

NLP Task 2: Sentiment Analysis (750 words)

Literature Review and Rationale: (150 words) 15% 174w

Sentiment analysis is an NLP process that aims to classify language data concerning underlying attitudes, emotions, or sentiments (Valdivia, Luzon, & Herrera, 2017). A statement can be classified along a continuum from positive to negative sentiment. These data can then be used to evaluate opinions towards any number of topics. Sentiment analysis can be applied to individuals, but it is beneficial for evaluating public sentiment on a given topic. User-generated content websites are an ideal source of sentiment data, and the results can inform decision-making in several situations.

The interest in sentiment analysis has increased significantly due to the large amount of stored text in web applications and the importance of online customer opinions. As a result, more than 1 million research papers contain the term "sentiment analysis," and various start-ups have been created to analyse sentiments in social media companies (Valdivia, Luzon, & Herrera, Sentiment Analysis on TripAdvisor: Are There Inconsistencies in User Reviews?, 2017).

Multiple studies on TripAdvisor exist, but there is no complete analysis from the sentiment analysis viewpoint. This article proposes TripAdvisor as a source of data for sentiment analysis tasks.

Data: (200) 20% 695w

Wrangling

VADER (Valence Aware Dictionary and Sentiment Reasoner) was used to get sentimental scores of the user reviews and convert them into three categorical sentiments (positive, negative, and neutral). VADER is a lexicon and rule-based sentiment analysis tool specifically attuned to sentiments expressed in social media (Borg & Boldt, 2020). The following code was used to create two new columns that give the sentiment score and the sentiment category of the review:

```
# Creating sentimental polarity
analyzer = SentimentIntensityAnalyzer()
def compound_score(txt):
    return analyzer.polarity_scores(txt)["compound"]

# Sentiments
def sentiment(score):
    emotion = ""
    if score >= 0.5:
```

```

        emotion = "Positive"
    elif score <= -0.5:
        emotion = "Negative"
    else:
        emotion = "Neutral"
    return emotion

# Applying compound score
polarity_scores = df["review"].astype("str").apply(compound_score)
df["Sentiment_Score"] = polarity_scores

## Applying Sentiment
df["Sentiment"] = df["Sentiment_Score"].apply(sentiment)

```

Fig. 1 displays the output of the dataframe.

	hotelId	reviewId	stars	date	review	Sentiment_Score	Sentiment
0	0	0	5	November 2021	Just overnight and breakfast. The breakfast is...	0.5254	Positive
1	0	1	5	November 2021	Great stay and fantastic room view, the staff ...	0.9493	Positive
2	0	2	4	November 2021	Central, and able to walk to most things in To...	0.9451	Positive
3	0	3	5	November 2021	Was there for work. It was close to everything...	0.9538	Positive
4	0	4	5	November 2021	The room was very clean. Great location close ...	0.9728	Positive

Fig 1. Screenshot of top five rows in 'hotelReviews' data frame.

It was determined that no missing data were present in the data frame using the `isna().sum()` function.

Summary of data for NLP task

After checking our data for class balance, this report found that most of the user reviews were positive (Fig. 2).

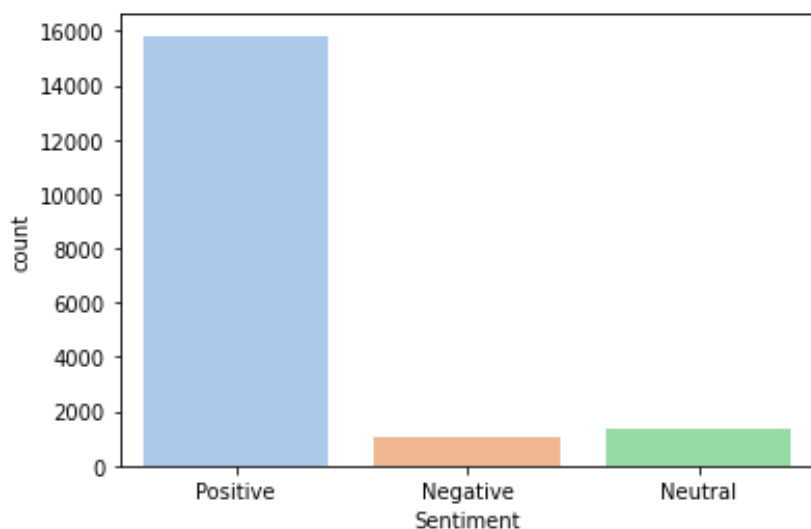


Fig 2. Distribution of sentiment class in 'hotelReviews' data frame.

This study performed down-sampling to balance the classes and created a subset of the data which contained the top 1000 rows from each of the three 'Sentiment' categories (Fig. 3). This was achieved via the following code:

```
# Function to retrieve top few numbers of each category
def get_top_data(top_n = 1000):
    top_data_df_positive = df[df['Sentiment'] == 'Positive'].head(top_n)
    top_data_df_negative = df[df['Sentiment'] == 'Negative'].head(top_n)
    top_data_df_neutral = df[df['Sentiment'] == 'Neutral'].head(top_n)
    top_data_df_small = pd.concat([top_data_df_positive, top_data_df_negative,
top_data_df_neutral])
    return top_data_df_small

# Function call to get the top 1000 from each sentiment
df = get_top_data(top_n=1000)

# After selecting top few samples of each sentiment
print("After segregating and taking equal number of rows for each sentiment:")
print(df['Sentiment'].value_counts())
df.head(10)
```

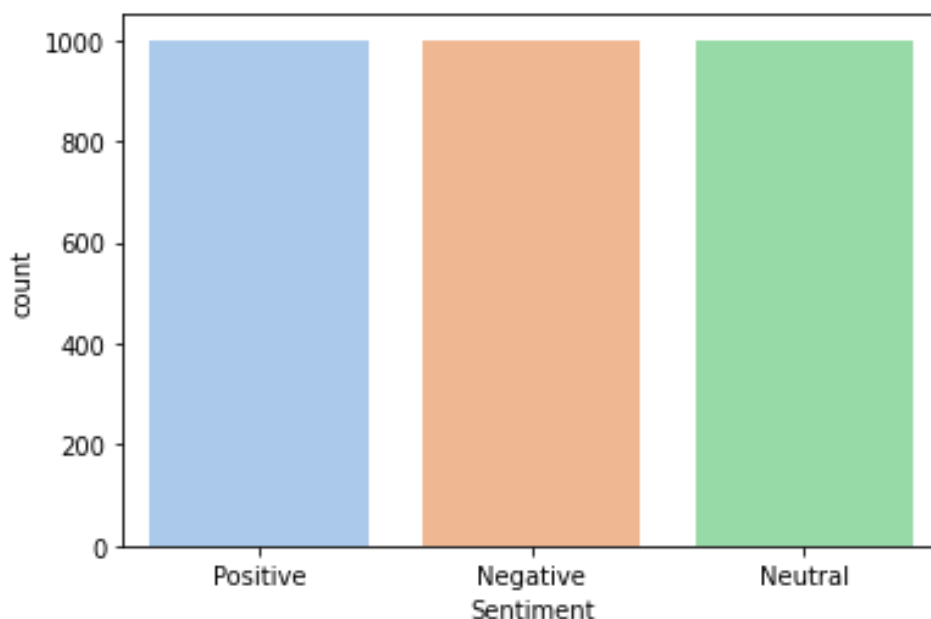


Fig 3. Distribution of sentiment class in 'hotelReviews' data frame after downsampling.

Visualisation of data

Fig. 4 illustrates that people with 5-star ratings have the highest positive sentiment. Whereas at lower ratings, the sentiment is primarily negative.

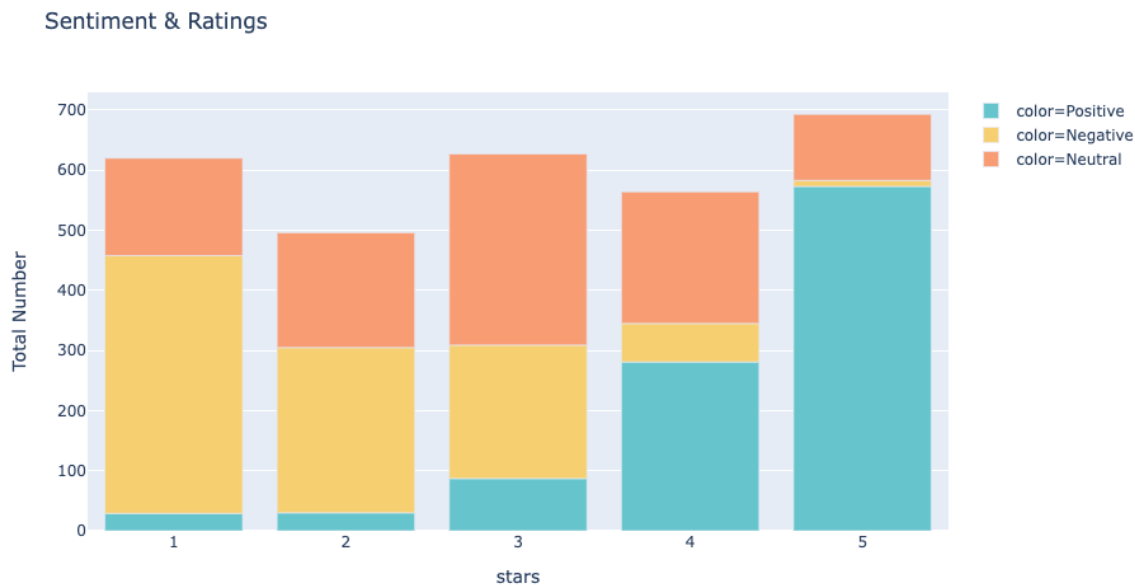


Fig 4. Distribution of sentiment class and customer ratings.

A pie chart using the **Plotly** library (Fig. 5) displays the distribution of different ratings. The range is between 16.5% (2-star review) and 23.1% (5-star review), which demonstrates the distribution is well spread across the ratings.

Rating Distribution of the data

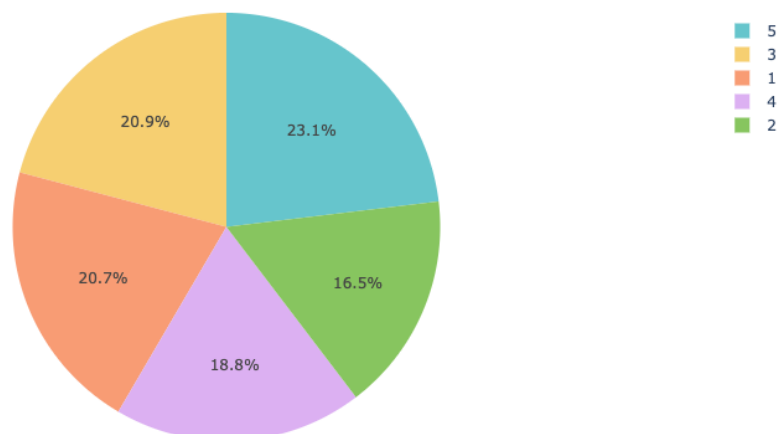


Fig 8. Distribution of top ten keywords.

Most used words were removed to improve the performance and increase accuracy of the model.

```
words = ["hotel", "room", "rooms", "hotels"]
for x in words:
    data["review"] = data["review"].astype(str).str.replace(x, "")
```

In computing, stop words are words that are filtered out before or after the natural language data (text) are processed. An example of stop words are: 'a', 'an', 'the', 'this', 'not'. There is no universal set of stop words, and the list of stop words varies between packages. Removal of stop words removes the necessary words required to get the sentiment, and sometimes it can change the sentence's meaning (Arumugam & Shanmugamani, 2018). For this reason, the study skipped the removal of stop words for the sentiment analysis.

Tokenisation is when the sentence/text is split into an array of words called tokens (Reese & Bhatia, 2018). This helps to transform each word separately and is also required to transform words into numbers. The **simple_preprocess** function converts text to lower case and remove punctuations. It has min and max length parameters, which help filter out rare words that will fall in that range of lengths (Fig. 9). **Simple_preprocess** was used to get the tokens for the data frame as it does most of the pre-processing under the hood. The following code was used to tokenise the review text.

```
# Tokenisation

from gensim.utils import simple_preprocess
# Tokenize the text column to get the new column 'tokenized_text'
df['tokenized_text'] = [simple_preprocess(line, deacc=True) for line in df['review']]
```

```
0      [just, overnight, and, breakfast, the, breakfa...
1      [great, stay, and, fantastic, view, the, staff...
2      [central, and, able, to, walk, to, most, thing...
3      [was, there, for, work, it, was, close, to, ev...
4      [the, was, very, clean, great, location, close...
6      [booked, nights, in, deluxe, suite, ocean, vie...
7      [before, they, turned, the, palm, house, to, t...
8      [amazing, with, amazing, team, members, very, ...
9      [the, quarterdeck, bar, is, great, setting, ov...
11     [it, was, pleasurable, stay, and, short, walk,...
Name: tokenized_text, dtype: object
```

Fig 9. Output following Tokenisation of customer reviews text.

The Stemming process reduces the words to its' root word. Unlike Lemmatization which uses grammar rules and dictionaries for mapping words to root form, stemming removes suffixes/prefixes (Fig. 10). Stemming is widely used in the application of search engine optimisation (SEOs), web search results, and information retrieval. If the root word matches in the text somewhere, it will help retrieve all the related documents in the search. The following code was used to stem the tokens in 'tokenized_text':

```
# Stemming
from gensim.parsing.porter import PorterStemmer
porter_stemmer = PorterStemmer()

# Get the stemmed_tokens
df['stemmed_tokens'] = [[porter_stemmer.stem(word) for word in tokens] for tokens in
df['tokenized_text']] ]
```

```
0    [just, overnight, and, breakfast, the, breakfa...
1    [great, stai, and, fantast, view, the, staff, ...
2    [central, and, abl, to, walk, to, most, thing,...
3    [wa, there, for, work, it, wa, close, to, ever...
4    [the, wa, veri, clean, great, locat, close, to...
6    [book, night, in, delux, suit, ocean, view, th...
7    [befor, thei, turn, the, palm, hous, to, the, ...
8    [amaz, with, amaz, team, member, veri, friendl...
9    [the, quarterdeck, bar, is, great, set, overlo...
11   [it, wa, pleasur, stai, and, short, walk, to, ...
Name: stemmed_tokens, dtype: object
```

Fig 10. Output following Stemming of tokenised text.

A unique id identified each unique word in the dictionary object. This was used for creating representations of texts. A Bag-of-Words (BOW) corpus was created using this method to build the TF-IDF model. The list of words created a dictionary. This allowed the sentences to be converted to a list of words and then fed to the **corpora. Dictionary** as a parameter:

```
# Building a dictionary

from gensim import corpora
# Build the dictionary
mydict = corpora.Dictionary(df['stemmed_tokens'])
print("Total unique words:")
print(len(mydict.token2id))
print("\nSample data from dictionary:")
i = 0
# Print top 4 (word, id) tuples
for key in mydict.token2id.keys():
    print("Word: {} - ID: {}".format(key, mydict.token2id[key]))
    if i == 3:
```

```
        break
    i += 1
```

After building the dictionary, the report identified a total of 7156 unique words in the 'stemmed_tokens' column (Fig.11).

```
Total unique words:
7156

Sample data from dictionary:
Word: across - ID: 0
Word: and - ID: 1
Word: breakfast - ID: 2
Word: choic - ID: 3
```

Fig 11. The output of the dictionary.

NLP Task Implementation: (200 words) 35% 208 w

Feature normalisation

The data was split into training (70%) and test (30%) sets to train the model on a separate data set to the validation data set. The test data is what the model will predict the classes on, and it will be compared with the original labels to check the accuracy and other model performance metrics.

```
# Splitting into train and test sets

from sklearn.model_selection import train_test_split
# Train Test Split Function
def split_train_test(df, test_size=0.3, shuffle_state=True):
    X_train, X_test, Y_train, Y_test = train_test_split(df[['stars', 'stemmed_tokens']],
                                                         df['Sentiment'],
                                                         shuffle = shuffle_state,
                                                         test_size = test_size,
                                                         random_state = 15)

    print("Value counts for Train sentiments")
    print(Y_train.value_counts())
    print("Value counts for Test sentiments")
    print(Y_test.value_counts())
    print(type(X_train))
    print(type(Y_train))
    X_train = X_train.reset_index()
    X_test = X_test.reset_index()
    Y_train = Y_train.to_frame()
    Y_train = Y_train.reset_index()
    Y_test = Y_test.to_frame()
    Y_test = Y_test.reset_index()
    print(X_train.head())
    return X_train, X_test, Y_train, Y_test
```



```
# Call the train_test_split
X_train, X_test, Y_train, Y_test = split_train_test(df)
```

As seen below, the train and test sets both have an even distribution of all three sentiment categories (Fig. 12). This was crucial to prevent model bias.

```
Value counts for Train sentiments
Neutral      707
Negative     700
Positive     693
Name: Sentiment, dtype: int64
Value counts for Test sentiments
Positive      307
Negative      300
Neutral       293
Name: Sentiment, dtype: int64
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.series.Series'>
   index  ...  stemmed_tokens
0  12648  ...  [we, have, visit, casino, in, sydney, gold, co...
1   5044  ...  [hard, to, find, on, the, strand, for, easter,...
2    865  ...  [the, locat, is, perfect, the, price, is, veri...
3    826  ...  [thi, properti, ha, everyth, go, for, it, and,...
4   3756  ...  [have, stai, in, mani, of, the, motel, in, the...
```

Fig 12. Output after splitting the data into test and training sets.

TF-IDF is computed by multiplying a local component like term frequency (TF) with a global component, inverse document frequency (IDF), and optionally normalising the result to unit length. The Gensim package was used to create the TF-IDF model. The TF-IDF model was trained, and TF-IDF vectors were generated using the following code:

```
# Created TFIDF model

# Train the tfidf Model
from gensim.models import TfidfModel
# Make sure the dictionary is created from the previous block
# BOW corpus is required for tfidf model
corpus = [mydict.doc2bow(line) for line in df['stemmed_tokens']]

# TF-IDF Model
tfidf_model = TfidfModel(corpus)
# Generating TFIDF vectors
start_time = time.time()
OUTPUT_FOLDER = '/content/drive/My Drive/A3/'

tfidf_filename = OUTPUT_FOLDER + 'train_review_tfidf.csv'
# Storing the tfidf vectors for training data in a file
vocab_len = len(mydict.token2id)
with open(tfidf_filename, 'w+') as tfidf_file:
```

```

for index, row in X_train.iterrows():
    doc = mydict.doc2bow(row['stemmed_tokens'])
    features = gensim.matutils.corpus2csc([tfidf_model[doc]],
num_terms=vocab_len).toarray()[:,0]
    if index == 0:
        header = ",".join(str(mydict[ele]) for ele in range(vocab_len))
        print(header)
        print(tfidf_model[doc])
        tfidf_file.write(header)
        tfidf_file.write("\n")
    line1 = ",".join( [str(vector_element) for vector_element in features] )
    tfidf_file.write(line1)
    tfidf_file.write('\n')
print("Time taken to create tfidf for :" + str(time.time() - start_time))

```

ML algorithm

A decision tree classification model was used to perform the sentiment classification. A decision tree classifier is a supervised machine learning algorithm for classification problems. The **scikit-learn** package was used for implementing the decision tree classifier. The fit function is used to fit the input feature vectors against the sentiments in the train data set. The following code shows how this was achieved:

```

# Training sentiment classification model using TF-IDF vectors

from sklearn.tree import DecisionTreeClassifier
import time
start_time = time.time()

# Read the TFIDF vectors
tfidf_df = pd.read_csv('/content/drive/My Drive/A3/train_review_tfidf.csv')

# Initialize the model
clf_decision_tfidf = DecisionTreeClassifier(random_state=2)

# Fit the model
clf_decision_tfidf.fit(tfidf_df, Y_train['Sentiment'])
print("Time to taken to fit the TF-IDF as input for classifier: " + str(time.time() -
start_time))

```

Hyper-parameters

Quality metrics

The `feature_importances_` attribute was used to get the important features of the model. The value for each feature was produced with the higher the value suggesting more importance (Fig. 13). Getting the important features was implemented via the following code:

```
# Find out the important features from the TFIDF classification model

importances = list(clf_decision_tfidf.feature_importances_)
feature_importances = [(feature, round(importance, 10)) for feature, importance in
zip(tfidf_df.columns, importances)]
feature_importances = sorted(feature_importances, key = lambda x: x[1], reverse = True)
# print(feature_importances)
top_i = 0
for pair in feature_importances:
    print('Variable: {:10} Importance: {}'.format(*pair))
    if top_i == 10:
        break
    top_i += 1
```

Variable: not	Importance: 0.0743225813
Variable: love	Importance: 0.037668821
Variable: friendli	Importance: 0.0375842673
Variable: great	Importance: 0.0375575124
Variable: no	Importance: 0.0280894298
Variable: comfort	Importance: 0.0216480563
Variable: and	Importance: 0.0196681452
Variable: locat	Importance: 0.0147433007
Variable: on	Importance: 0.0144248843
Variable: the	Importance: 0.0128809058
Variable: recommend	Importance: 0.0125819778

Fig 13. Ten most important features of the TFIDF classification model.

Fig. 13 displays the three most important features: 'not', 'love' and 'friendly', respectively.

Summary of outputs

As shown below, the decision tree classifier generates an overall accuracy of 56% (Fig 14). The model appears to predict the 'positive' sentiment class more accurately. This may be due to the slightly higher number of observations.

```
# Testing the model

from sklearn.metrics import classification_report
test_features_tfidf = []
import time
start_time = time.time()
for index, row in X_test.iterrows():
    doc = mydict.doc2bow(row['stemmed_tokens'])
    features = gensim.matutils.corpus2csc([tfidf_model[doc]],
num_terms=vocab_len).toarray()[:,0]
    test_features_tfidf.append(features)
test_predictions_tfidf = clf_decision_tfidf.predict(test_features_tfidf)
print(classification_report(Y_test['Sentiment'],test_predictions_tfidf))
print("Time taken to predict using TF-IDF:" + str(time.time() - start_time))
```

	precision	recall	f1-score	support
Negative	0.56	0.55	0.55	300
Neutral	0.44	0.44	0.44	293
Positive	0.68	0.69	0.68	307
accuracy			0.56	900
macro avg	0.56	0.56	0.56	900
weighted avg	0.56	0.56	0.56	900

Fig 14. Summary of model performance.

The model's precision, also known as positive pred value, is 68% for the 'positive' class. This means the model correctly predicted 68% of the 'positive' class out of all the predicted 'positive' class observations. The recall, also known as sensitivity, was 55% for the 'negative' class. This means the model correctly 55% of the 'negative' class out of the actual 'negative' class observations. The F1 score can be interpreted as the harmonic mean of precision and recall. The measure is between 0-1, with a score closer to 1 representing a higher performance. The F1 score for the 'positive' class is 68%. The accuracy of the model from the test data is 56%. This means the model correctly classified 56% of the total 3000 observations.

Visualise outputs

NLP task alignment with issue

By performing Sentiment Analysis, hotels can better understand the sentiment of the customer review without relying on the overall customer rating. The overall customer rating can be limiting in expressing the various aspects of a customer's experience. Hotels could improve their customer experience by using a sentiment analyser to expand their knowledge of customer reviews on TripAdvisor.