incomplete views of the data.

## ***10.2. Growing a Plant — The Multivariate Problem***

Consider what it takes to grow a plant. Five elements are required: sunlight, water, soil quality, temperature, and nutrients. Remove any one of them and the plant struggles or dies. Provide all five in the right proportions and it thrives. The growth of the plant is not the result of any single input — it is the product of all five working together, each one influencing the others.

Business growth works the same way. Revenue does not grow because of one thing. Customer retention does not improve because of one change. Product adoption does not accelerate because a single feature was released. Growth is the result of multiple variables acting simultaneously — pricing, onboarding experience, feature adoption sequence, support quality, platform reliability, sales motion, account size — each one interacting with the others in ways that compound, cancel out, or amplify depending on the customer segment and the moment in the lifecycle.

The problem this creates is both mathematical and organizational. Mathematically, isolating the contribution of any single variable when multiple variables are changing simultaneously requires calculus —

specifically, the kind of multivariate analysis that examines partial derivatives: how much does outcome Y change when variable X changes, holding all other variables constant? Organizationally, the instinct to identify "the one thing that drove growth" is almost always wrong, and acting on that instinct leads to optimizing a single lever while leaving the others misaligned.

The goal of advanced analytics is not to find the one variable that explains everything. It is to map the contribution of each variable, understand their interactions, and identify the combination of conditions under which growth becomes reliable and repeatable. That requires statistical methods that go beyond what any dashboard can show.

### ***10.3. Navigating Complexity Without Fear***

The mathematics behind these methods — regression models, Bayesian inference, gradient boosting, multivariate calculus — can sound intimidating to leaders who are not data scientists. It should not. The role of a business leader is not to build or validate these models. It is to understand what questions each method is designed to answer, when to commission the analysis, and how to interpret the output in terms that connect to decisions.

The methods described in this chapter are tools, not obstacles. Each one was designed to answer a specific class of question that descriptive analytics cannot — whether an observed difference is real or random, which customers are most likely to churn, how much revenue the business will generate next quarter, which combination of product behaviors predicts long-term retention. Understanding that these tools exist, and

knowing which business question maps to which tool, is the analytical literacy that separates leaders who can drive a data-informed organization from those who are merely consumers of what the data team produces.

A deeper treatment of these methods — including worked examples using the case study data introduced in this book — will be covered in the next edition. The goal here is to give you a working vocabulary: enough to ask the right questions, commission the right analyses, and interpret the results when they land on your desk.

### ***10.4. Statistical Inference — Is What You Observed Real?***

Statistical inference answers a question that descriptive analytics cannot: is the pattern I am seeing in the data real, or is it a product of chance? At small sample sizes and short time windows, many things that look like signals are noise. Without a formal method for distinguishing between the two, organizations make expensive decisions based on patterns that will not hold.

#### **A/B Testing**

**Definition:** A controlled experiment in which two versions of a product, feature, or experience — version A and version B — are exposed to separate, randomly assigned groups of users simultaneously, and their outcomes are compared.

A/B testing removes the confounding effect of time. When a team ships a change and then compares metrics before and after, they are not measuring the impact of the change — they are measuring the combined effect of the change, the passage of time, seasonality, external events, and

everything else that shifted between the two periods. A/B testing holds time constant by running both versions in parallel, so the only systematic difference between the two groups is the thing being tested.

In practice, A/B testing has been applied to product UI decisions — comparing two designs for an onboarding flow to determine which one drives higher activation rates — and to pricing experiments, where different customer segments are offered different price points to understand the revenue and retention implications before a change is rolled out to the full base.

The output of an A/B test is not just which version "won." It is the magnitude of the difference, the confidence level at which that difference can be attributed to the change rather than chance, and the estimated effect size — how much impact the change is likely to have if deployed at scale.

### **Confidence Intervals**

**Definition:** A range of values within which the true value of a metric is estimated to fall, given the observed data, at a specified probability — typically 95%.

A single metric number — "our weekly retention rate is 94.3%" — tells you what was observed. A confidence interval tells you how much to trust it. A 95% confidence interval of [93.1%, 95.5%] means that if the same measurement were repeated many times under the same conditions, 95% of those measurements would fall within that range. Narrow intervals indicate high precision. Wide intervals indicate that the underlying data is noisy, the sample size is too small, or the metric is inherently variable.

Confidence intervals have been used to draw inferences from datasets where the full population cannot be observed directly — for example, estimating the true churn rate across a customer segment based on a sample, or understanding whether a measured improvement in customer satisfaction is large enough to be meaningful or within the margin of natural variation. They turn a single observation into a defensible claim about the underlying reality.

### ***10.5. Predictive Analysis — What Will Happen Next?***

Predictive analysis uses historical patterns to forecast future outcomes and score future probabilities. Where statistical inference looks backward to validate what already happened, predictive analysis looks forward to inform decisions that have not yet been made.

#### **Customer Segmentation — Clustering**

**Definition:** An unsupervised machine learning technique that groups customers into segments based on the similarity of their behavior, without pre-defining what those groups should be.

Traditional segmentation is rule-based: Enterprise accounts, Growth accounts, Starter accounts — defined by contract value or employee count. Clustering lets the data define the segments. When applied to behavioral data — usage frequency, feature adoption patterns, support ticket history, time to activation, API call volume — clustering surfaces natural groupings that rule-based segmentation misses.

In practice, clustering has been used for customer segmentation to identify behavioral cohorts that cut across pricing tiers. The insight it consistently

produces is that accounts within the same pricing tier often behave more differently from each other than accounts across tiers — and that behavioral similarity predicts retention and expansion far better than tier alone. That finding reshapes how customer success resources are allocated and which accounts get proactive outreach.

## **Revenue Forecasting**

**Definition:** A time-series or regression-based model that projects future revenue based on historical patterns, seasonality, growth trends, and leading indicators.

Revenue forecasting translates the current state of the business — active accounts, expansion pipeline, known churn, seasonal patterns — into a probability distribution of future outcomes. It answers the question every leadership team asks but rarely has a rigorous answer to: given where we are today, what revenue range should we expect next quarter, and what assumptions is that forecast most sensitive to?

In practice, revenue forecasting has been used to project quarterly and annual revenue using a combination of historical weekly revenue trends, cohort-level retention curves, and pipeline data. The value is not the point estimate — the single number — but the range and the sensitivity analysis: which inputs, if they change, move the forecast the most. That sensitivity analysis is where forecasting connects back to strategy. If the forecast is highly sensitive to retention rate changes, retention becomes the strategic priority. If it is more sensitive to new account acquisition, the growth investment should be weighted toward top-of-funnel.

## **Churn Score Prediction**

**Definition:** A classification model that assigns each active account a probability of churning within a defined time window — typically 30, 60, or 90 days — based on behavioral signals observed in the data.

Churn scoring converts the lagging signal of actual churn into a leading signal of churn risk. By the time an account churns, the window for intervention has already closed. A churn score, updated weekly or daily, gives customer success teams the ability to intervene while there is still time — prioritizing outreach toward accounts whose behavioral patterns most closely resemble accounts that churned in the past.

In practice, churn score prediction has been built using a combination of product usage signals — declining API call volume, increasing error rates, support ticket frequency, time since last active session — weighted by their historical correlation with churn outcomes. The model does not tell the CS team why an account is at risk. It tells them which accounts to call first. The conversation that follows — which is a human one — determines the underlying reason and the right intervention.

### **Buyer Propensity Scoring**

**Definition:** A model that scores the likelihood of a prospect or existing account taking a desired commercial action — upgrading a tier, expanding usage, or converting from trial to paid — within a defined time window.

Where churn scoring identifies risk, propensity scoring identifies opportunity. The same behavioral signals that predict churn in one direction often predict expansion in the other — accounts that are increasing usage, adopting new features, and generating low support ticket volumes are both less likely to churn and more likely to expand. Propensity

scoring formalizes that intuition into a ranked list that sales and customer success teams can act on systematically.

In practice, buyer propensity scoring has been applied to identify accounts most likely to upgrade from Growth to Enterprise tier, using a combination of usage growth rate, feature breadth, account size, and industry vertical as input signals. The output is a weekly ranked list that allows the commercial team to prioritize conversations with the accounts that the data suggests are ready — rather than relying on relationship intuition alone.

### ***10.6. The Common Thread***

Across all of these methods — A/B testing, confidence intervals, clustering, forecasting, churn scoring, propensity scoring — the common thread is the same: they are all responses to the same fundamental limitation of descriptive analytics at scale. They exist because the data is too large, too multidimensional, and too noisy for any human to interpret directly. They are the tools that make it possible to navigate complexity without being paralyzed by it.

The plant still needs sunlight, water, soil, temperature, and nutrients. Advanced analytics does not reduce the complexity of growth to a single variable. What it does is tell you, with measurable confidence, how much each variable is contributing, which combination of conditions produces the most reliable outcomes, and — most practically — where to focus first when resources are limited and the business cannot optimize everything at once.

That is the promise of statistical inference and predictive analysis: not certainty, but better-informed decisions in the presence of complexity. And

in a business with millions of users, hundreds of variables, and outcomes that compound over years, better-informed decisions are the most valuable asset the analytics function can produce.

### ***10.7. Anomaly Detection — Catching What You Were Not Looking For***

Statistical inference confirms whether a known pattern is real. Predictive analysis forecasts known outcomes. Anomaly detection does something different: it surfaces unexpected patterns automatically, before anyone thought to look for them.

**Definition:** An automated method for identifying data points, events, or patterns that deviate significantly from expected behavior — flagging them for investigation before they surface as incidents, missed targets, or customer complaints.

The value of anomaly detection is in what it catches between the metrics you are actively watching. A weekly dashboard reviews thirty metrics. An anomaly detection system monitors thousands — every combination of account, feature, time window, and usage dimension — and alerts when something moves outside its expected range. The analyst does not need to know which signal to watch. The system tells them which one changed.

In practice, anomaly detection has been applied in two configurations. The first is operational: monitoring error rates, latency percentiles, and throttle rates continuously, and triggering alerts when any of them cross a threshold defined by historical variance rather than a manually set number. This catches platform degradation events earlier than dashboard

reviews — often within minutes of onset rather than at the next reporting cycle. The second is behavioral: monitoring per-account usage patterns and flagging accounts whose behavior diverges from their historical baseline. A sudden drop in API call volume, an unusual spike in error rate for a single account, or a feature usage pattern that stops — each of these can appear in the anomaly detection system days before the account raises a support ticket or shows up on a churn risk list.

**Illustrative example:** A platform monitors weekly API call volume per account. Account 0147 has averaged 2.3M calls per week for eight consecutive weeks. In week nine, volume drops to 800K. No support ticket has been raised. No CS contact has occurred. The anomaly detection system flags the deviation — a 65% drop against the account’s expected range — and routes an alert to the assigned CS manager. The CS manager contacts the account and discovers the primary technical contact left the company two weeks ago. The account’s integration has partially broken. With three weeks remaining before the renewal date, the intervention is timely. Without anomaly detection, the signal would have appeared in a churn score update five days later — and on the renewal risk list two weeks after that.

Anomaly detection does not replace the metrics in this chapter. It amplifies them. The three-layer framework defines what to measure. Anomaly detection ensures that unexpected movement within those measurements gets caught — not only when someone checks the dashboard, but the moment it happens.

## ***10.8. What to Commission First: A Practical Guide***

The methods in this chapter are not equally urgent at every stage of business growth. The right question is not which method is most sophisticated — it is which method answers the question the business most needs to answer right now.

The following guide maps business stage and pressing question to the appropriate first commission. It is not a sequence to follow in order. It is a decision framework for prioritizing analytical investment when resources are limited.

<b>Business Stage</b>	<b>The Question That Matters Most</b>	<b>Commission First</b>
Early growth ( $< 200$ accounts)	Are our retention numbers real, or too small to trust?	Confidence intervals on retention and adoption rates. Establish whether observed metrics are statistically meaningful before building strategy on them.
Scaling (200–1,000 accounts)	Which accounts are most at risk before the renewal cycle?	Churn score prediction. Build the model on the behavioral signals available. Even a simple model with three to five input features outperforms gut-feel prioritization at this account volume.
Scaling (200–1,000 accounts)	Did the product change we shipped actually improve activation?	A/B testing on the specific change. Run it for long enough to reach statistical significance before drawing conclusions.

<b>Business Stage</b>	<b>The Question That Matters Most</b>	<b>Commission First</b>
Growth (1,000+ accounts)	Are our customers behaviorally different from each other in ways our tier structure does not capture?	Customer segmentation via clustering. Apply it to behavioral data — usage frequency, feature adoption sequence, support history. The output reshapes CS resource allocation.
Growth (1,000+ accounts)	Where will we be at end of quarter, and what is the forecast most sensitive to?	Revenue forecasting with sensitivity analysis. The point estimate matters less than the input sensitivity — it tells you which lever to pull.
At scale (any stage)	Are we catching platform and account-level anomalies before customers report them?	Anomaly detection on operational metrics and per-account usage patterns. Commission this in parallel with other methods — it is the earliest warning system available and has a continuous payoff.
At scale (any stage)	Which accounts are ready to expand before the commercial team calls them?	Buyer propensity scoring. Combine with churn scores to build a full account health view: risk on one axis, opportunity on the other.

The sequence that emerges from this guide for most enterprise SaaS businesses at the scaling stage is: confidence intervals first, to validate that the metrics being tracked are trustworthy; churn scoring second, to protect the account base; revenue forecasting third, to align the business on where it is going; clustering fourth, to understand the account base at a behavioral

level; and anomaly detection and propensity scoring as continuous capabilities that run alongside everything else.

The methods are the tools. The guide is the sequencing. What drives the decision is which question the business cannot currently answer — and which gap, left unanswered, costs the most.

# SECTION III: FRAMEWORK TO ACTION

This section (Chapters 11–15) turns the framework into practical playbooks for building and delivering analytics in production. A scalable, trustworthy data foundation is not optional—it is what makes the framework work. At the same time, a strong foundation alone is not sufficient without a clear connection to business outcomes.

These chapters address both: constructing the data foundation, implementing metric pipelines, and building dashboards for reliable reporting, while keeping the focus on outcomes rather than outputs. The analysis chapters teach exploratory methods, hypothesis testing, and how to package insights so they lead to concrete actions. The section closes by showing how to operationalize outcomes and measure downstream business impact.

## Chapter 11. Building the Data Foundation

When two teams look at the same data and get different answers, the disagreement rarely stems from a person — it stems from the foundation on which the data was built. That foundation has four layers: Bronze captures everything as it arrives for traceability; Silver cleans, standardizes, and connects data across sources into a single coherent view; Gold applies business rules to produce business-ready models that reflect how the company actually operates; and the Semantic layer ensures every metric has a single agreed-upon definition that means the same thing to everyone, every time. Together, these four layers transform raw, messy data into decisions that can be trusted — and understanding them does not require writing code; it requires knowing which questions to ask when a number looks wrong. The metrics leaders rely on are only as trustworthy as the layers beneath them.

These chapters are about how modern data teams turn raw, messy information into clear, reliable answers that leaders can trust. They follow the journey of data from the moment it lands in your systems, through the stages where it is cleaned and connected, into the layers where your business rules and definitions are applied, and finally into the metrics you see in dashboards and reports. You will see how each layer has a specific job: one focuses on capturing everything safely, another on making it consistent, another on reflecting how your business actually works, and another on locking in shared definitions for things like revenue, churn, and active users. Together, these layers explain why two teams can look at "the same data" and still get different answers — and how to prevent that.

## **11.1. The Bronze Layer**

The bronze layer is where your data journey truly begins. Before you can clean, model, or analyze anything, you first need a dependable place where data from all your systems lands in a consistent way. This includes data from your CRM, billing and finance tools, product event logs, marketing platforms, support systems, and various third-party tools. At this stage, the goal is simple but crucial: capture what is happening in your business without losing or "fixing" anything too early. You are recording reality as it is, not as you wish it looked.

A helpful way to picture the bronze layer is as a warehouse receiving dock. Trucks show up from many suppliers, carrying boxes of different sizes, shapes, and labels. The workers on the dock do not immediately unpack and organize everything into neat aisles. Instead, they focus on a few basic but important tasks: check that the shipment arrived, log who sent it and when, and make sure the boxes are stored safely. Only later, in other parts of the warehouse, will people open the boxes, inspect the contents, and arrange everything for easy access.

In data terms, the bronze layer stores raw files and tables more or less exactly as they arrive. It captures key metadata — which system the data came from, when it was loaded, what the original file name was, and which version or batch it belongs to. It runs simple quality checks: did the data arrive on schedule, is the file readable at all, and does the structure roughly match what's expected?

For business and engineering leaders, the bronze layer matters because it provides traceability. When a dashboard number looks suspicious, the

team needs a way to trace it back through each step to the original source. Without a solid bronze layer, you might only see the cleaned and transformed version, with no easy way to verify whether the original input was already broken. With a good bronze layer, you can always go back and answer questions like, "Did the CRM actually send us that record?" or "What did this event look like before we changed the schema?"

You can think of the bronze layer as the "black box recorder" of your data platform. In aviation, the black box captures everything that happens on a flight so investigators can reconstruct events if something goes wrong. In your data stack, the bronze layer serves a similar purpose. If a metric suddenly drops or jumps, or if two teams see different results, the bronze layer is where your data experts go to investigate what really happened in the source systems at that time.

On its own, the bronze layer will not answer high-level questions like "What is our churn rate?" or "Which campaign delivered the best ROI?" However, without a strong bronze layer, any answers you get later are less trustworthy.

## ***11.2. The Silver Layer***

Once data is safely stored in the bronze layer, the next step is to make it truly usable. If bronze is about "capture everything," silver is about "make sense of it." Here, your data team builds pipelines that clean, standardize, and connect data from different systems so it starts to tell one coherent story about your business instead of many disconnected fragments.

At this stage, the team moves from "What did we receive?" to "What does this actually represent in the real world?" They focus on questions like: Are

customer IDs consistent across systems, or does the same customer show up with three different identifiers? Do we know which records refer to the same person or company, even if names are slightly different or there are typos? Are dates, currencies, and country codes stored in a consistent format?

Returning to the warehouse analogy, if the bronze layer is the receiving dock, the silver layer is the organized storeroom. Boxes are opened, items are inspected, labeled, and grouped by type. Broken items are set aside or flagged. Similar items from different suppliers are placed on the same shelf, with a common label.

In data terms, the silver layer fixes obvious errors, standardizes formats across sources so tables can be reliably joined, builds core tables like Customers, Accounts, Orders, Products, and Events by combining multiple raw inputs, and applies light, widely agreed rules such as removing exact duplicate records or filtering out clearly bad test data.

A strong silver layer has two major benefits. First, it reduces rework — instead of each analyst reinventing their own version of "customers" or "orders" for every report, they all build on the same, trusted tables. Second, it prepares the ground for the gold and semantic layers where you will encode more detailed business rules and define metrics.

### ***11.3. The Gold Layer***

The gold layer is where your data finally starts to look and feel like your business. Up to this point, earlier layers focused on capturing data reliably (bronze) and organizing it into clean, consistent building blocks (silver). In the gold layer, your team takes those building blocks and shapes them

using your specific business rules.

The key idea in the gold layer is **business logic** — the rules, definitions, and edge-case decisions that describe how your company actually operates. How do you define an "active" customer? How do you treat free trials, coupons, and refunds when calculating revenue? How do you classify accounts into segments such as small, mid-market, and enterprise? How do you define terms like "marketing-qualified lead" or a "healthy" subscription?

Different companies, and even different teams inside the same company, can make different choices here even when they are using the same raw data. That is why the gold layer is so important for leaders. It is the place where those choices are made explicit, written down in data models, and ideally agreed upon across functions. Instead of "finance has one version, marketing has another," the goal is to have shared definitions that everyone can see and challenge when needed.

If the silver layer is the organized storeroom, the gold layer is the set of curated shelves arranged for specific groups of people — "Executive KPIs," "Marketing Funnel," "Sales Pipeline," "Product Engagement." Each gold-layer object aims to be easy to use for reporting and dashboards, aligned with how the business measures success, and consistent across departments.

For non-data professionals, the gold layer is often the first part of the data platform that truly feels familiar. The tables and metrics carry the same words you use in meetings: "active customers," "net revenue," "churn," "upsell rate." The difference is that, in the gold layer, these terms are

backed by transparent, shared rules in code, not by individual interpretations hidden in someone's workbook.

## ***11.4. The Semantic Layer***

The semantic layer sits on top of all the previous layers and focuses on one thing: giving your metrics clear, shared meanings that everyone understands in the same way. Bronze, silver, and gold layers handle how data is stored, cleaned, and shaped with business rules. The semantic layer answers a different question: "When this dashboard says 'MRR' or 'Active Users,' what exactly does that mean, and is that the same definition everyone else is using?"

In many organizations, marketing, finance, and product all answer metric questions differently. That leads to recurring meetings where people argue about whose number is correct instead of what action to take.

The semantic layer solves this by defining metrics once, in a central and transparent way, then reusing those definitions everywhere. Instead of each analyst writing their own formula for churn or MRR inside a spreadsheet or dashboard, the calculation lives in a shared semantic model. When someone drags "MRR" into a chart in a BI tool, they are using the same logic that finance approved. When a product manager filters on "Active Accounts," they are using the same definition sales uses. This consistency is what turns metrics from debate topics into decision tools.

For business and engineering leaders, the semantic layer is like a glossary that runs on code. A normal glossary might say, "Revenue: total money earned from customers." The semantic layer goes further and encodes the exact formula and data sources behind that definition — which tables are

used, which dates matter, how to handle refunds, and how to convert currencies.

This trust is what unlocks actionable insights. When leaders know that metrics are consistent, they can focus on interpreting patterns and deciding what to do next. If "Active Users" are dropping in a region, you can dig into why — without first spending half the meeting checking that everyone's numbers match. If "Net Revenue Retention" falls below a target, you can ask whether to adjust pricing, improve onboarding, or invest in customer success, instead of asking whether finance and sales reports are aligned.

The semantic layer is the final step that turns raw data into a reliable, shared source of truth for decision-making. It reduces confusion, speeds up alignment across teams, and makes it easier to move from dashboards to concrete decisions.

The metrics leaders rely on are only as trustworthy as the layers beneath them.

## Chapter 12. Business Metrics in Action

The metrics defined in Chapter 8 become useful the moment they are expressed as SQL. A formula answers the question "what does this measure?" — SQL answers the question "how do we calculate it from the data we actually have?" Together they form the operational definition of a metric: precise enough to build, consistent enough to trust, and clear enough to explain to any stakeholder who questions the number.

This chapter works through all three metric layers — Financial, Customer, and Product Performance — with a formula and working SQL for each. Every query uses the same synthetic dataset: four consecutive weeks in January 2026 across a small set of accounts and products. The dataset is deliberately simple — the goal is to make the SQL readable, not to simulate production volume.

### 12.1. The Dataset

All metrics in this chapter derive from four tables. The same tables power metrics across all three layers — accounts and revenue for financial and customer metrics, and platform events for product performance metrics.

```
-- accounts: one row per customer account
-- Used by: Active Accounts, Adoption Rate, Retention Rate, ARPA, CLTV
WITH accounts AS (
  SELECT * FROM (VALUES
    (1, 'Acme Corp',      'Enterprise', DATE '2025-10-01', TRUE),
    (2, 'Beta Systems',  'Growth',     DATE '2025-11-15', TRUE),
    (3, 'Crest Digital', 'Enterprise', DATE '2025-12-01', TRUE),
    (4, 'Delta Works',   'Starter',    DATE '2026-01-06', TRUE),
    (5, 'Echo Labs',     'Growth',    DATE '2025-09-01', TRUE),
    (6, 'Foxtrot Inc',   'Enterprise', DATE '2025-08-01', TRUE),
```

```

    (7, 'Gamma Tech',      'Starter',    DATE '2026-01-13', TRUE),
    (8, 'Horizon AI',     'Growth',    DATE '2025-10-15', TRUE),
    (9, 'Iris Cloud',     'Enterprise', DATE '2025-11-01', FALSE), -- churned
    (10, 'Jade Networks', 'Growth',    DATE '2025-12-15', TRUE)
  ) AS t(account_id, account_name, tier, signup_date, is_active)
),

-- weekly_revenue: one row per account per week per product
-- Used by: YTD Revenue, Weekly Revenue, Growth Rate, CWGR, ARR, ARPA, CLTV
weekly_revenue AS (
  SELECT * FROM (VALUES
    -- acct product      week_start      mrr
    (1, 'Product A', DATE '2026-01-05', 8500.00),
    (2, 'Product A', DATE '2026-01-05', 3200.00),
    (3, 'Product B', DATE '2026-01-05', 12000.00),
    (4, 'Product A', DATE '2026-01-05', 950.00),
    (5, 'Product A', DATE '2026-01-05', 4100.00),
    (6, 'Product B', DATE '2026-01-05', 11500.00),
    (8, 'Product A', DATE '2026-01-05', 3800.00),
    (10, 'Product B', DATE '2026-01-05', 5200.00),
    (1, 'Product A', DATE '2026-01-12', 8800.00),
    (2, 'Product A', DATE '2026-01-12', 3400.00),
    (3, 'Product B', DATE '2026-01-12', 12500.00),
    (4, 'Product A', DATE '2026-01-12', 1100.00),
    (5, 'Product A', DATE '2026-01-12', 4300.00),
    (6, 'Product B', DATE '2026-01-12', 11800.00),
    (7, 'Product A', DATE '2026-01-12', 750.00),
    (8, 'Product A', DATE '2026-01-12', 4000.00),
    (10, 'Product B', DATE '2026-01-12', 5500.00),
    (1, 'Product A', DATE '2026-01-19', 9100.00),
    (2, 'Product A', DATE '2026-01-19', 3300.00),
    (3, 'Product B', DATE '2026-01-19', 12800.00),
    (4, 'Product A', DATE '2026-01-19', 1050.00),
    (5, 'Product A', DATE '2026-01-19', 4500.00),
    (6, 'Product B', DATE '2026-01-19', 12000.00),
    (7, 'Product A', DATE '2026-01-19', 800.00),
    (8, 'Product A', DATE '2026-01-19', 4100.00),
    (10, 'Product B', DATE '2026-01-19', 5600.00),
    (1, 'Product A', DATE '2026-01-26', 9400.00),
    (2, 'Product A', DATE '2026-01-26', 3500.00),

```

```

    (3, 'Product B', DATE '2026-01-26', 13000.00),
    (4, 'Product A', DATE '2026-01-26', 1100.00),
    (5, 'Product A', DATE '2026-01-26', 4700.00),
    (6, 'Product B', DATE '2026-01-26', 12200.00),
    (7, 'Product A', DATE '2026-01-26', 900.00),
    (8, 'Product A', DATE '2026-01-26', 4200.00),
    (10, 'Product B', DATE '2026-01-26', 5800.00)
  ) AS t(account_id, product, week_start, weekly_revenue)
),

-- usage_events: one row per API request
-- Used by: Adoption Rate, Usage Retention Rate, Error Rate, Latency, Time to Value
usage_events AS (
  SELECT * FROM (VALUES
    -- req_id acct week_start          tokens latency_ms is_error event_ts
    (1001, 1,  DATE '2026-01-05',    4500,  320, FALSE, TIMESTAMP '2026-01-
05 09:00:00'),
    (1002, 2,  DATE '2026-01-05',    2100,  410, FALSE, TIMESTAMP '2026-01-
05 10:30:00'),
    (1003, 3,  DATE '2026-01-05',    8200,  290, FALSE, TIMESTAMP '2026-01-
05 11:00:00'),
    (1004, 4,  DATE '2026-01-06',    1200,  380, FALSE, TIMESTAMP '2026-01-
06 09:15:00'),
    (1005, 5,  DATE '2026-01-05',    3600,  350, FALSE, TIMESTAMP '2026-01-
05 14:00:00'),
    (1006, 6,  DATE '2026-01-05',    9100,  300, FALSE, TIMESTAMP '2026-01-
05 15:30:00'),
    (1007, 8,  DATE '2026-01-05',    2800,  420, FALSE, TIMESTAMP '2026-01-
05 16:00:00'),
    (1008, 10, DATE '2026-01-05',    5500,  310, FALSE, TIMESTAMP '2026-01-
05 17:00:00'),
    (1009, 1,  DATE '2026-01-12',    4800,  330, FALSE, TIMESTAMP '2026-01-
12 09:00:00'),
    (1010, 2,  DATE '2026-01-12',    2300,  390, FALSE, TIMESTAMP '2026-01-
12 10:00:00'),
    (1011, 3,  DATE '2026-01-12',    8600,  280, FALSE, TIMESTAMP '2026-01-
12 11:30:00'),
    (1012, 4,  DATE '2026-01-12',    1400, 3800, TRUE,  TIMESTAMP '2026-01-
13 14:00:00'),
    (1013, 5,  DATE '2026-01-12',    3900,  360, FALSE, TIMESTAMP '2026-01-

```

```

12 15:00:00'),
    (1014, 6, DATE '2026-01-12', 9400, 295, FALSE, TIMESTAMP '2026-01-
12 16:00:00'),
    (1015, 7, DATE '2026-01-13', 1800, 440, FALSE, TIMESTAMP '2026-01-
13 09:30:00'),
    (1016, 8, DATE '2026-01-12', 3100, 410, FALSE, TIMESTAMP '2026-01-
12 17:00:00'),
    (1017, 10, DATE '2026-01-12', 5800, 320, FALSE, TIMESTAMP '2026-01-
12 18:00:00'),
    (1018, 1, DATE '2026-01-19', 5100, 340, FALSE, TIMESTAMP '2026-01-
19 09:00:00'),
    (1019, 2, DATE '2026-01-19', 2200, 400, FALSE, TIMESTAMP '2026-01-
19 10:00:00'),
    (1020, 3, DATE '2026-01-19', 9000, 4200, TRUE, TIMESTAMP '2026-01-
19 11:00:00'),
    (1021, 4, DATE '2026-01-19', 1500, 360, FALSE, TIMESTAMP '2026-01-
19 12:00:00'),
    (1022, 5, DATE '2026-01-19', 4200, 345, FALSE, TIMESTAMP '2026-01-
19 13:00:00'),
    (1023, 6, DATE '2026-01-19', 9700, 305, FALSE, TIMESTAMP '2026-01-
19 14:00:00'),
    (1024, 7, DATE '2026-01-19', 2000, 430, FALSE, TIMESTAMP '2026-01-
19 15:00:00'),
    (1025, 8, DATE '2026-01-19', 3400, 415, FALSE, TIMESTAMP '2026-01-
19 16:00:00'),
    (1026, 10, DATE '2026-01-19', 6100, 315, FALSE, TIMESTAMP '2026-01-
19 17:00:00'),
    (1027, 1, DATE '2026-01-26', 5400, 325, FALSE, TIMESTAMP '2026-01-
26 09:00:00'),
    (1028, 2, DATE '2026-01-26', 2500, 385, FALSE, TIMESTAMP '2026-01-
26 10:00:00'),
    (1029, 3, DATE '2026-01-26', 9300, 275, FALSE, TIMESTAMP '2026-01-
26 11:00:00'),
    (1030, 4, DATE '2026-01-26', 1600, 370, FALSE, TIMESTAMP '2026-01-
26 12:00:00'),
    (1031, 5, DATE '2026-01-26', 4500, 355, FALSE, TIMESTAMP '2026-01-
26 13:00:00'),
    (1032, 6, DATE '2026-01-26', 10100, 298, FALSE, TIMESTAMP '2026-01-
26 14:00:00'),
    (1033, 7, DATE '2026-01-26', 2200, 425, FALSE, TIMESTAMP '2026-01-

```

```

26 15:00:00'),
      (1034, 8, DATE '2026-01-26', 3600, 405, FALSE, TIMESTAMP '2026-01-
26 16:00:00'),
      (1035, 10, DATE '2026-01-26', 6400, 310, FALSE, TIMESTAMP '2026-01-
26 17:00:00')
    ) AS t(request_id, account_id, week_start, tokens,
           latency_ms, is_error, event_ts)
),

-- incidents: one row per platform incident
-- Used by: Uptime, MTTR
incidents AS (
  SELECT * FROM (VALUES
    -- inc_id week_start          started_at          resolved_at
downtime_mins
    (1, DATE '2026-01-12',    TIMESTAMP '2026-01-13 02:15:00', TIMESTAMP
'2026-01-13 03:40:00', 85),
    (2, DATE '2026-01-19',    TIMESTAMP '2026-01-20 14:00:00', TIMESTAMP
'2026-01-20 14:35:00', 35),
    (3, DATE '2026-01-19',    TIMESTAMP '2026-01-22 22:00:00', TIMESTAMP
'2026-01-23 01:30:00', 210)
  ) AS t(incident_id, week_start, started_at, resolved_at, downtime_minutes)
),

-- support_tickets: one row per ticket
-- Used by: Ticket Volume, Resolution Time
support_tickets AS (
  SELECT * FROM (VALUES
    -- ticket acct week_start      created_at
resolved_at          category
    (201, 4, DATE '2026-01-12',    TIMESTAMP '2026-01-13 14:30:00', TIMESTAMP
'2026-01-14 09:00:00', 'error'),
    (202, 3, DATE '2026-01-19',    TIMESTAMP '2026-01-19 11:45:00', TIMESTAMP
'2026-01-20 10:00:00', 'latency'),
    (203, 1, DATE '2026-01-19',    TIMESTAMP '2026-01-20 08:00:00', TIMESTAMP
'2026-01-20 16:30:00', 'billing'),
    (204, 5, DATE '2026-01-19',    TIMESTAMP '2026-01-21 09:00:00', TIMESTAMP
'2026-01-22 11:00:00', 'onboarding'),
    (205, 7, DATE '2026-01-19',    TIMESTAMP '2026-01-22 14:00:00', TIMESTAMP
'2026-01-23 09:30:00', 'latency'),

```

```

        (206, 2, DATE '2026-01-26', TIMESTAMP '2026-01-26 10:00:00', TIMESTAMP
'2026-01-26 15:00:00', 'billing'),
        (207, 8, DATE '2026-01-26', TIMESTAMP '2026-01-27 09:30:00', TIMESTAMP
'2026-01-28 14:00:00', 'onboarding')
    ) AS t(ticket_id, account_id, week_start, created_at, resolved_at, category)
)

```

The queries that follow reference these CTEs inline. In production, each CTE would be replaced by a `ref()` call to the corresponding dbt model in the Gold layer.

## 12.2. Financial Metrics

### 12.2.1. R.1 Year-to-Date Revenue

**Formula:** YTD Revenue = SUM(revenue) WHERE week\_start >= fiscal\_year\_start AND week\_start <= today

```

SELECT
    SUM(weekly_revenue) AS ytd_revenue,
    -- $209,650 – sum of all revenue from Jan 5 to Jan 26
    COUNT(DISTINCT account_id) AS contributing_accounts,
    -- 9 – accounts with at least one revenue record in January
    ROUND(SUM(weekly_revenue) / 20000000.0 * 100, 1) AS pct_of_annual_target
    -- 1.0% – $209,650 of a $20M annual target after 4 weeks
FROM weekly_revenue
WHERE week_start >= DATE '2026-01-01' -- fiscal year start
    AND week_start <= CURRENT_DATE;

```

### 12.2.2. R.2 Weekly Revenue

**Formula:** Weekly Revenue = SUM(revenue) WHERE week\_start = [week]

```

SELECT
    week_start,

```

```

SUM(weekly_revenue)           AS weekly_revenue,
-- Week 1 (Jan 05): $49,250
-- Week 2 (Jan 12): $52,150
-- Week 3 (Jan 19): $53,150
-- Week 4 (Jan 26): $54,800
COUNT(DISTINCT account_id)   AS active_accounts
FROM weekly_revenue
GROUP BY week_start
ORDER BY week_start;

```

### 12.2.3. R.3 Weekly Growth Rate

**Formula:**  $\text{Weekly Growth Rate} = \frac{(\text{this\_week\_revenue} - \text{last\_week\_revenue})}{\text{last\_week\_revenue}}$

```

WITH weekly AS (
  SELECT
    week_start,
    SUM(weekly_revenue) AS revenue
  FROM weekly_revenue
  GROUP BY week_start
)
SELECT
  week_start,
  revenue,
  LAG(revenue) OVER (ORDER BY week_start) AS prev_week_revenue,
  ROUND(
    (revenue - LAG(revenue) OVER (ORDER BY week_start))
    / LAG(revenue) OVER (ORDER BY week_start) * 100
  , 1) AS weekly_growth_rate_pct
-- Week 1 (Jan 05): NULL (no prior week)
-- Week 2 (Jan 12): +5.9% ($49,250 → $52,150)
-- Week 3 (Jan 19): +1.9% ($52,150 → $53,150)
-- Week 4 (Jan 26): +3.1% ($53,150 → $54,800)
FROM weekly
ORDER BY week_start;

```

### 12.2.4. R.4 Cumulative Weekly Growth Rate (CWGR)

**Formula:**  $CWGR = (last\_week\_revenue / first\_week\_revenue)^{(1 / (n\_weeks - 1))} - 1$

```

WITH weekly AS (
  SELECT
    week_start,
    SUM(weekly_revenue) AS revenue,
    ROW_NUMBER() OVER (ORDER BY week_start) AS week_num
  FROM weekly_revenue
  GROUP BY week_start
),
bounds AS (
  SELECT
    MIN(revenue) FILTER (WHERE week_num = 1)           AS first_week,
    MAX(revenue) FILTER (WHERE week_num = (SELECT MAX(week_num) FROM weekly))
                                                    AS last_week,
    MAX(week_num) - 1                               AS periods
  FROM weekly
)
SELECT
  first_week,
  last_week,
  periods,
  ROUND(
    (POWER(last_week / first_week, 1.0 / periods) - 1) * 100
  , 2)
  AS cwgr_pct
-- first_week: $49,250  last_week: $54,800  periods: 3
-- CWGR: 3.6% per week – the average weekly growth rate across January
FROM bounds;

```

### 12.2.5. R.5 Catch-up CWGR

**Formula:**  $Catch\text{-}up\ CWGR = (annual\_target / ytd\_revenue)^{(1 / remaining\_weeks)} - 1$

```

WITH ytd AS (
  SELECT SUM(weekly_revenue) AS ytd_revenue
  FROM weekly_revenue

```

```

WHERE week_start >= DATE '2026-01-01'
)
SELECT
  ytd_revenue,
  20000000.00                                AS annual_target,
  48                                           AS remaining_weeks,
  -- 52 weeks in year minus 4 weeks elapsed = 48 remaining
  ROUND(
    (POWER(20000000.00 / ytd_revenue, 1.0 / 48) - 1) * 100
  , 2)                                         AS catchup_cwgr_pct
  -- ytd_revenue: $209,650   annual_target: $20,000,000
  -- Catch-up CWGR: 21.6% per week – significantly above current 3.6% CWGR
  -- This gap signals the annual target requires a step change in growth
FROM ytd;

```

### 12.2.6. R.6 Annual Recurring Revenue (ARR)

**Formula:**  $ARR = MRR \times 12$

```

WITH latest_week AS (
  SELECT
    account_id,
    SUM(weekly_revenue) * 4.33                AS mrr
    -- multiply weekly revenue by 4.33 to approximate monthly revenue
    -- in production this would use actual monthly billing records
  FROM weekly_revenue
  WHERE week_start = DATE '2026-01-26'       -- most recent week as proxy for
current MRR
  GROUP BY account_id
)
SELECT
  SUM(mrr)                                    AS total_mrr,
  SUM(mrr) * 12                               AS arr
  -- total_mrr: ~$237,284   arr: ~$2,847,408
  -- This is the annualised revenue run rate based on the most recent week
FROM latest_week;

```

### 12.2.7. R.7 Average Revenue Per Account (ARPA)

**Formula:**  $ARPA = \text{Total Revenue} / \text{Number of Active Accounts}$

```

WITH weekly_totals AS (
  SELECT
    week_start,
    SUM(weekly_revenue)           AS total_revenue,
    COUNT(DISTINCT account_id)   AS active_accounts
  FROM weekly_revenue
  GROUP BY week_start
)
SELECT
  week_start,
  total_revenue,
  active_accounts,
  ROUND(total_revenue / active_accounts, 0) AS arpa
-- Week 1 (Jan 05): $49,250 / 8 accounts = $6,156
-- Week 2 (Jan 12): $52,150 / 9 accounts = $5,794
-- Week 3 (Jan 19): $53,150 / 9 accounts = $5,906
-- Week 4 (Jan 26): $54,800 / 9 accounts = $6,089
-- ARPA dip in week 2 reflects new lower-value accounts joining (Delta, Gamma)
FROM weekly_totals
ORDER BY week_start;

```

### 12.2.8. R.8 Customer Lifetime Value (CLTV)

**Formula:**  $CLTV = ARPA \times \text{Average Account Lifetime (months)}$

```

WITH account_tenure AS (
  SELECT
    account_id,
    account_name,
    tier,
    -- months since signup to end of January 2026
    EXTRACT(YEAR FROM AGE(DATE '2026-01-31', signup_date)) * 12
    + EXTRACT(MONTH FROM AGE(DATE '2026-01-31', signup_date)) AS tenure_months
  FROM accounts
  WHERE is_active = TRUE
),

```

```

current_arpa AS (
  SELECT
    account_id,
    SUM(weekly_revenue) * 4.33           AS mrr
  FROM weekly_revenue
  WHERE week_start = DATE '2026-01-26'
  GROUP BY account_id
)
SELECT
  a.account_name,
  a.tier,
  a.tenure_months,
  ROUND(c.mrr, 0)                       AS current_mrr,
  ROUND(c.mrr * 24, 0)                  AS cltv_24mo
  -- assumes 24-month average account lifetime – replace with actual retention
  data
  -- Acme Corp:  tenure 4mo, MRR $9,400, CLTV $225,600
  -- Foxtrot Inc: tenure 6mo, MRR $12,200, CLTV $292,800
FROM account_tenure a
JOIN current_arpa c ON a.account_id = c.account_id
ORDER BY cltv_24mo DESC;

```

## 12.3. Customer Metrics

### 12.3.1. C.1 Active Accounts

**Formula:** Active Accounts = COUNT(accounts WHERE is\_active = TRUE AND week\_start = [week])

```

SELECT
  wr.week_start,
  COUNT(DISTINCT wr.account_id)         AS active_accounts
  -- Week 1 (Jan 05): 8 – Delta and Gamma not yet signed up
  -- Week 2 (Jan 12): 9 – Delta joins Jan 6, Gamma joins Jan 13
  -- Week 3 (Jan 19): 9
  -- Week 4 (Jan 26): 9
  -- Account 9 (Iris Cloud) churned – never appears in revenue data

```

```

FROM weekly_revenue wr
JOIN accounts a ON wr.account_id = a.account_id
WHERE a.is_active = TRUE
GROUP BY wr.week_start
ORDER BY wr.week_start;

```

### 12.3.2. C.2 Weekly Sign-up

**Formula:** Weekly Sign-ups = COUNT(accounts WHERE signup\_date BETWEEN week\_start AND week\_end)

```

SELECT
    DATE_TRUNC('week', signup_date)::date      AS signup_week,
    COUNT(*)                                  AS new_signups,
    -- Week of Jan 05: 0 (no new accounts this week)
    -- Week of Jan 06: 1 (Delta Works signed up Jan 6)
    -- Week of Jan 13: 1 (Gamma Tech signed up Jan 13)
    -- All prior accounts signed up before January 2026
    STRING_AGG(account_name, ', ')           AS new_accounts
FROM accounts
WHERE signup_date >= DATE '2026-01-01'
GROUP BY DATE_TRUNC('week', signup_date)
ORDER BY signup_week;

```

### 12.3.3. C.3 Adoption Rate

**Formula:** Adoption Rate = Accounts reaching usage threshold / Total new accounts in cohort

```

-- Threshold: at least 1,000 tokens consumed within 13 weeks of signup
WITH new_accounts AS (
    SELECT account_id, signup_date
    FROM accounts
    WHERE signup_date >= DATE '2026-01-01'      -- January 2026 cohort
),
usage_check AS (
    SELECT
        na.account_id,

```

```

        na.signup_date,
        SUM(ue.tokens)                AS total_tokens
FROM new_accounts na
LEFT JOIN usage_events ue
    ON ue.account_id = na.account_id
    AND ue.event_ts <= na.signup_date + INTERVAL '91 days'
GROUP BY na.account_id, na.signup_date
)
SELECT
    COUNT(*)                        AS cohort_size,
    -- 2 accounts: Delta Works, Gamma Tech
    COUNT(*) FILTER (WHERE total_tokens >= 1000) AS adopted,
    -- 2 – both accounts made API calls within January
    ROUND(
        COUNT(*) FILTER (WHERE total_tokens >= 1000)::numeric
        / COUNT(*) * 100
    , 1)                            AS adoption_rate_pct
    -- 100.0% – both January sign-ups adopted within the observation window
FROM usage_check;

```

#### 12.3.4. C.4 Retention Rate

**Formula:** Retention Rate = Accounts active at end of period / Accounts active at start of period

```

WITH week1_accounts AS (
    -- accounts present in week 1
    SELECT DISTINCT account_id FROM weekly_revenue
    WHERE week_start = DATE '2026-01-05'
),
week4_accounts AS (
    -- accounts still present in week 4
    SELECT DISTINCT account_id FROM weekly_revenue
    WHERE week_start = DATE '2026-01-26'
)
SELECT
    COUNT(DISTINCT w1.account_id)    AS start_of_period_accounts,
    -- 8 accounts in week 1
    COUNT(DISTINCT w4.account_id)    AS end_of_period_accounts,

```

```

-- 9 accounts in week 4 (includes new signups)
COUNT(DISTINCT w1.account_id)
  FILTER (WHERE w4.account_id IS NOT NULL) AS retained_accounts,
-- 8 – all original accounts were retained (Iris Cloud never had revenue)
ROUND(
  COUNT(DISTINCT w1.account_id)
    FILTER (WHERE w4.account_id IS NOT NULL)::numeric
  / COUNT(DISTINCT w1.account_id) * 100
, 1) AS retention_rate_pct
-- 100.0% – no accounts from week 1 were lost by week 4
FROM week1_accounts w1
LEFT JOIN week4_accounts w4 ON w1.account_id = w4.account_id;

```

### 12.3.5. C.5 Usage Retention Rate

**Formula:** Usage Retention Rate = Period 2 tokens / Period 1 tokens (existing accounts only)

```

WITH existing_accounts AS (
  -- accounts active before January 2026 – excludes new signups
  SELECT account_id FROM accounts
  WHERE signup_date < DATE '2026-01-01'
    AND is_active = TRUE
),
week1_usage AS (
  SELECT SUM(tokens) AS tokens
  FROM usage_events
  WHERE week_start = DATE '2026-01-05'
    AND account_id IN (SELECT account_id FROM existing_accounts)
),
week4_usage AS (
  SELECT SUM(tokens) AS tokens
  FROM usage_events
  WHERE week_start = DATE '2026-01-26'
    AND account_id IN (SELECT account_id FROM existing_accounts)
)
SELECT
  w1.tokens AS week1_tokens,
  -- 35,800 tokens from existing accounts in week 1

```

```

w4.tokens AS week4_tokens,
-- 43,500 tokens from existing accounts in week 4
ROUND(w4.tokens::numeric / w1.tokens * 100, 1) AS usage_retention_rate_pct
-- 121.5% – existing accounts grew token consumption by 21.5%
-- above 100% means the base is expanding, not just holding steady
FROM week1_usage w1, week4_usage w4;

```

### 12.3.6. C.6 Customer Satisfaction Score

**Formula:** CS Score = weighted average of response score, ticket resolution time, and uptime (0-100 scale)

```

-- CS Score is a composite – this query approximates it from available signals:
-- component 1: uptime score (40% weight)
-- component 2: ticket resolution score (30% weight)
-- component 3: error rate score (30% weight)
WITH uptime_score AS (
  SELECT
    week_start,
    -- 100 = perfect, deduct points for downtime
    GREATEST(0, 100 - (SUM(downtime_minutes)::numeric / 10080 * 100) * 10) AS
score
    -- 10080 = minutes in a week; scale impact to a 0-100 point deduction
  FROM incidents
  GROUP BY week_start
),
error_score AS (
  SELECT
    week_start,
    GREATEST(0, 100 - (COUNT(*) FILTER (WHERE is_error)::numeric
/ COUNT(*) * 100) * 5) AS score
  FROM usage_events
  GROUP BY week_start
),
ticket_score AS (
  SELECT
    week_start,
    GREATEST(0, 100 - (AVG(EXTRACT(EPOCH FROM (resolved_at - created_at))
/ 3600))::numeric) AS score

```

```

        -- deduct 1 point per hour of average resolution time
    FROM support_tickets
    GROUP BY week_start
)
SELECT
    e.week_start,
    ROUND(
        COALESCE(u.score, 100) * 0.4
        + COALESCE(t.score, 100) * 0.3
        + e.score * 0.3
    , 1) AS cs_score
-- Week 1 (Jan 05): ~100.0 – no incidents or tickets
-- Week 2 (Jan 12): ~91.4 – one incident (85 min), one error, one ticket
-- Week 3 (Jan 19): ~78.2 – two incidents (245 min combined), one error, three
tickets
-- Week 4 (Jan 26): ~95.8 – no incidents, two tickets with fast resolution
FROM error_score e
LEFT JOIN uptime_score u ON e.week_start = u.week_start
LEFT JOIN ticket_score t ON e.week_start = t.week_start
ORDER BY e.week_start;

```

## 12.4. Product Performance Metrics

### 12.4.1. P.1 Time to Value

**Formula:** TTV = PERCENTILE(days from signup\_date to first API call) at p50, p90, p100

```

WITH first_usage AS (
    SELECT
        a.account_id,
        a.account_name,
        a.signup_date,
        MIN(ue.event_ts)::date AS first_api_call_date,
        MIN(ue.event_ts)::date - a.signup_date AS days_to_value
    FROM accounts a
    JOIN usage_events ue ON a.account_id = ue.account_id

```

```

GROUP BY a.account_id, a.account_name, a.signup_date
)
SELECT
  PERCENTILE_CONT(0.50) WITHIN GROUP (ORDER BY days_to_value) AS p50_days,
  -- p50: 0 days – half of accounts made their first call on day of signup
  PERCENTILE_CONT(0.90) WITHIN GROUP (ORDER BY days_to_value) AS p90_days,
  -- p90: 7 days
  MAX(days_to_value) AS p100_days
  -- p100: 7 days (Gamma Tech, newest account)
FROM first_usage;

```

### 12.4.2. P.2 Error Rate

**Formula:** Error Rate = error requests / total requests

```

SELECT
  week_start,
  COUNT(*) AS total_requests,
  COUNT(*) FILTER (WHERE is_error) AS error_requests,
  ROUND(
    COUNT(*) FILTER (WHERE is_error)::numeric
    / COUNT(*) * 100
  , 2) AS error_rate_pct
  -- Week 1 (Jan 05): 0.00% – 8 requests, 0 errors
  -- Week 2 (Jan 12): 9.09% – 11 requests, 1 error (account 4, latency spike)
  -- Week 3 (Jan 19): 9.09% – 11 requests, 1 error (account 3, latency spike)
  -- Week 4 (Jan 26): 0.00% – 9 requests, 0 errors
FROM usage_events
GROUP BY week_start
ORDER BY week_start;

```

### 12.4.3. P.3 Latency

**Formula:** Latency = PERCENTILE(request\_latency\_ms) at p50, p90, p99

```

SELECT
  week_start,

```

```

PERCENTILE_CONT(0.50) WITHIN GROUP
  (ORDER BY latency_ms)::int          AS p50_ms,
PERCENTILE_CONT(0.90) WITHIN GROUP
  (ORDER BY latency_ms)::int          AS p90_ms,
PERCENTILE_CONT(0.99) WITHIN GROUP
  (ORDER BY latency_ms)::int          AS p99_ms
-- Week 1: p50=355ms p90=432ms p99=432ms – all healthy
-- Week 2: p50=390ms p90=448ms p99=3800ms – p99 spike from account 4 error
event
-- Week 3: p50=355ms p90=430ms p99=4200ms – p99 spike from account 3 error
event
-- Week 4: p50=355ms p90=425ms p99=425ms – recovered
FROM usage_events
GROUP BY week_start
ORDER BY week_start;

```

#### 12.4.4. P.4 Uptime and Availability

**Formula:** 
$$\text{Uptime \%} = \frac{(\text{total\_minutes\_in\_period} - \text{downtime\_minutes})}{\text{total\_minutes\_in\_period}} \times 100$$

```

WITH all_weeks AS (
  SELECT DISTINCT week_start FROM usage_events
),
week_downtime AS (
  SELECT
    week_start,
    COALESCE(SUM(downtime_minutes), 0) AS downtime_mins
  FROM incidents
  GROUP BY week_start
)
SELECT
  w.week_start,
  COALESCE(d.downtime_mins, 0) AS downtime_minutes,
  10080 AS total_minutes_in_week,
  -- 7 days x 24 hours x 60 minutes = 10,080 minutes per week
  ROUND(
    (10080 - COALESCE(d.downtime_mins, 0))::numeric
    / 10080 * 100

```

```

, 3) AS uptime_pct
-- Week 1 (Jan 05): 100.000% – no incidents
-- Week 2 (Jan 12): 99.157% – 85 min downtime (incident 1)
-- Week 3 (Jan 19): 97.569% – 245 min downtime (incidents 2 and 3 combined)
-- Week 4 (Jan 26): 100.000% – no incidents
FROM all_weeks w
LEFT JOIN week_downtime d ON w.week_start = d.week_start
ORDER BY w.week_start;

```

#### 12.4.5. P.5 Mean Time to Recovery (MTTR)

**Formula:**  $MTTR = \text{AVG}(\text{resolved\_at} - \text{started\_at})$  across all incidents in period

```

SELECT
  COUNT(*) AS total_incidents,
  -- 3 incidents across January
  ROUND(AVG(downtime_minutes), 0) AS mttr_minutes,
  -- MTTR: 110 minutes (85 + 35 + 210) / 3
  MIN(downtime_minutes) AS fastest_recovery_mins,
  -- 35 minutes (incident 2)
  MAX(downtime_minutes) AS slowest_recovery_mins,
  -- 210 minutes (incident 3 – late-night extended outage)
  ROUND(SUM(downtime_minutes)::numeric / 60, 1) AS total_downtime_hours
  -- 5.5 hours total downtime in January
FROM incidents;

```

#### 12.4.6. P.6 Support Ticket Volume

**Formula:**  $\text{Ticket Volume} = \text{COUNT}(\text{tickets WHERE created\_at BETWEEN week\_start AND week\_end})$

```

SELECT
  week_start,
  COUNT(*) AS ticket_volume,
  -- Week 1 (Jan 05): 0
  -- Week 2 (Jan 12): 1 (error-related)
  -- Week 3 (Jan 19): 3 (1 latency + 1 billing + 1 onboarding)

```

```

-- Week 4 (Jan 26): 2 (1 billing + 1 onboarding)
COUNT(*) FILTER (WHERE category IN ('error', 'latency')) AS platform_tickets,
-- tickets directly related to platform reliability issues
COUNT(*) FILTER (WHERE category IN ('billing', 'onboarding')) AS other_tickets
FROM support_tickets
GROUP BY week_start
ORDER BY week_start;

```

### 12.4.7. P.7 Support Ticket Resolution Time

**Formula:** Resolution Time = AVG(resolved\_at - created\_at) in hours

```

SELECT
  week_start,
  COUNT(*) AS tickets,
  ROUND(AVG(
    EXTRACT(EPOCH FROM (resolved_at - created_at)) / 3600
  ), 1) AS avg_resolution_hours,
  -- Week 2 (Jan 12): 18.5 hours (1 ticket)
  -- Week 3 (Jan 19): 26.8 hours (3 tickets – onboarding ticket slowest)
  -- Week 4 (Jan 26): 14.3 hours (2 tickets – fastest week)
  ROUND(MIN(
    EXTRACT(EPOCH FROM (resolved_at - created_at)) / 3600
  ), 1) AS fastest_resolution_hours,
  ROUND(MAX(
    EXTRACT(EPOCH FROM (resolved_at - created_at)) / 3600
  ), 1) AS slowest_resolution_hours
FROM support_tickets
GROUP BY week_start
ORDER BY week_start;

```

## 12.5. Reading the Metrics Together

The real value of these metrics emerges when they are read as a connected story rather than independent numbers. The January data tells a specific one.

Weeks one and two are clean — revenue grows steadily, all accounts are retained, no incidents, no tickets. Week three changes the picture. Two incidents combine for 245 minutes of downtime, error rate spikes to 9%, p99 latency jumps to 4,200ms, three support tickets arrive, and CS score drops from the high nineties to 78. Week four shows recovery — no incidents, error rate returns to zero, resolution time improves — but revenue growth slows to 3.1%, its weakest week of the month.

The sequence is the lesson. The platform events in week three are visible before their impact reaches revenue. An analytics team reading these metrics in real time would have seen the incident cluster, flagged the CS score drop, and reached out to the three affected accounts before the revenue line softened. That is what the three-layer framework is designed to enable — not explanation after the fact, but intervention while there is still time to act.

## Chapter 13. The Complete Reporting Framework

An effective reporting framework typically includes these core elements:

1. **Objectives and Scope:** Defines the "why" and boundaries of the report (e.g., time period, audience)
2. **Data Collection and Sources:** Identifies where data comes from and how it is gathered
3. **Standards and Guidelines:** Specific rules followed (e.g., IFRS, GRI)
4. **Metrics and KPIs:** The quantitative and qualitative measures used to track performance
5. **Analysis and Interpretation:** Methods used to turn raw data into insights (e.g., trend analysis)
6. **Tools and Technology:** Software used for generation, such as Power BI or Tableau
7. **Presentation:** The format (e.g., dashboard, written report) and distribution method
8. **Review and Assurance:** Processes to verify accuracy, often via internal or external audits
9. **Compliance and Ethics:** Ensures legal requirements and transparency standards are met
10. **Feedback and Improvement:** Mechanisms for stakeholders to provide input for future reporting cycles

### ***13.1. The Four Types of Reports***

**Analytical (The Storyteller):** These are designed for deep exploration

using large volumes of data. They use filters and drill-downs to allow analysts to "connect the dots" and build a narrative about causes, patterns, and future predictions.

**Strategic (The Summary):** These provide the "high-level plot" for executives by showing whether the business is hitting its long-term goals. They tell a simplified story of overall health (e.g., "Are we winning or losing?") without getting bogged down in details.

**Tactical (The Chapter):** These focus on specific departmental projects or mid-term goals. They tell the story of a specific initiative's progress, like a marketing campaign's journey from launch to lead generation.

**Operational (The Status Report):** These are built for real-time monitoring and often lack a complex narrative. They act as an "early warning system," telling you what is happening right now so you can take immediate action.

## Chapter 14. Analysis and Insights

Metrics tell you what happened. Analysis tells you why — and what to do about it. The numbers produced by the queries in Chapter 12 are not insights by themselves. They are raw material. Turning that material into a recommendation someone will act on requires three distinct moves: exploration, which surfaces what is interesting in the data; hypothesis-driven investigation, which tests whether what looks interesting is real; and insight packaging, which translates what is real into a recommendation a stakeholder can act on. Each move requires a different question, a different method, and a different output. Together they form the analysis workflow that bridges the data foundation to the decisions that drive business outcomes.

### ***14.1. From Metrics to Questions***

The most common misuse of a dashboard is to open it and look at it, waiting for something to stand out. The problem with that approach is that dashboards are designed to confirm what you already know. They organize the metrics you defined, in the layout you chose, for the audience you had in mind when you built them. What they are not designed to do is surface what you did not think to measure, question assumptions you built into the metric definitions, or reveal patterns that span multiple metrics at once.

Analysis starts not with a number but with a question. A question has a specific shape: it names a metric, a time period, and a suspected cause or comparison. "Why did the Customer Satisfaction Score drop in week three?" is a question. "How do our dashboards look?" is not.

Good analysis questions come from three sources. The first is anomalies — a metric that moved outside its expected range, as discussed in Chapter 9. The second is decisions that have not yet been made — a pricing change under consideration, a CS engagement model being debated, a reliability investment requiring prioritization. The third is patterns observed in the framework’s driver structure, introduced in Chapter 10: when a success metric moves, the factors and drivers that connect to it are the natural starting point for the investigation.

The January 2026 data from Chapter 12 contains a clear anomaly that generates a good analysis question. Customer Satisfaction Score drops from the high nineties in weeks one and two to 78.2 in week three, and revenue growth slows to its lowest rate of the month in week four. The question that follows naturally is: what drove the CS score drop, and what is the risk of a sustained revenue impact if the underlying cause is not addressed?

That is the question this chapter answers — not as an exercise in the data, but as an illustration of how the analysis workflow moves from numbers to action.

## ***14.2. Exploratory Analysis***

Exploratory analysis is the practice of examining data without a predetermined conclusion. Its purpose is not to confirm a hypothesis. Its purpose is to find one. The analyst starts with the anomaly and works outward, asking a sequence of questions that progressively narrow the search.

There are five fundamental moves in exploratory analysis, and they apply regardless of the tool, the size of the dataset, or the specific metric under

investigation.

**Compare:** What changed between the period that looks normal and the period that looks abnormal? Comparison is the most basic move. In the January data, week three stands out immediately when compared to weeks one and two: two incidents instead of zero, 245 minutes of downtime instead of none, an error rate of 9% instead of 0%, and three support tickets instead of one. Comparison reveals the magnitude of the change and confirms that something different happened.

**Segment:** Does the pattern hold across all segments, or is it concentrated in one? Segmentation tests whether the observed effect is widespread or isolated. In the January data, the two incidents in week three are separate events — one at 2:00 AM on Tuesday and one at 10:00 PM on Thursday — suggesting platform-wide exposure rather than a single account issue. However, the error events are concentrated in accounts 3 and 4, not distributed evenly. That segmentation matters: it distinguishes a broad reliability issue from a set of account-specific integration problems.

**Trend:** Is this a new pattern or a continuation of an existing one? Trend analysis places the current observation in historical context. If CS scores had been declining for five consecutive weeks, a drop in week three would be part of a sustained degradation. If weeks one and two were clean and week three was the first disruption, the event in week three is the probable cause. In the January data, the trend is clear: the drop is discontinuous, not gradual, which points to a specific triggering event rather than a slow erosion.

**Correlate:** Do multiple metrics move together in a way that suggests a

common cause? Correlation surfaces relationships between variables that may not be visible in a single metric view. In the January data, three signals move simultaneously in week three: incident count rises, error rate rises, and support ticket volume rises. CS score falls. Revenue growth rate decelerates in week four. The co-movement across these variables — platform reliability metrics leading customer satisfaction metrics, which in turn leads financial metrics — is consistent with the causal sequence the three-layer framework predicts: product performance metrics are the earliest signal, customer metrics follow, and financial metrics lag.

**Rank:** Which instances, accounts, or dimensions are contributing the most to the observed pattern? Ranking identifies where to look next. In the January data, ranking support tickets by category reveals that the three tickets in week three span different problem types — one latency issue, one billing question, one onboarding question — suggesting that the platform reliability problem is creating secondary effects beyond the direct technical impact.

These five moves do not need to be applied in sequence, and no single application will exhaust the analytical value of a dataset. The goal of exploratory analysis is to generate a hypothesis: a specific, testable claim about what caused the pattern and what the implications are. Once the exploration produces a plausible explanation, the next step is to test it.

### ***14.3. Hypothesis-Driven Investigation***

A hypothesis has a specific structure: it connects an observed pattern in the data to a claimed cause, and it implies a prediction that can be checked. "The CS score drop in week three was caused by the two platform incidents,

and accounts directly affected by those incidents will show the largest decline in subsequent usage" is a hypothesis. "The CS score dropped because of reliability issues" is a supposition. The difference is precision: a hypothesis names the mechanism, identifies the affected population, and predicts a measurable outcome.

Hypothesis-driven investigation follows three steps: form the hypothesis, define what evidence would confirm or challenge it, and then go look for that evidence.

**Forming the hypothesis.** In the January data, the exploratory analysis identified three co-moving signals in week three: platform incidents, rising error rate, and rising support ticket volume. The natural hypothesis is: the two platform incidents in week three caused the CS score drop by degrading platform reliability for affected accounts, and the downstream impact will be visible in week four usage volume and revenue growth.

**Defining confirming and challenging evidence.** If the hypothesis is correct, accounts that experienced direct platform impact — those generating support tickets related to latency or errors — should show lower API usage in week four relative to their week two baseline. The revenue deceleration in week four (3.1% growth versus prior weeks) should be partially attributable to those accounts specifically, not distributed evenly across the base. If the hypothesis is wrong — if week four usage is flat across all accounts, or if the impacted accounts actually increase their usage — a different explanation is needed.

**Looking for that evidence.** The following query applies the hypothesis directly to the January dataset, comparing week four API token usage for

## incident-affected accounts against the unaffected base:

```

WITH affected_accounts AS (
  -- Accounts with latency- or error-related support tickets in week 3
  SELECT DISTINCT account_id
  FROM support_tickets
  WHERE week_start = DATE '2026-01-19'
      AND category IN ('error', 'latency')
),
usage_by_week AS (
  SELECT
    ue.account_id,
    ue.week_start,
    SUM(ue.tokens) AS tokens,
    CASE WHEN aa.account_id IS NOT NULL THEN 'affected' ELSE 'unaffected' END AS
group_label
  FROM usage_events ue
  LEFT JOIN affected_accounts aa ON ue.account_id = aa.account_id
  GROUP BY ue.account_id, ue.week_start, aa.account_id
),
pivoted AS (
  SELECT
    account_id,
    group_label,
    MAX(tokens) FILTER (WHERE week_start = DATE '2026-01-12') AS wk2_tokens,
    MAX(tokens) FILTER (WHERE week_start = DATE '2026-01-26') AS wk4_tokens
  FROM usage_by_week
  GROUP BY account_id, group_label
)
SELECT
  group_label,
  COUNT(*) AS accounts,
  ROUND(AVG(wk2_tokens), 0) AS avg_wk2_tokens,
  ROUND(AVG(wk4_tokens), 0) AS avg_wk4_tokens,
  ROUND(
    (AVG(wk4_tokens) - AVG(wk2_tokens)) / AVG(wk2_tokens) * 100
  , 1) AS usage_change_pct
  -- If hypothesis holds: affected accounts show usage decline or flat growth
  -- Unaffected accounts should show continued growth

```

```
FROM pivoted
WHERE wk2_tokens IS NOT NULL AND wk4_tokens IS NOT NULL
GROUP BY group_label
ORDER BY group_label;
```

The output either supports the hypothesis or challenges it. If the affected accounts show weaker week-over-week usage growth than the unaffected base, the hypothesis survives. If both groups show similar growth rates, the incident impact on behavior is not visible in the usage data — which suggests either that the affected accounts recovered quickly, or that the revenue deceleration has a different cause entirely.

**Two types of hypotheses.** The example above is a **diagnostic hypothesis** — it attempts to explain a pattern that has already occurred. A second type is a **predictive hypothesis**: given the signals visible now, what will happen in the next period if no intervention occurs? In the January context, a predictive hypothesis would state: if the week three incidents are not followed by proactive CS outreach to the three directly affected accounts, at least one of them will show declining token usage by week six, consistent with early disengagement.

Predictive hypotheses are more valuable to the business because they create a window for intervention. They also require more care: a predictive claim based on three data points is not statistically robust. It is directionally useful — enough to justify a targeted action — but not a reliable forecast. The distinction matters when communicating the finding to a stakeholder.

#### **14.4. Packaging the Insight**

A finding and an insight are not the same thing. A finding is what the data shows. An insight is what it means and what to do about it. The gap

between them is the gap between analysis and decision — and it is where most analytical work gets lost.

An insight has three components: what happened, why it matters, and what action it implies. All three must be present. A finding without the "why it matters" becomes a data dump. A finding with context but no recommended action becomes a discussion that ends without resolution. An insight that includes all three gives the recipient everything they need to make a decision.

**What happened.** State the observation in plain terms, anchored to specific metrics and time periods. "Customer Satisfaction Score dropped from 97 to 78 in week three following two platform incidents totaling 245 minutes of downtime, a 9% error rate, and three inbound support tickets."

**Why it matters.** Connect the observation to a business outcome, using the causal chain the framework makes visible. "CS score is a leading indicator of retention rate — the framework predicts that a sustained drop in CS score is followed by a decline in renewal probability within one to two quarters. The week four revenue deceleration is consistent with the early signal of customer friction. If the impacted accounts reduce usage, the ARR impact over the next quarter will reflect that contraction."

**What action it implies.** Name a specific action, an owner, and a timeframe. "Customer Success should initiate direct outreach to accounts 3 and 4 before the week four renewal pipeline review. The goal is to confirm platform stability, address any residual frustration, and identify whether any integration issues remain unresolved. The outcome to track is usage volume recovery to the week-two baseline by week six."

This three-part structure scales to any audience. For an executive, compress it to two sentences and lead with business impact. For a team lead, include the data chain. For an engineer, lead with the technical signal and the customer impact as context. The structure stays the same — what changes is the level of detail and the entry point.

**The insight brief.** For findings that require cross-functional action — which is most findings that matter — a brief format reduces the friction between analysis and decision. The brief has four fields:

Field	Content
Observation	What the data shows, in one to two sentences, with specific metrics cited.
Implication	Why it matters — the connection to a success metric, a driver, or a downstream outcome.
Recommended Action	A specific action, an owner, and a timeframe. Avoid passive constructions ("it may be worth exploring") in favor of direct ones ("CS should contact accounts 3 and 4 before Thursday").
Measurement	How you will know whether the action worked. Name the metric, the direction it should move, and the window in which you expect to see the change.

The brief is not a report. It is a decision prompt. Its purpose is to reduce the time between the finding and the action — and to make the action trackable, so the next analysis cycle can measure whether it worked.

That connection between action and measurement is what Chapter 15 addresses.

## Chapter 15. From Insight to Business Impact

An insight that does not change a decision has no business value. The test of analytics work is not whether the findings are statistically valid — it is whether they led to action, and whether that action moved the metric it was designed to move. Most analytics functions are good at producing findings. Far fewer close the loop: tracking whether the action was taken, measuring whether the downstream outcome moved, and feeding that result back into the next cycle of analysis. This last mile — from insight to business impact — is not a technical problem. It is an accountability problem. Solving it requires three things: a clear handoff between the insight and the action owner, a measurement framework that captures whether the action worked, and a regular review cycle that keeps findings connected to outcomes over time.

### ***15.1. The Insight-to-Action Gap***

The most common failure in analytics is not bad data or weak analysis. It is an insight that lands without an owner.

An insight brief reaches a leadership meeting, generates agreement that the finding is important, and then enters the organization without anyone whose job it is to track what happens next. The meeting ends. People return to their queues. The analysis is filed in a shared folder. Four weeks later, someone asks what happened with the week three CS score drop — and the honest answer is that no one followed up.

This is not a failure of intention. It is a structural gap. The analytics function has a clear mandate to find and communicate findings. The

stakeholder functions — CS, Engineering, GTM — have clear mandates to manage their teams and outcomes. Nobody has a mandate to track the connection between a specific finding and the specific action it triggered. That gap is where business impact gets lost.

Three structural problems create it.

**Insight without ownership.** A finding communicated to a team without naming a single decision-maker is a suggestion, not an action trigger. "CS should consider reaching out to affected accounts" diffuses accountability across everyone on the CS team, which in practice means no one acts. An insight with a named owner — "CS Manager responsible for accounts 3 and 4 initiates outreach by Thursday" — has a single point of accountability.

**Ownership without a measurement window.** An action without a defined timeframe and success metric is not a commitment — it is an intention. "We will reach out to the impacted accounts" leaves open when, and whether the outreach moved anything. "We will confirm usage recovery to week-two baseline by week six" sets a standard against which the action can be evaluated.

**Measurement without a feedback loop.** Even when actions are owned and measured, the results rarely return to the analytics team in a structured way. The CS manager knows whether the outreach worked. The analytics team does not. The learning from that outcome — whether the causal hypothesis was confirmed, whether the action was sufficient, whether the metric responded as predicted — never feeds back into the framework. The organization acts, but does not learn systematically.

Closing the insight-to-action gap requires addressing all three failures. The

mechanism that does this is the outcome brief.

## **15.2. Operationalizing an Outcome**

An outcome becomes operational the moment it has four properties: a metric that will move, an owner who is accountable for moving it, an action that is expected to cause the movement, and a measurement window that defines when the result will be assessed.

The insight brief introduced at the end of Chapter 13 provides the starting structure. The measurement field in that brief is what makes an insight operational. Without it, the insight is complete. With it, the insight becomes a commitment — a claim that if a specific action is taken, a specific metric will move in a specific direction within a specific time.

The January scenario makes this concrete. The week three incident cluster triggered a CS score drop to 78 and a week four revenue deceleration. The insight brief recommended CS outreach to accounts 3 and 4 before the week four renewal review. The operational outcome looks like this:

<b>Field</b>	<b>Content</b>
Metric	API token usage per week for accounts 3 and 4 (current: below week-two baseline)
Owner	CS Manager for Enterprise accounts
Action	Direct outreach to accounts 3 and 4 to confirm platform stability, address residual friction, and resolve any open integration issues
Success Threshold	Token usage returns to within 10% of week-two baseline for both accounts by week six (DATE '2026-02-09')

Field	Content
Fallback Trigger	If usage has not recovered by week five, escalate to a technical review call with Engineering and the account's primary technical contact

This brief does three things. It converts an analytical finding into a managed action. It creates a success standard that does not require judgment at assessment time — either the metric moved to within the threshold or it did not. And it includes a fallback trigger, which acknowledges that first-order actions sometimes need to be reinforced, and defines in advance when escalation is appropriate.

The outcome brief is not a guarantee of business impact. It is a structure that makes impact measurable and accountability visible. Whether the action worked is a separate question — the one the measurement window answers.

### ***15.3. Measuring Downstream Business Impact***

Measuring whether an action produced business impact requires a comparison: what the outcome metric did after the action, versus what it would have done without it. That counterfactual — what would have happened without the intervention — is the central challenge of impact measurement.

Three approaches are available, each appropriate at different levels of data maturity and organizational complexity.

**Before-and-after comparison.** The simplest approach compares the metric in the period immediately before the action to the metric in the period immediately after. For the January scenario, this means comparing

the weekly token usage of accounts 3 and 4 in weeks three and four — the pre-intervention period — to weeks five and six — the post-intervention period.

```

WITH account_usage AS (
  SELECT
    account_id,
    week_start,
    SUM(tokens) AS tokens
  FROM usage_events
  WHERE account_id IN (3, 4)
  GROUP BY account_id, week_start
)
SELECT
  account_id,
  MAX(tokens) FILTER (WHERE week_start = DATE '2026-01-12') AS baseline_wk2,
  MAX(tokens) FILTER (WHERE week_start = DATE '2026-01-19') AS incident_wk3,
  MAX(tokens) FILTER (WHERE week_start = DATE '2026-01-26') AS post_incident_wk4
  -- After intervention: compare week 5 and 6 values once available
  -- Recovery target: return to within 10% of baseline_wk2
FROM account_usage
GROUP BY account_id;

```

Before-and-after comparison is easy to compute and easy to communicate. Its limitation is that other factors may have changed between the two periods — a product update, a pricing change, a market event — making it difficult to attribute the metric movement to the specific action rather than the broader environment.

**Matched comparison.** A more rigorous approach compares the treated accounts — those that received the intervention — to a set of comparable accounts that did not. If accounts 3 and 4 received CS outreach and recovered their usage by week six, and accounts that experienced similar incidents in a prior period without CS outreach did not recover at the same

rate, the difference in recovery rate provides an estimate of the intervention's contribution.

Matched comparison requires a larger dataset and a historical baseline of similar events. It is the appropriate method when the business has been running long enough to accumulate comparable cases and when the impact being measured is large enough to justify the analytical investment.

**Contribution framing.** In most operational contexts, the right question is not whether the action caused the outcome in a strict causal sense — it is whether the action contributed to an outcome that would not have occurred without it. Contribution framing acknowledges that business outcomes have multiple causes and focuses on whether the specific action accelerated, prevented, or amplified a particular outcome.

For the January scenario, the appropriate contribution framing is: did the CS outreach to accounts 3 and 4 prevent a renewal risk that would otherwise have materialized in the Q1 renewal cycle? The answer to that question does not require statistical proof of causation. It requires a clear chain of evidence: the accounts had degraded signals, the outreach occurred, the signals recovered, and the renewal conversation proceeded without escalation. That chain — visible in the data and confirmed in the account record — is sufficient to attribute the outcome to the action.

Contribution framing is the most practical approach for operational analytics work, because most actions taken in response to insights are not controlled experiments. They are operational decisions made under uncertainty, by people with multiple competing priorities, in environments where running a true control group is not feasible. Acknowledging that

limitation does not reduce the value of the measurement. It makes the measurement honest.

### ***15.4. The Analytics Review Cycle***

Impact measurement is not a one-time activity. It is a recurring process that connects the current cycle of analysis to the previous one — closing the loop between what was recommended, what was done, and what happened as a result.

An effective analytics review cycle operates at three tempos.

**Weekly:** Review current metric status against the framework. Identify anomalies that have emerged since the previous cycle. Confirm that actions from the previous week's insight briefs have been initiated. Flag any outcome thresholds that will close in the coming week. This review is operational — its purpose is to maintain visibility into in-flight actions and ensure nothing has slipped.

**Monthly:** Review whether the actions taken over the past month moved the metrics they were designed to move. Compare actual outcomes against the success thresholds set in outcome briefs. Where thresholds were met, document what worked and why. Where thresholds were missed, investigate whether the action was taken as specified, whether the timeframe was realistic, or whether the causal hypothesis underlying the insight was wrong. This review is diagnostic — its purpose is to separate effective interventions from ineffective ones.

**Quarterly:** Review the overall performance of the analytics function against the business outcomes it was designed to support. Which success

metrics moved? Which driver-level actions had the most impact? Which findings generated the most significant business outcomes? Where did the framework reveal connections that were not visible before? This review is strategic — its purpose is to assess whether the analytics investment is producing returns that justify its cost, and to identify where the next cycle of investment should go.

The quarterly review is also where the framework itself is updated. If a driver turned out to predict a success metric more reliably than expected, that connection should be strengthened. If a metric consistently fails to move in response to the actions designed to influence it, the causal model underlying that connection should be reconsidered. The framework is not static. It is a working model of how the business creates value — and like any model, it improves when it is tested against outcomes and updated when the evidence requires it.

The most valuable output of a mature analytics function is not a dashboard or a report. It is a track record: a documented history of insights that triggered actions, actions that moved metrics, and metrics that produced business outcomes. That track record is the evidence that the gap between data investment and business value has been closed — not in theory, but in practice, quarter after quarter.

That is what the framework is built to produce.

# SECTION IV: TECHNICAL REFERENCE

This section (Chapters 16-18) provides guidance and practical references to accelerate daily analytics work. Chapter 16 provides foundational SQL and Python knowledge for anyone getting started with data analysis. Chapter 17 provides advanced SQL references and Python snippets for data modeling, automation, and orchestration. Chapter 18 is a consolidated dbt reference covering the core concepts and commands used in the case studies in the next section.

## Chapter 16. SQL and Python Foundations

SQL and Python are the two languages data teams speak every day — and while writing them is not a requirement for business leaders, understanding their fundamentals creates better conversations with the data team, sets realistic expectations, and sharpens judgment about what is and is not possible with data.

SQL's core building blocks — joining tables, filtering results, and advanced analytics functions — are business logic expressed in code, and the language of choice for querying, aggregating, and transforming structured data at scale. Python extends that capability into data science and predictive analytics — its core data structures and libraries are the foundation for the statistical models, forecasts, and machine learning applications that take analysis beyond what happened and into what will happen next. Together, they form a complete analytical toolkit — SQL addresses descriptive and diagnostic questions, while Python addresses predictive and prescriptive ones.

### ***16.1. The Core Concept: Tables Are Sets***

Every table in a database is a collection of records — essentially a spreadsheet. The Customers table holds customer information. The Orders table contains purchase history. The Products table stores product details. The power comes from connecting these tables to answer business questions.

Imagine two overlapping circles in a Venn diagram. One circle contains all California customers. The other contains all customers who bought Product

A. The overlap — customers in both circles — represents California customers who bought Product A. That overlap is what data people call an "intersection" or "inner join."

This Venn diagram mental model is how database queries work. Every time an analyst talks about joining tables, they're essentially combining circles to answer specific questions.

## ***16.2. The Three Types of Joins***

**Inner joins** find only the overlap — records that match in both tables. "Show me customers in California who bought Product A" returns only those customers who meet both criteria. This is the most common join for focused analysis.

**Left joins** keep everyone from the first table and add optional information from the second. "Show me all California customers and whether they bought Product A" returns all California customers — some with Product A purchase data and others without. This is critical for understanding who doesn't have something — inactive customers, products without sales, regions without coverage.

**Full outer joins** grab everyone from both tables. "Show me customers in California or customers who bought Product A" returns everyone in either circle. Less common, but useful for the complete picture across multiple conditions.

When an analyst asks which join to use, they're really asking: only matches, everyone from one side, or everyone from both sides?

### **16.3. Grouping and Aggregating: Turning Data into Insights**

Grouping organizes data into categories for calculating summaries. "How many customers by state?" requires grouping all customers by their state and counting each group. "What's total revenue by product?" requires grouping all orders by product and summing the revenue.

This is where data becomes insight. Raw data shows individual transactions. Grouped data shows patterns: which products sell best, which regions are growing, which customer segments are most valuable.

Teams use aggregate functions constantly:

- **COUNT** reveals how many (customers, orders, products)
- **SUM** adds things up (total revenue, total quantity sold)
- **AVERAGE** finds the middle (average order value, average customer lifetime)
- **MAX/MIN** find extremes (largest order, smallest purchase, newest customer)

### **16.4. Filtering: Getting to the Data That Matters**

Filters narrow results to what actually matters. "California customers" filters by state. "Orders over one hundred dollars" filters by amount. "Customers who joined last quarter" filters by date.

The trick is knowing when to filter before grouping versus after. Filtering before grouping is like deciding which ingredients to use. Filtering after grouping is like tasting the dish and deciding if it's spicy enough. Both matter, but they happen at different stages.

Teams also use conditional logic to categorize data on the fly. CASE statements assign each customer to the right tier based on spending — like an automated sorting hat for data.

## **16.5. SQL Reference Functions**

### **Date and Time Functions**

- Converting UTC Timestamps
- Extracting Date Parts (Day, Week, Month, Quarter, Year)
- Timezone Conversions
- Date Arithmetic

### **Aggregation Functions**

- SUM, COUNT, AVG, MIN, MAX
- GROUP BY and HAVING

### **Window Functions**

- ROW\_NUMBER, RANK, DENSE\_RANK
- Running Totals and Moving Averages

### **String Functions**

- CONCAT, SUBSTRING, UPPER, LOWER
- Pattern Matching

### **Conditional Logic**

- CASE Statements

- COALESCE and NULL Handling

## **16.6. SQL Best Practices**

Writing SQL that works is one thing. Writing SQL that performs well and scales with data growth is another. These practices come from watching queries that ran fine at gigabyte scale grind to a halt at terabyte scale.

**Read What You Need:** The easiest way to speed up queries is reducing data touched. In big data environments, queries can time out if tables have 800 columns — `SELECT *` becomes impossible. Filter early with meaningful `WHERE` conditions that eliminate irrelevant data before processing. Only select needed columns, not everything. Use partition filters like date ranges to skip unnecessary data segments entirely.

**Make Joins Work for You:** Joins often account for the heaviest SQL work. Use the correct join type — don't use `LEFT JOIN` when `INNER JOIN` suffices. Ensure join keys are indexed or well-distributed. Apply filters before joins to shrink tables. Pre-aggregate large tables before joining when possible. Smaller, cleaner join inputs create dramatically faster queries.

**Keep Aggregations Lean:** Group only on columns that matter — avoid high-cardinality fields unless necessary. Use approximate functions when exact results aren't crucial. Perform early summarization to reduce final result size. These choices reduce compute time and memory usage, keeping queries performant as data grows.

**Use Window Functions Wisely:** Window functions are powerful but can slow execution if used blindly. Keep partitions small by filtering before applying window operations. Reuse window definitions rather than

repeating logic. Filter or aggregate data before window operations when feasible.

**Balance Clarity and Speed with CTEs:** Common Table Expressions improve clarity but may be reprocessed multiple times depending on the database engine. Use CTEs for readability, but consider temporary tables for repeated, expensive logic. Balance readability with performance — understandable queries matter, but not if they take hours to run.

**Give the Optimizer What It Needs:** Databases rely on metadata to choose execution plans. Keep table statistics current so the optimizer knows sizes, distributions, and patterns. Use composite indexes for multi-column filters. Drop unused indexes — they slow writes without query benefits.

**Work in Sets, Not Rows:** SQL engines excel at set-based operations, not row-by-row processing. Avoid cursor-style patterns that process one record at a time. Use CASE, COALESCE, and set operations instead of procedural logic.

**Watch Performance as Data Grows:** Check execution plans regularly to understand how queries actually run. Monitor performance over time — data growth impacts previously fast queries. Use database-specific features like clustering keys, materialized views, or caching when available.

SQL best practices aren't about perfect queries — they're about queries that stay fast, remain understandable, and adapt as data grows. The difference between a query that works and one that works well often determines whether analytics delivers value quickly or becomes a bottleneck.

## ***16.7. Python Foundations for Data Pipelines***

Python is the second language data teams speak every day. Where SQL queries and transforms data inside a database, Python moves data between systems — reading files, calling APIs, inserting records, and orchestrating what runs when. Understanding its fundamentals creates better conversations about what pipelines actually do, where they can fail, and why they are built the way they are.

### *16.7.1. The Core Concept: Code Runs Top to Bottom*

Python executes line by line from top to bottom. That sequential nature is exactly why it works well for pipelines — read the file, then transform the records, then write to the database, in that order. When something goes wrong, the error message tells you which line failed, making problems straightforward to locate.

Two things make Python feel different from SQL. First, indentation is not cosmetic — it is syntax. Everything indented inside a function or loop belongs to it. Change the indentation and you change what the code does. Second, Python uses libraries — collections of pre-built tools you import at the top of the file. The same way SQL has built-in functions like `SUM` and `DATE_TRUNC`, Python has libraries like `json` for reading JSON files and `psycopg2` for connecting to PostgreSQL.

### *16.7.2. Reading Files and Connecting to Databases*

Every bronze layer DAG in this stack reads from a file and writes to a database. The pattern is always the same — open the file, parse the contents, open a database connection, insert the records, close everything cleanly.

```

import json
import psycopg2
from pathlib import Path

# Path() builds file paths that work on any operating system
# On a Mac: /opt/airflow/config-files/model_configuration.json
CONFIG_DIR = Path('/opt/airflow/config-files')
filepath = CONFIG_DIR / 'model_configuration.json'

# Open and read a JSON file
# 'with' automatically closes the file when the block ends – even if an error occurs
with open(filepath, 'r') as f:
    data = json.load(f) # parses the JSON text into a Python dictionary

# data is now a Python dictionary – navigate it with keys
models = data.get('models', [])
# .get('models', []) means: give me data['models'], and if it doesn't exist return
# []
# safer than data['models'] which would crash if the key is missing

print(f"Found {len(models)} models")
# f-strings embed variables directly into text – len(models) is replaced by the
# count

# Connect to PostgreSQL
DB_CONFIG = {
    'host': 'postgres', # container name on the Docker network
    'port': '5432',
    'database': 'resource_utilization',
    'user': 'postgres',
    'password': 'postgres'
}

conn = psycopg2.connect(**DB_CONFIG) # ** unpacks the dictionary as keyword
arguments
cursor = conn.cursor() # cursor is what executes SQL statements

# Always close the connection when done – open connections accumulate
# and eventually exhaust the database's connection limit

```

```

cursor.close()
conn.close()

```

### 16.7.3. Building Records and Inserting Data

Once a file is read, the data gets shaped into records — a list of tuples, one tuple per row — before being written to the database. This is the pattern used in every bronze DAG in the stack.

```

from psycopg2.extras import execute_values
from datetime import datetime

# Build records as a list of tuples
# Each tuple matches the column order in the INSERT statement
records = []
snapshot_date = datetime.now().date() # today's date - used as a watermark

for model in models:
    record = (
        model['publisher_name'],      # first column
        model['model_display_name'],  # second column
        model['model_variant'],       # third column
        model['is_open_source'],      # fourth column
        snapshot_date,                # fifth column - added by the pipeline, not
the file
        filepath.name                  # sixth column - which file this came from
    )
    records.append(record)

# Insert all records in one call
# execute_values is dramatically faster than inserting one row at a time
# It sends all rows to the database in a single round-trip
insert_sql = """
    INSERT INTO raw_bronze.config_model_dimensions (
        publisher_name, model_display_name, model_variant,
        is_open_source, snapshot_date, source_file
    ) VALUES %s
    ON CONFLICT (model_variant, snapshot_date) DO NOTHING

```

```

"""
# ON CONFLICT DO NOTHING: if this model + date already exists, skip it silently
# This makes the DAG safe to re-run – running it twice won't create duplicates

execute_values(cursor, insert_sql, records)
inserted = cursor.rowcount    # how many rows were actually written

conn.commit()    # make the changes permanent – without this, nothing is saved
cursor.close()
conn.close()

print(f"Inserted {inserted} records")

# xcom_push shares a value with other tasks in the same DAG run
# The validate task can then read this count to confirm what was loaded
context['task_instance'].xcom_push(key='model_config_count', value=inserted)

```

#### 16.7.4. Validating Data in the Pipeline

Every DAG in this stack has a validation task that runs after loading. Validation is not optional — it is the pipeline’s own quality check, catching problems before downstream models run against bad data.

```

def validate_bronze_data(**context):
    """Check data quality after loading – called by the validate task"""

    conn = psycopg2.connect(**DB_CONFIG)
    cursor = conn.cursor()
    issues = []    # collect problems – if this list is empty at the end, all checks
    passed

    # Check 1: every model loaded today has at least one region assigned
    # If a model has no regions it cannot be deployed – that is a data problem
    cursor.execute("""
        SELECT cmd.model_variant
        FROM raw_bronze.config_model_dimensions cmd
        WHERE snapshot_date = CURRENT_DATE
        AND NOT EXISTS (

```

```

        SELECT 1
        FROM raw_bronze.config_model_region_availability mra
        WHERE mra.model_variant = cmd.model_variant
              AND mra.snapshot_date = CURRENT_DATE
    )
    """)
    orphan_models = cursor.fetchall() # fetchall returns a list of tuples
    if orphan_models:
        issues.append(f"{len(orphan_models)} models have no region availability")

    # Check 2: confirm records were actually loaded today
    cursor.execute("""
        SELECT COUNT(*) FROM raw_bronze.config_model_dimensions
        WHERE snapshot_date = CURRENT_DATE
    """)
    count = cursor.fetchone()[0] # fetchone returns one tuple - [0] gets the
    first value

    cursor.close()
    conn.close()

    # Print results - Airflow captures this output in the task logs
    print(f"Records loaded today: {count}")

    if issues:
        for issue in issues:
            print(f"WARNING: {issue}")
    else:
        print("All validation checks passed")

    return len(issues) == 0 # True = all good, False = problems found

```

### 16.7.5. Python Reference Functions

#### Imports — what each library does

- `import json` — read and parse JSON files into Python dictionaries
- `import psycopg2` — connect to PostgreSQL and execute SQL from Python

- `from psycopg2.extras import execute_values` — fast bulk inserts (all rows in one call)
- `from pathlib import Path` — file paths that work on Mac, Linux, and Windows
- `from datetime import datetime, timedelta` — dates, times, and arithmetic on them

## File and data operations

- `Path('/opt/airflow') / 'file.json'` — build a file path safely
- `open(path, 'r') as f` — open a file for reading; closes automatically when done
- `json.load(f)` — parse JSON file into a Python dictionary
- `data.get('key', [])` — read a dictionary key safely; returns default if missing
- `len(list)` — count items in a list
- `list.append(item)` — add one item to the end of a list
- `f"text {variable}"` — embed a variable's value directly into a string

## Database operations

- `psycopg2.connect(**DB_CONFIG)` — open a connection to PostgreSQL
- `conn.cursor()` — create a cursor to execute SQL statements
- `cursor.execute(sql)` — run a SQL query
- `cursor.fetchone()` — get one result row as a tuple
- `cursor.fetchall()` — get all result rows as a list of tuples

- `execute_values(cursor, sql, records)` — bulk insert a list of tuples
- `cursor.rowcount` — how many rows the last statement affected
- `conn.commit()` — make all changes permanent
- `cursor.close() / conn.close()` — release the connection back to the pool

### 16.7.6. YAML Foundations

YAML is the configuration language of the modern data stack. dbt uses it to define models, tests, sources, and connection profiles. Docker uses it to define which containers to run. Airflow uses it for environment settings. The ability to read and edit YAML confidently — without accidentally breaking indentation or misplacing a dash — is a practical skill that saves hours of debugging.

YAML is built on three rules. Indentation defines hierarchy — a key indented under another key belongs to it. Colons separate keys from values. Dashes mark list items. Everything else follows from these three.

### 16.7.7. Reading Indentation and Hierarchy

```
# This is a comment – ignored by all tools
# YAML uses 2-space indentation to show that something belongs to something else
# Tabs are not allowed – use spaces only

# A simple key: value pair
name: resource_utilization

# A nested key – 'dev' belongs to 'outputs', 'outputs' belongs to
# 'resource_utilization'
resource_utilization:
  outputs:
    dev:
      type: postgres
```

```

    host: postgres
    port: 5432

# The same structure as a sentence:
# "resource_utilization has outputs, one of which is dev,
# which is a postgres connection on host postgres port 5432"

# A list – each item starts with a dash
clean-targets:
  - target
  - dbt_packages

# A list of objects – each item starts with a dash followed by key: value pairs
models:
  - name: stg_customer_details      # first item in the list
    description: "Customer data"
    columns:
      - name: account_id           # nested list inside the first item
        tests:
          - not_null
          - unique

  - name: stg_token_usage           # second item in the list
    description: "Token usage data"

```

### 16.7.8. Reading a `profiles.yml`

The `profiles.yml` tells dbt how to connect to the database. It defines one or more connection environments and which one to use by default. This is the file you edit when changing a host, password, or schema.

```

resource_utilization:      # project name – must match the name in dbt_project.yml

  outputs:                 # all available connection environments live here

    dev:                   # environment name – used when running: dbt run --target dev
    type: postgres         # database type

```

```

host: postgres      # 'postgres' = the Docker container name on the network
                   # change to 'localhost' when running dbt from your Mac
port: 5432
user: postgres
password: postgres
dbname: resource_utilization
schema: dbt_dev    # base schema prefix – dbt appends model layer names to
this
threads: 4        # how many models dbt builds in parallel

docker:           # second environment – used when running dbt inside a
container
  type: postgres
  host: postgres
  port: 5432
  user: postgres
  password: postgres
  dbname: resource_utilization
  schema: dbt_dev
  threads: 4

target: dev       # which environment is active by default

```

### 16.7.9. Reading a `sources.yml`

The `sources.yml` declares the raw Bronze tables that staging models read from. Every table a staging model references with `{{ source('raw_bronze', 'table_name') }}` must be declared here. It also defines freshness thresholds — how old data can be before dbt raises a warning or error.

```

version: 2        # always 2 for modern dbt

sources:
  - name: raw_bronze    # the source group name – first argument to {{ source()
}}
    schema: raw_bronze  # the actual PostgreSQL schema
    description: "Bronze layer – raw ingested data from source systems"

```

```

freshness:          # default freshness applied to all tables in this
source
  warn_after: {count: 2, period: day}  # warn if data is older than 2 days
  error_after: {count: 3, period: day}  # fail if data is older than 3 days

loaded_at_field: loaded_at  # which column dbt checks to determine freshness

tables:
  - name: customer_details          # table name – second argument to {{
source() }}
  description: "Customer records"
  freshness:                        # overrides the source-level freshness for
this table
    warn_after: {count: 30, period: day}
    error_after: {count: 60, period: day}
  columns:
    - name: account_id
      tests:
        - not_null                  # dbt generates: SELECT COUNT(*) WHERE
account_id IS NULL
        - unique                    # dbt generates: SELECT account_id, COUNT(*)
HAVING COUNT(*) > 1

    - name: inference_user_token_usage_proprietary
      description: "Token usage events for proprietary models"
      freshness:
        warn_after: {count: 2, period: hour}  # high-frequency table – warn
after 2 hours
        error_after: {count: 4, period: hour}
      columns:
        - name: request_id
          tests:
            - not_null
            - unique

```

### 16.7.10. Reading a `schema.yml`

The `schema.yml` documents staging models and defines column-level data

quality tests. Think of it as a contract — "this column must never be null, this column only accepts these values." When `dbt test` runs, `dbt` generates SQL from these declarations and runs it against the database.

```

version: 2

models:
  - name: stg_customer_details      # must match the .sql file name exactly
    description: "Cleaned customer account data. One row per account."

    columns:
      - name: account_id
        description: "Unique account identifier"
        tests:
          - unique                  # no two rows share the same account_id
          - not_null                # account_id is never empty

      - name: segment
        tests:
          - not_null
          - accepted_values:        # only these values are valid in this column
              values: ['Strategic', 'Commercial', 'Unknown']

      - name: company_name
        tests:
          - not_null:
              config:
                severity: warn      # warn instead of failing the run
                                   # use warn when the issue is important but
                                   # not blocking

```

### 16.7.11. Reading a `dbt_project.yml`

The `dbt_project.yml` is the master configuration for the entire `dbt` project. It defines where files live, how each layer materializes, and what runs after each model builds. The `+` prefix on a key means "apply this to all models in

this folder and every folder underneath it."

```

name: 'resource_utilization'    # project name – must match profiles.yml
version: '1.0.0'
config-version: 2

profile: 'resource_utilization' # which profile to use from profiles.yml

# Where dbt looks for each type of file
model-paths:    ["models"]
seed-paths:     ["seeds"]
test-paths:     ["tests"]
macro-paths:    ["macros"]
snapshot-paths: ["snapshots"]

# Folders that dbt cleans when you run: dbt clean
clean-targets:
  - "target"      # compiled SQL and run artifacts
  - "dbt_packages" # installed packages – regenerated by: dbt deps

models:
  resource_utilization:      # must match the project name above

  staging:
    +materialized: view      # + means: apply to ALL models in the staging folder
    +schema: staging_silver # models build into the staging_silver schema in
                             PostgreSQL

  intermediate:
    +materialized: table
    +schema: intermediate_silver
    +post-hook:
      - "ANALYZE {{ this }}" # SQL that runs after each model builds
                              # {{ this }} is replaced by the actual table name at
runtime

  marts:
    +materialized: table
    +schema: mart_gold

```

```

+grants:
  select: ['reporter', 'analyst', 'data_engineer'] # who can query mart
tables
+post-hook:
  - "ANALYZE {{ this }}"

seeds:
  resource_utilization:
+schema: seeds          # seed tables land in the seeds schema

```

### 16.7.12. Reading a *docker-compose.yml*

Docker Compose uses YAML to define the entire local data stack — which containers run, how they connect to each other, and where data is stored. You will read this file when adding a service, changing a port, or debugging a connection.

```

version: '3.8'      # Docker Compose file format version

services:          # each key under services is one container

  postgres:       # service name – other containers connect to this using
'postgres' as the host
  image: postgres:13      # which Docker image to run
  container_name: postgres14 # name shown in: docker ps
  environment:          # environment variables passed into the container
    POSTGRES_USER: postgres
    POSTGRES_PASSWORD: postgres
    POSTGRES_DB: resource_utilization
  ports:
    - "5432:5432"      # mac_port:container_port – maps your Mac's port to the
container's port
  volumes:
    - postgres-data:/var/lib/postgresql/data # named volume – data persists
across restarts
    - ./init-db.sql:/docker-entrypoint-initdb.d/init-db.sql # file from your Mac
mounted in

```

```

networks:
  - data-stack-network    # containers on the same network can reach each other
by service name
  healthcheck:            # how Docker knows the container is ready to accept
connections
  test: ["CMD-SHELL", "pg_isready -U postgres"]
  interval: 10s
  retries: 5

airflow-scheduler:
  build:
    context: ./airflow    # build a custom image from the Dockerfile in
./airflow/
    dockerfile: Dockerfile
  container_name: airflow-scheduler
  environment:
    AIRFLOW__CORE__EXECUTOR: CeleryExecutor
    AIRFLOW__DATABASE__SQL_ALCHEMY_CONN:
postgresql+psycopg2://postgres:postgres@postgres/airflow
    # AIRFLOW__CORE__EXECUTOR is an Airflow config key written as an environment
variable
    # double underscores (__) replace the dots and section separators in
airflow.cfg
  volumes:
    - ./airflow/dags:/opt/airflow/dags    # your DAG files on Mac → container's
dags folder
    - ./dbt:/opt/airflow/dbt            # dbt project → accessible from
inside Airflow
  networks:
    - data-stack-network
  depends_on:
    - postgres    # Docker starts postgres before this container

volumes:    # named volumes – Docker manages storage, data survives
container restarts
  postgres-data:
    driver: local

networks:
  data-stack-network:

```

```
driver: bridge # bridge = containers can communicate; isolated from your Mac's
network
```

### 16.7.13. YAML Best Practices

**Use spaces, never tabs.** YAML parsers reject tabs entirely. If a YAML file throws an unexpected error, tab characters hiding as spaces are the first thing to check. Most editors show invisible characters — enable that setting when editing YAML.

**Keep indentation consistent.** Two spaces per level is the universal standard across dbt, Docker, and Airflow configuration files. Mixing two and four spaces in the same file causes parse errors that can be difficult to locate.

**Quote strings that contain special characters.** Colons, hash signs, curly braces, and square brackets have special meaning in YAML. When a value contains any of these, wrap it in single or double quotes: `password: 'db:secret'` not `password: db:secret`. Airflow environment variable values and dbt descriptions frequently need quoting.

**Every list item starts with a dash.** If you are defining multiple tables, multiple columns, or multiple tests, each item must start with ``- ``. A missing dash merges the second item into the first silently — the file parses without error but produces wrong behaviour.

**Comments are free — use them.** Anything after `#` is ignored. The dbt and Docker files in this stack use comments heavily to explain non-obvious values. Annotate freshness thresholds, explain why a severity is set to warn, note which port maps to which tool. The next person reading the file

— often you, six months later — will be grateful.

#### 16.7.14. Airflow DAG Anatomy

An Airflow DAG is a Python file that defines a workflow — which tasks to run, in what order, and on what schedule. Every DAG in this stack follows the same four-part structure, and recognising that structure makes reading or debugging any DAG straightforward.

```

from airflow import DAG
from airflow.operators.bash import BashOperator
from airflow.operators.python import PythonOperator
from datetime import datetime, timedelta

# — Part 1: Default arguments


---


# These apply to every task in the DAG unless a task overrides them
default_args = {
    'owner': 'data_engineering',
    'depends_on_past': False,      # do not wait for yesterday's run to succeed
    'email_on_failure': True,     # send an alert email if a task fails
    'email_on_retry': False,     # do not email on each retry attempt
    'retries': 2,                # retry a failed task twice before giving up
    'retry_delay': timedelta(minutes=5), # wait 5 minutes between retries
}

# — Part 2: DAG definition


---


# The DAG object holds all metadata about the pipeline
dag = DAG(
    'silver_layer_staging',      # unique pipeline name – shown in the
    Airflow UI
    default_args=default_args,
    description='Run dbt staging models every 6 hours',
    schedule_interval='0 */6 * * *', # cron expression: at minute 0, every 6 hours

```

```

    start_date=datetime(2026, 2, 17), # when the pipeline became active
    catchup=False,                    # do not run missed historical runs on
startup
    tags=['silver', 'dbt', 'staging'], # labels in the Airflow UI – useful for
filtering
)

# Path constants defined once at the top – change them here, they update everywhere
DBT_PROJECT_DIR = '/opt/airflow/dbt'
DBT_PROFILES_DIR = '/opt/airflow/dbt'
DB_CONN         = 'postgresql://postgres:postgres@postgres/resource_utilization'

# ——— Part 3: Task definitions
—————

# BashOperator runs a shell command inside the Airflow container
# Used for dbt commands – dbt is a command-line tool, not a Python library
run_dbt_staging = BashOperator(
    task_id='run_dbt_staging',          # unique name within this DAG
    bash_command=f"""
        cd {DBT_PROJECT_DIR} && \\\
        /usr/local/airflow/dbt_venv/bin/dbt run \\\
        --select stg_customer_details \\\
            stg_token_usage_proprietary \\\
            stg_token_usage_open_source \\\
        --profiles-dir {DBT_PROFILES_DIR} || exit 1
    """,
    # || exit 1 means: if dbt returns an error, fail this task immediately
    # Without it, Airflow might not detect the failure
    dag=dag,
)

run_dbt_tests = BashOperator(
    task_id='run_dbt_tests',
    bash_command=f"""
        cd {DBT_PROJECT_DIR} && \\\
        /usr/local/airflow/dbt_venv/bin/dbt test \\\
        --select staging \\\
        --profiles-dir {DBT_PROFILES_DIR} || exit 1
    """
)

```

```

    """
    dag=dag,
)

validate_staging = BashOperator(
    task_id='validate_staging',
    bash_command=f"""
        psql {DB_CONN} -c "
            SELECT table_name
            FROM information_schema.tables
            WHERE table_schema = 'staging_silver'
            ORDER BY table_name;
        """
    dag=dag,
)

# PythonOperator runs a Python function – used when logic is too complex for a shell
# command
def load_configuration(**context):
    # **context gives this function access to Airflow metadata
    # context['task_instance'].xcom_push() shares data with other tasks
    pass

load_task = PythonOperator(
    task_id='load_configuration',
    python_callable=load_configuration, # the function to call
    provide_context=True,               # pass the context dictionary into the
    function                             function
    dag=dag,
)

# ——— Part 4: Task dependencies

—————

# >> means "this must succeed before that starts"
# run_dbt_staging must pass before run_dbt_tests starts
# run_dbt_tests must pass before validate_staging starts
run_dbt_staging >> run_dbt_tests >> validate_staging

```

```
# Tasks can also fan out – run two tasks in parallel after one completes:
# check_files >> [load_config, load_regions] >> validate
# load_config and load_regions run at the same time; validate waits for both
```

## Cron schedule quick reference

<code>schedule_interval='0 */6 * * *'</code>	every 6 hours – midnight, 6am, noon, 6pm
<code>schedule_interval='0 1 * * *'</code>	daily at 1am
<code>schedule_interval='0 2 * * *'</code>	daily at 2am
<code>schedule_interval='0 3 * * 1'</code>	every Monday at 3am
<code>schedule_interval='0 4 1 * *'</code>	1st of every month at 4am
<code>schedule_interval=None</code>	manual trigger only – no automatic schedule

### 16.7.15. Python Best Practices for Data Pipelines

**Always close database connections.** Every `psycopg2.connect()` call must be followed by `conn.close()`. Connections left open accumulate and eventually exhaust the database’s connection limit, causing failures across all pipelines simultaneously. Close the cursor first, then the connection.

**Use `execute_values` for bulk inserts.** Inserting one row at a time with a loop sends one database round-trip per row — thousands of rows means thousands of round-trips. `execute_values` sends all rows in a single call. The difference is seconds versus minutes at scale.

**Commit only after all inserts succeed.** `conn.commit()` makes changes permanent. Committing after each row means partial data survives if the script crashes halfway through. Commit once at the end when everything has succeeded. If the script fails before commit, no data is written and the DAG can be re-run safely.

**One function per task.** Every task in the bronze DAG has exactly one function — `load_model_configuration`, `load_model_region_availability`,

`validate_bronze_data`. Each function does one thing. When a task fails, the function name in the error log tells you immediately what went wrong without reading the entire DAG file.

**Define path constants at the top.** `DBT_PROJECT_DIR`, `DB_CONN`, and `CONFIG_DIR` are defined once at the top of each DAG file. Changing a path or connection string means editing one line, not hunting through every `BashOperator` command.

**Print progress inside task functions.** Airflow captures everything printed to stdout and displays it in the task logs. `print(f"Inserted {inserted} records")` gives you visibility into what the pipeline actually did when you review a completed run. Silent pipelines are difficult to audit.

## Chapter 17. Advanced SQL and Python Reference

Where Chapter 16 builds conceptual fluency, this reference chapter provides the hands-on patterns that data teams use to solve real analytical problems — date and time manipulation, string cleaning and parsing, and window functions for rankings, running totals, and period-over-period comparisons. These are not academic exercises — every pattern is drawn from the same SaaS business context used throughout the book, making them immediately applicable to the metrics and analyses the framework depends on.

Window functions in particular unlock a category of analysis that simple aggregations cannot — comparing each account’s current usage to its own history, ranking accounts within tiers, and calculating growth rates across periods — all without leaving SQL. Treat this as a working reference to return to whenever a business question requires a pattern you have not used before — the goal is not memorization but recognition.

*This reference guide provides practical SQL examples for common data transformation tasks in PostgreSQL syntax. To use this reference in any other SQL version, upload the code to any AI tool and ask to convert. Each section includes working queries with inline comments explaining datatypes and outputs.*

---

### 17.1. Date and Time Functions

**NOTE**

When to use these functions

Date and time functions are needed at every layer of the

medallion architecture — but the use cases differ by layer.

- **Bronze layer:** parse raw timestamps from API responses or event streams into usable date types.
- **Silver layer:** truncate timestamps to the right grain (day, week, month) and calculate date differences for SLA checks and late-arrival detection.
- **Gold layer:** extract year, quarter, and month for time-based aggregations and period-over-period comparisons in dashboards.

Use this code snippet to work with timestamps, dates, and time dimensions for reporting and analytics. This code converts a single timestamp into multiple time formats and dimensions. Use this to understand how to extract dates, weeks, months, quarters, and years from timestamps.

```
WITH time_dimensions AS (
  SELECT
    '2025-12-31T13:01:49.459Z'::timestampz AS ts_utc,
    1735651309.459 AS epoch_seconds,
    'Europe/London' AS timezone,
    '2025-12-31 13:01:49' AS ts_string
)
SELECT
  -- varchar: '2025-12-31T13:01:49.459Z'
  ts_utc,

  -- double: 1735651309.459
  epoch_seconds,

  -- timestamp: 2025-12-31 13:01:49+00
  TO_TIMESTAMP(epoch_seconds) AS ts_from_epoch,

  -- timestampz: '2025-12-31 13:01:49+00'
```

```

DATE_TRUNC('second', ts_utc)                AS ts_truncated,

-- date: '2025-12-31'
ts_utc::date                                AS utc_date,

-- date: '2025-12-29' (Monday of the ISO week)
DATE_TRUNC('week', ts_utc)::date            AS week_start_date,

-- date: '2026-01-04' (Sunday of the ISO week)
(DATE_TRUNC('week', ts_utc)
 + INTERVAL '6 days')::date                 AS week_end_date,

-- date: '2025-12-01'
DATE_TRUNC('month', ts_utc)::date           AS month_start_date,

-- date: '2025-12-31'
(DATE_TRUNC('month', ts_utc + INTERVAL '1 month')
 - INTERVAL '1 day')::date                 AS month_end_date,

-- date: '2025-10-01'
DATE_TRUNC('quarter', ts_utc)::date        AS quarter_start_date,

-- int: 2025
EXTRACT(YEAR FROM ts_utc)::int             AS utc_year,

-- date: '2026-01-30'
ts_utc::date + INTERVAL '30 days'         AS date_plus_30,

-- integer: days since 2025-12-31
CURRENT_DATE - ts_utc::date               AS days_since,

-- timestampz: '2025-12-31 13:01:49+00'
TO_TIMESTAMP(ts_string,
 'YYYY-MM-DD HH24:MI:SS')                AS parsed_timestamp,

-- varchar: '2025-12-31'
TO_CHAR(ts_utc, 'YYYY-MM-DD')             AS formatted_date,

-- varchar: '2025-12-31 13:01:49'
TO_CHAR(ts_utc, 'YYYY-MM-DD HH24:MI:SS')  AS formatted_datetime

```

```
FROM time_dimensions;
```

## 17.2. Date Dimension Calendar

When to use a date dimension table

A date dimension is a lookup table with one row per calendar day, pre-calculated with week, month, quarter, and year boundaries. Ideally it lives permanently in your data warehouse, and all fact tables join to it.

### NOTE

- **Use the warehouse date dimension** when one already exists — joining to it is faster than computing date boundaries inline on every query.
- **Generate it with GENERATE\_SERIES** when you do not have a warehouse date dimension, or when you need a specific rolling window (e.g. the last 365 days only) rather than a full calendar.
- **Common pattern in dbt:** create a `date_spine` model in the seeds or staging layer once, then `ref()` it everywhere else. Repeating `GENERATE_SERIES` in multiple models violates the DRY principle — generate once, reference everywhere.

Generates a row for every day in the last 365 days with week and month boundaries. Use this as a lookup table to join with fact tables for time-based reporting and aggregations.

```
WITH date_spine AS (
  SELECT
    GENERATE_SERIES(
```

```

        CURRENT_DATE - INTERVAL '365 days',
        CURRENT_DATE,
        INTERVAL '1 day'
    )::date AS date_column
)
SELECT
    -- date: varies (last 365 days)
    date_column                                AS date,

    -- date: Monday of the ISO week
    DATE_TRUNC('week', date_column)::date     AS week_start_date,

    -- date: Sunday of the ISO week
    (DATE_TRUNC('week', date_column)
     + INTERVAL '6 days')::date              AS week_end_date,

    -- date: first day of the month
    DATE_TRUNC('month', date_column)::date    AS month_start_date,

    -- date: last day of the month
    (DATE_TRUNC('month',
        date_column + INTERVAL '1 month')
     - INTERVAL '1 day')::date               AS month_end_date

FROM date_spine
ORDER BY date_column;
```

### 17.3. String Functions

When to use string functions

**NOTE**

String functions are primarily a concern of the Silver layer. Raw data arrives messy — inconsistent casing, extra whitespace, embedded delimiters, concatenated fields — and string functions are the tool to fix it before data reaches the Gold layer.

- **Bronze layer:** load raw values as-is. Avoid transforming in bronze — keep it a faithful copy of the source.
- **Silver / Staging layer:** this is where string functions belong. Trim whitespace, normalize case, split composite fields, validate formats with REGEXP, and fill or flag unexpected values. Your goal is clean, typed, consistent columns.
- **Gold layer:** cleaning string functions should rarely appear here. If they do, it is a signal that something was missed in staging. Presentation string functions (CONCAT, REPLACE for display labels) are appropriate in Gold for formatting clean data for dashboards.

Demonstrates text manipulation, pattern matching, and string parsing. Use this to clean, transform, and extract data from text fields like names, emails, URLs, and product codes.

### Cleaning Functions — Silver / Staging Layer

These functions fix raw data. They belong in your staging models — run once so every downstream model gets clean, consistent columns.

```
WITH customer_data AS (
  SELECT * FROM (VALUES
    (1, 'quantum.nebula@turmeric-tech.com', 'Premium Plan - Annual'),
    (2, 'stellar.phoenix@cardamom.co', 'STARTER plan (monthly)'),
    (3, 'Cosmic.Thunder@SAFFRON.IO', 'Enterprise Plan - Quarterly'),
    (4, 'atomic_lightning@coriander-labs.net', 'growth PLAN - annual'),
    (5, 'digital.vortex@peppercorn.com', 'Premium Plan - Monthly')
  ) AS t(customer_id, email, plan_description)
)
SELECT
  -- integer: 1-5
```

```

customer_id,

-- Whitespace and case normalisation -----

-- varchar: 'quantum.nebula@turmeric-tech.com' (whitespace removed)
TRIM(email) AS email_cleaned,

-- varchar: 'quantum.nebula@turmeric-tech.com' (forced lowercase)
LOWER(TRIM(email)) AS email_lower,

-- varchar: 'QUANTUM.NEBULA@TURMERIC-TECH.COM'
UPPER(email) AS email_upper,

-- Splitting and extracting fields -----

-- varchar: 'quantum.nebula' (username from email)
SPLIT_PART(TRIM(email), '@', 1) AS username,

-- varchar: 'quantum.nebula' (same result using regex)
SUBSTRING(TRIM(email)
  FROM '^([^@]+)') AS username_regex,

-- varchar: 'quantum.nebula' (using POSITION to find @ boundary)
SUBSTRING(TRIM(email), 1,
  POSITION('@' IN TRIM(email)) - 1) AS username_substr,

-- Pattern matching and replacement -----

-- varchar: 'premium_plan_annual' (slug-safe string)
REGEXP_REPLACE(LOWER(plan_description),
  '^[^a-z0-9]+', '_', 'g') AS plan_slug,

-- boolean: true for rows 1 and 4 (filter or flag annual plans)
LOWER(plan_description) LIKE '%annual%' AS is_annual,

-- Validation helpers -----

-- integer: 34 (use to flag suspiciously short or long values)
LENGTH(TRIM(email)) AS email_length,

```

```

-- integer: 15 (use to validate @ exists in expected position)
POSITION('@' IN TRIM(email))           AS at_sign_position

FROM customer_data
ORDER BY customer_id;

```

## Presentation Functions — Gold Layer

These functions format clean data for display. They belong in mart models where data is being shaped for dashboards and reports — not in staging where data is still being cleaned.

```

WITH customer_data AS (
  SELECT * FROM (VALUES
    (1, 'quantum.nebula@turmeric-tech.com', 'Premium Plan - Annual'),
    (2, 'stellar.phoenix@cardamom.co',     'STARTER plan (monthly)'),
    (3, 'cosmic.thunder@saffron.io',      'Enterprise Plan - Quarterly'),
    (4, 'atomic_lightning@coriander-labs.net', 'growth PLAN - annual'),
    (5, 'digital.vortex@peppercorn.com',   'Premium Plan - Monthly')
  ) AS t(customer_id, email, plan_description)
)
SELECT
  -- Combining fields for display -----
  -- varchar: 'Customer: 1 - quantum.nebula@turmeric-tech.com'
  CONCAT('Customer: ', customer_id::text,
    ' - ', email)                               AS customer_label,

  -- Formatting for display -----
  -- varchar: 'Premium Plan - Annual' (em-dash for display)
  REPLACE(plan_description, '-', ' ')          AS plan_display,

  -- varchar: 'turmeric-tech.com' (show domain only in UI)
  SPLIT_PART(email, '@', 2)                    AS email_domain,

  -- varchar: 'quantum.nebula' (show username only in UI)
  SPLIT_PART(email, '@', 1)                    AS display_name

```

```
FROM customer_data
ORDER BY customer_id;
```

## 17.4. Window Functions

Window functions in the medallion architecture

Advanced window functions (such as `ROW_NUMBER()`, `RANK()`, `LEAD/LAG`, and `SUM() OVER()`) are used to transform, clean, and analyze data as it moves from the Silver to the Gold layer.

### NOTE

- **Silver Layer (Cleansing and Transformation):** window functions are used to deduplicate data, handle late-arriving events, and normalize data from multiple sources.
- **Gold Layer (Reporting and Analytics):** window functions compute complex business metrics such as running totals and trends, creating pre-aggregated tables optimized for dashboard performance.

Silver Layer — Advanced window functions for data cleansing

### NOTE

- **Deduplication (`ROW_NUMBER()`):** when ingesting data from sources like Kafka or IoT, duplicates often appear. `ROW_NUMBER() OVER(PARTITION BY id ORDER BY timestamp DESC)` identifies and keeps only the latest record for each entity.
- **Handling late-arriving data:** use `LAG` to check if a specific event arrived out of order compared to the previous record in the timestamp sequence.
- **Data quality checks:** rank records to find invalid or null

values within specific partitions.

Gold Layer — Window functions for aggregation and reporting

**NOTE**

- **Running totals and moving averages (SUM...OVER):** calculate KPIs like 7-day average sales or customer lifetime spend to optimize dashboard performance.
- **Ranking top performers (RANK() / DENSE\_RANK()):** identify the top 10 products sold or top customers per region.
- **Time series analysis (LEAD/LAG):** determine the time difference between user clicks or transactional steps in a customer journey.

Demonstrates ranking, aggregation, and offset functions for B2B usage-based pricing. Use this to understand how to perform running totals, rankings, comparisons, and analytics across ordered data.

```
WITH usage_billing_data AS (
  SELECT * FROM (VALUES
    (1, 'Turmeric Technologies', 'Enterprise', DATE '2025-01-31', 45000,
    2250.00),
    (2, 'Cardamom Systems', 'Growth', DATE '2025-01-31', 18000,
    900.00),
    (3, 'Saffron Solutions', 'Enterprise', DATE '2025-01-31', 52000,
    2600.00),
    (4, 'Cumin Innovations', 'Starter', DATE '2025-01-31', 3500,
    175.00),
    (5, 'Turmeric Technologies', 'Enterprise', DATE '2025-02-28', 48000,
    2400.00),
    (6, 'Coriander Labs', 'Growth', DATE '2025-02-28', 22000,
    1100.00),
    (7, 'Cardamom Systems', 'Growth', DATE '2025-02-28', 16500,
    825.00),
```

```

        (8, 'Peppercorn Group',      'Enterprise', DATE '2025-02-28', 61000,
3050.00),
        (9, 'Fenugreek Digital',    'Growth',     DATE '2025-02-28', 25000,
1250.00),
        (10, 'Mustard Seed Tech',   'Starter',   DATE '2025-02-28', 4200,
210.00)
    ) AS t(invoice_id, customer_name, plan_tier,
          billing_period_end, api_calls, amount_billed)
)
SELECT
    -- integer: 1-10
    invoice_id,
    -- varchar: 'Turmeric Technologies', 'Cardamom Systems', etc.
    customer_name,
    -- varchar: 'Enterprise', 'Growth', 'Starter'
    plan_tier,
    -- date: end of billing period
    billing_period_end,
    -- integer: total API calls (3500 to 61000)
    api_calls,
    -- numeric: amount billed ($175 to $3050)
    amount_billed,

    -- Ranking -----
    -- integer: 1, 2, 3, ... sequential invoice number
    ROW_NUMBER() OVER (
        ORDER BY billing_period_end, customer_name) AS row_num,

    -- integer: rank by API call volume (1=highest)
    -- ties share rank, next rank skips
    RANK() OVER (
        ORDER BY api_calls DESC) AS rank_by_usage,

    -- integer: dense rank by revenue, no gaps between ranks
    DENSE_RANK() OVER (
        ORDER BY amount_billed DESC) AS dense_rank_by_revenue,

    -- integer: usage ranking restarted within each plan tier
    ROW_NUMBER() OVER (

```

```

PARTITION BY plan_tier
ORDER BY api_calls DESC)                AS usage_rank_within_tier,

-- integer: 1=Low, 2=Med-Low, 3=Med-High, 4=High quartile
NTILE(4) OVER (
  ORDER BY api_calls)                    AS usage_quartile,

-- Running totals and rolling windows -----

-- numeric: running total of revenue (2250, 3150, 5750, ...)
SUM(amount_billed) OVER (
  ORDER BY billing_period_end, customer_name
  ROWS BETWEEN UNBOUNDED PRECEDING
  AND CURRENT ROW)                      AS cumulative_revenue,

-- numeric: 3-month rolling average API calls per customer
AVG(api_calls) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end
  ROWS BETWEEN 2 PRECEDING
  AND CURRENT ROW)                      AS rolling_avg_usage_3mo,

-- Partition aggregates -----

-- integer: total invoices for each plan tier
COUNT(*) OVER (
  PARTITION BY plan_tier)                AS invoices_per_tier,

-- integer: total API calls for the whole tier
SUM(api_calls) OVER (
  PARTITION BY plan_tier)                AS total_usage_by_tier,

-- double: this customer's share of tier usage (0.35 = 35%)
api_calls::double precision /
  NULLIF(SUM(api_calls) OVER (
    PARTITION BY plan_tier), 0)          AS pct_of_tier_usage,

-- integer: highest single-month usage for this customer
MAX(api_calls) OVER (
  PARTITION BY customer_name)            AS customer_peak_usage,

```

```

-- LAG and LEAD -----

-- integer: previous month API calls (NULL for first month)
LAG(api_calls, 1) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end)          AS prev_month_usage,

-- integer: next month API calls (NULL for last month)
LEAD(api_calls, 1) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end)         AS next_month_usage,

-- integer: month-over-month change (+3000 = 3000 more calls)
api_calls - LAG(api_calls, 1) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end)        AS usage_change_mom,

-- double: month-over-month growth rate (0.067 = 6.7%)
(api_calls - LAG(api_calls, 1) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end))::double precision /
NULLIF(LAG(api_calls, 1) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end), 0)    AS usage_growth_rate,

-- numeric: previous month's billed amount
LAG(amount_billed, 1) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end)       AS prev_month_revenue,

-- numeric: month-over-month revenue change
amount_billed - LAG(amount_billed, 1) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end)       AS revenue_change_mom,

-- Value functions -----

-- integer: first month usage for this tier (baseline)
FIRST_VALUE(api_calls) OVER (

```

```

PARTITION BY plan_tier
ORDER BY billing_period_end) AS tier_baseline_usage,

-- integer: most recent usage for this tier
-- requires explicit frame or LAST_VALUE returns current row
LAST_VALUE(api_calls) OVER (
  PARTITION BY plan_tier
  ORDER BY billing_period_end
  ROWS BETWEEN UNBOUNDED PRECEDING
    AND UNBOUNDED FOLLOWING) AS tier_latest_usage,

-- numeric: second highest billed amount in this tier
NTH_VALUE(amount_billed, 2) OVER (
  PARTITION BY plan_tier
  ORDER BY amount_billed DESC) AS
second_highest_revenue_in_tier,

-- text[]: ordered list of plan tiers per customer over time
ARRAY_AGG(plan_tier) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end
  ROWS BETWEEN UNBOUNDED PRECEDING
    AND CURRENT ROW) AS plan_history,

-- Advanced window frames -----

-- double: symmetric 3-row window (1 before + current + 1 after)
AVG(api_calls) OVER (
  ORDER BY billing_period_end
  ROWS BETWEEN 1 PRECEDING
    AND 1 FOLLOWING) AS centered_moving_avg,

-- numeric: reverse running total (complement of cumulative_revenue)
SUM(amount_billed) OVER (
  ORDER BY billing_period_end, customer_name
  ROWS BETWEEN CURRENT ROW
    AND UNBOUNDED FOLLOWING) AS remaining_revenue,

-- integer: asymmetric window (2 before + current + 1 after)
SUM(api_calls) OVER (

```

```

PARTITION BY customer_name
ORDER BY billing_period_end
ROWS BETWEEN 2 PRECEDING
        AND 1 FOLLOWING)                AS usage_4mo_window,

-- double: RANGE groups tied ORDER BY values together
SUM(amount_billed) OVER (
    ORDER BY billing_period_end
    RANGE BETWEEN UNBOUNDED PRECEDING
        AND CURRENT ROW)                AS range_cumulative_revenue,

-- double: value-based band – averages rows within ±5000 of current
AVG(api_calls) OVER (
    PARTITION BY plan_tier
    ORDER BY api_calls
    RANGE BETWEEN 5000 PRECEDING
        AND 5000 FOLLOWING)            AS range_band_avg_usage,

-- Percentile ranking -----

-- double: 0.0-1.0 relative rank: (rank-1)/(total_rows-1)
-- lowest=0.0, highest=1.0
PERCENT_RANK() OVER (
    ORDER BY api_calls)                AS pct_rank_usage,

-- double: same as pct_rank_usage, restarted within each tier
PERCENT_RANK() OVER (
    PARTITION BY plan_tier
    ORDER BY api_calls)                AS pct_rank_within_tier,

-- double: fraction of rows <= current row. Min=1/n, never 0.0
CUME_DIST() OVER (
    ORDER BY amount_billed)            AS cumulative_distribution,

-- double: same as cumulative_distribution, scoped to tier
CUME_DIST() OVER (
    PARTITION BY plan_tier
    ORDER BY amount_billed)            AS cume_dist_within_tier,

-- LAG/LEAD with defaults -----

```

```

-- integer: returns 0 instead of NULL for first month
LAG(api_calls, 1, 0) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end)          AS
prev_month_usage_default_zero,

-- numeric: returns 0.0 instead of NULL for last month
LEAD(amount_billed, 1, 0.0) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end)          AS
next_month_revenue_default_zero,

-- integer: api_calls from 2 months ago, 0 if fewer than 2 prior
LAG(api_calls, 2, 0) OVER (
  PARTITION BY customer_name
  ORDER BY billing_period_end)          AS two_months_ago_usage,

-- Peer comparison and conditional aggregation -----

-- double: difference between this row and the tier average
AVG(api_calls::double precision) OVER (
  PARTITION BY plan_tier)
- api_calls                             AS usage_diff_from_tier_avg,

-- numeric: total tier revenue minus this row's contribution
SUM(amount_billed) OVER (
  PARTITION BY plan_tier)
- amount_billed                         AS tier_revenue_excl_self,

-- integer: count of >20000 call invoices in this tier
SUM(CASE WHEN api_calls > 20000 THEN 1 ELSE 0 END) OVER (
  PARTITION BY plan_tier)                AS
high_usage_invoices_in_tier,

-- numeric: revenue from high-usage months only within tier
SUM(CASE WHEN api_calls > 20000
  THEN amount_billed ELSE 0.0 END) OVER (
  PARTITION BY plan_tier)                AS high_usage_revenue_in_tier

```

```
FROM usage_billing_data  
ORDER BY billing_period_end, customer_name;
```

*All examples use PostgreSQL 14+ syntax and include practical B2B SaaS use cases with fictitious company names for illustration purposes.*

## Chapter 18. dbt Command Reference

dbt — data build tool — is the transformation layer of the modern data stack. It takes raw data that has already been loaded into your warehouse and transforms it into business-ready models using SQL. It brings software engineering practices — version control, testing, documentation, and modular design — directly into analytics work. This chapter is a reference for every command, concept, and pattern you will use working with dbt in practice.

### 18.1. The dbt Project Structure

A dbt project is a folder. Everything dbt needs lives inside it — models, tests, seeds, macros, snapshots, and configuration files. Understanding what goes where makes every command easier to reason about.

**models/** is where all transformation SQL lives. Every `.sql` file in this folder is a model — a transformation that reads from a source or another model and writes an output. Models are organized into subfolders by layer: `staging/` for cleaning raw data, `intermediate/` for heavy aggregations, and `marts/` for business-ready outputs.

**tests/** holds singular tests — `.sql` files that define custom business rules. A singular test returns rows when something is wrong. Zero rows means the test passes. Any rows returned means it fails.

**seeds/** holds CSV files that dbt loads into the database as lookup tables. Region codes, product names, and other reference data that changes rarely belong here.

**macros/** holds reusable Jinja functions. Write a macro once and call it from any model. Custom generic tests also live here — named `test_<name>.sql` and applied in YAML like built-in tests.

**snapshots/** holds snapshot definitions. Snapshots track how a record changes over time by recording a new row with a timestamp whenever something changes — giving you a full history of slowly changing data.

**analyses/** holds SQL files you want to version control and compile with dbt but never materialize into the database. Useful for ad hoc queries and audits that need access to `ref()` and `source()`.

**target/** is generated — never commit it. Every dbt command writes compiled SQL and run artifacts here. `target/compiled/` contains the plain SQL dbt will run. `target/manifest.json` and `target/run_results.json` are the dbt artifacts used for state comparison and CI.

**dbt\_packages/** is generated — never commit it. Installed by dbt `deps` from `packages.yml`.

The key files that govern a dbt project:

**dbt\_project.yml** is the main configuration file. It sets the project name, tells dbt where to find each type of file, configures how each model layer materializes, and defines post-hooks for grants. Every dbt project must have exactly one.

**profiles.yml** holds database connection credentials. It lives outside the project folder at `~/.dbt/profiles.yml` — never inside the project — so it cannot be accidentally committed to Git. It defines one or more targets: `personal`, `dev`, `ci`, and `prod`.

**packages.yml** lists the dbt packages to install — `dbt_utils`, `dbt_expectations`, `astronomer-cosmos`, and others. Run `dbt deps` after changing it.

**sources.yml** declares the raw tables that dbt models read from but did not create. Without it, dbt has no visibility into those tables — no lineage, no freshness checks, no compile-time validation.

**schema.yml** documents and tests model output columns. One per folder is the convention. This is where `not_null`, `unique`, `accepted_values`, and custom tests are applied, and where model contracts and access modifiers are declared.

## 18.2. Project Setup Commands

`dbt init` creates a new dbt project from scratch — generating the folder structure, `dbt_project.yml`, and a sample `profiles.yml`. Run it once when starting a new project.

```
dbt init project_name
```

`dbt debug` tests the database connection and validates all configuration files. Run this whenever a connection issue occurs — it checks the profile, the target, the adapter version, and whether the database is reachable.

```
dbt debug
```

### What dbt debug checks.

#### NOTE

It validates four things in sequence: that `dbt_project.yml` and `profiles.yml` parse correctly, that the profile name in `dbt_project.yml` matches an entry in `profiles.yml`, that the database adapter is installed, and that the connection

credentials actually reach the database. If any check fails it stops and tells you exactly which one.

`dbt deps` installs packages listed in `packages.yml` into `dbt_packages/`. Run it once after cloning a project and again whenever `packages.yml` changes.

```
dbt deps

# Upgrade all packages to latest compatible versions
dbt deps --upgrade
```

## 18.3. Development Commands

### 18.3.1. *dbt compile*

`dbt compile` translates all Jinja templating into plain SQL and writes it to `target/compiled/`. Nothing runs in the database. Use it to inspect what SQL dbt will actually execute before running it, to debug Jinja expressions, or to see which schema a `ref()` resolves to in each target.

```
# Compile the entire project
dbt compile

# Compile one model and inspect the output
dbt compile --select model_name

# Force full reparse – clears cached parse state
dbt compile --no-partial-parse
```

#### NOTE

#### Why compile before run?

Your dbt SQL files contain Jinja — `{{ ref('stg_orders') }}`, `{% if is_incremental() %}`, `{{ source('raw_bronze', 'customer_details') }}`. `dbt compile` resolves all of that into the plain SQL that

PostgreSQL actually receives. Copying the compiled SQL from `target/compiled/` and running it directly in pgAdmin is the fastest way to debug a model without waiting for a full run.

### 18.3.2. *dbt seed*

`dbt seed` loads CSV files from the `seeds/` folder into the database as tables. Seeds are for small, static lookup tables — region codes, product names, configuration values. Run it once when setting up the project and again whenever a CSV changes.

```
# Load all seeds
dbt seed

# Load one specific seed
dbt seed --select seed_name

# Full refresh – drop and recreate even if unchanged
dbt seed --full-refresh
```

### 18.3.3. *dbt run*

`dbt run` executes the SQL models and builds views and tables in the database. It does not run tests.

```
# Run all models
dbt run

# Run one model
dbt run --select model_name

# Run a folder of models
dbt run --select staging
```

```
# Run a model and all its downstream dependents
dbt run --select model_name+

# Run only models you changed plus their downstream dependents
dbt run --select state:modified+ --state ./prod-state/

# For any model not selected, read from prod instead of rebuilding
dbt run --select state:modified+ --defer --state ./prod-state/
```

### 18.3.4. *dbt build*

`dbt build` is the main command for running the full pipeline. It runs seeds, models, snapshots, and tests in dependency order — in one command. This is what CI and Airflow run in production.

```
# Full build – all seeds, models, tests, snapshots
dbt build

# Build in personal schema
dbt build --target personal

# Build in dev schema
dbt build --target dev

# Build only modified models plus downstream – the standard development loop
dbt build --select state:modified+ --defer --state ./prod-state/

# Stop immediately on first failure
dbt build --select state:modified+ --fail-fast --state ./prod-state/

# Build schema only with zero rows – fastest SQL syntax check
dbt build --select state:modified+ --empty --state ./prod-state/
```

#### **dbt build vs dbt run.**

#### **NOTE**

`dbt run` executes models only — no tests. `dbt build` runs everything: seeds, models, tests, and snapshots. Use `dbt run`

during fast iterative development when you just want to see if the SQL executes. Use `dbt build` before opening a Pull Request and in all automated pipelines — it is the complete check.

### 18.3.5. Model Selection Syntax

Model selection is how you tell `dbt` which models to include in any command. The same syntax works with `dbt run`, `dbt build`, `dbt test`, and `dbt ls`.

Selector	What it selects
<code>model_name</code>	That specific model only
<code>model_name+</code>	That model and all its downstream dependents
<code>+model_name</code>	That model and all its upstream dependencies
<code>model_name</code>	The model, all upstream, and all downstream
<code>staging</code>	All models in the staging folder
<code>tag:revenue</code>	All models tagged with <code>revenue</code>
<code>source:raw_bronze</code>	All models that read directly from the <code>raw_bronze</code> source
<code>state:modified+</code>	Models whose SQL changed vs the reference manifest, plus downstream
<code>state:new</code>	Brand new models that did not exist in the reference manifest
<code>result:error+</code>	Models that errored in the last run, plus downstream
<code>result:skipped+</code>	Models that were skipped in the last run, plus downstream
<code>path:models/marts</code>	All models in the <code>marts</code> directory path
<code>exposure:dashboard_name</code>	All models feeding a specific exposure

Combine selectors with space (union) or `--exclude`:

```
# Union – run modified models AND new models
dbt build --select state:modified+ state:new --state ./prod-state/
```

```
# Exclude – run all staging except one model
dbt run --select staging --exclude stg_customer_details
```

## 18.4. Testing

### 18.4.1. dbt test

`dbt test` runs all data quality tests against the models currently in the database. Tests must be defined in `schema.yml` or as `.sql` files in `dbt/tests/` before `dbt test` can run them.

```
# Run all tests
dbt test

# Run tests on one model
dbt test --select model_name

# Run only source tests
dbt test --select source:raw_bronze

# Run only generic tests (not_null, unique, accepted_values, custom generics)
dbt test --select test_type:generic

# Run only singular tests (SQL files in dbt/tests/)
dbt test --select test_type:singular

# Run only unit tests (unit_tests: blocks in schema.yml)
dbt test --select test_type:unit

# Run a specific singular test by name
dbt test --select test_revenue_always_positive

# Test sources then build models – fail-fast pattern
dbt test --select source:raw_bronze && dbt build
```

### 18.4.2. The Three Test Types

dbt has three distinct test types. Each serves a different purpose and dbt infers the type automatically from location and structure — no explicit declaration needed.

**Generic tests** validate column-level constraints. They are declared in `schema.yml` and apply the same check to any column in any model. Built-in generic tests are `not_null`, `unique`, `accepted_values`, and `relationships`. Custom generic tests are written as macros in `dbt/macros/test_<name>.sql` and applied the same way as built-ins. Use a generic test when the check is simple, column-level, and reusable.

**Singular tests** validate business rules specific to a dataset. They are `.sql` files in `dbt/tests/`. Zero rows returned means the test passes — any rows returned means it fails. Use a singular test when the check requires joining multiple models, involves business arithmetic, or is too specific to express in a YAML declaration.

**Unit tests** validate transformation logic using mock data — no database connection needed. They are defined as `unit_tests:` blocks in `schema.yml`. Given controlled mock inputs, the output must match the expected rows exactly. Unit tests catch logic bugs before they reach real data. They are the only test type that runs entirely in memory.

	<b>Generic</b>	<b>Singular</b>	<b>Unit</b>
Validates	Column constraints	Business rules across dataset	Transformation logic
Input	Real data	Real data	Mock data in YAML
Written in	Macro + YAML	Pure SQL	YAML

	Generic	Singular	Unit
Requires DB	Yes	Yes	No

**NOTE**

**Unit tests do not replace data tests — they complement them.**

Unit tests catch logic bugs. Data tests catch bad data. Both are necessary. A unit test confirms that your CASE WHEN logic handles nulls correctly. A data test confirms that no nulls actually slipped through in production.

### 18.4.3. Test Configuration

Tests can be configured individually for severity, failure thresholds, and failure storage:

```
# Run tests – store failing rows to the database for inspection
dbt test --store-failures

# Promote all warnings to errors – useful in CI
dbt build --warn-error
```

Key test configuration options applied in `schema.yml`:

**severity** — `error` stops the run, `warn` logs and continues. Default is `error`.

**error\_if / warn\_if** — threshold-based severity. `error_if: ">10"` only errors if more than 10 rows fail. `warn_if: ">0"` warns if any rows fail but does not error.

**store\_failures** — persists failing rows to a database table for inspection in pgAdmin. Useful for debugging which specific rows violated a test.

**limit** — caps the number of failing rows returned. Useful for large tables

where returning all failures would be slow.

## 18.5. Data and History

### 18.5.1. dbt snapshot

`dbt snapshot` runs all snapshot definitions. Snapshots track how records change over time by recording a new row with a timestamp whenever a value changes — giving you a full history of slowly changing data. Bronze tables are append-only and never update existing rows. Snapshots bridge that gap for source data that does change.

```
# Run all snapshots
dbt snapshot

# Run one specific snapshot
dbt snapshot --select snapshot_name
```

#### NOTE

#### When to use snapshots.

Use snapshots for source data that changes over time and where you need the history — a customer’s segment changing from Commercial to Strategic, a model going inactive in a region, an account’s contract tier being upgraded. Without a snapshot, you only ever see the current state.

### 18.5.2. dbt source freshness

`dbt source freshness` checks how old the raw data in your source tables is. It reads the most recent timestamp in each source table and compares it against the thresholds defined in `sources.yml`. If the data is older than the threshold, dbt warns or errors before any models are built on top of stale

data.

```
# Check all sources
dbt source freshness

# Check one specific source table
dbt source freshness --select source:raw_bronze.revenue_account_daily

# Check sources feeding a specific exposure
dbt source freshness --select +exposure:revenue_dashboard
```

**NOTE**

**Run freshness before build in production.**

In a production pipeline, run `dbt source freshness` as a separate step before `dbt build`. If any source errors, the build step never starts — preventing models from being built on top of stale data and silently serving wrong numbers in dashboards.

## 18.6. Production

### 18.6.1. Running Against Targets

Every `dbt` command runs against a target defined in `profiles.yml`. The default target is `personal` — your private development schema. Use `--target` to specify a different one.

```
# Run against personal schema (default)
dbt build

# Run against dev schema
dbt build --target dev

# Run against prod schema
dbt build --target prod
```

### 18.6.2. State-Based Selection

State-based selection is how dbt knows what changed between runs. It compares your current code against a reference `manifest.json` — the prod manifest — and selects only what changed and its downstream dependents.

```
# Save the prod manifest after a successful prod build
mkdir -p prod-state
cp target/manifest.json prod-state/manifest.json

# Preview what changed vs prod manifest
dbt ls --select state:modified+ --state ./prod-state/

# Run only what changed and its downstream dependents
dbt build --select state:modified+ --state ./prod-state/ --target prod

# Update the manifest after every successful prod run
cp target/manifest.json prod-state/manifest.json
```

#### NOTE

#### Keep prod-state/ out of Git.

The prod manifest is a local reference file, not source code. Add `prod-state/` to `.gitignore`. In CI, the manifest is downloaded from the artifact store — not committed.

### 18.6.3. --defer

`--defer` tells dbt: for any model not selected, read from prod instead of rebuilding it. This means you only rebuild the models you changed — unchanged upstream dependencies are read directly from production tables.

```
# The standard CI and development command
# Only changed models rebuild – everything else reads from prod
dbt build \
```

```
--select state:modified+ state:new \  
--defer \  
--state ./prod-state/
```

**NOTE**

**--defer does not copy prod data into your schema.**

Your changed models still write to your personal or CI schema. What changes is what your model reads as input — it reads from prod tables instead of rebuilding them first. Your personal schema gets one new table. It does not become a copy of prod.

#### 18.6.4. State and Result Selectors

Two selector families answer two different questions:

**State selectors** ask "what changed compared to a previous run?" They compare your current code against a `manifest.json` snapshot.

**Result selectors** ask "what happened in the last run?" They read `run_results.json` from the last run to find what failed, errored, or was skipped.

```
# State selectors  
dbt ls --select state:modified+ # changed models + downstream  
dbt ls --select state:new      # new models not in prod manifest  
  
# Result selectors – after a failed run  
dbt ls --select result:error+  # errored models + downstream  
dbt ls --select result:skipped+ # skipped models + downstream  
  
# Retry only what failed and what was skipped  
dbt build --select result:error+ result:skipped+ --state ./prod-state/  
  
# Or simply – re-execute the last command from the point of failure  
dbt retry
```

### 18.6.5. *dbt retry*

`dbt retry` re-executes the exact last `dbt` command from the point of failure. It reads `run_results.json` automatically — no selectors needed. Use it when you fix a bug and want to rerun what failed without typing the full command again.

```
# Fix the bug, then:
dbt retry
```

## 18.7. *Inspection and Utilities*

### 18.7.1. *dbt docs*

`dbt docs generate` builds the documentation website — a full lineage graph from raw sources through staging, intermediate, and marts, with all model descriptions, column descriptions, tests, and source freshness status. `dbt docs serve` opens it in your browser.

```
# Generate docs
dbt docs generate

# Serve on default port 8080
dbt docs serve

# Serve on a different port
dbt docs serve --port 8081
```

#### NOTE

#### What the lineage graph shows.

The lineage graph is the most valuable output of `dbt docs generate`. It shows every model as a node, every dependency as an edge, and makes it immediately visible when a model is

pulling directly from a raw source instead of going through staging — a DAG design violation. Exposures appear as orange nodes downstream of your marts, making dashboard dependencies visible in the same graph.

### 18.7.2. *dbt ls*

`dbt ls` lists resources in the project without running anything. Use it to preview what a selector will match before running it.

```
# List all models
dbt ls

# List all tests
dbt ls --resource-type test

# Preview what state:modified+ will select
dbt ls --select state:modified+ --state ./prod-state/

# List all models with a specific tag
dbt ls --select tag:revenue

# List all versions of a versioned model
dbt ls --select customer_revenue_monthly
```

### 18.7.3. *dbt clean*

`dbt clean` deletes the `target/` and `dbt_packages/` folders. Use it when you want a completely fresh start — no cached compiled SQL, no installed packages.

```
dbt clean
```

### 18.7.4. *dbt parse*

`dbt parse` parses the entire project — validates all YAML files, resolves all `ref()` and `source()` calls, and checks the dependency graph — without running anything in the database. Faster than `dbt compile` for catching configuration errors.

```
dbt parse

# Force full reparse - ignore cached parse state
dbt compile --no-partial-parse
```

## 18.8. Jinja and Macros

dbt uses Jinja — a templating language — to make SQL dynamic. Jinja runs at compile time on your machine, before any SQL reaches the database. This is a critical distinction: Jinja cannot access row data — it only knows about project configuration, targets, and variables.

### 18.8.1. Built-in dbt Functions

`ref('model_name')` references another dbt model. It resolves to the correct schema for whatever target you are running against — personal, dev, or prod. It also builds the dependency graph so dbt knows the execution order. Always use `ref()` to reference dbt models — never hardcode schema and table names.

`source('source_name', 'table_name')` references a raw table that dbt did not create. It resolves to the correct schema, validates the table exists at compile time, enables source freshness checks, and makes the lineage visible in dbt docs. Always use `source()` for bronze tables — never `ref()`.

#### NOTE

**`ref()` vs `source()` — the mental model.**

`ref()` means "dbt built this, dbt manages this." `source()` means "something else built this, dbt just reads it." Use `ref()` for everything in `models/`. Use `source()` for everything in `raw_bronze`.

`config()` sets model configuration inline — overriding anything set in `dbt_project.yml` or `schema.yml`. The most specific config always wins: SQL file config > `schema.yml` config > folder config in `dbt_project.yml`.

`this` refers to the current model's fully qualified table name. Used in incremental models to reference the existing table when deciding which rows to process. Resolves to the correct schema for the current target automatically.

`is_incremental()` returns `True` when dbt is running an incremental model and the target table already exists. Use it to filter for only new rows. Returns `False` on a full refresh or the first run.

### 18.8.2. Variables and Environment

`var('variable_name')` reads a project variable defined in `dbt_project.yml` or passed at runtime with `--vars`. Use it for configurable values like date ranges or feature flags.

```
# Override a variable at runtime
dbt build --vars '{"date_spine_start": "2025-01-01"}'
```

`env_var('ENV_VAR_NAME')` reads an environment variable from your shell. Used in `profiles.yml` to keep credentials out of the file — `POSTGRES_PASSWORD`, `DBT_SCHEMA`. Case-sensitive — `DBT_SCHEMA` and `dbt_schema` are different variables.

### 18.8.3. Macros

A macro is a reusable Jinja function defined in `dbt/macros/`. Write it once, call it from any model. Macros generate SQL at compile time — they are not functions that run in the database.

The primary benefit of macros is DRY — Don't Repeat Yourself. If the same window function pattern, the same type casting logic, or the same conditional expression appears in multiple models, it belongs in a macro.

```
# Run a macro directly – useful for one-off operations and code generation
dbt run-operation macro_name --args '{"arg1": "value1"}'

# Generate schema YAML for an existing model
dbt run-operation generate_model_yaml --args '{"model_names":
["model_revenue_monthly"]}'
```

**NOTE****When macros are worth the complexity.**

The primary concern with overusing macros is that compiled SQL becomes harder to read and debug. A macro that generates 50 lines of SQL from 2 lines of Jinja makes the compiled output dense and difficult to inspect. Use a macro when the same pattern repeats across five or more models and the compiled output remains readable. If you cannot read the compiled SQL and immediately understand it, the macro is too abstract.

*18.8.4. Selectors File*

Reusable selection logic can be saved in `dbt/selectors.yml` and referenced by name instead of typing the full selector every time.

```
# Use a named selector
dbt build --selector ci_modified
dbt build --selector revenue_only --target prod
```

```
dbt build --selector failed_retry --state ./prod-state/
```

## 18.9. Scheduling with Airflow

dbt Core is a command-line tool — it runs when you run it and stops when it finishes. It has no built-in scheduler. For automated, recurring pipeline runs in production, you need an orchestrator. Apache Airflow is the standard choice.

### 18.9.1. How Airflow Works

Airflow runs workflows defined as DAGs — Directed Acyclic Graphs. A DAG is a Python file that defines a set of tasks, the order they run in, and when they are triggered. Each task in the DAG is a unit of work — running a dbt command, checking source freshness, or sending an alert.

**Directed** — tasks flow in one direction. Task A must finish before Task B starts. **Acyclic** — no loops. A task can never depend on itself or create a cycle. **Graph** — the tasks and their dependencies form a visual graph visible in the Airflow UI at `localhost:8080`.

Airflow does not run dbt directly. It runs shell commands inside the Airflow container — the same `dbt build` commands you run manually, but triggered on a schedule or by a CI/CD event.

### 18.9.2. The Two Production DAGs

A complete setup has two DAGs:

`pipeline_daily` runs in production on a schedule — every day at 2:00am. It runs `dbt source freshness` first — if any source errors, the build never starts.

Then it runs `dbt build --target prod` against the production database. Triggered also by CI/CD on merge to main.

`pipeline_dev` runs in the dev environment on a schedule — every day at midnight. Same pattern: freshness check, then `dbt build --target dev`. Triggered also by CI/CD on merge to dev.

Neither DAG is triggered manually in production. Both run on schedule and on CI/CD events.

### 18.9.3. Why Plain `dbt build` Instead of Cosmos

Cosmos is a package that converts each dbt model into its own Airflow task — giving you model-level visibility and retry in the Airflow UI. However, Cosmos breaks dbt's natural test ordering. With Cosmos, a staging model's tests can fire before the mart models it feeds are built, causing `relation does not exist` errors on relationship tests.

Plain `dbt build` lets dbt handle ordering internally — it builds models in dependency order and runs tests at the right point. Use Cosmos when your pipeline takes longer than 10 minutes and you need model-level retry. For most projects, plain `dbt build` is the right choice.

### 18.9.4. CI/CD — Where Airflow Fits

| Environment | Who builds | Trigger | Database |

Personal	You manually	<code>dbt build</code>	Dev database	
Dev	Airflow <code>pipeline_dev</code>	Daily + on merge to dev	Dev database	

CI	GitHub Actions	On every Pull Request	Dev database	
Production	Airflow pipeline_daily	Daily + on merge to main	Prod database	

The efficient development loop: develop in personal schema, open a Pull Request, CI runs `dbt build --select state:modified+ --defer --state ./prod-state/` automatically, merge to dev, Airflow builds the dev schema overnight, merge to main, Airflow builds production.

## 18.10. Troubleshooting

**dbt debug fails with "connection refused"** — PostgreSQL is not running. Check `docker compose ps`. If postgres shows Exit — run `docker compose start postgres`.

**dbt compile fails with "Could not find relation"** — a `ref()` or `source()` is pointing to a model or table that does not exist. Check the spelling. Run `dbt parse` to see the full error with the model name.

**State selectors return nothing** — the prod manifest is missing or in the wrong location. Confirm `prod-state/manifest.json` exists. Run `dbt ls --select state:modified+ --state ./prod-state/` to preview before building.

**env\_var() fails with "is not set"** — the environment variable is not loaded in your shell. Run `source ~/.zshrc` and check with `env | grep DBT`. Remember `env_var()` is case-sensitive.

**Deprecation warnings from installed packages** — dbt 1.9+ requires the arguments: `wrapper` and `severity` inside `config:`. Try upgrading packages first with `dbt deps --upgrade`. Use `dbt test --show-all-deprecations` to find all

instances.

**Compiled SQL not updating after a schema change** — run `dbt compile --no-partial-parse` to force a full reparse and clear the cached parse state.

**DAG does not appear in Airflow after 60 seconds** — check for a Python syntax error in the DAG file. Run `docker exec btg-airflow-scheduler airflow dags list-import-errors`. Force a reserialize with `docker exec btg-airflow-scheduler airflow dags reserialize`.

# SECTION V: CASE STUDIES

This final section (Chapters 19-20) puts the framework to work in realistic business settings. Chapter 19 walks through setting up the data stack end to end — this edition covers the open source stack, with AWS and other cloud environments covered in the next edition. Chapter 20 presents two case studies: resource utilization and customer lifecycle. The resource utilization study demonstrates capacity planning, forecasting, and cost optimization using operational data. The customer lifecycle case covers cohort analysis, activation funnels, and retention levers that connect analytics directly to growth. A marketing attribution case study — covering lift measurement, channel optimization, and experiment-driven reporting — will be added in the next edition.

## Chapter 19. Building Your Analytics Environment

The case studies that follow share a common foundation. Before any analysis is possible, before any metric can be defined or any dashboard built, the data infrastructure that makes analysis possible has to exist. This chapter describes that foundation: what was built, why each component was chosen, and how the pattern is designed to be portable beyond this edition.

### ***19.1. The Local Stack Setup***

#### *19.1.1. What Was Built*

The stack used across all case studies is a local, open-source modern data stack running on a Mac. Every component is containerized using Docker, which means the entire environment — the database, the orchestration layer, the transformation tool, and the visualization layer — starts with a single command and runs identically on any machine.

The components and their roles:

**PostgreSQL** serves as the data warehouse. It stores raw source data in a bronze layer, cleaned and modelled data in silver and gold layers, and provides the query engine that powers every metric and dashboard in the case studies. PostgreSQL was chosen for this edition because it runs locally without a cloud account, costs nothing, and behaves consistently with cloud data warehouses in terms of SQL syntax and data modelling patterns.

**Apache Airflow** handles orchestration — the scheduling, sequencing, and monitoring of data pipelines. Each case study has its own set of DAGs:

bronze layer ingestion runs on demand when source data changes, silver layer staging runs every six hours, and gold layer mart builds run nightly. Airflow makes the pipeline visible — every run is logged, every failure is surfaced, and every dependency between tasks is explicit.

**dbt** handles transformation — the SQL models that move data from raw bronze tables through cleaned silver views into gold mart tables ready for analysis. dbt enforces the medallion architecture, runs data quality tests on every model, and generates documentation that makes the data lineage visible to anyone on the team. It runs inside the Airflow container through a virtual environment for scheduled automation.

**Metabase** provides the visualization layer — the dashboards and reports that turn gold layer data into something a business leader can read without writing SQL. Each case study has a corresponding Metabase dashboard built on top of the mart tables dbt produces.

**pgAdmin** provides direct database access for setup, validation, and troubleshooting — useful during the setup steps and whenever a query needs to be run directly against the database without going through dbt or Metabase.

### *19.1.2. Why This Pattern*

Every component in this stack has a direct equivalent in every major cloud data platform. The pattern — ingest to bronze, clean to silver, aggregate to gold, orchestrate with a scheduler, transform with a SQL-based tool, visualize with a BI layer — is the same whether the underlying infrastructure is a local PostgreSQL database or a cloud data warehouse running at petabyte scale.

This is deliberate. The goal of this edition is to make the pattern learnable on a laptop, without a cloud account, without infrastructure costs, and without organizational approvals. The same framework, the same dbt models, the same Airflow DAGs, and the same analytical approach transfer directly to production environments on any cloud platform.

### *19.1.3. Stack Setup Runbook — Mac Edition*

The following runbook walks through every step required to build and start the data stack on a Mac. Follow it from top to bottom — the stack will be running by the end.

### *19.1.4. Before You Begin*

#### **Before you start — what you need**

- A Mac (Intel or Apple Silicon)
- An internet connection
- About 1-2 hours
- 10 GB of free disk space — check: Apple menu → **About This Mac** → **Storage**

### *19.1.5. Decisions to Make First*

Before touching a single command, there are two things you need to decide. These are hard to change later — especially the database name.

#### **WARNING**

#### **Decide your database name now — before you start.**

The database name is set when Docker first starts PostgreSQL. It gets baked into the Airflow connection string,

the Metabase connection, the dbt profile, and your Git history. Changing it later means updating all of those places and recreating the database from scratch.

In this runbook we use: `btg_resource_utilization`

**Rules for naming:** Use underscores not hyphens — `btg_resource_utilization` not `btg-resource-utilization`. Hyphens require double quotes every time you reference the database in SQL. Underscores never do.

### Decide your project folder name now.

In this runbook:

`~/Documents/btg-case-studies-with-dbt/` — the Git repository

`~/Documents/btg-case-studies-with-dbt/single-dbt-opensource/` — this specific case study

If you use different names, substitute them everywhere you see these paths in this runbook.

#### WARNING

### 19.1.6. Step 1 — Check & Install Git

```
git --version
```

Already installed → `git version 2.x.x` → move to Step 2.

Not installed → a popup appears, click **Install**, wait, then re-run `git --version`.

You should see: `git version 2.x.x`

### 19.1.7. Step 2 — Check & Install Docker

#### NOTE

#### What is Docker?

Docker runs each tool in its own sealed container. Your Mac stays clean — nothing conflicts. Think of it as a lunchbox where each food has its own compartment.

```
docker --version
```

Already installed → confirm the Docker whale icon is in your menu bar → move to Step 3.

Not installed:

1. Go to [docker.com/products/docker-desktop](https://docker.com/products/docker-desktop)
2. Click **Download for Mac**
  - Apple M1/M2/M3 → **Apple Silicon**
  - Intel → **Intel Chip**
3. Open the `.dmg` → drag Docker to Applications
4. Open Docker → enter password if asked
5. Wait for the whale icon in menu bar to stop animating
6. Re-run `docker --version`

You should see: `Docker version 29.x.x` and the Docker whale icon in menu bar

### 19.1.8. Step 3 — Check & Install VS Code

#### What is VS Code?

##### NOTE

VS Code is a free code editor made by Microsoft. It is where you will read and write SQL, Python, and config files.

```
code --version
```

Already installed → move to installing extensions.

Not installed:

1. Go to [code.visualstudio.com](https://code.visualstudio.com)
2. Click **Download for Mac**
3. Open the `.zip` → drag VS Code to Applications
4. Open VS Code → press `Cmd + Shift + P` → type `shell command` → click **Install 'code' command in PATH**
5. Close and reopen Terminal → re-run `code --version`

Install VS Code extensions — open VS Code, click the Extensions icon, search and install:

- **dbt Power User** by Altimate AI
- **Python** by Microsoft
- **Docker** by Microsoft
- **Remote - SSH** by Microsoft
- **Claude Code** by Anthropic

You should see: `1.9x.x`

#### 19.1.9. Step 4 — Check & Install GitHub CLI

**NOTE**

**What is the GitHub CLI?**

A command-line tool for interacting with GitHub — cloning repositories, pushing code, and managing branches — without leaving Terminal.

```
gh --version
```

Already installed → move to Step 5.

Not installed:

```
brew install gh
```

Log in to GitHub:

```
gh auth login
```

Choose: `GitHub.com` → `HTTPS` → `Login with a web browser` → follow the prompts.

You should see: `Logged in to github.com account yourname`

### 19.1.10. Step 5 — Check & Install uv

#### NOTE

#### What is uv?

uv is a fast Python package manager used to install dbt on your Mac for local development. It is significantly faster than pip and handles virtual environments cleanly.

```
uv --version
```

Already installed → move to Step 6.

Not installed:

```
brew install uv
```

You should see: `uv 0.x.x`

### 19.1.11. Step 6 — Install dbt

#### NOTE

#### Why install dbt on your Mac AND inside Docker?

dbt runs natively on your Mac for development — writing

models, running tests, and checking output interactively. It also runs inside the Airflow container for scheduled automation. Both are needed and serve different purposes.

```
cd ~
uv venv .dbt-venv
source .dbt-venv/bin/activate
uv pip install dbt-core dbt-postgres
dbt --version
```

Activate the virtual environment every time you start working on a dbt project:

```
source ~/.dbt-venv/bin/activate
```

You should see: Core: 1.x.x

### 19.1.12. Step 7 — Configure Git

Set your Git identity — required for every commit:

```
git config --global init.defaultBranch main
git config --global user.name "Your Name"
git config --global user.email "you@example.com"
```

#### NOTE

#### Why `init.defaultBranch main`?

Git defaults to calling your first branch `master` — an old convention. GitHub changed their default to `main` in 2020. Running this once means every repo you create from now on starts on `main`.

#### Why `user.name` and `user.email`?

Every commit in Git is stamped with this name and email — who made the change and when. This is what appears in your

GitHub commit history.

**WARNING**

If you skip this step, your first commit will fail with "Author identity unknown."

### 19.1.13. Step 8 — Clone the Project Repository

**Two ways to use this repo.**

**NOTE**

**Mode 1 — Reference implementation:** Clone and run — everything is already built and working. Good for experienced practitioners who want to explore a working codebase.

**Mode 2 — Learn from scratch:** Clone, delete the `dbt/` folder, and follow the runbooks step by step to build everything yourself.

```
cd ~/Documents/btg-case-studies-with-dbt
git clone https://github.com/btg-case-studies-with-dbt/single-dbt-opensource
cd single-dbt-opensource
git checkout dev
git pull origin dev
```

Confirm you are in the right folder:

```
cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource && pwd
```

You should see: `.../single-dbt-opensource`

### 19.1.14. Step 9 — Create `.gitignore`

**NOTE**

**What is `.gitignore`?**

A file that tells Git which files to never track or commit. Secrets like passwords, generated files like compiled SQL, and operating

system files like `.DS_Store` should never go to GitHub. The `.gitignore` file handles all of this automatically.

```
cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource
cat > .gitignore << 'EOF'
# Secrets – never commit these
.env

# dbt – credentials
dbt/profiles.yml

# dbt – generated files
dbt/target/
dbt/dbt_packages/
dbt/logs/

# Airflow – runtime logs
airflow/logs/

# Python
__pycache__/*
*.pyc
.venv/
venv/

# OS
.DS_Store
*.swp
*.swo

# IDE
.vscode/
.idea/

# Docker
*.log
EOF

git add .gitignore
```

```
git commit -m "add .gitignore"
git push origin dev
```

You should see: `git status` shows `.gitignore` committed, no other files tracked.

### 19.1.15. Step 10 — Create Environment File

#### What is a `.env` file?

Holds private settings — passwords, port numbers — that Docker needs to start the stack. Lives only on your Mac, blocked by `.gitignore` so it never goes to GitHub.

#### NOTE

**What is `.env.example`?** A safe template with the same variable names but no real passwords. Committed to Git so teammates know exactly what to fill in.

```
cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource && pwd
```

Create the example template:

```
cat > .env.example << 'EOF'
# PostgreSQL – your data warehouse
POSTGRES_USER=mds_user
POSTGRES_PASSWORD=mds_password
POSTGRES_DB=btg_resource_utilization
POSTGRES_PORT=5432

# Airflow – your pipeline scheduler
AIRFLOW_UID=50000
AIRFLOW__WEBSERVER__SECRET_KEY=changeme123
EOF
```

#### What does each variable mean?

#### NOTE

`POSTGRES_USER` — the username `dbt`, Airflow, and Metabase use to connect to PostgreSQL. `mds_user` stands for Modern Data Stack user.

`POSTGRES_PASSWORD` — the password for that user.

`POSTGRES_DB` — the database name PostgreSQL creates on first start. Hard to change later — choose carefully.

`POSTGRES_PORT` — port 5432 is the PostgreSQL standard.

`AIRFLOW_UID` — tells Airflow which system user to run as inside Docker. Prevents permission errors.

`AIRFLOWWEBSERVERSECRET_KEY` — a private key Airflow uses to secure browser sessions. Any random string works for local development.

Copy it to create the real `.env`:

```
cp .env.example .env
```

Commit the template — not the real `.env`:

```
git add .env.example
git commit -m "add environment variable template"
git push origin dev
```

You should see: `git status` — `.env` does not appear anywhere. Only `.env.example` is tracked.

### 19.1.16. Step 11 — Create Dockerfile

#### What is a Dockerfile?

A recipe for building a custom Docker image. The official Airflow image has Airflow but not dbt. We install dbt inside the Airflow container so Cosmos can run dbt models as Airflow tasks automatically. dbt is installed here for Airflow automation only — you run dbt manually using the native installation from Step 6.

#### NOTE

```

cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource
cat > Dockerfile << 'EOF'
FROM apache/airflow:2.11.2

USER root
RUN python -m venv /usr/local/airflow/dbt_venv && \
    /usr/local/airflow/dbt_venv/bin/pip install --no-cache-dir \
    dbt-core \
    dbt-postgres

USER airflow

COPY airflow/requirements.txt .
RUN pip install --no-cache-dir -r requirements.txt
EOF

```

### 19.1.17. Step 12 — Create `airflow/requirements.txt`

#### What is `requirements.txt`?

A list of Python packages for Airflow. Installs everything at once when Docker builds the image.

#### NOTE

`apache-airflow-providers-postgres` — lets Airflow connect to PostgreSQL. Without this, Airflow cannot talk to your database.

`astronomer-cosmos` — connects dbt and Airflow, running dbt models as native Airflow tasks so you can schedule dbt through Airflow automatically.

```

cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource
mkdir -p airflow/dags airflow/logs airflow/plugins

cat > airflow/requirements.txt << 'EOF'
apache-airflow-providers-postgres
astronomer-cosmos
EOF

```

### 19.1.18. Step 13 — Create `docker-compose.yml`

#### NOTE

#### What is `docker-compose.yml`?

The file that defines all the services that make up your stack and how they connect to each other. One file. One command. The entire stack starts.

```
cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource
cat > docker-compose.yml << 'EOF'
services:

  postgres:
    image: postgres:17
    container_name: postgres
    restart: unless-stopped
    environment:
      POSTGRES_USER: ${POSTGRES_USER}
      POSTGRES_PASSWORD: ${POSTGRES_PASSWORD}
      POSTGRES_DB: ${POSTGRES_DB}
    ports:
      - "${POSTGRES_PORT}:5432"
    volumes:
      - postgres_data:/var/lib/postgresql/data

  airflow-init:
    build:
      context: .
      dockerfile: Dockerfile
    container_name: btg-airflow-init
    depends_on:
      postgres:
        condition: service_started
    environment:
      AIRFLOW__CORE__EXECUTOR: LocalExecutor
      AIRFLOW__DATABASE__SQL_ALCHEMY_CONN:
postgresql+psycopy2://mds_user:mds_password@postgres:5432/btg_resource_utilization
      AIRFLOW__WEBSERVER__SECRET_KEY: changeme123
      _AIRFLOW_DB_MIGRATE: 'true'
```

```

_AIRFLOW_WWW_USER_CREATE: 'true'
_AIRFLOW_WWW_USER_USERNAME: admin
_AIRFLOW_WWW_USER_PASSWORD: admin
command: version
restart: on-failure

airflow-webserver:
  build:
    context: .
    dockerfile: Dockerfile
  container_name: btg-airflow-webserver
  restart: unless-stopped
  depends_on:
    postgres:
      condition: service_started
  environment:
    AIRFLOW__CORE__EXECUTOR: LocalExecutor
    AIRFLOW__DATABASE__SQL_ALCHEMY_CONN:
postgressql+psycopy2://mds_user:mds_password@postgres:5432/btg_resource_utilization
    AIRFLOW__WEBSERVER__SECRET_KEY: changeme123
  ports:
    - "8080:8080"
  volumes:
    - ./airflow/dags:/opt/airflow/dags
    - ./airflow/logs:/opt/airflow/logs
    - ./airflow/plugins:/opt/airflow/plugins
    - ./dbt:/opt/airflow/dbt
    - ${HOME}/Documents/btg-case-studies-with-
dbt/common/database_scripts:/opt/airflow/database_scripts
  command: webserver

airflow-scheduler:
  build:
    context: .
    dockerfile: Dockerfile
  container_name: btg-airflow-scheduler
  restart: unless-stopped
  depends_on:
    postgres:
      condition: service_started

```

```

environment:
  AIRFLOW__CORE__EXECUTOR: LocalExecutor
  AIRFLOW__DATABASE__SQL_ALCHEMY_CONN:
postgresql+psycopy2://mds_user:mds_password@postgres:5432/btg_resource_utilization
  AIRFLOW__WEBSERVER__SECRET_KEY: changeme123
volumes:
  - ./airflow/dags:/opt/airflow/dags
  - ./airflow/logs:/opt/airflow/logs
  - ./airflow/plugins:/opt/airflow/plugins
  - ./dbt:/opt/airflow/dbt
  - ${HOME}/Documents/btg-case-studies-with-
dbt/common/database_scripts:/opt/airflow/database_scripts
command: scheduler

metabase:
  image: metabase/metabase:v0.59.2
  container_name: metabase
  restart: unless-stopped
  depends_on:
    postgres:
      condition: service_started
  ports:
    - "3000:3000"
  environment:
    MB_DB_TYPE: postgres
    MB_DB_DBNAME: metabase
    MB_DB_PORT: 5432
    MB_DB_USER: ${POSTGRES_USER}
    MB_DB_PASS: ${POSTGRES_PASSWORD}
    MB_DB_HOST: postgres
  volumes:
    - metabase_data:/metabase-data

pgadmin:
  image: dpage/pgadmin4:latest
  container_name: btg-pgadmin
  restart: unless-stopped
  depends_on:
    postgres:
      condition: service_started

```

```

environment:
  PGADMIN_DEFAULT_EMAIL: "admin@admin.com"
  PGADMIN_DEFAULT_PASSWORD: "admin"
ports:
  - "5050:80"
volumes:
  - pgadmin_data:/var/lib/pgadmin

volumes:
  postgres_data:
  metabase_data:
  pgadmin_data:
EOF

```

**WARNING**

Email addresses in YAML must be quoted — the @ symbol confuses the YAML parser.

**WARNING**

Airflow env vars with double underscores like AIRFLOW COREEXECUTOR conflict with how Docker Compose reads `.env` files. Write these directly in the compose file, not in `.env`.

```

git add .
git commit -m "add Dockerfile, docker-compose, requirements"
git push origin dev

```

**19.1.19. Step 14 — Build and Start the Stack**

Two commands — build then start. Because Airflow uses a custom Dockerfile, Docker must build the image first.

```

cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource
docker compose build

```

**NOTE**

Takes 3-5 minutes. Docker downloads Airflow and installs dbt inside it. Wait for the prompt to return before continuing.

```
docker compose up -d
```

**NOTE**

First run takes 5-10 minutes pulling PostgreSQL, Metabase, and pgAdmin images.

If you see "Found orphan containers":

```
docker compose up -d --remove-orphans
```

**Restarting after a machine restart:**

```
open -a Docker
```

Wait for the Docker whale icon to stop animating, then:

```
cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource
docker compose up -d
```

*19.1.20. Step 15 — Verify Everything Is Running*

```
cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource
docker compose ps
```

Every container must show Up:

NAME	STATUS
btg-airflow-scheduler	Up
btg-airflow-webserver	Up 0.0.0.0:8080->8080/tcp
btg-pgadmin	Up 0.0.0.0:5050->80/tcp
metabase	Up 0.0.0.0:3000->3000/tcp
postgres	Up 0.0.0.0:5432->5432/tcp

**NOTE**

btg-airflow-init not in the list is normal — it runs once on startup then exits cleanly.

If any container shows Exit or Restarting:

```
docker compose logs airflow-webserver --tail=30
```

Replace `airflow-webserver` with the failing service name.

Verify `dbt` is installed inside the Airflow container:

```
docker exec btg-airflow-scheduler \
  /usr/local/airflow/dbt_venv/bin/dbt --version
```

You should see: Core: 1.x.x

Create the Metabase database:

```
docker exec postgres psql -U mds_user \
  -d btg_resource_utilization \
  -c "CREATE DATABASE metabase;"

docker compose restart metabase
```

### What does this command do?

`docker exec postgres` runs a command inside the postgres container.

#### NOTE

`psql -U mds_user -d btg_resource_utilization` connects using your credentials.

`-c "CREATE DATABASE metabase;"` creates an empty database for Metabase to store its settings.

You should see: After 1-2 minutes, `docker compose ps` shows metabase as `Up`.

### 19.1.21. Step 16 — Open Airflow and Metabase

Tool	URL	Login
Airflow	localhost:8080	admin / admin
Metabase	localhost:3000	set up on first visit
pgAdmin	localhost:5050	<a href="#">admin@admin.com</a> / admin

**What is localhost?****NOTE**

"Localhost" means "this computer." Port numbers (:8080, :3000, :5050) are like apartment numbers — each tool lives on its own port.

**WARNING**

Metabase takes 1-2 minutes on first start. If you see a loading screen, wait and refresh.

**Register your database in pgAdmin**

1. Open localhost:5050 and log in
2. Right-click **Servers** → **Register** → **Server**
3. **General** tab: Name = `btg-local`
4. **Connection** tab:

Field	Value
Host	postgres
Port	5432
Maintenance database	btg_resource_utilization
Username	mds_user
Password	mds_password

1. Check **Save password** → Click **Save**

**Why is the host postgres and not localhost?****NOTE**

Inside Docker, containers find each other by container name — not by `localhost`. Our PostgreSQL container is named `postgres`.

You should see: In pgAdmin — **btg-local** → **Databases** → **btg\_resource\_utilization**.

```
git add .
```

```
git commit -m "stack running - all services confirmed up"
git push
```

**TIP**

You did it. A fully working modern data stack is running on your Mac. PostgreSQL, Airflow, Metabase, and pgAdmin — all confirmed up. This is the same pattern used at real companies, running at scale in the cloud.

### 19.1.22. Stopping and Starting the Stack

```
# Stop (data is saved)
cd ~/Documents/btg-case-studies-with-dbt/single-dbt-opensource
docker compose stop

# Start again
docker compose start
```

**CAUTION**

**Last resort only — wipes ALL data including the database**

```
docker compose down -v
```

You will need to re-run all data population scripts. Use this only if you want to start completely from scratch.

### 19.1.23. Troubleshooting

#### A container shows "Exit" instead of "Up"

```
docker compose logs airflow-webserver --tail=30
```

Replace `airflow-webserver` with the service that failed. The last few log lines almost always say exactly what went wrong.

#### Port already in use

```
docker compose stop && docker compose up -d
```

### Metabase stuck on loading screen

```
docker compose logs metabase --tail=20
```

If you see `database "metabase" does not exist` — you skipped the `CREATE DATABASE` step in Step 15. Run it now:

```
docker exec postgres psql -U mds_user -d btg_resource_utilization \
  -c "CREATE DATABASE metabase;"
docker compose restart metabase
```

### "cannot use SQLite with LocalExecutor" in Airflow logs

Check that `AIRFLOWDATABASESQL_ALCHEMY_CONN` is hardcoded in `docker-compose.yml` — not read from `.env`.

### dbt command not found

```
source ~/.zshrc
```

### "container name already in use"

```
docker rm container-name-from-error
```

Then retry `docker compose up -d`.

### Git says "Author identity unknown" on first commit

```
git config --global user.name "Your Name"
git config --global user.email "you@example.com"
```

Then retry your commit.

Continue to **Chapter 18 — Project Setup and dbt Configuration** to load

data and connect dbt to the stack.

## ***19.2. The AWS Stack — Next Edition***

The next edition will provide a complete setup guide for running the same analytical framework on AWS. The stack will cover data ingestion with AWS Glue, transformation with dbt, query execution with Amazon Athena, pipeline orchestration with MWAA, and visualization with Amazon QuickSight — designed to mirror the local stack pattern and transfer directly to a production cloud environment.

## Chapter 20. The Framework in Practice

Any framework has three elements: a clear input, a defined process, and a measurable output. For the product analytics framework proposed in this book, the input is business metrics, specifically success metrics — the definition of what good looks like for a product or initiative. The process is the eight stages of the analytics lifecycle — Define, Collect, Ingest, Transform, Implement, Visualize, Predict, and Iterate. The output is four types of analytics: descriptive, diagnostic, predictive, and prescriptive — delivered as dashboards that different audiences can act on. The following case studies put the framework to work. Each one walks through all eight stages — from defining success metrics to delivering dashboards. The stack and the architecture remain constant. What changes is the defined success metrics. This edition covers the resource utilization case study in full. Customer journey milestones and marketing attribution will be covered in the next edition.

### ***20.1. Case Study 1: Resource Utilization***

Viewed through the O2R (Observability-to-Revenue) Framework, this case study starts from Product Performance metrics — utilization rates, queue depth, processing latency, and failure counts — and traces how they shape customer behavior and, ultimately, revenue from the platform. By explicitly modeling how resource saturation shows up first as slower workflows, then as lower adoption and expansion, and finally as softer revenue, the team turns capacity planning from an infrastructure exercise into a product and revenue decision.

Resource utilization serves as the first case study — a concrete, data-rich

context for demonstrating how the Product Analytics Framework works end to end in practice. The case study walks through all four analytics types applied to a single infrastructure problem:

Resource utilization is the first case study. The complete data foundation — dbt models, tests, pipelines, and CI configuration — is available as a working GitHub repository [here](#).

### *20.1.1. Stage 1 — Define Success Metrics*

### *20.1.2. Stage 2 — Collect Usage Data*

The resource utilization project ingests data from multiple bronze source tables — each representing a different signal about how AI models are being deployed and consumed. Before any transformation can begin, the database must be set up, the bronze tables created, and the raw data loaded.

The project setup runbook covers: creating the PostgreSQL schemas and roles that enforce the medallion architecture, creating the bronze tables that hold raw infrastructure and usage data, loading the configuration tables via an Airflow DAG, and populating the remaining bronze tables with synthetic data using the data population scripts.

Key bronze tables for this case study:

- `config_model_dimensions` — AI model specifications: publisher, accelerator type, replica count, performance metrics
- `config_model_region_availability` — which models are deployed in which regions

- `resource_accelerator_inventory` — GPU and TPU accelerator inventory snapshots
- `resource_model_utilization` — utilization rates by model and region
- `resource_model_instance_allocation` — instance allocation counts
- `inference_user_token_usage_open_source` — per-request token usage for open source models
- `inference_user_token_usage_proprietary` — per-request token usage for proprietary models
- `revenue_account_daily` — daily revenue at account, model, and region grain
- `quota_default_rate_limits` — default rate limits per model per region
- `quota_customer_rate_limit_adjustments` — customer-specific rate limit overrides

### 20.1.3. Stage 3 — Ingest Data

With the right signals in place, the ingestion layer moves data reliably from source systems into the data platform. Two paths exist: telemetry data is ingested directly by data engineers raw and as-is, while metering data passes through a centralized billing system before being made available for analysis.

For this project, ingestion is handled in two ways. Static configuration data — model specifications and region availability — is loaded once via a manually triggered Airflow DAG. The remaining bronze data is loaded through SQL population scripts that simulate the streaming ingestion pattern used in production.

In production, the ingestion pipeline runs on a schedule: bronze layer ingestion on demand when source data changes, silver layer staging every six hours, and gold layer mart builds nightly. The Airflow setup covers the full production pipeline — two DAGs, one for the dev environment running at midnight and one for production running at 2:00am.

#### *20.1.4. Stage 4 — Transform Data*

Raw data in the bronze layer can answer ad-hoc questions but cannot power automated reporting at scale. The transformation layer uses dbt to move data through the medallion architecture — from raw bronze tables through cleaned silver views into business-ready gold mart tables.

The transformation runbooks cover the complete pipeline in sequence:

**Project setup and first run** — creating the dbt folder structure, configuring profiles for multiple environments, installing packages, registering all bronze tables as sources, running the full pipeline for the first time in the personal schema, running the first production build, and saving the production manifest for state-based selection.

**Staging models** — 12 views that clean and standardize each bronze table one to one. Timestamps converted, strings standardized, nulls handled, segments coalesced.

**Intermediate models** — 2 incremental tables that perform heavy aggregations between staging and marts. Token usage aggregated per minute. Revenue aggregated per day. Both use the delete+insert incremental strategy with a 7-day lookback window to handle late-arriving data.

**Mart models** — 10 tables that produce business-ready outputs across two domains: token usage analytics by model and customer at weekly, monthly, and YTD granularity; and revenue analytics by model and customer at the same granularity.

**Advanced transformation patterns** — ephemeral models for shared filter logic, model contracts for schema enforcement, model versioning for non-breaking evolution, access modifiers and groups for ownership boundaries, custom macros for reusable window function patterns.

#### *20.1.5. Stage 5 — Implement Metrics*

This is where success metrics from Stage 1 become operational — translated into precise, reusable calculations that power every report, dashboard, and analysis across the organization.

For the resource utilization project, metric implementation happens at two levels. At the model level, dbt enforces contracts — every output column is declared with its data type, and dbt validates the schema on every run. At the semantic level, metric definitions are standardized across environments so that Finance, Product, and GTM all pull the same numbers from the same source of truth.

The advanced runbooks cover the patterns that make metrics trustworthy at scale: singular tests for business rules, custom generic tests for reusable column-level validation, unit tests for transformation logic using mock data, and source freshness checks that prevent stale data from reaching the marts silently.

#### *20.1.6. Stage 6 — Visualize Metrics and Deep Dive*

### *20.1.7. Stage 7 — Predict and Prescribe*

### *20.1.8. Stage 8 — Act and Iterate*

The value proposition for the analytics framework is right here. Recommendations become decisions — adjusting capacity allocation, reallocating infrastructure spend, modifying provisioning policies. Since every stage before is fully automated, incremental changes feed back into the pipeline and surface in dashboards within the next scheduled run.

The CI/CD runbook covers the automation that makes iteration fast and safe: GitHub Actions running the modified model selection on every Pull Request, Airflow triggering the production pipeline on merge to main, and the development loop that keeps personal, dev, and production schemas in sync. This is the stage that closes the loop — making the framework self-sustaining and continuously improving as the business evolves.

## **20.2. Case Study 2: Customer Journey Milestones**

Through an O2R lens, the customer journey in this case study is not just a sequence of milestones; it is a chain of observability signals that either accelerate or stall progress. Time-to-first-value, step-level error rates, and drop-off points in the funnel become Product Performance metrics that predict whether customers will reach key milestones, renew, and expand — long before those outcomes appear in the revenue reports.

Understanding how customers move through a product — where they succeed, where they hesitate, and where they drop off — is one of the most valuable and underutilized sources of insight available to product and business teams.

This case study will walk through all four analytics types applied to the customer journey:

**Descriptive analysis** to map milestone completion rates across the onboarding sequence — from first sign-up through first meaningful usage to full adoption.

**Diagnostic analysis** to pinpoint drop-off causes — which steps in the journey are losing customers, which segments struggle most, and what behavioral patterns precede disengagement.

**Predictive analysis** to identify customers at risk of disengaging before the signal becomes visible in retention metrics — using early usage patterns to score adoption likelihood.

**Prescriptive analysis** to recommend the product and onboarding changes that improve activation and retention — which interventions work at each stage of the journey and which customer segments respond to which approaches.

When journey data is used well, the business impact is tangible — activation improves, churn drops, and feature roadmap decisions shift from opinion to evidence.

*This case study will be covered in the next edition.*

### **20.3. Case Study 3: Marketing Attribution and Growth**

In the context of O2R, marketing attribution is not only about which channel sourced the opportunity, but also about how reliably the end-to-end experience performs for each segment. Differences in latency,

availability, and onboarding friction across channels become observability signals that help explain why superficially similar campaigns produce different patterns of activation, retention, and revenue growth.

For a B2B SaaS business relying on marketing campaigns to drive subscription growth, knowing which campaigns work and which do not is not a reporting exercise — it is a revenue decision where every dollar of marketing spend has to earn its place.

This case study will walk through all four analytics types applied to campaign performance:

**Descriptive analysis** to measure performance across channels and campaigns — reach, conversion rates, cost per acquisition, and downstream revenue by source.

**Diagnostic analysis** to identify which messages, segments, and timing drove the strongest results — and why campaigns that looked similar on paper produced different outcomes.

**Propensity modelling** to score leads by their likelihood to convert, enabling sales and marketing teams to prioritize effort toward the opportunities most likely to close.

**A/B testing** to validate what actually moves the needle before scaling spend — confirming that observed differences between campaigns are real and not the product of chance or timing.

When campaign data is connected to downstream subscription metrics — activation, retention, and revenue — marketing decisions shift from gut feel to evidence and budget flows to what demonstrably drives growth.

*This case study will be covered in the next edition.*

## AUTHOR'S NOTE

---

This book evolved from my experience transforming data into business value across a variety of teams and organizations, specializing in SaaS Enterprise. These are lessons learned from successes and challenges across various data roles. It is a collection of patterns you can adapt to make better decisions in your own environment.

It offers practical approaches using the modern data stack that help leaders deliver measurable business impact. It introduces the data landscape, highlights key considerations at each stage of an analytics project, and provides guidance for building cross-functional teams and making data-informed decisions.

## KANJA SAHA

Data Professional · Analytics Leader · Author

A data professional with hands-on experience building AI-ready data and analytics foundations across startups and large organizations. Specializing in SaaS B2B products, Kanja leads cross-functional teams to align business metrics with data strategy, building scalable pipelines, robust models, and reporting frameworks that drive business impact.