

# Sentiment Analysis

For IMDB movie reviews

Team members:

Chong Xian Jun,  
Wang Tian,  
Li Yexin,  
Ahmed hashim taha salim

# Executive Summary

In this study, we conducted a sentiment analysis of movie reviews from IMDB to predict the overall sentiment of the audience towards a given movie. We used a combination of natural language processing techniques and machine learning algorithms to analyze the text of the reviews and classify them as either positive or negative.

We found that our model was able to accurately predict the sentiment of the reviews with a high degree of accuracy. Additionally, we also identified some common themes and words that were more prevalent in positive or negative reviews that becomes indicator of the audiences' sentiment.

Overall, our study provides valuable insights into the sentiment of movie audiences and can be used to inform the marketing and promotion of films.

# Content List

## Introduction and Dataset

- Background and Research Questions
- Introducing the iMDB Movie Review Dataset

## Methodology

- Overview
- Preprocessing
- Model Training
- Model Evaluation

## Results and Findings

- Models to Predict Sentiment
- Important Features

## Conclusion

# **Introduction**

# Background and Research Questions

Conventionally, marketing and promotional campaigns for movies are based on historical data and tend to focus on movies that are similar to previously successful ones. This approach can exclude new or growing markets that might have different preferences.

To better capture the evolving preferences of movie audiences, we made use of natural-language processing methods and machine learning model to perform sentiment analysis on IMDB movie review data to build the predictive model. This can be helpful by automating the extraction of audiences' sentiment on a movie efficiently, thereby informing the promotion and distribution of film.

The study aims to answer the following questions:

1. How do we predict the overall sentiment of the audiences on a particular movie with reviews?
2. What are the important themes that can indicate the audience sentiment about a movie?

# IMDB Movie Reviews Dataset

Source: [IMDB Dataset of 50K Movie Reviews | Kaggle](#)

- Consist of only two columns, “review” and “sentiment”.
- The sentiments are labelled as positive or negative, with each sentiment consist of 25000 samples.
- Contains many syntaxes, symbols and punctuations to be pruned off before analysing with NLP methods

	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

Dataset samples

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 2 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   review      50000 non-null  object
1   sentiment   50000 non-null  object
dtypes: object(2)
memory usage: 781.4+ KB

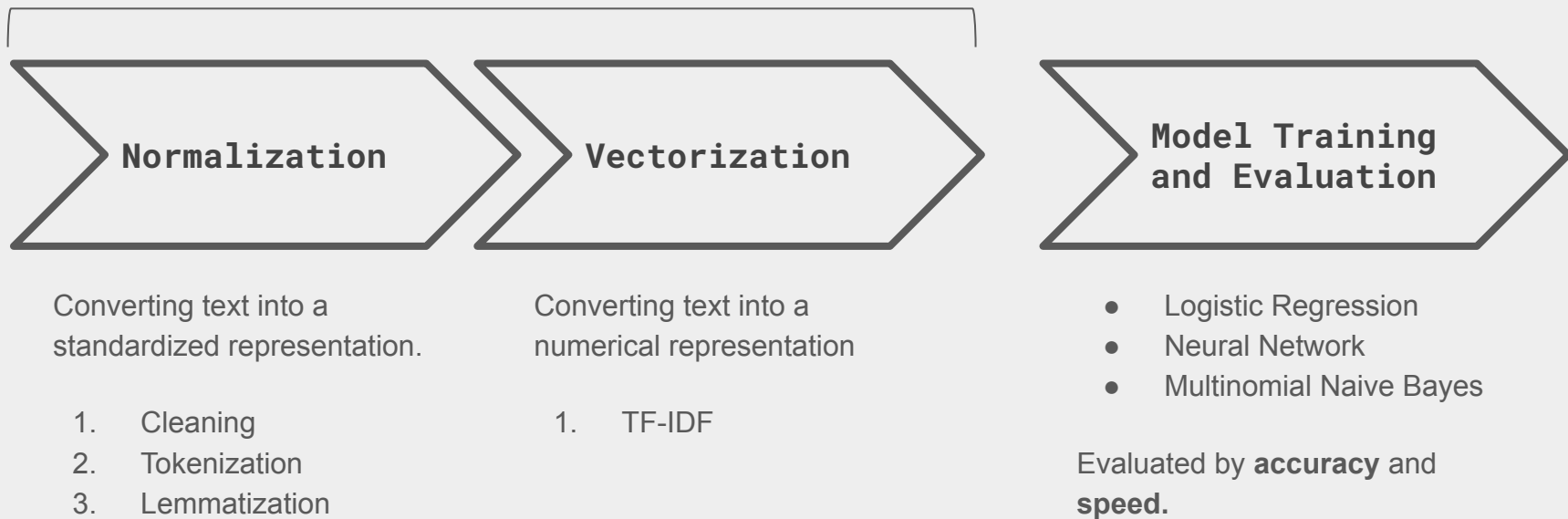
{'positive': 25000, 'negative': 25000}
```

Dataset info

# Methodology

# Methodology Overview

## Preprocessing





# Preprocessing

## Normalization (Cleaning)

The steps involved in cleaning vary according to data. For the movie reviews, we

1. remove html syntax
2. remove url
3. standardise letter casing
4. fix contractions
5. remove stopwords
6. remove all that are not words

Before actually cleaning, we write regex function to check if they are properly cleaned.

Call function:

`check\_if\_contain(data, x, pattern)`

```
1 # define function to check whether they consist of unwanted stuffs with regex
2 def search_review(target_str, pattern):
3     cpat = re.compile(pattern)
4     if cpat.search(target_str):
5         return True
6     else:
7         return False
8
9 def check_if_contain(data, x, pattern, contain=True):
10     # avoiding 'contains_{x}' columns accumulating, set up temp_df
11     temp_df = data.copy()
12     temp_df[f'contains_{x}'] = temp_df.apply(lambda x: search_review(x['review'], pattern), axis=1)
13     return temp_df[temp_df[f'contains_{x}'] == contain]
```

## 1. remove html syntax

- Use html-parser from 'BeautifulSoup' to get rid of the html syntaxes.

```
1 # clean html
2 cleaned_df.review = cleaned_df.review.apply(lambda string: BeautifulSoup(string, 'html.parser').get_text())
3
4 html_regex = r'<.*?>'
5 check_if_contain(cleaned_df, 'symbol', html_regex, True)
```



	review	sentiment	contains_symbol
1513	This all-but-ignored masterpiece is about the ...	positive	True
13736	>>> Great News there is a BBC DVD release sche...	positive	True
22164	In Truffaut book-length interview with Hitchco...	positive	True
27107	One minute into THE UNTOLD and it's already ri...	negative	True
33891	and quite frankly that just sums it up.It is a...	positive	True
42217	Well, I guess I'm emotionally attached to this...	positive	True
48617	<-----Minor Spoilers!----->A woman gets pregna...	positive	True

## 2. remove url

- remove url with regex pattern starting with https:// or www.

```
1 # clean url and check if there's still url
2 cleaned_df.review = cleaned_df.review.str.replace(r"https://\S+|www\.\S+", '', regex = True)
```

## 3. Standardise letter casing

- apply lower() function

```
✓ [9] 1 # standardise casing to lower case and check
1s 2 cleaned_df.review = cleaned_df.review.apply(lambda x: x.lower())
3
4 upper_case_regex = r'[A-Z]'
5 check_if_contain(cleaned_df, 'upper', upper_case_regex, True)
```

review sentiment contains\_upper



#### 4. Fix contractions

- contractions such as I'm, they're... are expanded with regex

```
4 replacement_patterns = [  
5     (r'won\\'t', 'will not'),  
6     (r'can\\'t', 'cannot'),  
7     (r'i\\'m', 'i am'),  
8     (r'ain\\'t', 'is not'),  
9     (r'(\w+)\\"'ll', '\\g<1> will'),  
10    (r'(\w+)n\\'t', '\\g<1> not'),  
11    (r'(\w+)\\"'ve', '\\g<1> have'),  
12    (r'(\w+)\\"'s', '\\g<1> is'),  
13    (r'(\w+)\\"'re', '\\g<1> are'),  
14    (r'(\w+)\\"'d', '\\g<1> would')  
15 ]
```

```
18 class RegexpReplacer(object):  
19     def __init__(self, patterns=replacement_patterns):  
20         self.patterns = [(re.compile(regex), repl) for (regex, repl) in patterns]  
21         # regex applies to regex pattern to be replaced  
22         # repl refers to the replacements  
23  
24     def replace(self, text):  
25         s = text  
26         for (pattern, repl) in self.patterns:  
27             s = re.sub(pattern, repl, s)  
28         return s
```

```
1 replacer = RegexpReplacer()  
2 cleaned_df.review = cleaned_df.review.apply(lambda x: replacer.replace(x))
```

#### 4. Fix contractions (cont)

- check if there is still contractions

```
1 contraction_regex = r"(?!\\")\\[A-Za-z]+('[A-Za-z]*)\\b"  
2 # \\b matches word boundary, ensuring regex matches complete word not only part of it  
3 # [A-Za-z]+ matches first part of the word  
4 # ('[A-Za-z]*) matches the contraction part  
5 # (?!\\") checks that the it does not detect a quotation mark (") at the beginning and take as contraction  
6  
7 check_if_contain(cleaned_df, "contraction", contraction_regex, True)
```

review sentiment contains\_contraction



## 5. Remove stopwords

- make use of nltk package to remove stop words such as "I'm, myself, the..."

```
1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import word_tokenize
```

```
1 nltk.download('stopwords')
2 nltk.download('punkt')
```

```
1 stopwords = nltk.corpus.stopwords.words('english')
2
3 # define a function to remove stopwords from a string
4 # not compulsory, it can be done by vectorization
5 def remove_stopwords(s):
6     # split the string into a list of words
7     words = word_tokenize(s)
8
9     # remove stopwords from the list of words
10    filtered_words = [w for w in words if w not in stopwords]
11
12    # join the filtered words back into a string
13    return " ".join(filtered_words)
```

```
1 cleaned_df.review = cleaned_df.review.apply(remove_stopwords)
```

## 5. Remove stopwords

- check the first review

one reviewers mentioned watching 1 oz episode hooked . right , exactly happened me.the first thing struck oz brutality unflinching scenes violence , set right word go . trust , show faint hearted timid . show pulls punches regards drugs , sex violence . hardcore , classic use word.it called oz nickname given oswald maximum security state penitentiary . focuses mainly emerald city , experimental section prison cells glass fronts face inwards , privacy high agenda . em city home many .. aryan , muslims , gangstas , latinos , christians , italians , irish .... scuffles , death stares , dodgy dealings shady agreements never far away.i would say main appeal show due fact goes shows would dare . forget pretty pictures painted mainstream audiences , forget charm , forget romance ... oz mess around . first episode ever saw struck nasty surreal , could say ready , watched , developed taste oz , got accustomed high levels graphic violence . violence , injustice ( crooked guards sold nickel , inmates kill order get away , well mannered , middle class inmates turned prison bitches due lack street skills prison experience ) watching oz , may become comfortable uncomfortable viewing .... thats get touch darker side



## 6. Remove all that are not words

- use regex to choose only words and white space

```
2 symbol_regex = r'^a-z|\s'  
3 cleaned_df.review = cleaned_df.review.str.replace(symbol_regex, '', regex = True)
```

The first review:

one reviewers mentioned watching oz episode hooked right exactly happened methe first thing struck oz brutality unflinching scenes violence set right word go trust show faint hearted timid show pulls punches regards drugs sex violence hardcore classic use wordit called oz nickname given oswald maximum security state penitentiary focuses mainly emerald city experimental section prison cells glass fronts face inwards privacy high agenda em city home many aryan muslims gangstas latinos christians italians irish scuffles death stares dodgy dealings shady agreements never far awayi would say main appeal show due fact goes shows would dare forget pretty pictures painted mainstream audiences forget charm forget romance oz mess around first episode ever saw struck nasty surreal could say ready watched developed taste oz got accustomed high levels graphic violence violence injustice crooked guards sold nickel inmates kill order get away well mannered middle class inmates turned prison bitches due lack street skills prison experience watching oz may become comfortable uncomfortable viewing thats get touch darker side

# Preprocessing

## Normalization (Tokenization + Lemmatization)

Tokenization is the process of extracting words as separate tokens.

Lemmatization is the process of returning words into its base form, e.g. “plays” -> “play”

Tokenization is necessary before lemmatization as it allows the lemmatization algorithm to operate on individual words rather than on the entire text. This makes the lemmatization process more efficient and accurate.

Several packages from nltk library are used:

```
1 # packages to lemmatize the dataset
2 from nltk.stem import WordNetLemmatizer
3 from nltk.tokenize import word_tokenize
4 nltk.download('wordnet')
5 nltk.download('omw-1.4')
```

## Normalization (Tokenization + Lemmatization)

```
1 def lemmatize_string(s):
2     lemmatizer = WordNetLemmatizer()
3     # split the string into a list of words (tokenization)
4     words = word_tokenize(s)
5
6     # lemmatize each word in the list
7     lemmas = [lemmatizer.lemmatize(w) for w in words]
8
9     # join the lemmas back into a string
10    return " ".join(lemmas)
```

```
1 cleaned_df['lemma_review']=cleaned_df['review'].apply(lambda z: lemmatize_string(z))
```

## Lemmatized review sample compared with that before lemmatization:

```
1 cleaned_df.lemma_review[0]
```

```
'one reviewer mentioned watching oz episode hooked right exactly happened me the first thing struck oz brutality unflinching scene violence set right word go tr  
ust show faint hearted timid show pull punch regard drug sex violence hardcore classic use word it called oz nickname given oswald maximum security state penite  
ntary focus mainly emerald city experimental section prison cell glass front face inwards privacy high agenda em city home many aryan muslim gangsta latino chri  
stian italian irish scuffle death stare dodgy dealing shady agreement never far away i would say main appeal show due fact go show would dare forget pretty pict  
ure painted mainstream audience forget charm forget romance oz mess around first episode ever saw struck nasty surreal could say ready watched developed taste o  
z got accustomed high level graphic violence violence injustice crooked guard sold nickel inmate kill order get away well mannered middle class inmate turned pr  
ison bitch due lack street skill prison experience watching oz may become comfortable uncomfortable viewing thats get touch darker side<'
```

```
1 cleaned_df.review[0]
```

```
'one reviewers mentioned watching oz episode hooked right exactly happened me the first thing struck oz brutality unflinching scenes violence set right  
word go trust show faint hearted timid show pulls punches regards drugs sex violence hardcore classic use word it called oz nickname given oswald ma  
ximum security state penitentiary focuses mainly emerald city experimental section prison cells glass fronts face inwards privacy high agenda em city hom  
e many arians muslims gangstas latinos christians italians irish scuffles death stares dodgy dealings shady agreements never far away i  
would say main appeal show due fact goes shows would dare forget pretty pictures painted mainstream audiences forget charm forget romance oz mess arou  
nd first episode ever saw struck nasty surreal could say ready watched developed taste oz got accustomed high levels graphic violence violence inj  
ustice crooked guards sold nickel inmates kill order get away well mannered middle class inmates turned prison bitches due lack street skills prison exp  
erience watching oz may become comfortable uncomfortable viewing thats get touch darker side <'
```

# Preprocessing

## Vectorization

Encode the texts into numerical form as before inputting the data into machine learning models.  
For this study, TF-IDF algorithm is used due to its ability to take into account the frequency and context of words.

Before vectorization, the normalized dataset is split into X\_train, X\_test, y\_train, y\_test.  
The following is the code for TF-IDF vectorization:

```
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 t_vect = TfidfVectorizer(max_df = 0.8, min_df = 3, tokenizer=word_tokenize,
6                          stop_words=frozenset(stop_words_tokens), max_features=5000)
7
8 X_tfidfvect_train = pd.DataFrame(t_vect.fit_transform(X_train).toarray(),
9                                  columns=t_vect.get_feature_names_out())
10 X_tfidfvect_test = pd.DataFrame(t_vect.transform(X_test).toarray(),
11                                 columns=t_vect.get_feature_names_out())
```

Vectorized output, `X_tfidfvect_train`

```
12 X_tfidfvect_train.head()
```

	<b>a</b>	<b>aaron</b>	<b>abandoned</b>	<b>abc</b>	<b>ability</b>	<b>able</b>	<b>absence</b>	<b>absent</b>	<b>absolute</b>	<b>absolutely</b>	<b>...</b>
<b>0</b>	0.053855	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...
<b>1</b>	0.000000	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...
<b>2</b>	0.000000	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...
<b>3</b>	0.000000	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...
<b>4</b>	0.000000	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...

```
X_tfidfvect_train.shape= (50000, 5000)
```

Vectorize y\_train and y\_test to numerical values as well. {'positive': 1, 'negative': 0}

```
1 def encode_y(y):
2     y[y == 'positive'] = 1
3     y[y == 'negative'] = 0
4     y.astype(int)
5     return y
6
7 y_train = encode_y(y_train).astype(int)
8 y_test = encode_y(y_test).astype(int)
```

# Model Training and Evaluation

Models trained are:

- Logistic Regression
- Neural Network
- Multinomial Naive Bayes

Using sklearn API, we are able to train and test different models easily. In our study, the models are evaluated by their accuracy as the data has equal class distribution. The speed of models are evaluated by using time() function as well to determine the best model.



- Logistic Regression

```
▼ # Logistic Regression
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()

t0 = time()    # start time
lr.fit(X_tfidfvect_train, y_train)
print ("training done in %fs" %(time() - t0))    # end of training

t1 = time()
lr_accuracy = lr.score(X_tfidfvect_test, y_test)
print ("testing done in %fs" %(time() - t1))
print('logistic regression accuracy:', lr_accuracy)
```

```
training done in 6.706810s
testing done in 0.146578s
logistic regression accuracy: 0.8864
```

## - Neural Network

```
▼ # neural network
from sklearn.neural_network import MLPClassifier
▼ mlp = MLPClassifier(solver='lbfgs', alpha=1e-5,
                      hidden_layer_sizes=(5, 2), random_state=1, max_iter=200)

t0 = time()    # start time
mlp.fit(X_tfidfvect_train, y_train)
print("training done in %fs" %(time() - t0))    # end of training

t1 = time()
mlp_accuracy = mlp.score(X_tfidfvect_test, y_test)
print("testing done in %fs" %(time() - t1))

print('mlp accuracy:', mlp_accuracy)
```

```
training done in 516.074917s
testing done in 0.197302s
mlp accuracy: 0.8795333333333333
```

- Multinomial Naive Bayes

```
▼ #Multinomial Naive Bayes
from sklearn.naive_bayes import MultinomialNB
mnb=MultinomialNB()

t0 = time()    # start time
mnb.fit(X_tfidfvec_train,y_train)
print ("training done in %fs" %(time() - t0))    # end of training

t1 = time()
mnb_accuracy = mnb.score(X_tfidfvec_test, y_test)
print ("testing done in %fs" %(time() - t1))

print('mnb accuracy:', mnb_accuracy)
```

```
training done in 0.396690s
testing done in 0.141874s
mnb accuracy: 0.8536666666666667
```

## **Results and Findings**

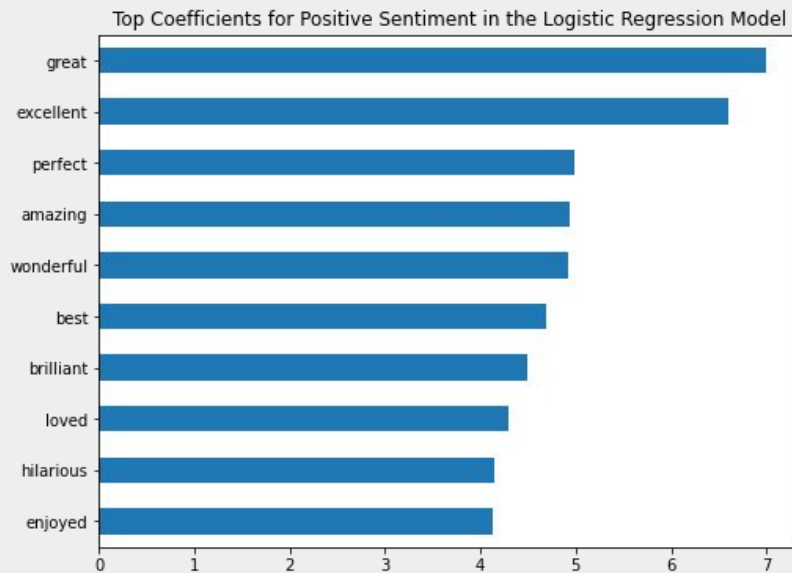
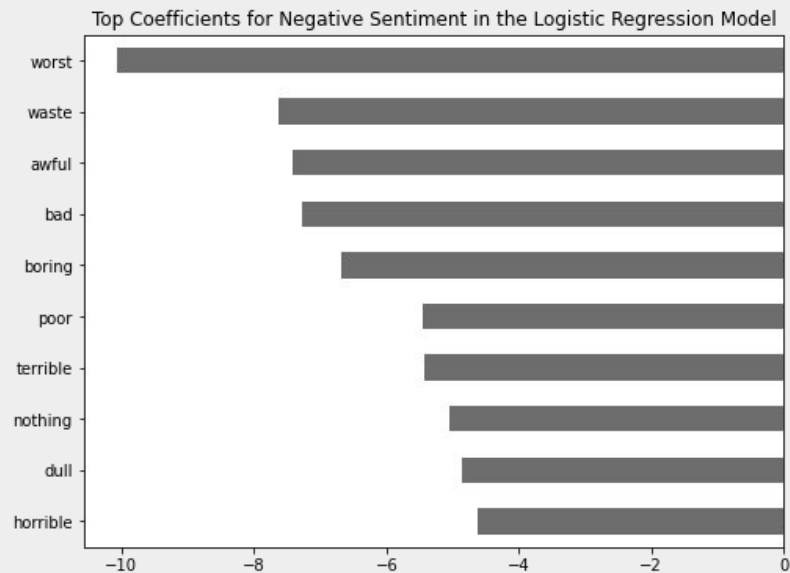
# Results and Findings

All in all, we can make use of predictive models to predict audience sentiment based on their review comments. For this purpose, **logistic regression** is the best performing model as it has the highest accuracy with fairly low training and testing time.

	Models		
	Logistic Regression	Multinomial Naive Bayes	Neural Network
Accuracy	0.8864	0.8537	0.8795
Training Time (s)	6.706810	0.396690	509.5316
Testing Time (s)	0.146578	0.141874	0.1730

# Results and Findings

From the logistic regression, we extract the most important features (words) that people associate strongly with positive or negative sentiments.



# Conclusion

# Conclusion

In this study, we conducted a sentiment analysis of movie reviews from IMDB to predict the sentiment of the audience towards a given movie based on unstructured review. It is important to note that thorough preprocessing involving text normalization and vectorization is needed before model training. Our findings showed that logistic regression was able to accurately predict the sentiment of the reviews with accuracy of 0.8864 and fairly efficient training time. The common themes and words that were indicative of positive or negative sentiment are extracted from the model to allow for interpretation.

These findings have important implications for the film industry as they can be used to inform the marketing and promotion of movies, helping the stakeholders to capture the new market with ever-evolving preferences on movies.