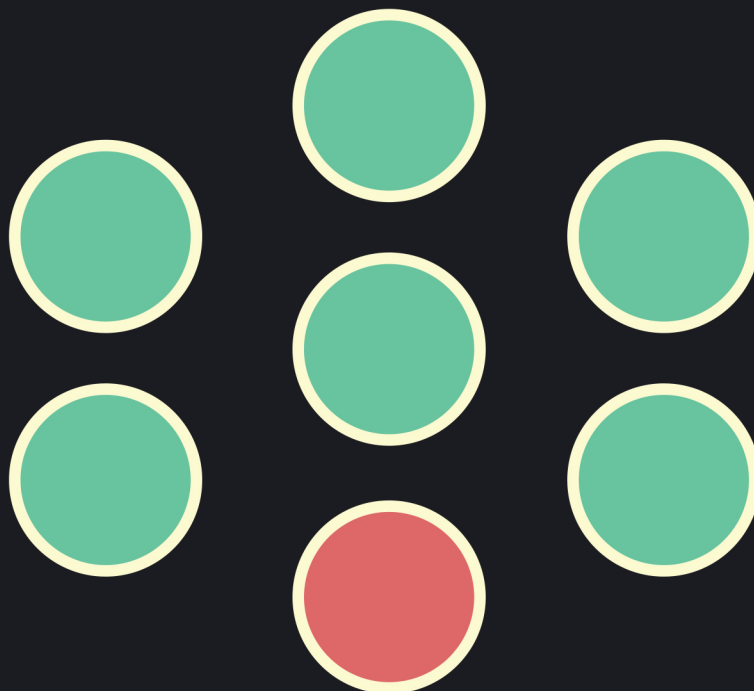


Venelin Valkov's

# Hands-On Machine Learning from Scratch



# Hands-On Machine Learning from Scratch

Develop a Deeper Understanding of Machine Learning Models by Implementing Them from Scratch in Python

Venelin Valkov

This book is for sale at <http://leanpub.com/hmls>

This version was published on 2019-06-29



This is a [Leanpub](#) book. Leanpub empowers authors and publishers with the Lean Publishing process. [Lean Publishing](#) is the act of publishing an in-progress ebook using lightweight tools and many iterations to get reader feedback, pivot until you have the right book and build traction once you do.

© 2019 Venelin Valkov

# **Tweet This Book!**

Please help Venelin Valkov by spreading the word about this book on [Twitter](#)!

The suggested hashtag for this book is [#hmls](#).

Find out what other people are saying about the book by clicking on this link to search for this hashtag on Twitter:

[#hmls](#)

# Contents

<b>Color palette extraction with K-means clustering</b> . . . . .	<b>1</b>
Unsupervised Learning . . . . .	1
The Data . . . . .	2
What is K-Means Clustering? . . . . .	2
Data Preprocessing . . . . .	3
Implementing K-Means clustering . . . . .	5
Evaluation . . . . .	6
Conclusion . . . . .	9
<b>Movie review sentiment analysis with Naive Bayes</b> . . . . .	<b>10</b>
Dealing with Text . . . . .	10
Naive Bayes . . . . .	13
Implementing Multinomial Naive Bayes . . . . .	14
Predicting sentiment . . . . .	17
Conclusion . . . . .	19

# Color palette extraction with K-means clustering

Choosing a color palette for your next big mobile app (re)design can be a daunting task, especially when you don't know what the heck you're doing. How can you make it easier (asking for a friend)?

One approach is to head over to a place where the *PROs* share their work. Pages like [Dribbble](https://dribbble.com/)<sup>1</sup>, [uplabs](https://www.uplabs.com/)<sup>2</sup> and [Behance](https://www.behance.net/)<sup>3</sup> have the goods.

After finding mockups you like, you might want to extract colors from them and use those. This might require opening specialized software, manually picking color with some tool(s) and other over-the-counter hacks. Let's use Machine Learning to make your life easier.

[Complete source code notebook](#)<sup>4</sup>

## Unsupervised Learning

Up until now we only looked at models that require training data in the form of features and labels. In other words, for each example, we need to have the correct answer, too.

Usually, such training data is hard to obtain and requires many hours of manual work done by humans (yes, we're already serving "The Terminators"). Can we skip all that?

Yes, at least for some problems we can use example data without knowing the correct answer for it. One such problem is *clustering*.

## What is clustering?

Given some vector of data points  $X$ , clustering methods allow you to put each point in a group. In other words, you can categorize a set of entities, based on their properties, automatically.

Yes, that is very useful in practice. Usually, you run the algorithm on a bunch of data points and specify how much groups you want. For example, your inbox contains two main groups of e-mail: spam and not-spam (were you waiting for something else?). You can let the clustering algorithm create two groups from your emails and use your beautiful brain to classify which is which.

More applications of clustering algorithms:

---

<sup>1</sup><https://dribbble.com/>

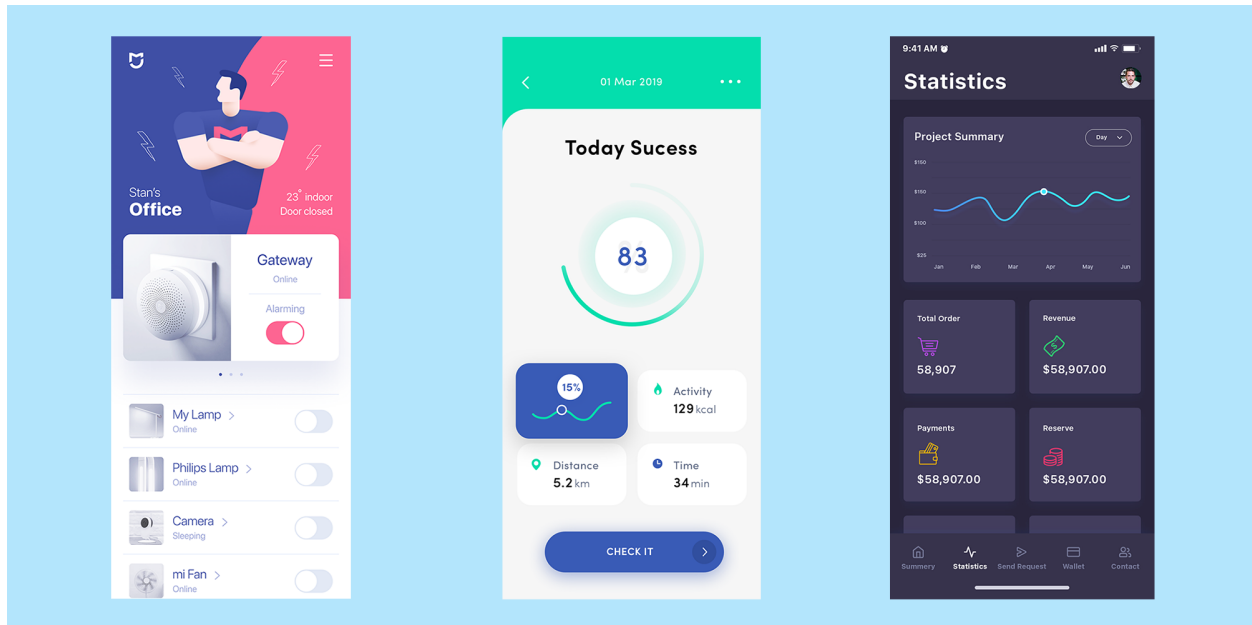
<sup>2</sup><https://www.uplabs.com/>

<sup>3</sup><https://www.behance.net/>

<sup>4</sup>[https://colab.research.google.com/drive/1\\_p1nptDfvJhcSvprmi1WtoIugsCI40Xa](https://colab.research.google.com/drive/1_p1nptDfvJhcSvprmi1WtoIugsCI40Xa)

- *Customer segmentation* - find groups of users that spend/behave the same way
- *Fraudulent transactions* - find bank transactions that belong to different clusters and identify them as fraudulent
- *Document analysis* - group similar documents

## The Data



source:

various authors on <https://www.uplabs.com/>

This time, our data doesn't come from some predefined or well-known dataset. Since Unsupervised learning does not require labeled data, the Internet can be your oyster.

Here, we'll use 3 mobile UI designs from various authors. Our model will run on each shot and try to extract the color palette for each.

## What is K-Means Clustering?

K-Means Clustering is defined by [Wikipedia](#)<sup>5</sup> as:

k-means clustering is a method of vector quantization, originally from signal processing, that is popular for cluster analysis in data mining. k-means clustering aims to partition  $n$  observations into  $k$  clusters in which each observation belongs to the cluster with the nearest mean, serving as a prototype of the cluster. This results in a partitioning of the data space into Voronoi cells.

<sup>5</sup>[https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering)

Wikipedia also tells us that solving K-Means clustering is difficult (in fact, [NP-hard](https://en.wikipedia.org/wiki/NP-hardness)<sup>6</sup>) but we can find local optimum solutions using some heuristics.

But how do *K-Means Clustering* works?

Let's say you have some vector  $X$  that contains  $n$  data points. Running our algorithm consists of the following steps:

1. Take random  $k$  points (called *centroids*) from  $X$
2. Assign every point to the closest centroid. The newly formed bunch of points is called *cluster*.
3. For each cluster, find new centroid by calculating a new center from the points
4. Repeat steps 2-3 until centroids stop changing

Let's see how can we use it to extract color palettes from mobile UI screenshots.

## Data Preprocessing

Given our data is stored in raw pixels (called images), we need a way to convert it to points that our clustering algorithm can use.

Let's first define two classes that represent a point and cluster:

```
1 class Point:
2
3     def __init__(self, coordinates):
4         self.coordinates = coordinates
```

Our Point is just a holder to the coordinates for each dimension in our space.

```
1 class Cluster:
2
3     def __init__(self, center, points):
4         self.center = center
5         self.points = points
```

The Cluster is defined by its center and all other points it contains.

Given a path to image file we can create the points as follows:

---

<sup>6</sup><https://en.wikipedia.org/wiki/NP-hardness>

```
1 def get_points(image_path):
2     img = Image.open(image_path)
3     img.thumbnail((200, 400))
4     img = img.convert("RGB")
5     w, h = img.size
6
7     points = []
8     for count, color in img.getcolors(w * h):
9         for _ in range(count):
10             points.append(Point(color))
11
12     return points
```

A couple of things are happening here:

- load the image into memory
- resize it to smaller image (mobile UX requires big elements on the screen, so we aren't losing much color information)
- drop the alpha (transparency) information

Note that we're creating a `Point` for each pixel in our image.

Alright! You can now extract points from an image. But how can we calculate the distance between points in our clusters?

## Distance function

Similar to the cost function in our supervised algorithm examples, we need a function that tells us how well we're doing. The objective of our algorithm is to minimize the distance between the points in each centroid.

One of the simplest distance functions we can use is the Euclidean distance, defined by:

$$d(p, q) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

where  $p$  and  $q$  are two points from our space.

Note that while Euclidean distance is simple to implement it might not be the best way to calculate the [color difference](https://en.wikipedia.org/wiki/Color_difference)<sup>7</sup>.

Here is the Python implementation:

---

<sup>7</sup>[https://en.wikipedia.org/wiki/Color\\_difference](https://en.wikipedia.org/wiki/Color_difference)

```
1 def euclidean(p, q):
2     n_dim = len(p.coordinates)
3     return sqrt(sum([
4         (p.coordinates[i] - q.coordinates[i]) ** 2 for i in range(n_dim)
5     ]))
```

## Implementing K-Means clustering

Now that you have all pieces of the puzzle you can implement the K-Means clustering algorithm. Let's start with the method that finds the center for a set of points:

```
1 def calculate_center(self, points):
2     n_dim = len(points[0].coordinates)
3     vals = [0.0 for i in range(n_dim)]
4     for p in points:
5         for i in range(n_dim):
6             vals[i] += p.coordinates[i]
7     coords = [(v / len(points)) for v in vals]
8     return Point(coords)
```

To find the center of a set of points, we add the values for each dimension and divide by the number of points.

Now for finding the actual clusters:

```
1 def assign_points(self, clusters, points):
2     plists = [[] for i in range(self.n_clusters)]
3
4     for p in points:
5         smallest_distance = float('inf')
6
7         for i in range(self.n_clusters):
8             distance = euclidean(p, clusters[i].center)
9             if distance < smallest_distance:
10                 smallest_distance = distance
11                 idx = i
12
13     plists[idx].append(p)
14
15     return plists
16
```

```

17 def fit(self, points):
18     clusters = [Cluster(center=p, points=[p]) for p in random.sample(points, self.n_clusters)]
19
20
21     while True:
22
23         plists = self.assign_points(clusters, points)
24
25         diff = 0
26
27         for i in range(self.n_clusters):
28             if not plists[i]:
29                 continue
30             old = clusters[i]
31             center = self.calculate_center(plists[i])
32             new = Cluster(center, plists[i])
33             clusters[i] = new
34             diff = max(diff, euclidean(old.center, new.center))
35
36         if diff < self.min_diff:
37             break
38
39     return clusters

```

The implementation follows the description of the algorithm given above. Note that we exit the training loop when the color difference is lower than the one set by us.

## Evaluation

Now that you have an implementation of K-Means clustering you can use it on the UI screenshots. We need a little more glue code to make it easier to extract color palettes:

```

1 def rgb_to_hex(rgb):
2     return '#%s' % ''.join('%02x' % p for p in rgb)
3
4 def get_colors(filename, n_colors=3):
5     points = get_points(filename)
6     clusters = KMeans(n_clusters=n_colors).fit(points)
7     clusters.sort(key=lambda c: len(c.points), reverse = True)
8     rgbs = [map(int, c.center.coordinates) for c in clusters]
9     return list(map(rgb_to_hex, rgbs))

```

The `get_colors` function takes a path to an image file and the number of colors you want to extract from the image. We sort the clusters obtained from our algorithm based on the points in each (in descending order). Finally, we convert the RGB colors to hexadecimal values.

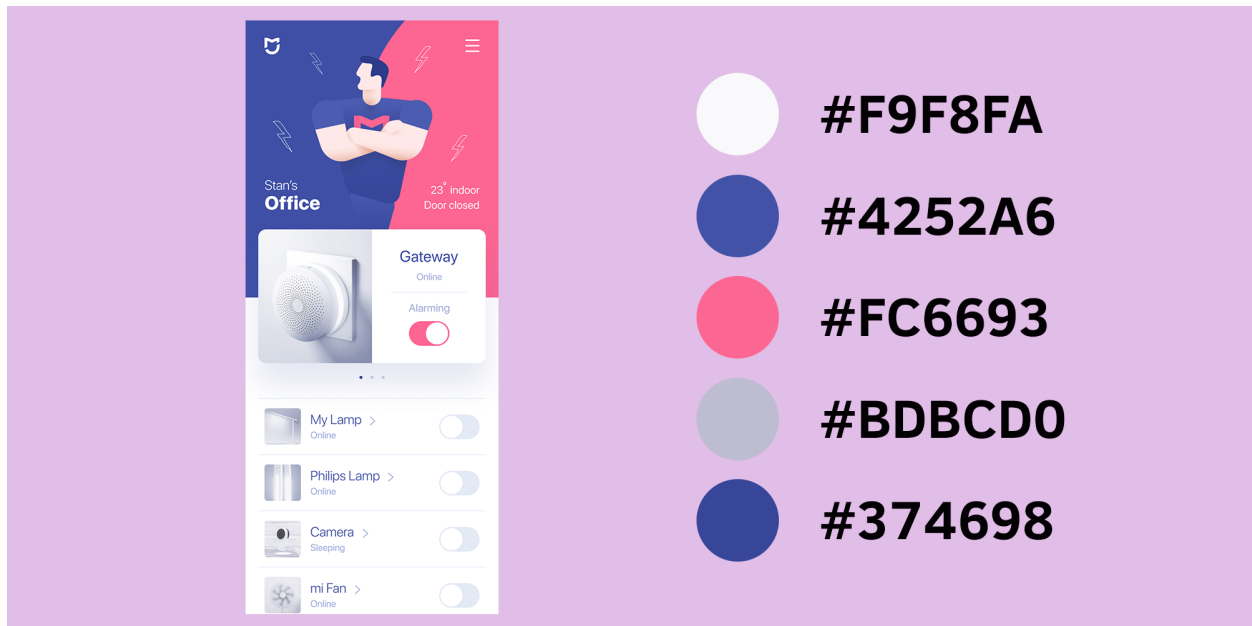
Let's extract the palette for the first UI screenshot from our data:

```
1 colors = get_colors('ui1.png', n_colors=5)
2 ", ".join(colors)
```

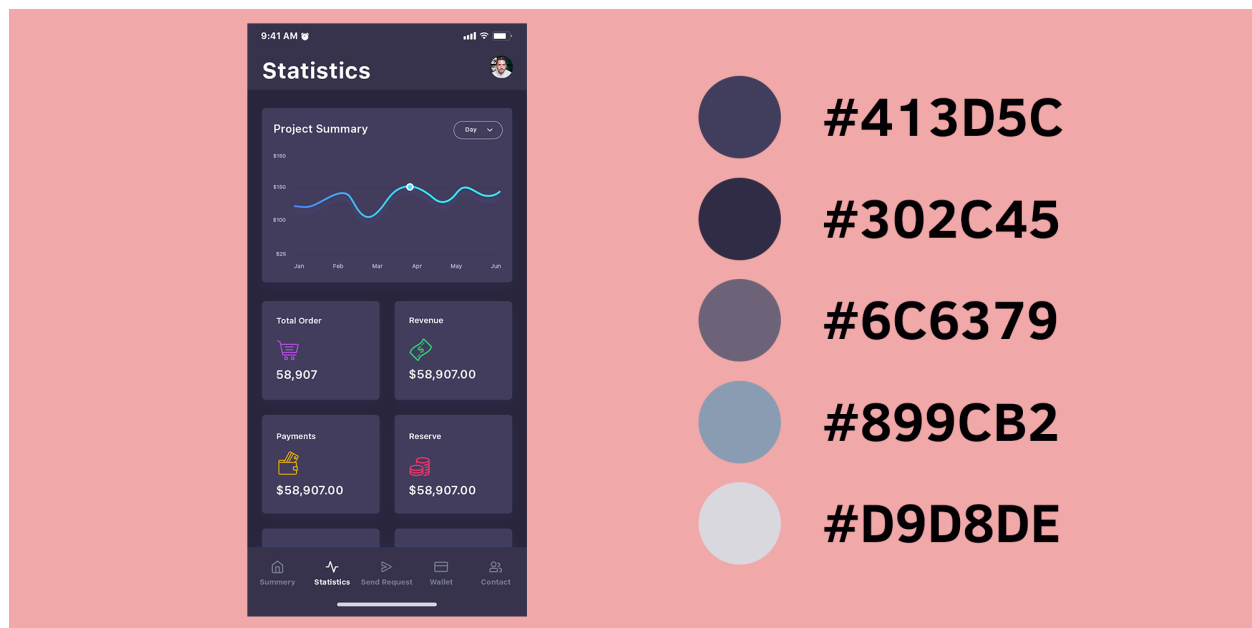
and here are the hexadecimal color values:

```
1 #f9f8fa, #4252a6, #fc6693, #bdbcd0, #374698
```

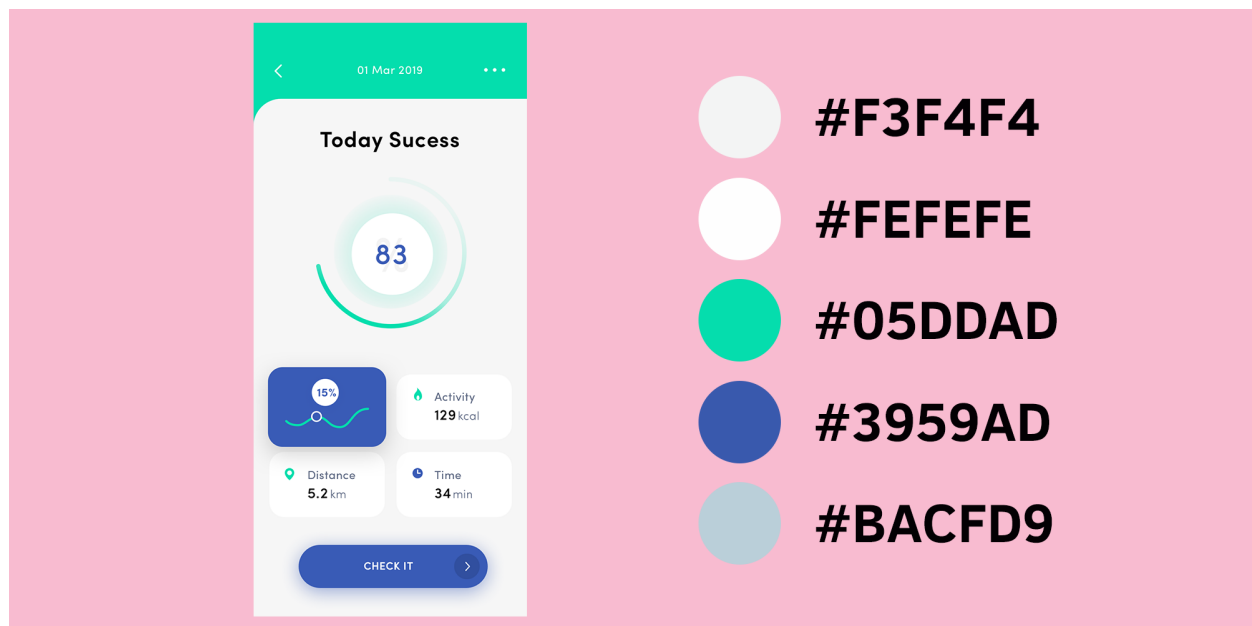
And for a visual representation of our results:



Running the clustering on the next two images, we obtain the following:



The last screenshot:



Well, that looks cool, right? Go on try it on your own images!

[Complete source code notebook on Google Colaboratory<sup>8</sup>](https://colab.research.google.com/drive/1_p1nptDfvJhcSvprmi1WtoIugsCI40Xa)

<sup>8</sup>[https://colab.research.google.com/drive/1\\_p1nptDfvJhcSvprmi1WtoIugsCI40Xa](https://colab.research.google.com/drive/1_p1nptDfvJhcSvprmi1WtoIugsCI40Xa)

## Conclusion

Congratulations, you just implemented your first unsupervised algorithm from scratch! It also appears that it obtains good results when trying to extract a color palette from images.

Next up, we're going to do sentiment analysis on short phrases and learn a bit more how we can handle text data.

# Movie review sentiment analysis with Naive Bayes

TL;DR Build Naive Bayes text classification model using Python from Scratch. Use the model to classify IMDB movie reviews as positive or negative.

Textual data dominates our world from the tweets you read to the timeless writings of Seneca. And while we're consuming images (looking at you Instagram) and videos at increasing rates, you still read Google search results multiple times per daily.

One frequently recurring problem with text data is [Sentiment analysis](#)<sup>9</sup> (classification). Imagine you're a big sugar + water beverage company. You want to know what people think of your products. You'll search for texts with some tags, logos or names. You can then use Sentiment analysis to figure out if the opinions are positive or negative. Of course, you'll send the negative ones to your highly underpaid support center in India to sort things out.

Here, we'll build a generic text classifier that puts movie review texts into one of two categories - negative or positive sentiment. We're going to have a brief look at the Bayes theorem and relax its requirements using the Naive assumption.

[Complete source code in Google Colaboratory Notebook](#)<sup>10</sup>

## Dealing with Text

Computers don't understand text data, though they do well with numbers. [Natural Language Processing \(NLP\)](#)<sup>11</sup> offers a set of approaches to solve text-related problems and represent text as numbers. While NLP is a vast field, we'll use some simple preprocessing techniques and [Bag of Words](#)<sup>12</sup> model.

## The Data

Our data comes from a Kaggle challenge - [“Bag of Words Meets Bags of Popcorn”](#)<sup>13</sup>.

We have 25,000 movie reviews from IMDB labeled as positive or negative. You might know that IMDB ratings are in the 0-10 range. An additional preprocessing step, done by the dataset authors, converts the rating to binary sentiment (<5 - negative ). Of course, a single movie can have multiple reviews, but no more than 30.

---

<sup>9</sup>[https://en.wikipedia.org/wiki/Sentiment\\_analysis](https://en.wikipedia.org/wiki/Sentiment_analysis)

<sup>10</sup><https://colab.research.google.com/drive/1OPQDDJTKy0b40pziZWSoBCmQV6HyXsm>

<sup>11</sup>[https://en.wikipedia.org/wiki/Natural\\_language\\_processing](https://en.wikipedia.org/wiki/Natural_language_processing)

<sup>12</sup>[https://en.wikipedia.org/wiki/Bag-of-words\\_model](https://en.wikipedia.org/wiki/Bag-of-words_model)

<sup>13</sup><https://www.kaggle.com/c/word2vec-nlp-tutorial>

## Reading the reviews

Let's load the training and test data in Pandas data frames:

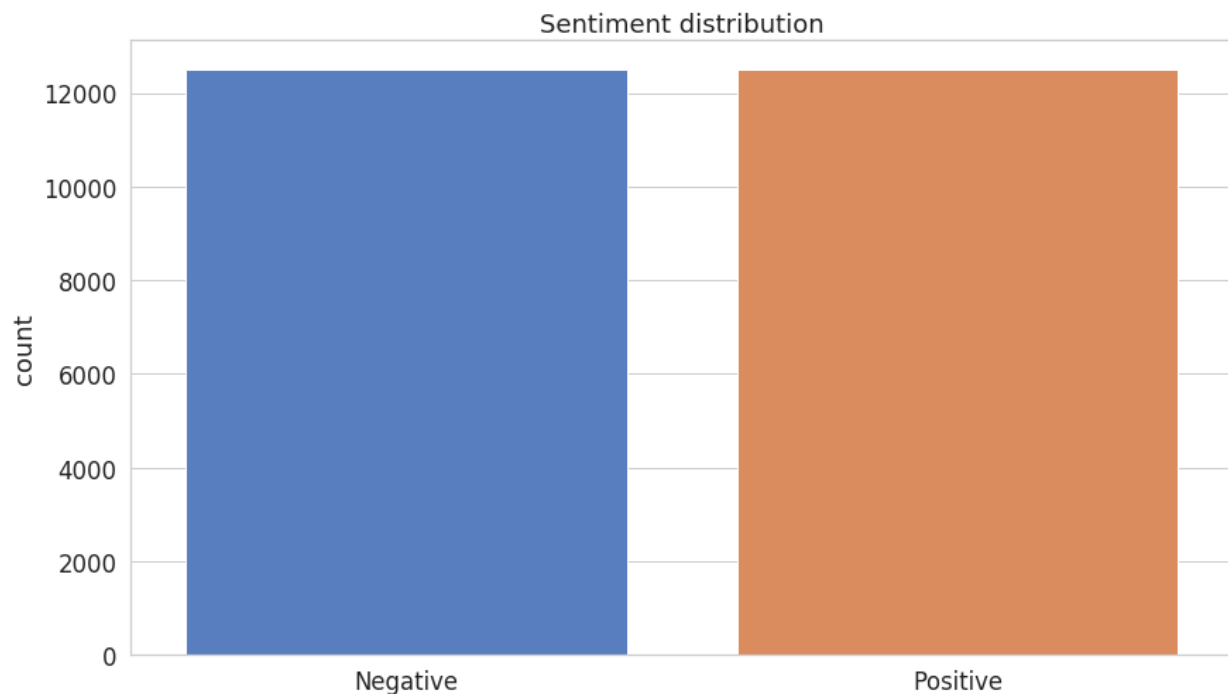
```
1 train = pd.read_csv("imdb_review_train.tsv", delimiter="\t")
2 test = pd.read_csv("imdb_review_test.tsv", delimiter="\t")
```

## Exploration

Let's get a feel for our data. Here are the first 5 rows of the training data:

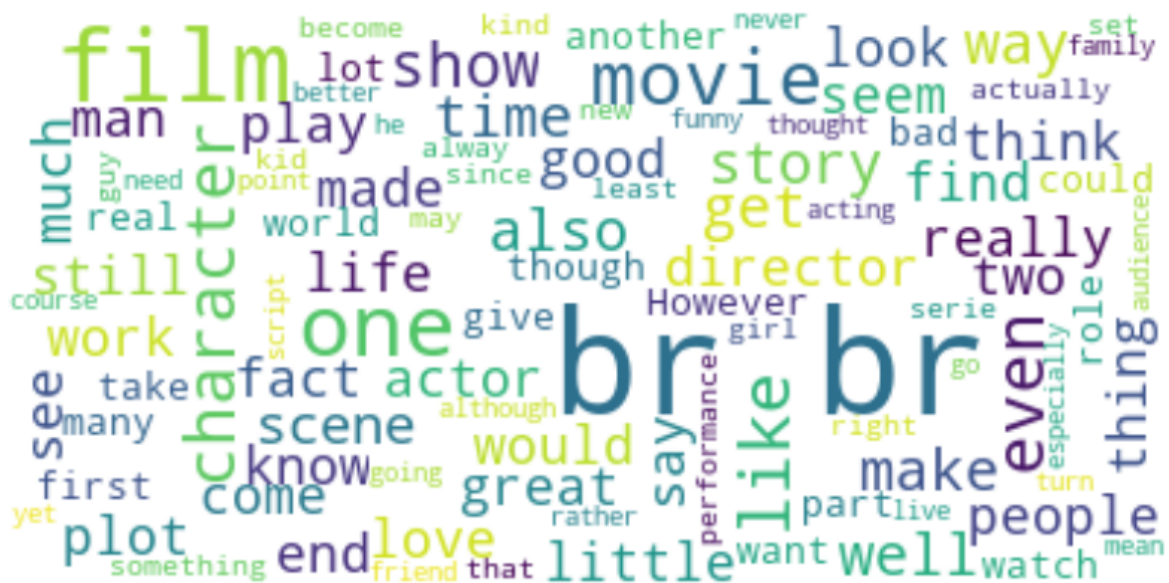
id	sentiment	review
5814_8	1	With all this stuff going down at the moment w...
2381_9	1	The Classic War of the Worlds" by Timothy Hi...
7759_3	0	The film starts with a manager (Nicholas Bell)...
3630_4	0	It must be assumed that those who praised this...
9495_8	1	Superbly trashy and wondrously unpretentious 8...

We're going to focus on the sentiment and review columns. The id column is combining the movie id with a unique number of a review. This might be a piece of important information in real-world scenarios, but we're going to keep it simple.



Both positive and negative sentiments have an equal presence. No need for additional gimmicks to fix that!

Here are the most common words in the training dataset reviews:



Hmm, that **br** looks weird, right?

## Cleaning

Real-world text data is really messy. It can contain excessive punctuation, HTML tags (including that br), multiple spaces, etc. We'll try to remove/normalize most of it.

Most of the cleaning we'll do using [Regular Expressions](#)<sup>14</sup>, but we'll use two libraries to handle HTML tags (surprisingly hard to remove) and removing common (stop) words:

```
1 def clean(self, text):
2     no_html = BeautifulSoup(text).get_text()
3     clean = re.sub("[^a-z\s]+", " ", no_html, flags=re.IGNORECASE)
4     return re.sub("(\\s+)", " ", clean)
```

First, we use [BeautifulSoup](#)<sup>15</sup> to remove HTML tags from our text. Second, we remove anything that is not a letter or space (note the ignoring of uppercase characters). Finally, we replace excessive spacing with a single one.

<sup>14</sup>[https://en.wikipedia.org/wiki/Regular\\_expression](https://en.wikipedia.org/wiki/Regular_expression)

<sup>15</sup><https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

## Tokenization

Now that our reviews are “clean”, we can further prepare them for our Bag of Words model. Let’s them to lowercase letters and split them into individual words. This process is known as [tokenization](#)<sup>16</sup>:

```
1 def tokenize(self, text):
2     clean = self.clean(text).lower()
3     stopwords_en = stopwords.words("english")
4     return [w for w in re.split("\W+", clean) if not w in stopwords_en]
```

The last step of our pre-processing is to remove [stop words](#)<sup>17</sup> using those defined in the [NLTK](#)<sup>18</sup> library. Stop words are usually frequently occurring words that might not significantly affect the meaning of the text. Some examples in English are: “is”, “the”, “and”.

An additional benefit of removing stop words is speeding up our models since we’re removing the amount of train/test data. However, in real-world scenarios, you should think about whether removing stop words can be justified.

We’ll place our clean and tokenize function in a class called *Tokenizer*.

## Naive Bayes

*Naive Bayes* models are probabilistic classifiers that use the [Bayes theorem](#)<sup>19</sup> and make a strong assumption that the features of the data are independent. For our case, this means that each word is independent of others.

Intuitively, this might sound like a dumb idea. You know that (even from reading) the prev word(s) influence the current and next ones. However, the assumption simplifies the math and [works really well in practice](#)<sup>20</sup>!

The Bayes theorem is defined as:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

where  $A$  and  $B$  are some events and  $P(.)$  is probability.

This equation gives us the conditional probability of event  $A$  occurring given  $B$  has happened. In order to find this, we need to calculate the probability of  $B$  happening given  $A$  has happened and

<sup>16</sup>[https://en.wikipedia.org/wiki/Lexical\\_analysis#Tokenization](https://en.wikipedia.org/wiki/Lexical_analysis#Tokenization)

<sup>17</sup>[https://en.wikipedia.org/wiki/Stop\\_words](https://en.wikipedia.org/wiki/Stop_words)

<sup>18</sup><https://www.nltk.org/>

<sup>19</sup>[https://en.wikipedia.org/wiki/Bayes%27\\_theorem](https://en.wikipedia.org/wiki/Bayes%27_theorem)

<sup>20</sup><https://www.cc.gatech.edu/~isbell/reading/papers/Rish.pdf>

multiply that by the probability of  $A$  (known as *Prior*) happening. All of this is divided by the probability of  $B$  happening on its own.

The naive assumption allows us to reformulate the Bayes theorem for our example as:

$$P(\textit{Sentiment}|w_1, \dots, w_n) = \frac{P(\textit{Sentiment}) \prod_{i=1}^n P(w_i|\textit{Sentiment})}{P(w_1, \dots, w_n)}$$

We don't really care about probabilities. We only want to know whether a given text has a positive or negative sentiment. We can skip the denominator entirely since it just scales the numerator:

$$P(\textit{Sentiment}|w_1, \dots, w_n) \propto P(\textit{Sentiment}) \prod_{i=1}^n P(w_i|\textit{Sentiment})$$

To choose the sentiment, we'll compare the scores for each sentiment and pick the one with a higher score.

## Implementing Multinomial Naive Bayes

As you might've guessed by now, we're classifying text into one of two groups/categories - positive and negative sentiment.

Multinomial Naive Bayes allows us to represent the features of the model as frequencies of their occurrences (how often some word is present in our review). In other words, it tells us that the probability distributions we're using are [multinomial](#)<sup>21</sup>.

### Note on numerical stability

Our model relies on multiplying many probabilities. Those might become increasingly small and our computer might round them down to zero. To combat this, we're going to use log probability by taking log of each side in our equation:

$$\log P(\textit{Sentiment}|w_1, \dots, w_n) = \log P(\textit{Sentiment}) + \log \prod_{i=1}^n P(w_i|\textit{Sentiment})$$

Note that we can still use the highest score of our model to predict the sentiment since log is [monotonic](#)<sup>22</sup>.

<sup>21</sup>[https://en.wikipedia.org/wiki/Multinomial\\_distribution](https://en.wikipedia.org/wiki/Multinomial_distribution)

<sup>22</sup>[https://en.wikipedia.org/wiki/Monotonic\\_function](https://en.wikipedia.org/wiki/Monotonic_function)

## Building our model

Finally, time to implement our model in Python. Let's start by defining some variables and group the data by class, so our training code is a bit tidier:

```
1 def fit(self, X, y):
2     self.n_class_items = {}
3     self.log_class_priors = {}
4     self.word_counts = {}
5     self.vocab = set()
6
7     n = len(X)
8     grouped_data = self.group_by_class(X, y)
9     ...
```

We're going to implement a generic text classifier that doesn't assume the number of classes. You can use it for news category prediction, sentiment analysis, email spam detection, etc.

For each class, we'll find the number of examples in it and the log probability (prior). We'll also record the number of occurrences of each word and create a vocabulary - set of all words we've seen in the training data:

```
1     for c, data in grouped_data.items():
2         self.n_class_items[c] = len(data)
3         self.log_class_priors[c] = math.log(self.n_class_items[c] / n)
4         self.word_counts[c] = defaultdict(lambda: 0)
5
6         for text in data:
7             counts = Counter(self.tokenizer.tokenize(text))
8             for word, count in counts.items():
9                 if word not in self.vocab:
10                     self.vocab.add(word)
11
12             self.word_counts[c][word] += count
```

Note that we use our *Tokenizer* and the built-in class [Counter](https://docs.python.org/3/library/collections.html#collections.Counter)<sup>23</sup> to convert a review to a bag of words.

In case you're interested, here's how `group_by_class` works:

---

<sup>23</sup><https://docs.python.org/3/library/collections.html#collections.Counter>

```
1 def group_by_class(self, X, y):
2     data = dict()
3     for c in self.classes:
4         data[c] = X[np.where(y == c)]
5     return data
```

## Making predictions

In order to predict sentiment from text data, we'll use our class priors and vocabulary:

```
1 def predict(self, X):
2     result = []
3     for text in X:
4
5         class_scores = {c: self.log_class_priors[c] for c in self.classes}
6         words = set(self.tokenizer.tokenize(text))
7
8         for word in words:
9             if word not in self.vocab: continue
10
11         for c in self.classes:
12
13             log_w_given_c = self.laplace_smoothing(word, c)
14             class_scores[c] += log_w_given_c
15
16         result.append(max(class_scores, key=class_scores.get))
17
18     return result
```

For each review, we use the log priors for positive and negative sentiment and tokenize the text. For each word, we check if it is in the vocabulary and skip it if it is not. Then we calculate the log probability of this word for each class. We add the class scores and pick the class with a max score as our prediction.

## Laplace smoothing

Note that  $\log(0)$  is *undefined* (and no, we're not using JavaScript here). It is entirely possible for a word in our vocabulary to be present in one class but not another - the probability of finding this word in that class will be 0! We can use [Laplace smoothing](https://en.wikipedia.org/wiki/Additive_smoothing)<sup>24</sup> to fix this problem. We'll simply add 1 to the numerator but also add the size of our vocabulary to the denominator:

---

<sup>24</sup>[https://en.wikipedia.org/wiki/Additive\\_smoothing](https://en.wikipedia.org/wiki/Additive_smoothing)

```
1 def laplace_smoothing(self, word, text_class):
2     num = self.word_counts[text_class][word] + 1
3     denom = self.n_class_items[text_class] + len(self.vocab)
4     return math.log(num / denom)
```

## Predicting sentiment

Now that your model can be trained and make predictions, let's use it to predict sentiment from movie reviews.

## Data preparation

As discussed previously, we'll use only the review and sentiment columns from the training data. Let split it for training and testing:

```
1 X = train['review'].values
2 y = train['sentiment'].values
3
4 X_train, X_test, y_train, y_test = train_test_split(
5     X, y,
6     test_size=0.2,
7     random_state=RANDOM_SEED
8 )
```

## Evaluation

We'll pack our fit and predict functions into a class called `MultinomialNaiveBayes`. Let's use it:

```
1 MNB = MultinomialNaiveBayes(
2     classes=np.unique(y),
3     tokenizer=Tokenizer()
4 ).fit(X_train, y_train)
```

Our classifier takes a list of possible classes and a `Tokenizer` as parameters. Also, the API is quite nice (thanks scikit-learn!)

```

1 y_hat = MNB.predict(X_test)
2 accuracy_score(y_test, y_hat)

```

```

1 0.8556

```

This looks nice we got an accuracy of ~86% on the test set.

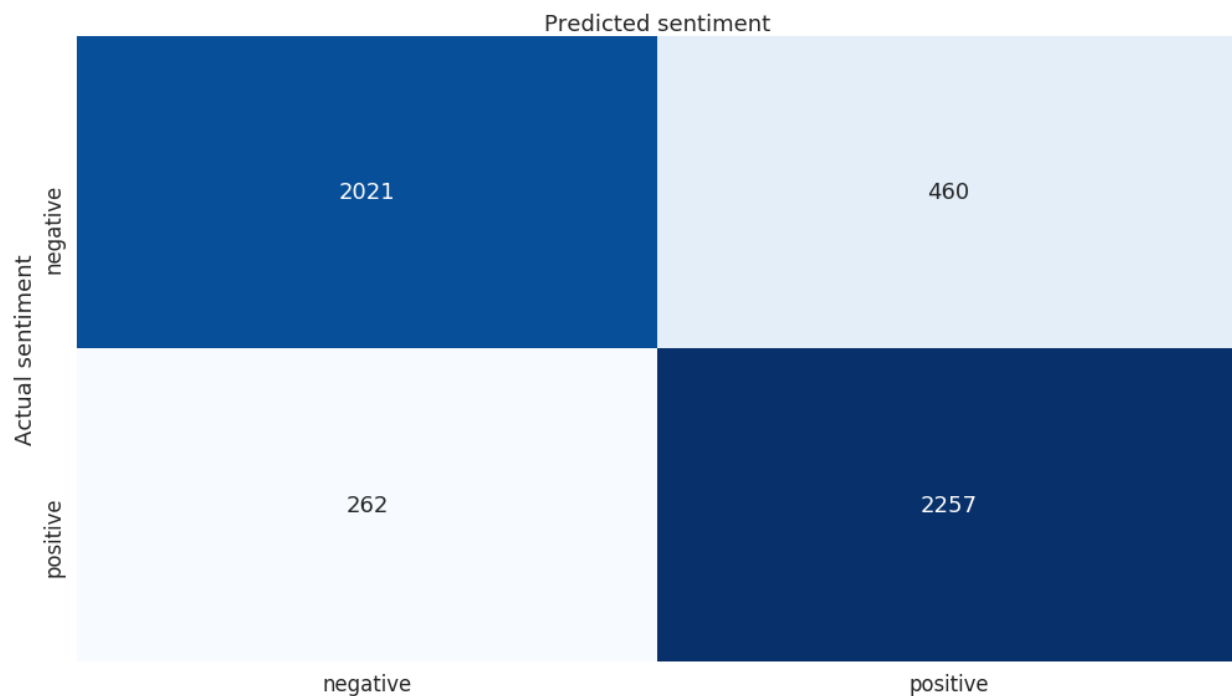
Here is the classification report:

```

1          precision    recall  f1-score   support
2
3     0           0.89      0.81      0.85     2481
4     1           0.83      0.90      0.86     2519
5
6  accuracy                   0.86     5000
7  macro avg                  0.86     5000
8  weighted avg               0.86     5000

```

And the confusion matrix:



Overall, our classifier does pretty well for himself. Submit the predictions to Kaggle and find out what place you'll get on the leaderboard.

## Conclusion

Well done! You just built a Multinomial Naive Bayes classifier that does pretty well on sentiment prediction. You also learned about Bayes theorem, text processing, and Laplace smoothing! Will another flavor of the Naive Bayes classifier perform better?

[Complete source code in Google Colaboratory Notebook<sup>25</sup>](#)

Next up, you'll build a recommender system from scratch!

---

<sup>25</sup><https://colab.research.google.com/drive/1OPQDDJTKy0b40pziZWSsoBCmQV6HyXsm>