

Objective

The main goal of the present project was to identify how a product is perceived in the market. This fall in the area of Sentimental Analysis, which is an important part of Natural Language Processing. This technology handles with the big challenge of make human language understandable to computers.

As complicated this might look, we chose a small and manageable objective: we wanted to create a model capable of classify comments on a product into positive and negative, in order to have a first insight in the consumers product perception. We thought that the best way to do this was to trough gathering social media comments in order to classify them into positive and negative, and then feed a machine learning classification model.

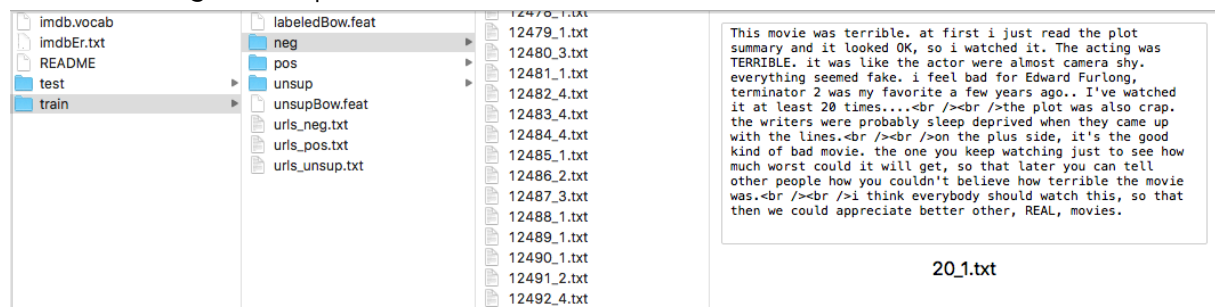
At first we tried to gather this data ourselves but soon it becomes evident that this was a major task out of our time constrains. So we decided to resolve a similar problem with existing data. Fortunately we found a well classified dataset of movie's comments posted on IMDB, who save us half of the work.

We thought that this was a good star, because movie industry is a multimillionaire business in which each project has a lot to lose if a production if fails to meet public interests and wishes. So it needs feedback of his production to meet public expectations. So the product we want to help to analyze were movies.

In this matter, Sentimental analysis is a particular useful machine learning method to help movie industry to meet public interests. On this respect we pretend to use Natural Language Processing model capable of identify the connotation of a movie comment, so in union with some other methods we can improve our understanding on public taste.

Database Construction

The data we used was compose of 50,000 comments, posted in IMDB, who were gathered for the Annual Meeting of the Association for Computational Linguistics: Human Language Technologies(A. Maas et.al, 2011). [Large Movie Review Dataset v1.0](#) This comments where composed by 30 comments maximum on one movie, and were already labeled as negative or positive.



The first problem we faced was that the data was provided in a multiple ".txt" files, so the

first task was to iterate throw it and put them inside a unified data frame. This task was best made with python pandas. Our data were provided in two folders differentiating them in positive and negative, so we have to add the correct label, and then apply a One-Hot Encoding in order to transform labels into numbers. This was our result:

		0	negative	positive
0	Story of a man who has unnatural feelings for ...		1.0	0.0
1	Airport '77 starts as a brand new luxury 747 p...		1.0	0.0
2	This film lacked something I couldn't put my f...		1.0	0.0
3	Sorry everyone,,, I know this is supposed to b...		1.0	0.0
4	When I was little my parents took me along to ...		1.0	0.0

Choosing a method

Determine if something is positive or negative is better manageable with a classifier model and binary data. At the end we chose a Support Vector Machine, because is the more straight forward classifier we know and was the default metaparameters provided by the sklearn library we finally use. Before reaching this point we have two options the Naive Bayes model from PySpark and using google Colab and Sklearn stochastic gradient descent (SGD) classifier. First we tried the Naive Bayes model which is explained in this image:

Naive Bayes: Is a **Probabilistic Classifier** algorithm based on Bayes' Theorem. Well know as a ML model for **Sentiment Analysis**.

Theorem

$$P(H|D) = \frac{P(D|H) P(H)}{P(D)}$$

Word: Great : 70 20
 + -
 75 25

what is the % of Probability that this message Could be + or -

70/90 = 0.7 %
 20/90 = 0.2 %

Inputs:
 50,000 labeled records as (+/-).
 70 % of records were used for model training, 30 % for testing.
 Data Cleanse process included: Drop Nulls, Stop Words, Steaming.
 Pyspark Lib was used to create NaiveBayes() model.
 GColab, spark, aws s3 storage tools were used.

Result

```
# Use the Class Evaluator for a cleaner description
from pyspark.ml.evaluation import MulticlassClassificationEvaluator

acc_eval = MulticlassClassificationEvaluator()
acc = acc_eval.evaluate(test_results)
print("Accuracy of model at predicting reviews was: %f" % acc)
```

Accuracy of model at predicting reviews was: 0.817837

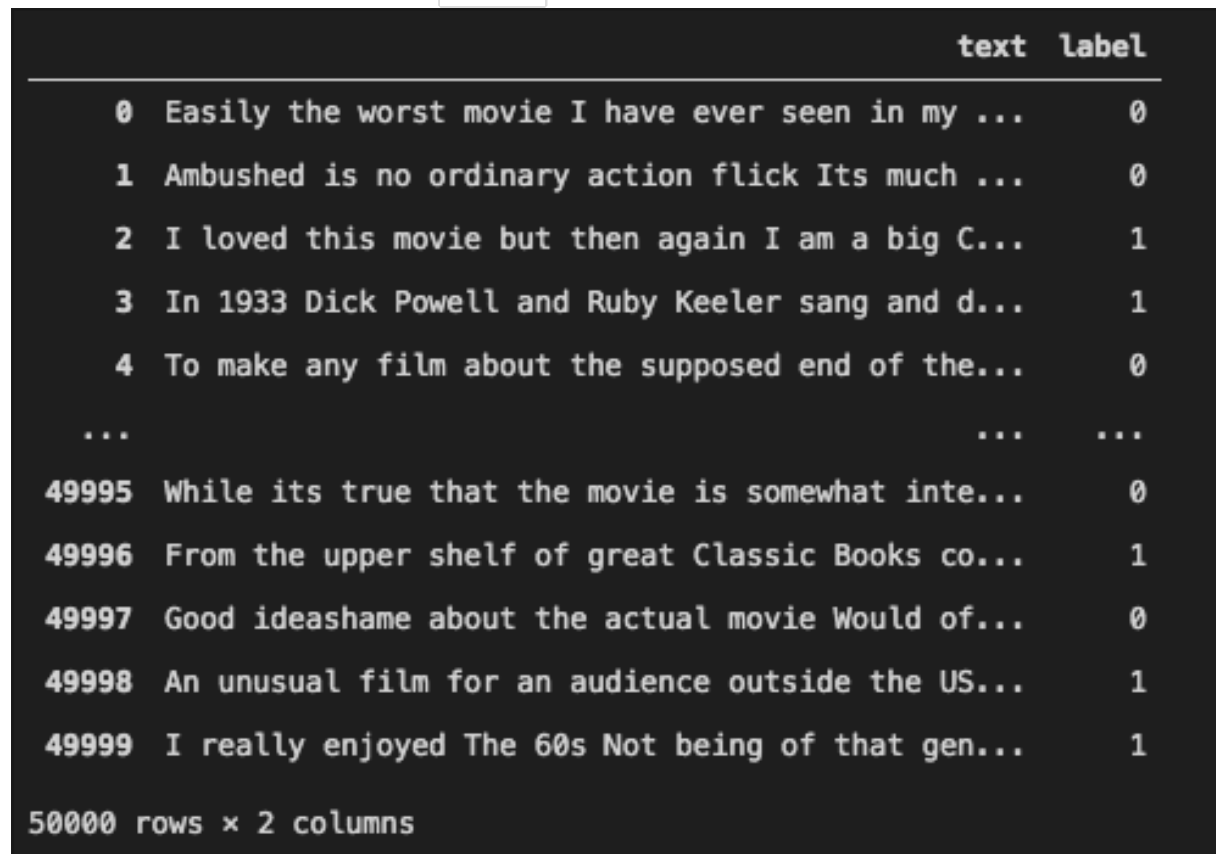
Findings:
 To Improve Model Accuracy we need:
 1.- More than 100 k rows must consider to Train the model.
 2.- A Hight Quality Labeled Dataset to train the model.
 3.- A robust Data Cleanse process.
 To choose a model:
 3.- We must try more than one model and compare results.
 To deliver and use a final solution take in account:
 1.-Cloud storage and data processing like spark for bigdata.

With this model we got an score of 81.56 accuracy as maximum. At the end we decided to change model, not because of the results, but because we found some limitation to deploy our model in a web app, so we decided to use sklearn instead, which luckily gave us better results as we will explain. But before use the sklearn model we need to meet

specific requirements of preprocessing.

Preprocessing

The first preprocessing task was to remove all punctuation to facilitate the recognition on every word. This was a decision on our own, that came with a cost because certain punctuation as "!" or "?" express certain emotions that could be lost. But putting this into a balance it was preferable to lose little accuracy than to get a duplicated word because a point or a comma, that can affect the weight our model can give to some words. This task was easily completed using the `string` library.



```

      text  label
0  Easily the worst movie I have ever seen in my ...    0
1  Ambushed is no ordinary action flick Its much ...    0
2  I loved this movie but then again I am a big C...    1
3  In 1933 Dick Powell and Ruby Keeler sang and d...    1
4  To make any film about the supposed end of the...    0
...
49995  While its true that the movie is somewhat inte...    0
49996  From the upper shelf of great Classic Books co...    1
49997  Good ideashame about the actual movie Would of...    0
49998  An unusual film for an audience outside the US...    1
49999  I really enjoyed The 60s Not being of that gen...    1
50000 rows x 2 columns
```

Vectorize

The second preprocessing task was something asked by sklearn model. We need to transform words into numbers via tokenization. To achieve this we had to vectorize our strings using `CountVectorizer` which transform our strings into a matrix of words count. For each word it find it create a column, and for each string it create a row. So we got a matrix who counts how many times a word appear in a string. If we would liked we could create more complex matrix not counting words but counting pair of words, we just had to specify `ngram_range=(2,2)` but this would had create a bigger matrix and made a heavier process, with interferes with our propose to create a light web app. In this image we use an example sentence to show up how vectorize works.

put a text box where the user can put their opinion and our model can tell if it is positive or negative.

When the text was clearly a bad or a good opinion we found great accuracy. But we found some complications when a comment were ambiguous, or had bad and good words combined. In this cases our model tend to fail, so this give us feedback, about maybe create a new category called neutral and train our model again.

Another problem we notice was that the stop words they have actually an emotional value, particularly thinking in words like "not", and our model tend to ignore this. We were worried about some false positives, but this seems not to affect to much as, we already tell, clear messages were well predicted.

Considering the little amount of time we had, we felt satisfied with the results we got. We think we have a good and working start point, to perfect our model. Maybe adding more data, adding more stop words for obvious irrelevant words as "movie".

References

=====

Maas, A., Daly, R., Pham, P., Huang, D., Ng, A. and Potts, C. (2011). Learning Word Vectors for Sentiment Analysis: Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies. [online] Portland, Oregon, USA: Association for Computational Linguistics, pp.142-150. Available at: <http://www.aclweb.org/anthology/P11-1015>.

Building a Sentiment Classifier using Scikit-Learn | by Dorian Lazar | Towards Data Science
