







DATA SCIENCE

MACHINE LEARNING

PROGRAMMING

VISUALIZATION

JOURNALISM

**EVENTS** 

SUBMIT

# Sentiment Analysis with Python (Part 1)

1)
Classifying IMDb Movie Reviews



Aaron Kub Follow

Jul 31, 2018 · 6 min read



Photo by Denise Jans on Unsplash

Sentiment Analysis is a common NLP task that Data Scientists need to perform. This is a straightforward guide to creating a barebones movie review classifier in Python. Future parts of this series will focus on improving the classifier.









DATA SCIENCE

MACHINE LEARNING

PROGRAMMING

VISUALIZATION

AI JOURNALISM

**EVENTS** 

SUBMIT

#### **Data Overview**

For this analysis we'll be using a dataset of 50,000 movie reviews taken from IMDb. The data was compiled by Andrew Maas and can be found here: IMDb Reviews.

The data is split evenly with 25k reviews intended for training and 25k for testing your classifier.

Moreover, each set has 12.5k positive and 12.5k negative reviews.

IMDb lets users rate movies on a scale from 1 to 10. To label these reviews the curator of the data labeled anything with  $\leq 4$  stars as negative and anything with  $\geq 7$  stars as positive. Reviews with 5 or 6 stars were left out.

• • •

#### **Step 1: Download and Combine Movie Reviews**

If you haven't yet, go to IMDb Reviews and click on "Large Movie Review Dataset v1.0". Once that is complete you'll have a file called  $aclindb_vl.tar.gz$  in your downloads folder.

**Shortcut:** If you want to get straight to the data analysis and/or aren't super comfortable with the terminal, I've put a tar file of the final directory that this step creates here: Merged Movie Data.

Double clicking this file should be sufficient to unpack it (at least on a Mac), otherwise gunzip -c

 $\verb"movie_data.tar.gz | [ \verb"tar xopf - in a terminal will do it."]$ 

# **Unpacking and Merging**

Follow these steps or run the shell script here: Preprocessing Script

- 1. Move the tar file to the directory where you want this data to be stored.
- 2. Open a terminal window and cd to the directory that you put aclimdb v1.tar.gz in.
- 3. gunzip -c aclImdb v1.tar.gz | tar xopf -
- 4. cd aclImdb && mkdir movie data
- 5. for split in train test; do for sentiment in pos neg; do for file in \$split/\$sentiment/\*; do cat \$file >> movie\_data/full\_\${split}.txt; echo >> movie\_data/full\_\${split}.txt; done; done; done;

. . .

#### Step 2: Read into Python

For most of what we want to do in this walkthrough we'll only need our reviews to be in a Python

```
list .Make sure to point_open_ to the directory where you put the movie data.

1    reviews_train = []

2    for line in open('../data/movie_data/full_train.txt', 'r'):

3        reviews_train.append(line.strip())

4
```









DATA SCIENCE

MACHINE LEARNING

PROGRAMMING

VISUALIZATION

I JOURNALISM

**EVENTS** 

SUBMIT

. .

### **Step 3: Clean and Preprocess**

The raw text is pretty messy for these reviews so before we can do any analytics we need to clean things up. Here's one example:

"This isn't the comedic Robin Williams, nor is it the quirky/insane Robin Williams of recent thriller fame. This is a hybrid of the classic drama without over-dramatization, mixed with Robin's new love of the thriller. But this isn't a thriller, per se. This is more a mystery/suspense vehicle through which Williams attempts to locate a sick boy and his keeper.<br/>
by />Also starring Sandra Oh and Rory Culkin, this Suspense Drama plays pretty much like a news report, until William's character gets close to achieving his goal.<br/>
br />I must say that I was highly entertained, though this movie fails to teach, guide, inspect, or amuse. It felt more like I was watching a guy (Williams), as he was actually performing the actions, from a third person perspective. In other words, it felt real, and I was able to subscribe to the premise of the story.<br/>
br />kbr />kbr />All in all, it's worth a watch, though it's definitely not Friday/Saturday night fare.<br/>
cbr />the Fiend :."

**Note:** Understanding and being able to use regular expressions is a prerequisite for doing any

Natural Language Processing task. If vou're unfamiliar with them perhaps start here: Regex Tutorial

```
import re

REPLACE_NO_SPACE = re.compile("[.;:!\'?,\"()\[\]]")

REPLACE_WITH_SPACE = re.compile("(<br\s*/><br\s*/>)|(\-)|(\/)")

def preprocess_reviews(reviews):
    reviews = [REPLACE_NO_SPACE.sub("", line.lower()) for line in reviews]
    reviews = [REPLACE_WITH_SPACE.sub(" ", line) for line in reviews]

return reviews

reviews_train_clean = preprocess_reviews(reviews_train)
    reviews_test_clean = preprocess_reviews(reviews_test)

clean_movie_reviews.py hosted with \( \rightarrow by GitHub

view raw
```

https://towardsdatascience.com/sentiment-analysis-with-python-part-1-5ce197074184









DATA SCIENCE MACHINE LEARNING PROGRAMMING VISUALIZATION AI JOURNALISM EVENTS SUBMIT

sick boy and his keeper also starring sandra oh and rory culkin this suspense drama plays pretty much like a news report until williams character gets close to achieving his goal i must say that i was highly entertained though this movie fails to teach guide inspect or amuse it felt more like i was watching a guy williams as he was actually performing the actions from a third person perspective in other words it felt real and i was able to subscribe to the premise of the story all in all its worth a watch though its definitely not friday saturday night fare it rates a from the fiend"

**Note:** There are a lot of different and more sophisticated ways to clean text data that would likely produce better results than what I've done here. I wanted part 1 of this tutorial to be as simple as possible. Also, I generally think it's best to get baseline predictions with the simplest possible solution before spending time doing potentially unnecessary transformations.

#### Vectorization

In order for this data to make sense to our machine learning algorithm we'll need to convert each review to a numeric representation, which we call *vectorization*.

The simplest form of this is to create one very large matrix with one column for every unique word in your corpus (where the corpus is all 50k reviews in our case). Then we transform each review into one row containing 0s and 1s, where 1 means that the word in the *corpus* corresponding to that column appears in that *review*. That being said, each row of the matrix will be very sparse (mostly

zeros). This process is also known as one hot encoding.

```
from sklearn.feature_extraction.text import CountVectorizer

cv = CountVectorizer(binary=True)

cv.fit(reviews_train_clean)

X = cv.transform(reviews_train_clean)

X_test = cv.transform(reviews_test_clean)

vectorize_reviews.py hosted with \(\sigma\) by GitHub

view raw
```

. . .

#### Step 4: Build Classifier

Now that we've transformed our dataset into a format suitable for modeling we can start building a classifier. Logistic Regression is a good baseline model for us to use for several reasons: (1) They're









SUBMIT

DATA SCIENCE MACHINE LEARNING PROGRAMMING VISUALIZATION AI JOURNALISM EVENTS

```
structured the same, where the first 12.5k are positive and the last 12.5k are negative.
       from sklearn.linear_model import LogisticRegression
       from sklearn.metrics import accuracy_score
       from sklearn.model selection import train test split
       target = [1 if i < 12500 else 0 for i in range(25000)]</pre>
      X_train, X_val, y_train, y_val = train_test_split(
           X, target, train_size = 0.75
  9
 10
       for c in [0.01, 0.05, 0.25, 0.5, 1]:
           lr = LogisticRegression(C=c)
 14
           lr.fit(X_train, y_train)
           print ("Accuracy for C=%s: %s"
                  % (c, accuracy_score(y_val, lr.predict(X_val))))
 17
 18
             Accuracy for C=0.01: 0.87472
 19
             Accuracy for C=0.05: 0.88368
             Accuracy for C=0.25: 0.88016
 20
             Accuracy for C=0.5: 0.87808
             Accuracy for C=1: 0.87648
 LR_review_classifier_part1.py hosted with \bigcirc by GitHub
                                                                                                 view raw
```

It looks like the value of C that gives us the highest accuracy is 0.05.

. . .

## **Train Final Model**

Now that we've found the optimal value for C, we should train a model using the entire training set









DATA SCIENCE MACHINE LEARNING PROGRAMMING VISUALIZATION AI JOURNALISM EVENTS SUBMIT

```
cv.get_reature_names(), Tinat_modef.coet_[v]
     for best_positive in sorted(
         feature_to_coef.items(),
         key=lambda x: x[1],
         reverse=True)[:5]:
         print (best_positive)
10
           ('excellent', 0.9288812418118644)
13
           ('perfect', 0.7934641227980576)
           ('great', 0.675040909917553)
           ('amazing', 0.6160398142631545)
           ('superb', 0.6063967799425831)
17
     for best_negative in sorted(
         feature_to_coef.items(),
20
         key=lambda x: x[1])[:5]:
         print (best_negative)
22
           ('worst', -1.367978497228895)
           ('waste', -1.1684451288279047)
24
           ('awful', -1.0277001734353677)
           ('poorly', -0.8748317895742782)
           ('boring', -0.8587249740682945)
best_predictors_Ir_part1.py hosted with ♥ by GitHub
                                                                                               view raw
```

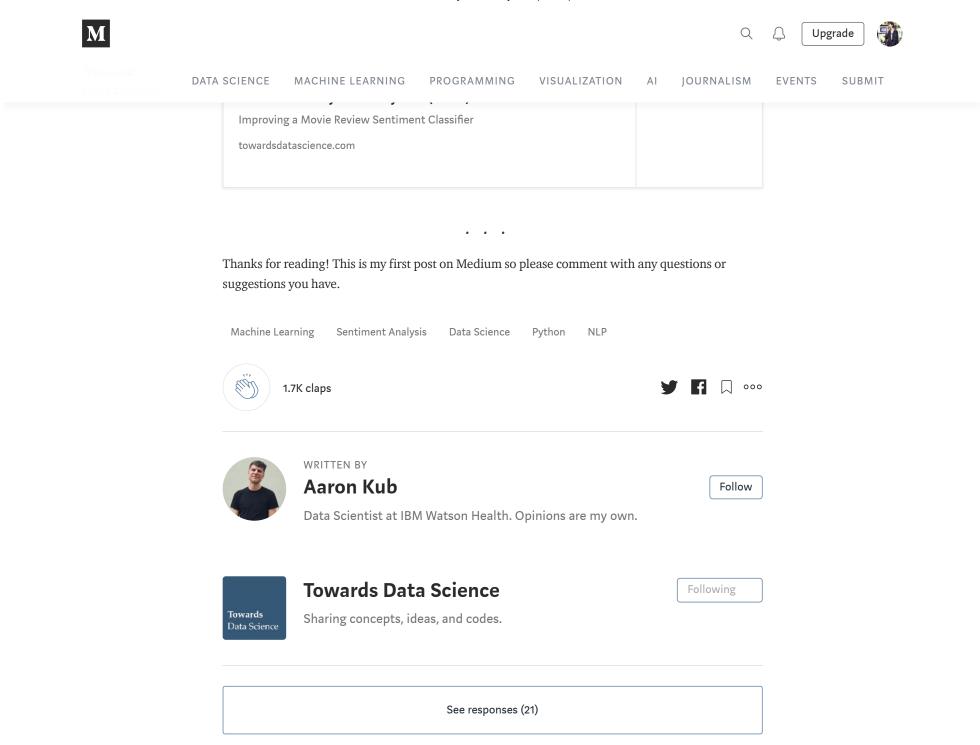
And there it is. A very simple classifier with pretty decent accuracy out of the box.

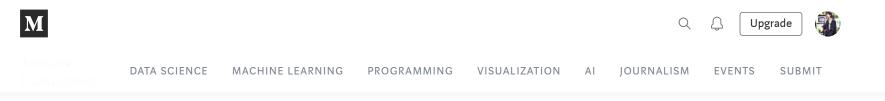
. . .

#### **Next Time**

In the next parts of this series we'll look at more sophisticated ways to get better performance out of our classifier.

- Text Processing: Stemming/Lemmatizing to convert different forms of each word into one.
- n-grams: Instead of just single-word tokens (1-gram/unigram) we can also include word pairs.
- **Representations**: Instead of simple, binary vectors we can use word counts or *TF-IDF* to transform those counts.





More from Towards Data Science

More from Towards Data Science

More from Towards Data Science

# 12 Things I Learned During My First Year as a Machine Learning Engineer





Hands On Bayesian Statistics with Python, PyMC3 & ArviZ



Susan Li in Towards Dat...
Jul 17 · 10 min read

586



P-values Explained By Data Scientist



