

DuckDB

空间数据管理指南

从 SQL 基础到高级地理空间分析



吴秋生

DuckDB 空间数据管 理指南

从 SQL 基础到高级地理空间分析

吴秋生
2026

Contents

前言	1
引言	3
本书适合谁	3
本书内容概览	4
如何充分利用本书	5
本书使用的约定	6
下载代码示例	7
视频教程和补充资源	7
社区和反馈	8
致谢	8
关于作者	8
许可和版权	9
I: DuckDB 基础知识	11
1. DuckDB 入门	13
1.1. 引言	13
1.2. 学习目标	13
1.3. DuckDB 与传统数据库的不同之处	13
1.4. 何时使用（以及何时不使用）DuckDB 进行空间工作	15
1.5. 安装 DuckDB CLI 并运行您的第一个查询	16
1.6. 安装 DuckDB Python 客户端	18
1.7. 安装 Visual Studio Code	21
1.8. 使用 DuckDB UI	24
1.9. 安装 DBeaver SQL IDE	25
1.10. 关键要点	29
1.11. 练习	30
2. 空间分析的 SQL 基础	32
2.1. 简介	32
2.2. 学习目标	32
2.3. 示例数据集	32
2.4. 环境设置	33
2.5. 连接到 DuckDB	34
2.6. 安装扩展	34
2.7. 从 URL 读取 CSV 文件	35
2.8. 创建表以提高性能	35
2.9. SQL SELECT 语句	37
2.10. 使用 WHERE 子句筛选数据	44
2.11. 使用 LIKE 进行模式匹配	45
2.12. IN 运算符	46
2.13. BETWEEN 运算符	47
2.14. 使用 SQL 连接组合数据	47
2.15. 查询计划和性能	53
2.16. 聚合数据以获取摘要统计	53
2.17. 条件语句	57

2.18. 保存结果	58
2.19. 空间分析的高级 SQL 特性	59
2.20. SQL 注释和文档	63
2.21. 关键要点	64
2.22. 练习	65
3. DuckDB Python 集成	68
3.1. 简介	68
3.2. 学习目标	68
3.3. 示例数据集	69
3.4. 安装和设置	69
3.5. 安装和加载扩展	70
3.6. 从多个来源读取数据	71
3.7. 与 Pandas DataFrame 的无缝集成	74
3.8. Polars 互操作性	77
3.9. 结果转换和输出格式	77
3.10. 将数据写入磁盘	80
3.11. 持久化存储和数据库文件	80
3.12. 预处理语句和参数	82
3.13. 关键要点	83
3.14. 练习	84
II: 空间数据操作	89
4. 加载空间数据格式	91
4.1. 简介	91
4.2. 学习目标	91
4.3. 示例数据集	92
4.4. 安装和设置	93
4.5. 安装和加载扩展	93
4.6. 下载示例数据	94
4.7. 加载带有坐标的 CSV 文件	94
4.8. 加载 JSON 文件	97
4.9. 直接查询 Pandas DataFrame	99
4.10. 加载 Parquet 文件以提高性能	100
4.11. 加载带有空间几何图形的 GeoJSON 文件	103
4.12. 将 Shapefile 加载到现代工作流程中	108
4.13. 加载 GeoParquet 用于云原生空间分析	110
4.14. 数据加载性能策略	112
4.15. 排查常见数据加载问题	112
4.16. 关键要点	113
4.17. 练习	114
5. 导出和转换空间数据	118
5.1. 简介	118
5.2. 学习目标	119
5.3. 示例数据集	119
5.4. 安装和设置	120
5.5. 安装和加载扩展	121

5.6. 加载示例数据	122
5.7. 导出到 Pandas DataFrame	123
5.8. 导出到 CSV 文件	128
5.9. 导出到 JSON 文件	131
5.10. 导出到 Excel 文件	133
5.11. 导出到 Parquet 文件	135
5.12. 导出到 GeoJSON 格式	137
5.13. 导出到 Shapefile 格式	141
5.14. 导出到 GeoPackage 格式	142
5.15. 核心要点	144
5.16. 练习	145
6. 几何操作和函数	150
6.1. 简介	150
6.2. 学习目标	150
6.3. 示例数据集	151
6.4. 安装和设置	152
6.5. 连接到 DuckDB 并加载扩展	152
6.6. 理解几何类型	154
6.7. 使用点	157
6.8. 使用线串	159
6.9. 使用多边形	161
6.10. 使用集合	163
6.11. 可视化纽约市空间数据	164
6.12. 几何处理	167
6.13. 几何有效性和鲁棒性	170
6.14. 关键要点	172
6.15. 练习	173
7. 空间查询和关系	177
7.1. 简介	177
7.2. 学习目标	177
7.3. 示例数据集	178
7.4. 安装和设置	178
7.5. 连接到 DuckDB 并加载扩展	179
7.6. 理解空间关系	180
7.7. 测试几何同一性	180
7.8. 拓扑关系	182
7.9. 基于距离的关系	190
7.10. 按接近阈值过滤	192
7.11. 最近邻查询	193
7.12. 要点总结	193
7.13. 练习	195
8. 高级空间连接	197
8.1. 简介	197
8.2. 学习目标	197
8.3. 示例数据集	198

8.4. 安装	198
8.5. 库导入和配置	198
8.6. 连接到 DuckDB	199
8.7. 相交连接	200
8.8. 距离范围内连接	204
8.9. 高级连接	207
8.10. 坐标系统转换	209
8.11. 空间关系函数参考	212
8.12. 关键要点	212
8.13. 练习	214
9. 交互式数据可视化	218
9.1. 简介	218
9.2. 学习目标	218
9.3. 示例数据集	219
9.4. 安装和设置	219
9.5. 下载示例数据	220
9.6. 连接到 DuckDB 并加载扩展	220
9.7. 可视化点数据	221
9.8. 可视化线数据	224
9.9. 可视化多边形数据	226
9.10. 关键要点	231
9.11. 练习	232
10. 矢量切片与 PMTiles	236
10.1. 简介	236
10.2. 学习目标	237
10.3. 示例数据集	237
10.4. 安装和设置	238
10.5. 直接从文件可视化矢量切片	239
10.6. 将矢量数据转换为 PMTiles	243
10.7. 可视化 PMTiles	244
10.8. 关键要点	254
10.9. 练习	256
III: 实际地理空间分析案例	261
11. 分析美国国家湿地清查数据	263
11.1. 引言	263
11.2. 学习目标	263
11.3. 本章使用的数据集	264
11.4. 理解数据集来源	266
11.5. 使用 DuckDB 访问湿地数据	268
11.6. 可视化湿地分布	271
11.7. 全国规模湿地分析	274
11.8. 关键要点	287
11.9. 练习	289
12. 分析全球建筑物轮廓数据	291
12.1. 引言	291

12.2. 学习目标	292
12.3. 关于数据集	292
12.4. 安装和设置	294
12.5. 安装和加载扩展	295
12.6. 探索可用数据	295
12.7. 使用边界框过滤的区域分析	298
12.8. 多区域比较	307
12.9. 数据源过滤和质量评估	311
12.10. 基于网格的空间聚合	316
12.11. 使用 H3 六边形网格聚合建筑物数据	320
12.12. 关键要点	330
12.13. 练习	331
13. 分析纽约市出租车数据	334
13.1. 引言	334
13.2. 学习目标	334
13.3. 关于数据集	335
13.4. 安装	337
13.5. 库导入	337
13.6. 安装和加载扩展	337
13.7. 加载出租车数据	337
13.8. 时间分析	341
13.9. 加载出租车区域查询数据	347
13.10. 空间分析	350
13.11. 行程流动分析	357
13.12. 支付和经济分析	361
13.13. 乘客行为分析	365
13.14. 多月分析	366
13.15. 可视化	368
13.16. 性能优化技巧	373
13.17. 关键要点	374
13.18. 练习	376
14. 使用 Voila 和 Solara 开发交互式仪表板	378
14.1. 简介	378
14.2. 学习目标	379
14.3. 安装 Voila 和 Solara	379
14.4. Hugging Face Spaces 简介	380
14.5. 创建基本的 Voila 应用程序	381
14.6. 使用 Solara 创建高级 Web 应用程序	391
14.7. 关键要点	398
14.8. 练习	398

前言

引言

在这个日益数据驱动的世界中，有效管理和分析空间信息的能力变得至关重要。从城市规划和环境监测到物流和个性化位置服务，地理空间数据构成了众多影响我们日常生活应用的基础。然而，处理空间数据通常被视为一个专业且复杂的领域，需要复杂的工具和陡峭的学习曲线。

DuckDB 的出现改变了这一切。

DuckDB 是一个创新的分析型数据库，专为高效性和易用性而设计。作为一个进程内 OLAP（联机分析处理）数据库，它直接在您的应用程序内运行，无需单独的服务器部署和繁琐的配置。这种嵌入式特性，加上其列式架构和向量化执行引擎，使 DuckDB 在处理大型数据集的分析查询时表现出色。DuckDB 最初因其通用数据处理能力而广受欢迎，其快速扩展的生态系统（尤其是空间数据扩展）代表了一次变革性的转变。

本书《*DuckDB 空间数据管理指南：从 SQL 基础到高级地理空间分析*》旨在揭开地理空间数据的神秘面纱，展示 DuckDB 如何以前所未有的简便性和速度，让每个人（从数据分析师和科学家到开发人员和 GIS 专业人员）都能利用其强大功能。我们相信，强大的空间分析应该是人人可及的，而不应局限于昂贵的专业软件或复杂的编程语言。借助 DuckDB，这种可及性成为现实。

我们的旅程从空间数据的基本概念及其表示方式开始，为使用 SQL 处理点、线和多边形奠定坚实的基础。随着学习的深入，您将发现 DuckDB 的原生地理空间功能（通过其兼容 PostGIS 的扩展增强）如何通过优雅的 SQL 查询实现复杂的操作，如空间连接、缓冲区分析和最近邻搜索。我们将探索各种实际应用，演示如何加载、转换、分析和可视化空间数据集，使您能够从地理信息中提取有意义的洞察。

无论您是想将空间分析集成到数据管道中、执行快速的临时地理空间查询，还是开发交互式的位置感知应用程序，本书都将作为您的全面指南。我们将涵盖从设置 DuckDB 环境和导入各种空间文件格式（如 Shapefile、GeoJSON 和 GeoParquet）到执行复杂分析任务和与可视化工具集成等主题。

我们的目标不仅仅是教您语法，而是培养您对这些工具和技术为何强大的理解。在本书结束时，您将能够熟练地使用 DuckDB 作为空间数据管理和分析的首选引擎，为您的项目解锁新的可能性，并使您能够做出明智的、具有空间意识的决策。

请与我们一起深入探索 DuckDB 分析能力与丰富地理空间数据世界的精彩交汇点。可及的空间分析的未来已经到来，它运行在 DuckDB 上。

本书适合谁

本书专为那些正在应对现代空间数据分析复杂性的人设计。如果您曾经花费数小时等待空间连接完成、难以将大型地理数据集加载到内存中，或者希望有一种更简单的方式将 SQL 的强大功能与空间操作结合起来，那么这本书正适合您。

如果您是以下人员，将获得最大价值

GIS 专业人员：对桌面软件处理大型数据集的局限性感到沮丧。您熟悉 QGIS 或 ArcGIS，但需要分析数百万个要素、处理大量 GPS 轨迹，或将空间分析集成到自动化工作流程中。

数据科学家或分析师：经常遇到位置数据。您熟悉 Python 和 pandas，但空间数据常常让您感到困惑。您希望将地理维度纳入分析，而无需深入学习复杂的 GIS 软件。

软件开发人员: 正在构建包含空间功能的应用程序。您需要快速的空间查询，希望避免繁重的数据基础设施，并且更喜欢使用熟悉的 SQL 而不是专门的空间库。

研究人员或学者: 从事地理学、环境科学或城市规划等领域的研究。您的研究涉及大型空间数据集，需要可重复、可扩展的分析方法，以适应不断增长的数据量。

商业智能专业人员: 处理基于位置的业务数据。无论是门店位置、配送路线、客户区域还是房地产投资组合，您都需要将业务指标与空间洞察相结合。

基本前提条件

您应该熟悉以下内容：

- **Python 编程**: 理解变量、函数以及如何导入库（不需要专家水平）
- **数据分析概念**: 筛选记录、聚合数据和连接表
- **SQL 基础**: 熟悉 SELECT、WHERE 和 GROUP BY 子句（我们将介绍空间方面的内容）
- **空间数据基础**: 理解数据具有位置属性（经纬度、投影）

有帮助的背景知识（非必需）

- 使用 pandas、GeoPandas 或 Jupyter notebooks 的经验
- 之前接触过数据库或数据仓库
- 熟悉 GIS 软件（QGIS、ArcGIS、PostGIS）
- 了解空间文件格式（GeoJSON、Shapefiles、Parquet）

如果您是 Python 编程新手

如果您是地理空间 Python 编程的新手，以下书籍为 GIS 基础概念和 Python 编程提供了很好的入门介绍：

Wu, Q. (2025). *Introduction to GIS Programming: A Practical Python Guide to Open Source Geospatial Tools*. Independently published. ISBN 979-8286979455. <https://www.amazon.com/dp/B0FFW34LL3>

本书内容概览

本书提供了从 SQL 基础到高级地理空间分析的系统化学习路径，通过实际案例帮助您掌握实用技能。每章都从简单查询逐步过渡到复杂的空间分析，逐步提升您在现代地理空间数据管理方面的专业能力。

第一部分：DuckDB 基础（第 1-3 章）

掌握所有后续内容的基本概念：

- **第 1 章：DuckDB 入门**: 安装、初始查询，以及了解 DuckDB 如何革新空间分析。
- **第 2 章：空间分析的 SQL 基础**: 针对空间数据的筛选、聚合、连接和查询优化的关键 SQL 模式。
- **第 3 章：DuckDB Python 集成**: 将 SQL 的强大功能与 pandas 的灵活性相结合，打造无缝的空间分析工作流程。

完成第一部分后，您将能够自信地查询空间数据集，并将 DuckDB 集成到任何基于 Python 的分析管道中。

第二部分：空间数据操作（第 4-10 章）

深入了解核心空间工具包，涵盖从数据加载到高级分析的所有内容：

- **第 4 章：加载空间数据格式：**导入各种格式，包括 CSV 坐标、来自 API 的 GeoJSON、大型 Shapefile 和云托管的 GeoParquet。
- **第 5 章：导出和转换空间数据：**将结果转换为利益相关者所需的任何格式。
- **第 6 章：几何操作和函数：**使用 SQL 函数创建、测量和转换空间要素。
- **第 7 章：空间查询和关系：**掌握表连接的空间等价物：包含、相交和邻近关系。
- **第 8 章：高级空间连接：**按位置而非 ID 组合数据集，这是空间分析的精髓。
- **第 9 章：交互式数据可视化：**生成引人注目的地图和图表，有效传达数据的空间叙事。
- **第 10 章：矢量切片与 PMTiles：**部署能够流畅处理数百万要素的交互式地图。

完成第二部分后，您将能够熟练管理任何空间数据格式，执行复杂操作，并创建专业级可视化。

第三部分：实际地理空间分析（第 11-14 章）

探索四个使用大规模真实数据集的综合案例研究：

- **第 11 章：分析美国国家湿地清查数据：**对全美 50 个州进行环境分析，处理数百万个湿地多边形。
- **第 12 章：分析全球建筑物轮廓数据：**使用 Overture Maps 的全球建筑物轮廓数据分析城市数据。
- **第 13 章：分析纽约市出租车行程数据：**从数亿次出租车行程中发现时空模式，揭示城市出行的洞察。
- **第 14 章：使用 Voilà 和 Solara 开发交互式仪表板：**构建和部署 Web 应用程序，使您的分析成果可供利益相关者访问。

完成第三部分后，您将拥有展示高级空间分析能力的优质项目作品集。

贯穿全书的主题

- **大规模性能：**无论处理数千还是数百万空间要素都有效的技术。
- **云原生工作流程：**直接从 S3 处理数据，与现代数据栈无缝集成。
- **可重复分析：**可版本控制并部署到生产环境的可共享代码和方法。
- **真实数据挑战：**解决投影问题、缺失值和数据质量问题。
- **集成模式：**将 DuckDB 与更广泛的 Python 地理空间生态系统相结合，以增强功能。

本书的独特之处

与理论讨论或特定工具教程不同，本书强调解决实际问题。每种技术都植根于实际的分析挑战，使用真实数据集进行演示，并清晰解释何时以及为何使用它。

如何充分利用本书

为了最大化您使用本书的学习体验，请考虑以下建议：

建立适当的开发环境：按照第 1 章的说明安装 Python 和所需的库。配置良好的环境将在整个学习过程中为您节省时间和减少挫折。考虑使用 conda 或 uv 来管理您的 Python 包，因为这简化了地理空间库的安装。

跟随代码示例练习：本书设计为互动式的。不要只是阅读代码；亲自输入、运行并尝试修改。理解来自实践，每个示例都在培养您以后需要的技能。

完成练习：每章都包含旨在巩固所学概念的练习。这些不是可选的附加内容；它们是学习过程的重要组成部分。从指导练习开始，然后用自己的项目挑战自己。

使用真实数据：虽然本书为示例和练习提供了数据集，但请尝试将这些技术应用于您自己领域或兴趣的数据。这将帮助您理解概念如何应用于实际场景，并建立对自己能力的信心。

构建项目：随着您在本书中的进展，考虑开展一个您感兴趣的个人项目。这可以是分析您研究中的数据、为您的社区创建地图，或解决您在工作中遇到的问题。

对自己有耐心：编程可能令人沮丧，尤其是在学习时。预期会遇到错误、花时间调试，偶尔会感到困惑。这是正常的，也是学习过程的一部分。需要时休息一下，记住专业技能是通过持续练习逐渐发展的。如果遇到困难，不要犹豫在本书的 GitHub 仓库中寻求帮助。

持续练习：本书中的技能需要定期练习来保持和发展。定期安排时间进行地理空间编程项目，即使是小项目也可以。

本书使用的约定

本书使用了几种约定来帮助您浏览内容并理解代码示例：

代码格式：所有 Python 代码都以等宽字体显示在代码块中。当代码出现在正文中时，格式为这样。文件和目录名也使用等宽字体格式。

代码示例：大多数代码示例都是完整且可运行的。它们包含解释所演示的关键概念和技术的注释。可能会包含行号以供正文参考。

```
# 这是一个代码块示例
import leafmap
m = leafmap.Map()
m.add_basemap("OpenTopoMap") # 向地图添加底图
m
```

SQL 风格指南：为了一致性和可读性，SQL 示例遵循以下模式：

- **关键字大写：** SELECT、 FROM、 WHERE、 JOIN
- **函数名保留大小写：** ST_Area()、 read_csv_auto()
- **表名和列名小写：** cities、 population
- **使用缩进提高可读性：**多行查询经过格式化以便清晰阅读

```
SELECT name, ST_Area(geometry) as area
FROM neighborhoods
WHERE borough = 'Manhattan'
ORDER BY area DESC;
```

命令行指令：在命令行或终端输入的命令以 \$ 提示符显示（不要输入 \$ 符号本身）：

```
$ pip install leafmap  
$ python script.py
```

下载代码示例

本书的所有代码示例、数据集和补充材料都可在 GitHub 上免费获取：

<https://github.com/giswqs/duckdb-spatial>

要下载材料，您可以使用以下方法之一：

- **克隆仓库**（如果您已安装 Git）：

```
$ git clone https://github.com/giswqs/duckdb-spatial.git
```

- **下载 ZIP 文件**（如果您不想使用 Git）：

- 访问 GitHub 仓库页面
 - 点击绿色的 **Code** 按钮
 - 选择 **Download ZIP**
 - 将文件解压到您喜欢的位置
- 如果您只需要特定示例，可以通过 GitHub 界面[在线浏览单个文件](#)

仓库会定期更新修正、改进和额外示例。请定期查看更新，或在 GitHub 上 **watch** 该仓库以获得变更通知。

如果您在代码中发现错误或有改进建议，请在 GitHub 上提交 issue 或 pull request。社区贡献有助于使这个资源对每个人都更好。

视频教程和补充资源

除了书面内容，本书还配有一系列全面的视频教程，讲解关键概念并提供额外示例：

<https://tinyurl.com/duckdb-spatial-videos>

视频旨在补充而非替代书面材料。它们对以下情况特别有帮助：

- 视觉型学习者，通过观看代码编写和执行获益
- 通过多种解释理解复杂概念
- 学习开发工作流程和最佳实践
- 了解如何处理问题和调试

播放列表按照本书的结构组织。您可以在阅读本书时按顺序观看，或根据需要跳转到特定主题。

这些视频是我在 2023 年秋季于田纳西大学教授**空间数据管理**¹ 课程时创建的。虽然课程已经结束，但视频仍然具有参考价值，可以作为本书的参考资料。未来还会添加更多视频。

¹<https://geog-414.gishub.org>

社区和反馈

我欢迎读者的反馈、问题和建议。您的意见有助于改进本书，使其对地理空间编程社区更有用。

书籍相关问题和讨论：

- GitHub Issues：<https://github.com/giswqs/duckdb-spatial/issues>
- GitHub Discussions：<https://github.com/giswqs/duckdb-spatial/discussions>

特别有帮助的反馈类型：

- 文本或代码中的错误或不清楚的解释
- 对额外示例或用例的建议
- 新主题或章节的想法
- 不同操作系统或库版本的兼容性问题报告
- 您如何应用本书技术的成功案例

致谢

本书的完成要感谢许多人的贡献以及更广泛的开源地理空间社区。

首先，我要感谢 DuckDB 开发团队 创建了如此出色的数据库系统。他们对性能、简洁性和开源原则的承诺使空间分析对更广泛的受众变得可及。特别感谢为 DuckDB 空间功能做出贡献的团队成员。

我感谢田纳西大学的 同事和学生们，他们通过 空间数据管理 课程提供了反馈、测试了示例，并帮助完善了内容。他们的问题和见解使本书更加完善。

开源地理空间社区 值得特别认可。像 GDAL、GeoPandas、Shapely 等项目以及无数其他项目构成了使现代空间分析成为可能的基础。这个社区的协作精神继续激励着我的工作。

感谢 早期读者 对草稿章节提供的反馈，并帮助确定需要澄清或改进的地方。他们多样化的视角（从经验丰富的 GIS 专业人员到数据科学新手）帮助确保本书服务于其目标受众。

我要感谢 我的家人 在我花费许多晚上和周末写作时的耐心和支持。他们的理解和鼓励使这个项目成为可能。

最后，感谢 您，读者，对空间数据分析和开源工具的兴趣。正是像您这样的从业者推动了创新，使地理空间领域成为一个令人兴奋的工作场所。

如果本书能帮助您更有效地解决空间问题，那么所有的努力都是值得的。

关于作者

吴秋生博士是田纳西大学诺克斯维尔分校地理与可持续发展系副教授，同时也是亚马逊学者。吴博士的研究专注于通过云计算和 GeoAI 推进开源地理空间分析。他是多个广泛使用的开源 Python 包的创建者和维护者，包括 Geemap²、Leafmap³、SAMGeo⁴ 和 GeoAI⁵，这些工具将基于云的地理空间平台与 AI 驱动的分析和可视化相集成。吴博士的工作连接了遥感、地球观测和人工智能，使大规

²<https://geemap.org>

³<https://leafmap.org>

⁴<https://samgeo.gishub.org>

⁵<https://opengeoai.org>

模地理空间数据对全球的研究人员、教育工作者和从业者更加可访问、可重复和智能化。他的开源项目可以在 GitHub 上找到：<https://github.com/opengeos>。

许可和版权

本书秉承开放科学和开放教育的原则。为了支持透明度、学习和再利用，本书中的 **代码示例** 采用 [知识共享署名 4.0 国际许可协议 \(CC BY 4.0\)](#) 发布。这意味着您可以自由复制、修改和分发代码，甚至用于商业目的，只要给予适当的署名。

请通过引用本书或链接到 GitHub 仓库来注明代码使用的归属：

Wu, Q. (2025). *Spatial Data Management with DuckDB: From SQL Basics to Advanced Geospatial Analytics*. Independently published. PDF edition ISBN 979-8993859705; Print edition ISBN 979-8274710572. <https://duckdb.gishub.org>

虽然代码可以免费获取，但本书中的 **文本、图表和图像受作者版权保护**，未经明确许可不得复制、再分发或修改。这包括所有书面内容、自定义图表和嵌入的可视化，除非另有说明。

如果您希望重新使用或改编本书中的任何非代码材料（例如用于教学、演示或出版），请联系作者申请许可。

这种双重许可方法有助于平衡学习材料的开放获取与原创作品的保护。感谢您尊重这些条款并支持开源地理空间社区。