

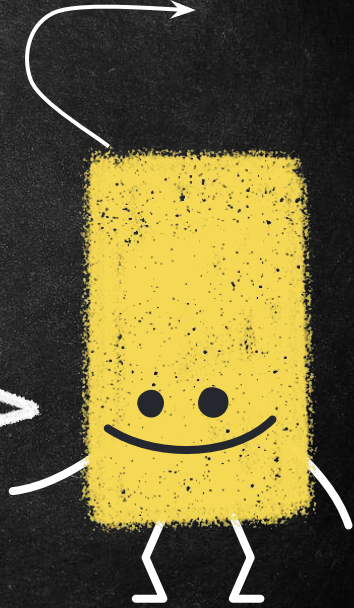
SENTIMENT ANALYSIS

ON SHORT MOVIE REVIEWS

EDDY



SIENNA



AGENDA

Introduction

1

DataSet

2

Our Solution

3

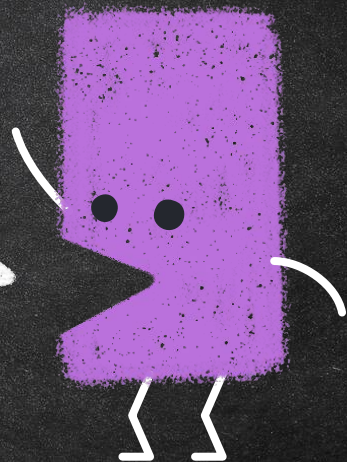
Future Work

4

“

PART 1

THE PROBLEM

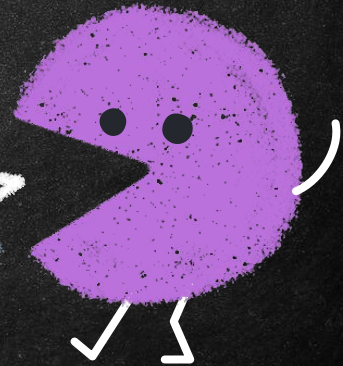


SENTIMENT ANALYSIS



BALABALABALABALAB
ALABALABALABALABA
LABALABALABALABAL
ABALABALABALA...
x&% ¥ #@.....

BINARY
CLASSIFICATION

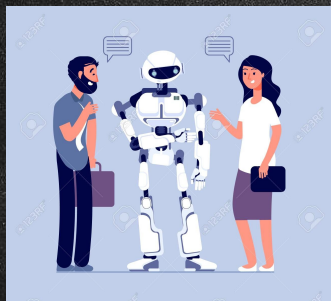


IMPORTANCE ! ! !

1.

Machines'
interaction with
humans

E.g: AI robot Siri



2

Efficiently decide
the emotion in
natural language

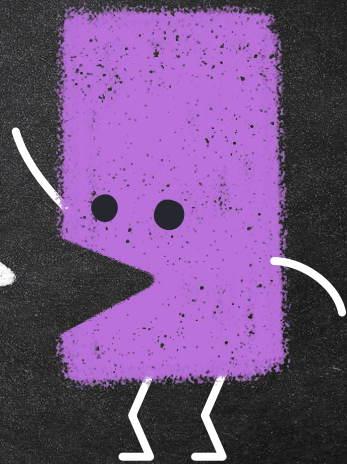
E.g. our dataset



“

PART 2

THE DATASET

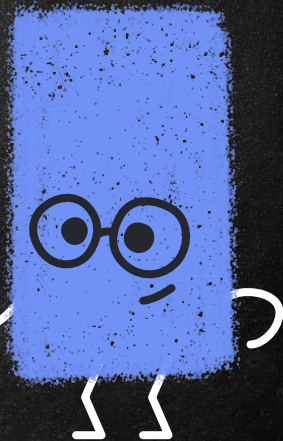


KAGGLE

 Dataset

IMDB Dataset of 50K Movie Reviews

Large Movie Review Dataset

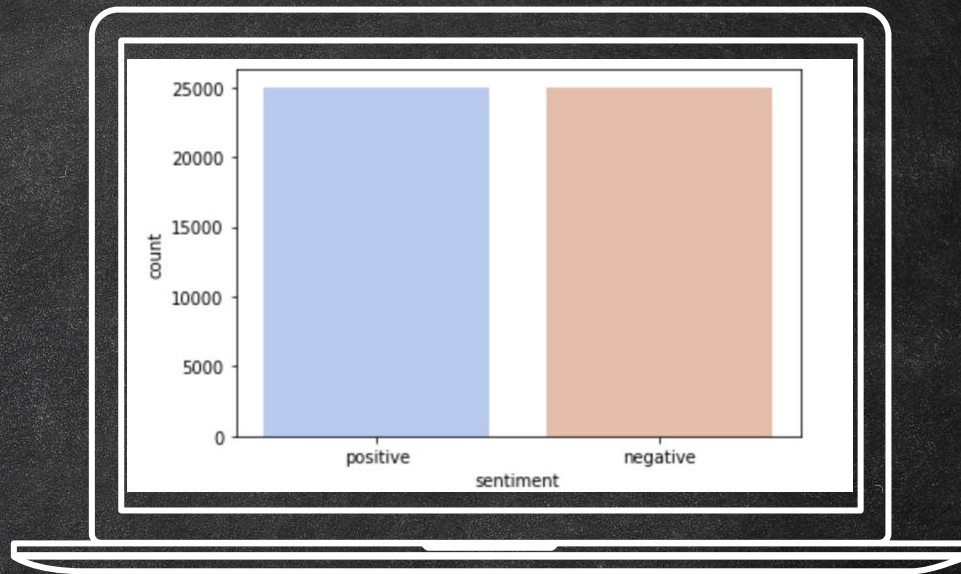


DATASET

8

DISTRIBUTION

POSITIVE:NEGATIVE = 1:1



DATASET

9

AROUND 231
WORDS IN EACH
COMMENT

AVERAGE
SENTENCE COUNT
26

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production. The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Mattei's "Love in the Time of Money" is...	positive

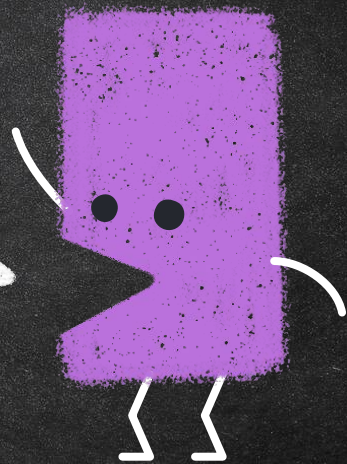
For data preprocessing, we do the basic text preprocessing in the very beginning.
(e.g. Remove stopwords and punctuation)

It turns out that the preprocessing of texts affects the performance largely(will explain later)

“

PART 3

OUR TRIAL



IN THIS PART...

1. First few trials
2. Final solution
3. Optimize the model
 - a. Change Text preprocessing
 - b. Hyperparameters

Sentiment
Analysis
Model

=

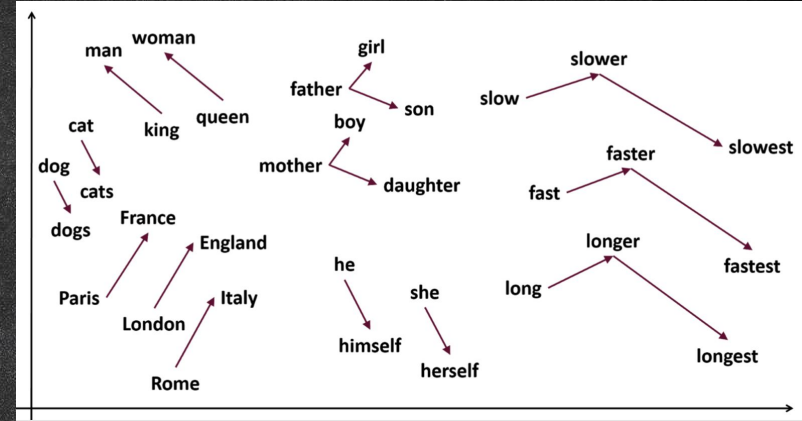
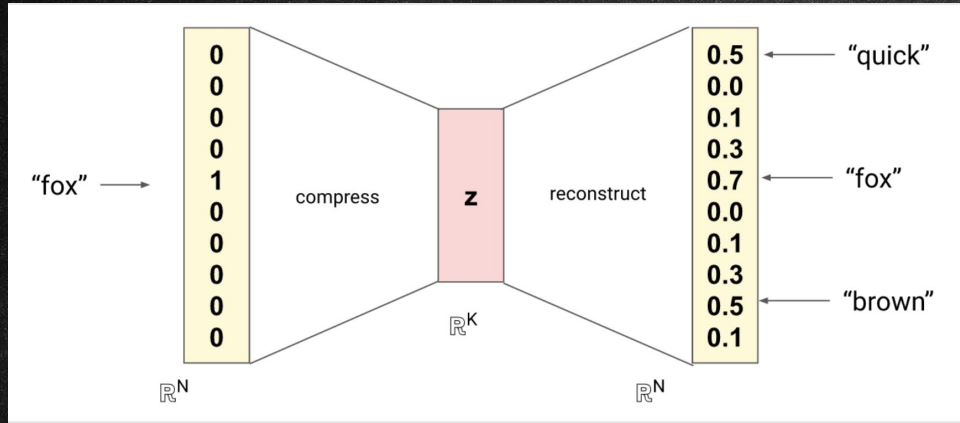
W_{ord} E_{mbedding}

+

C_{lassifier}

NLP

First Trial for Word Embedding: Word2Vec



- Convert each word to a **vector** (direction matters)
- If everything goes on correctly:

$$\text{King} - \text{Queen} = \text{Man} - \text{Woman}$$

Word2Vec + SVM (Gensim)

```
model = Word2Vec(review[:40000], vector_size=500, window=10, min_count=10)
```

```
clf = svm.SVC()  
clf.fit(X, Y)
```

```
SVC()
```

```
VX = val_X.apply(generate_vec).values  
Vy = val_y.values
```

```
error = 0  
for i in range(len(Vy)):  
    if clf.predict([VX[i]]) != Vy[i]:  
        error += 1  
print("The accuracy is", (1 - error/len(Vy))*100, "%")
```

```
The accuracy is 85.6 %
```



Accuracy:85.6%

However, surprisingly

ONE-HOT ENCODING + NAIVEBAYES

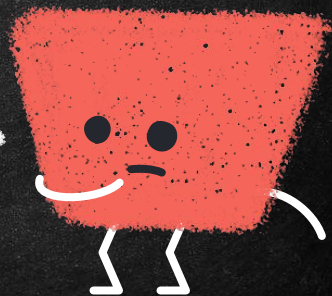
(The exactly the same method that we used in HW3, to identify spam emails)

Gives us a SIMILAR RESULT*

*Accuracy: 84%

RESULTS

Word Embedding	Classifier	Accuracy
W2V	SVM	85.6% (gensim) 84.8% (TF pretrained)
	Random Forest	84.2 %
	FCNN	84.6%
OneHot Encoding	Naive Bayes	84.5%



NEW CHOICE: TF-IDF



ABBREVIATION OF :

Term Frequency - Inverse Document Frequency

TFIDF (RELATIVE FREQUENCY)

$$TF_{w,d_i} = \frac{\text{count}(w)}{\sum_{t \in d_i} \text{count}(t)}$$

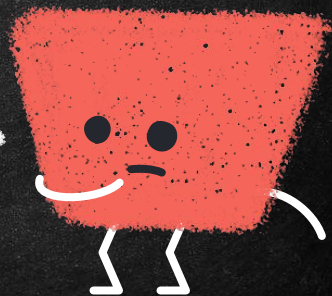
$$IDF_{w,D} = \log \frac{|D|}{1 + \sum_{i=1}^{|D|} I(w, d_i)}$$

$$TFIDF_{w,d_i} = TF_{w,d_i} * IDF_{w,D}$$

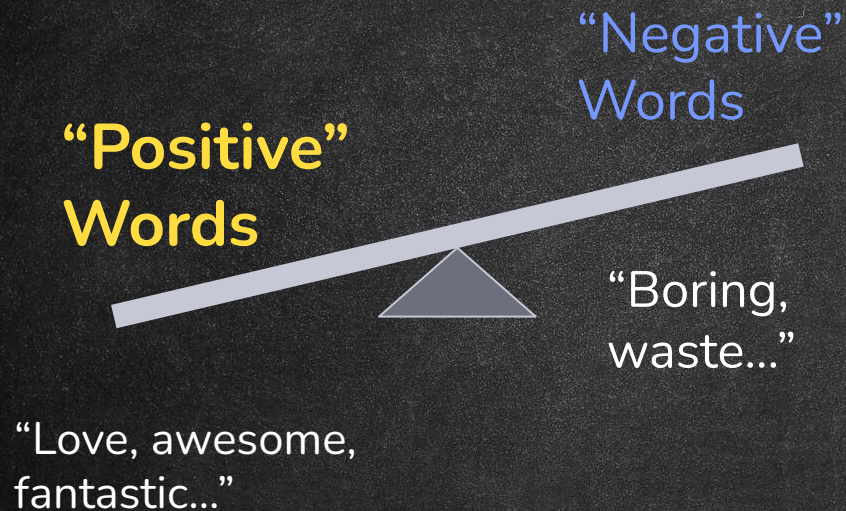
- If a term is frequently occur in one review, but not in the other reviews, its TFIDF value will be high.
- The logic:
Higher TFIDF => Informative
- In some way, it can be viewed as a enhanced OneHot encoding

RESULTS

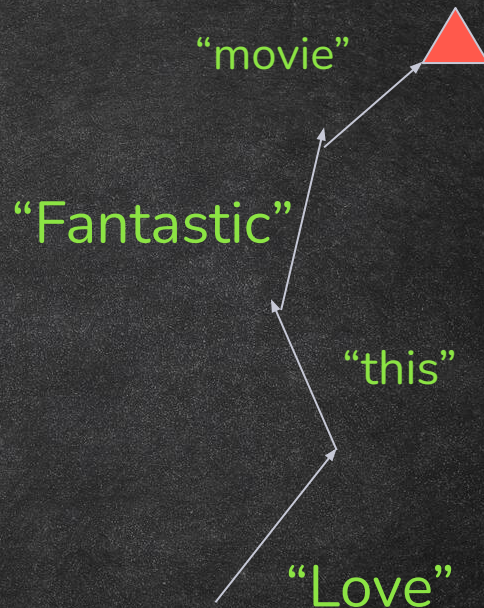
Word Embedding	Classifier	Accuracy
TFIDF	FCNN	Around 88.5%
	Random Forest	83.2%



TFIDF



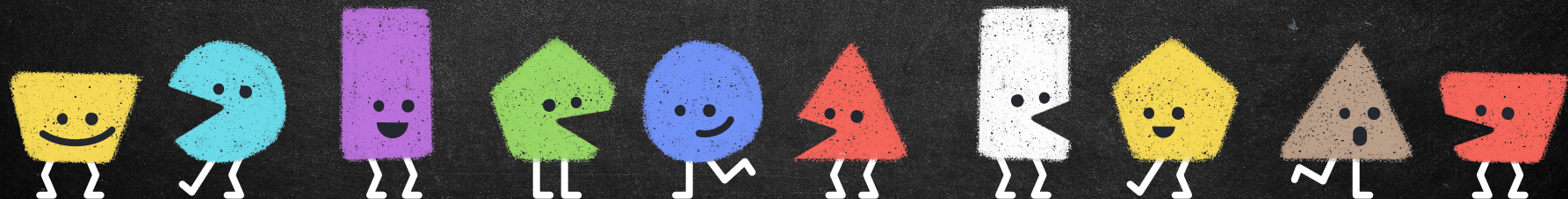
WORD2VEC



TFIDF + FCNN

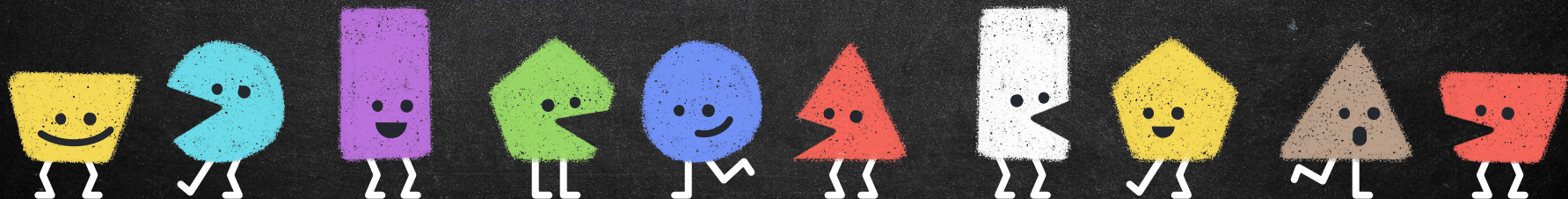
Activation function for output layer: sigmoid

Loss function: Binary Cross entropy loss

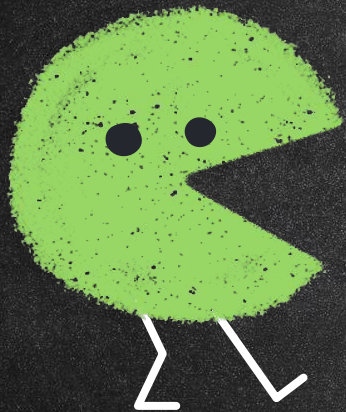


CAN WE DO BETTER?

Optimization of the model



AN UNEXPECTED DISCOVERY!



ONE DAY WHEN EDDY WAS TRAINING
THE MODEL, HE **FORGOT TO DO
THE TEXT PREPROCESSING.**

HOWEVER, THE PERFORMANCE OF
THE MODEL IS **EVEN BETTER**
THAN THE PREVIOUS MODEL!



TEXT PREPROCESSING

1. Removing all numbers
2. Convert all letters to lowercase
3. Removing punctuation
4. Tokenize
5. Remove stopwords
6. Lemmatize

TEXT PREPROCESSING



1. Removing all numbers
2. Convert all letters to lowercase
3. Removing punctuation
4. Tokenize
5. Remove stopwords
6. Lemmatize


IMPROVEMENTS



1. Removing all numbers
2. Convert all letters to lowercase(**weigh** uppercase more)
3. **Removing** punctuation
4. Tokenize **E.g. ABSOLUTELY GREAT**
5. Remove stopwords **Is counted as 2 * absolutely great**
6. Lemmatize

IMPROVEMENTS



1. Removing all numbers
2. Convert all letters to lowercase(**weigh** uppercase more)
3. Removing punctuation(**except ! ?**)
4. Tokenize  split
5. Remove stopwords
6. Lemmatize

IMPROVEMENTS



1. Removing all numbers
2. Convert all letters to lowercase(**weigh uppercase more**)
3. Removing punctuation(**except ! ?**)
4. Tokenize
5. **Remove stopwords**
6. Lemmatize



```
In [1]: from nltk.corpus import stopwords
```

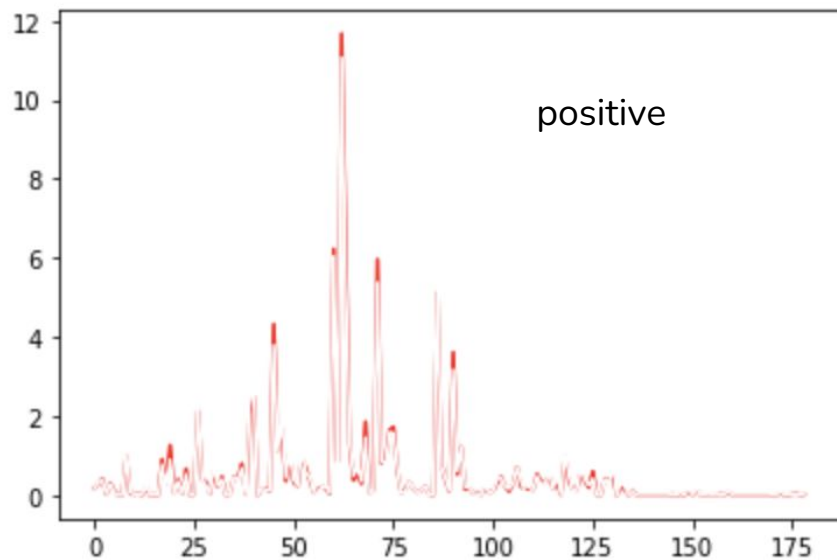
```
In [2]: print(stopwords.words("english"))
```

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your',  
'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's", 'her', 'hers', 'herself', 'it', "it's",  
'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that',  
"that'll", 'these', 'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',  
'do', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', 'at',  
'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', 'above', 'below',  
'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once', 'here',  
'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some', 'such',  
'no', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'don', "don't",  
'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', 've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't",  
'didn', "didn't", 'doesn', "doesn't", 'hadn', "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma',  
'mightn', "mightn't", 'mustn', "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',  
"wasn't", 'weren', "weren't", 'won', "won't", 'wouldn', "wouldn't"]
```


STOPWORDS ?

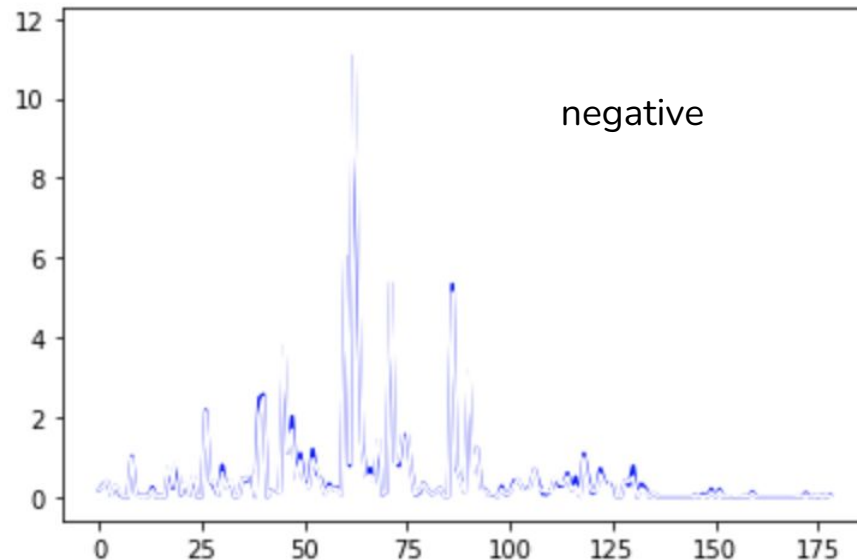


Average
occurrence



index of stopwords

Average
occurrence



index of stopwords

IMPROVEMENTS



1. Removing all numbers
2. Convert all letters to lowercase(**weigh** uppercase more)
3. Removing punctuation(**except ! ?**)
4. Tokenize
- ~~Remove stopwords~~
5. Lemmatize

```
-> rocks : rock  
-> corpora : corpus  
-> better : good
```


OTHER PROBLEM IN WORDLIST



im
you
Ll

- Expand slangs and contractions

a

ahh

ahhhhhh

Ahhhhhhhhh

haha

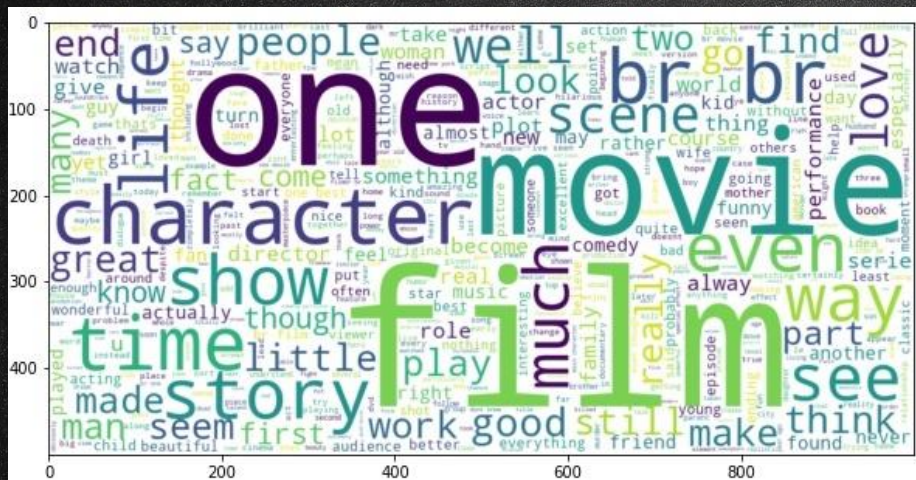
hahahaha

hahaaaahahahaha

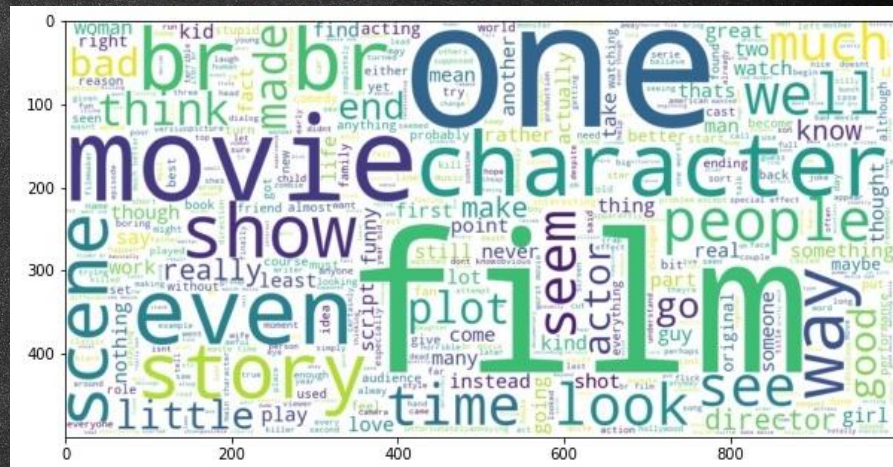
- Reweigh long interjection

E.G. "Hahaha" = 3* "ha"

WORDCLOUD



positive



negative

IMPROVEMENTS



- Removing all numbers
- Convert all letters to lowercase(**weigh** uppercase more)
- **Expand slangs and contractions**
- **Reweigh long interjection**
- Removing punctuation(**except ! ?**)
- Tokenize
- ~~Remove stopwords~~
- Lemmatize
- **Remove high occurrence words**

OTHER

HYPER-PARAMETERS

- The structure of the neural network and its hyperparameters (including dropout rate, regularization)
- The threshold for classification

AFTER TRYING SOME STRUCTURES, WE FINALLY GET A LOW-COST FCNN

Model: "sequential"

Layer (type)	Output Shape	Param #
=====	=====	=====
dropout (Dropout)	(None, 31690)	0
dense (Dense)	(None, 32)	1014112
dense_1 (Dense)	(None, 1)	33
=====	=====	=====

Total params: 1,014,145

Trainable params: 1,014,145

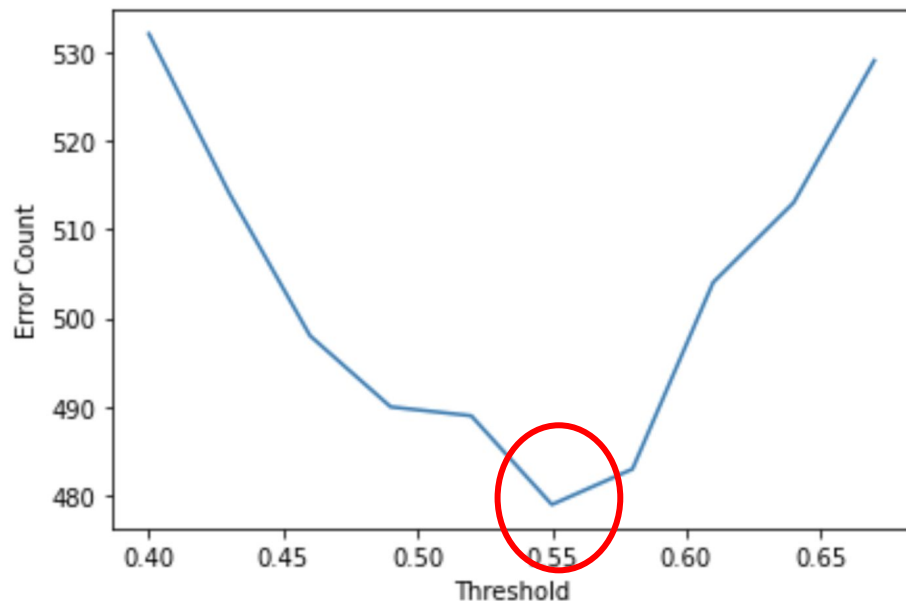
Non-trainable params: 0

Dropout rate = 0.2

GD: ADAM

Batch Size = 32

TUNE THE THRESHOLD ON VALIDATION SET



VALIDATION SIZE:
5000

OPTIMAL
THRESHOLD:
0.55



TFIDFvectorizer

Max_Df=0.95, no stop words removing

Vector dimension around 30000

Dropout=0.2

Hidden Layer (activation:relu)

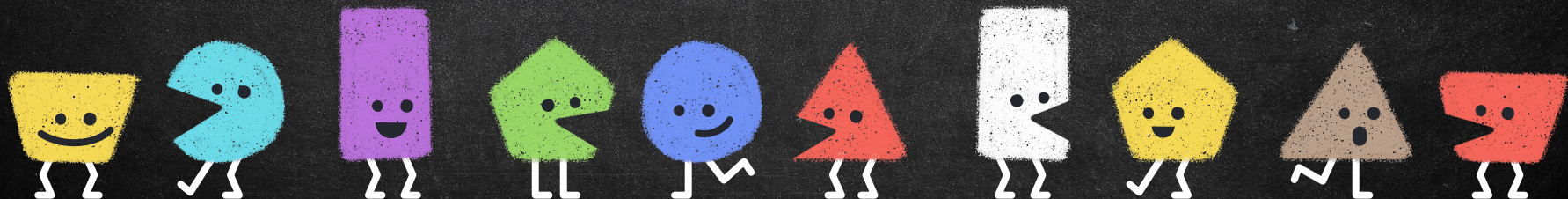
32 nodes in the hidden layer

Output
(sigmoid)

With threshold=0.55,

90.8%

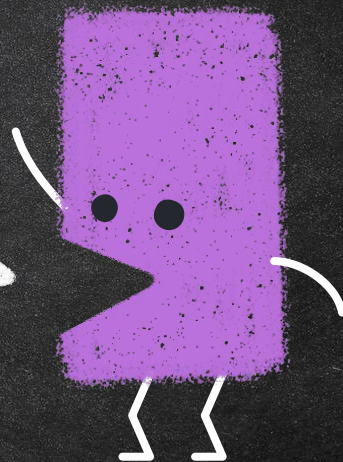
Whoa! That's a big number, aren't you proud of us?



“

PART 4

FUTURE WORK



IF WE HAVE INFINITE TIME...

- FCNN WITH BETTER PERFORMANCE
- TFIDF LOSES THE INFORMATION OF CONTEXT,
WE WANT TO TAKE THAT INTO
CONSIDERATION
- WE WANT TO MAKE OUR WORD EMBEDDING
PROCESS TRAINABLE



THANK YOU FOR
YOUR LISTENING

