

```
In [1]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns
from sklearn.linear_model import LogisticRegression, LogisticRegressionCV
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfTransformer
from nltk.corpus import stopwords
import nltk
from nltk.corpus import stopwords
import re
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
import pickle
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
from sklearn.pipeline import Pipeline
from sklearn.metrics import roc_curve
from sklearn.metrics import roc_auc_score
from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator
from sklearn.decomposition import TruncatedSVD

In [95]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib as mpl
from sklearn.linear_model import LogisticRegression, LogisticRegressionCV
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfTransformer
from nltk.corpus import stopwords
import nltk
from nltk.corpus import stopwords
import re
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
import pickle
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
from sklearn.pipeline import Pipeline
from nltk.stem import PorterStemmer
from nltk.stem import LancasterStemmer
```

Data Collection

```
In [2]: review_json_path = 'yelp_academic_dataset_review.json'
business_json_path = 'yelp_academic_dataset_business.json'
checkin_json_path = 'yelp_academic_dataset_review.json'
```

Reading Business Data

```
In [3]: df_b = pd.read_json(business_json_path, lines=True)

Out [22]: df_b.head(1)
```

	business_id	name	city	state	stars	review_count	is_open	attributes	categories
0	A1dDB8S4Mbm0w7JnUjVA	Veggie House	Las Vegas	NV	4.5	1142	1	['Restaurants', 'Ice Cream Parlor', '2', 'BikeParking', ...]	Restaurants, Specialty Food, Japanese, Sushi B...

Cleaning Data

Dropping Irrelevant Columns

```
In [5]: df_b = df_b.drop(columns = ['postal_code', 'address', 'latitude', 'longitude', 'hours'])
```

Filtering

```
In [6]: df_b = df_b.loc[df_b['categories'].str.contains("Restaurants", na=False)]

In [7]: df_b = df_b.loc[df_b['city'] == "Las Vegas"]

In [8]: df_b = df_b.loc[df_b['categories'].str.contains("Japanese", na=False)]
```

NA Values

```
In [9]: df_b.isna().sum()

Out [9]: business_id      0
name                  0
city                  0
state                 0
stars                 0
review_count         0
is_open              0
attributes            2
categories            0
dtypes: int64
```

```
In [10]: df_b = df_b[df_b['attributes'].notna()]

In [11]: business_data_for_edu = df_b.copy()
```

Reading Review Data

```
In [12]: reviews_json = 'yelp_academic_dataset_review.json'
reviews = pd.read_json(reviews_json, lines=True,
                        dtype={'review_id':'str', 'user_id':'str',
                               'business_id':'str', 'stars':'int',
                               'useful':'int', 'funny':'int', 'cool':'int',
                               'text':'str'},
                        chunksize=size)
```

Joining the Data

```
In [13]: chunk_list = []
for chunk_reviews in reviews:
    chunk_review = chunk_review.rename(columns={'stars': 'review_stars'})
    # Inner merge with edited business file so only reviews related to the business remain
    chunk_merged = pd.merge(df_b, chunk_review, on='business_id', how='inner')
    # Show progress
    print(f'{chunk_merged.shape[0]} out of {size}, {related_reviews}')
    chunk_list.append(chunk_merged)
# After trimming down the review file, concatenate all relevant data back to one dataframe
df_b = pd.concat(chunk_list, ignore_index=True, join='outer', axis=0)
```

13793 out of 500,000 related reviews
12887 out of 500,000 related reviews
9511 out of 500,000 related reviews
7024 out of 500,000 related reviews
6014 out of 500,000 related reviews
6175 out of 500,000 related reviews
9005 out of 500,000 related reviews
10734 out of 500,000 related reviews
8803 out of 500,000 related reviews
7187 out of 500,000 related reviews
9327 out of 500,000 related reviews
8578 out of 500,000 related reviews
8681 out of 500,000 related reviews
10684 out of 500,000 related reviews
12022 out of 500,000 related reviews
377 out of 500,000 related reviews

Descriptive Analytics and EDA

```
In [26]: business_data_for_edu.describe()
```

	stars	review_count	is_open
count	437.000000	437.000000	437.000000
mean	3.896314	328.407283	0.581286
std	0.612503	447.283049	0.493922
min	1.000000	3.000000	0.000000
25%	3.500000	44.000000	0.000000
50%	4.000000	161.000000	1.000000
75%	4.500000	429.000000	1.000000
max	5.000000	3512.000000	1.000000

```
In [110]: business_data_for_edu.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 437, entries: 2148, dtypes: object
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   stars       437 non-null    object
 1   name        437 non-null    object
 2   city        437 non-null    object
 3   stars       437 non-null    object
 4   stars       437 non-null    float64
 5   review_count 437 non-null    int64
 6   is_open     437 non-null    int64
 7   attributes  437 non-null    object
 8   categories  437 non-null    object
dtypes: float64(1), int64(2), object(6)
memory usage: 34.1+ KB
```

Star Ratings Distribution

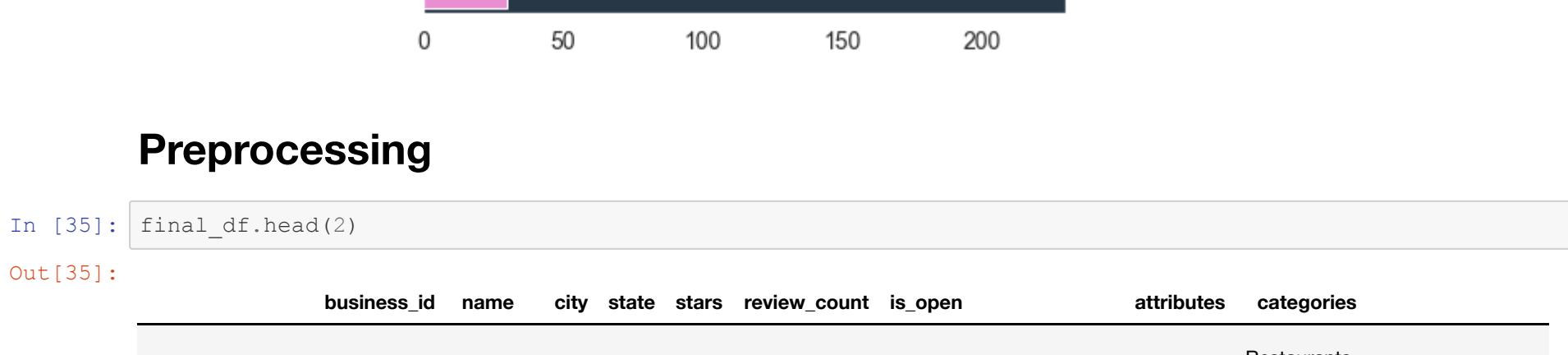
```
In [27]: x = business_data_for_edu['stars'].value_counts()
x = x.sort_index()
```

```
In [163]: plt.figure(figsize=(22,5))
sns.barplot(x=index,x.values)
plt.xlabel("Review Stars")
sns.set_palette("cubehelix")
plt.savefig("stars.png")
```



Distribution of Review Count

```
In [120]: plt.figure(figsize=(6,5))
sns.distplot(business_data_for_edu['review_count'], bins=10)
plt.show()
sns.set_palette("cubehelix")
plt.savefig("histo.png")
```



<Figure size 432x288 with 0 Axes>

Top 10 Restaurants

```
In [31]: rest = business_data_for_edu.sort_values('review_count', ascending=False)
rest = rest.head(10)
```

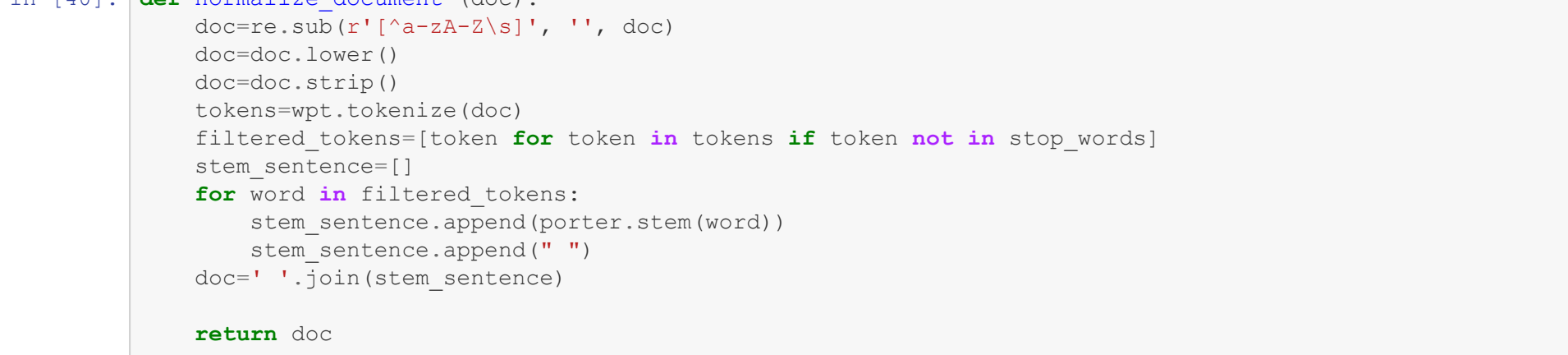
```
In [121]: plt.figure(figsize=(30,10))
sns.barplot(x=rest['name'], y=rest['review_count'])
plt.xlabel("Name", fontsize=15)
plt.savefig("top_rest.png")
sns.set_palette("cubehelix")
```



Top Categories

```
In [33]: business_data = ''.join(business_data_for_edu['categories'].astype('str'))
cats=pd.DataFrame(business_data.split(','),columns=['categories'])
x=cats.categories.value_counts()
x=x.sort_values(ascending=False)
x=x.sort(1)
```

```
In [122]: plt.figure(figsize=(7,5))
sns.barplot(x=index,x.values)
plt.xlabel("Category", fontsize=15)
plt.ylabel("Count", fontsize=15)
plt.savefig("categories.png")
sns.set_palette("cubehelix")
```



Preprocessing

```
In [35]: final_df.head(2)
```

```
Out [35]: business_id name city state stars review_count is_open attributes categories
0 A1dDB8S4Mbm0w7JnUjVA Veggie House Las Vegas NV 4.5 1142 1 ['Restaurants','Ice Cream Parlor', '2', 'BikeParking', ...] Restaurants, Specialty Food, Japanese, Sushi B...
```

Label Sentiment of Reviews

```
In [36]: final_df['sentiment'] = final_df['review_stars'].apply(lambda rating : '1' if rating > 3 else 0)
```

```
In [37]: plt.figure(figsize=(6,5))
sns.set(rc={'axes.facecolor':'#283747', 'axes.grid':False, 'xtick.labelsize':14, 'ytick.labelsize':14})
ax = sns.countplot(x=final_df['sentiment'])
plt.savefig("sentiment_counts.png")
```



```
In [38]: wpt=nltk.WordFuncTokenizier()
stopwords=nltk.corpus.stopwords.words('english')
```

```
In [39]: from nltk.stem import PorterStemmer
from nltk.stem import LancasterStemmer
porter=PorterStemmer()
```

Preprocessing Pipeline

```
In [40]: def normalize_document(doc):
    doc_tokens = [token.lower() for token in doc.split()]
    tokens=wpt.tokenize(doc)
    stem_sentence=[]
    for word in filtered_tokens:
        stem_sentence.append(porter.stem(word))
        stem_sentence.append(' ')
    doc=" ".join(stem_sentence)
    return doc

In [41]: text_column = final_df['text']
normalize_corpus=np.vectorize(normalize_document)
norm_corpus = normalize_corpus(text_column)
```

```
Out [42]: type(norm_corpus)

In [42]: norm_corpus

Out [42]: numpy.ndarray
```

```
In [43]: final_df['text_preprocessed'] = norm_corpus

In [57]: X = final_df['text_preprocessed']
y = final_df['sentiment']
```

```
In [45]: X.shape, y.shape

Out [45]: ((148224,)), (148224,))
```

```
In [46]: norm_corpus = X.tolist()
y = y.values
```

```
In [47]: type(norm_corpus), type(y)

Out [47]: (list, numpy.ndarray)
```

Saving Normalized Data

```
In [48]: with open('X.pickle','wb') as f:#wb, write-byte
    pickle.dump(X,f)
```

```
In [49]: with open('y.pickle','wb') as f:#wb, write-byte
    pickle.dump(y,f)
```

```
In [50]: with open('norm_corpus.pickle','wb') as f:
    pickle.dump(norm_corpus,f)
```

```
In [51]: #final_df = final_df[['text', 'sentiment']]
```

Wordclouds

```
In [52]: negative = final_df.loc[final_df['review_stars'] <= 2.0]
positive = final_df.loc[final_df['review_stars'] > 2.0]
```

Negative Sentiment

```
In [58]: stopwords = set(STOPWORDS)
negative_rev = " ".join(review for review in negative.text_preprocessed)
```

```
In [62]: wordcloud_neg = WordCloud(stopwords=STOPWORDS,
                                background_color = 'black',
                                width = 1200,
                                height = 1000,
                                collocation_threshold = 5
                                ).generate(negative_rev)
```

```
In [67]: plt.figure(figsize=(6,5))
plt.imshow(wordcloud_neg,interpolation="bilinear")
plt.axis('off')
plt.show()
plt.savefig('neg_wordcloud.png')
```

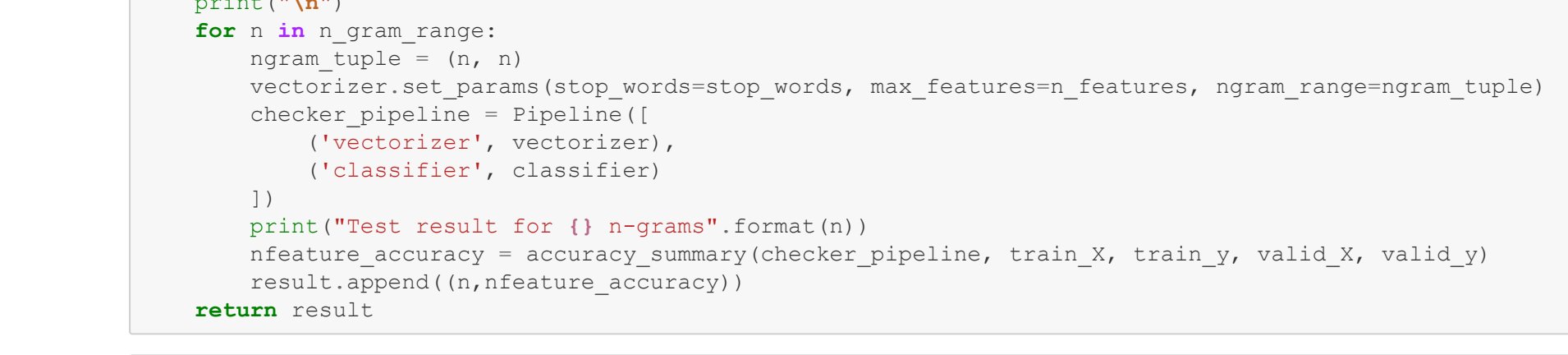


<Figure size 432x288 with 0 Axes>

Positive Sentiment

```
In [68]: positive_rev = " ".join(review for review in positive.text_preprocessed)
```

```
In [69]: wordcloud_pos = WordCloud(stopwords=STOPWORDS,
                                background_color = 'black',
                                width = 1200,
                                height = 1000,
                                collocation_threshold = 5
                                ).generate(positive_rev)
```



<Figure size 432x288 with 0 Axes>

Feature Selection

```
In [71]: with open('norm_corpus.pickle','rb') as f:
    norm_corpus = pickle.load(f)
```

```
In [72]: train_X, valid_X, train_y, valid_y = train_test_split(norm_corpus, y, random_state = 0, test_size=0.2)
```

Pipeline

```
In [73]: def accuracy_summary(pipeline, train_X, train_y, valid_X, valid_y):
    sentiment_fit = pipeline.fit(train_X, train_y)
    y_pred = sentiment_fit.predict(valid_X)
    accuracy = accuracy_score(valid_y, y_pred)
    print("Accuracy score: {0:.2f}%".format(accuracy*100))
    return accuracy
```

```
In [74]: cv = CountVectorizer()
lg = LogisticRegression(max_iter=500)
n_features = np.arange(5000,15000,2500)
```

```
In [75]: def nfeature_accuracy_checker(vectorizer=cv, n_features=n_features, stop_words=None, ngram_range=(1, 3), classifier=lg):
    result = []
    print(classifier)
    print("n")
    for i in n_features:
        vectorizer.set_params(stop_words=stop_words, max_features=n_features, ngram_range=ngram_range)
        checker_pipeline = Pipeline([
            ('vectorizer', vectorizer),
            ('classifier', classifier),
        ])
        print("Test result for {} n-grams".format(i))
        nfeature_accuracy = accuracy_summary(checker_pipeline, train_X, train_y, valid_X, valid_y)
        result.append((i,nfeature_accuracy))
    return result
```

```
In [76]: tfidf = TfidfVectorizer()
print("Result for trigram with stop words (Tfidf)")
feature_result_tfidf = nfeature_accuracy_checker(vectorizer=tfidf, ngram_range=(1, 3))
Result for trigram with stop words (Tfidf)
LogisticRegression(max_iter=500)
```

Test result for 5000 features
accuracy score: 92.18%
Test result for 10000 features
accuracy score: 92.23%
Test result for 12500 features
accuracy score: 92.33%
Test result for 15000 features
accuracy score: 92.30%

N-gram Selection

Pipeline

```
In [77]: train_X, valid_X, train_y, valid_y = train_test_split(norm_corpus, y, random_state = 0, test_size=0.2)
```

```
In [78]: def accuracy_summary(pipeline, train_X, train_y, valid_X, valid_y):
    sentiment_fit = pipeline.fit(train_X, train_y)
    y_pred = sentiment_fit.predict(valid_X)
    accuracy = accuracy_score(valid_y, y_pred)
    print("Accuracy score: {0:.2f}%".format(accuracy*100))
    return accuracy
```

```
In [79]: cv = CountVectorizer()
lg = LogisticRegression(max_iter=500)
n_features = np.arange(1,4,1)
```

```
In [80]: def nfeature_accuracy_checker(vectorizer=cv, n_features=n_features, stop_words=None, ngram_range=ngram_range, classifier=lg):
    result = []
    print(classifier)
    print("n")
    for i in ngram_range:
        vectorizer.set_params(stop_words=stop_words, max_features=n_features, ngram_range=ngram_range)
        checker_pipeline = Pipeline([
            ('vectorizer', vectorizer),
            ('classifier', classifier),
        ])
        print("Test result for {} n-grams".format(i))
        nfeature_accuracy = accuracy_summary(checker_pipeline, train_X, train_y, valid_X, valid_y)
        result.append((i,nfeature_accuracy))
    return result
```

```
In [81]: tfidf = TfidfVectorizer()
print("Result for N-grams with stop words (Tfidf)\n")
feature_result_tfidf = nfeature_accuracy_checker(vectorizer=tfidf)
Result for N-grams with stop words (Tfidf)
LogisticRegression(max_iter=500)
```

Test result for 1 n-grams
accuracy score: 91.41%
Test result for 2 n-grams
accuracy score: 89.66%
Test result for 3 n-grams
accuracy score: 82.76%

LSA and Topic Modeling

```
In [123]: vectorizer = TfidfVectorizer(ngram_range=(2,2),max_features=12500)
X = vectorizer.fit_transform(norm_corpus)
```

```
In [83]: print(X[0])

(0, 8843) 0.1881276654629726
(0, 2393) 0.16630467617742986
(0, 9493) 0.1672654300662116
(0, 4563) 0.18564291899443825
(0, 8339) 0.1459470188553208
(0, 3493) 0.1694436325039138
(0, 4649) 0.16776055186072664
(0, 7943) 0.18978249203128297
(0, 10113) 0.181550435017228
(0, 3933) 0.193840302449804
(0, 10060) 0.17729801259287284
(0, 4417) 0.1371868474298085
(0, 2387) 0.1392816978925348
(0, 7368) 0.1761204154307408
(0, 9690) 0.18399520648254186
(0, 285) 0.14798044432496385
(0, 10767) 0.20159820548591126
(0, 7074) 0.1683299716810083
(0, 6489) 0.1851204534279217
(0, 12249) 0.18577465105045105
(0, 6376) 0.1893937123435476
(0, 5184) 0.2063570448339135
(0, 2324) 0.188686286444025
(0, 11648) 0.1881276654629726
(0, 11347) 0.1829169945236995
(0, 5197) 0.1376313527006373
(0, 11763) 0.1757634082139892
(0, 6044) 0.11412466713140795
(0, 1311) 0.1593970235666776
(0, 10802) 0.19206732758593884
(0, 2837) 0.18538198732048225
(0, 7751) 0.16696107137527108
(0, 6918) 0.15958318741877893
```

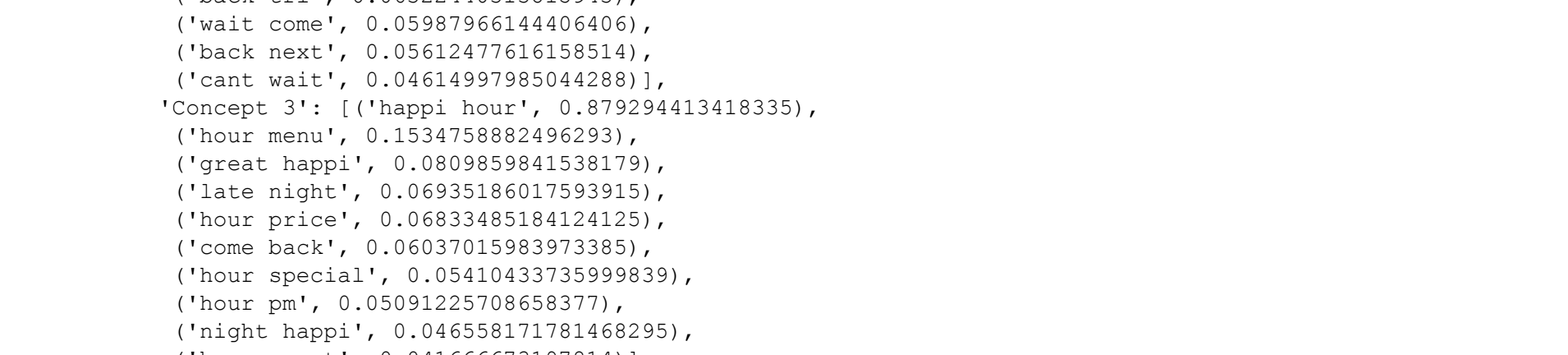
```
In [152]: lsa = TruncatedSVD(n_components=5, n_iter=100)#n_component number of concepts we want to find
lsa.fit(X)
```

```
Out [152]: TruncatedSVD(n_components=5, n_iter=100)
```

```
In [153]: Sigma = lsa.singular_values
Sigma
```

```
Out [153]: array([25.40447251, 17.73228187, 17.35564355, 16.83238397, 15.75597973])
```

```
In [154]: ax = sns.barplot(x=list(range(len(Sigma))), y = Sigma)
ax.set_xlabel('Components', y_label='Relative Importance')
plt.savefig('lsa.png')
```



```
In [155]: row1 = lsa.components_[0] #the row 1 row1

Out [155]: array([0.00245209, 0.00384008, 0.00201168, ..., 0.0036876 , 0.00173308,
0.00252203])
```

```
In [156]: terms = vectorizer.get_feature_names()
terms[:10]
```

```
Out [156]: ['abl accomod',
'abl eat',
'abl enjoy',
'abl find',
'abl finish',
'abl get',
'abl make',
'abl order',
'abl seat',
'abl see',
'abl sit',
'abl try',
'absolut amaz',
'absolut best',
'absolut delici',
'absolut fantast',
'absolut favorit',
'absolut great',
'absolut horribl',
'absolut incred']
```

```
In [157]: concept_words = []

for i,comp in enumerate(lsa.components_):#enumerate return index and row, a list of tuples
    component_terms = zip(terms,comp) #we use zip to combine values and terms
    #sort component terms, by concept value,lambda x (x corresponds tuples) and X[i] the value
    sortedTerms = sortedTerms, key=lambda xx:[1], reverse=True) #ascending order
    sortedTerms = sortedTerms[:10] #select 10 most imp. terms in a specific concept
    concept_words["concept " + str(i)] = sortedTerms #all concepts mapped with list of tuples
```

```
In [158]: concept_words
```

```
Out [158]: {'concept 0': ['(come back', 0.2610881420989711),
'sushi place', 0.174676828390287],
'concept 1': ['great service', 0.1585246599128933],
'concept 2': ['great time', 0.1631607480158179],
'concept 3': ['happy hour', 0.1683718919550834],
'concept 4': ['sushi bar', 0.1643542457895137],
'concept 5': ['great food', 0.1643542457895137],
'concept 6': ['great service', 0.1446263466971937],
'concept 7': ['great time', 0.1385246599128933],
'concept 8': ['great time', 0.1385246599128933],
'concept 9': ['great time', 0.1385246599128933],
'concept 10': ['great time', 0.1385246599128933],
'concept 11': ['great time', 0.1385246599128933],
'concept 12': ['great time', 0.1385246599128933],
'concept 13': ['great time', 0.1385246599128933],
'concept 14': ['great time', 0.1385246599128933],
'concept 15': ['great time', 0.1385246599128933],
'concept 16': ['great time', 0.1385246599128933],
'concept 17': ['great time', 0.1385246599128933],
'concept 18': ['great time', 0.1385246599128933],
'concept 19': ['great time', 0.1385246599128933],
'concept 20': ['great time', 0.1385246599128933],
'concept 21': ['great time', 0.1385246599128933],
'concept 22': ['great time', 0.1385246599128933],
'concept 23': ['great time', 0.1385246599128933],
'concept 24': ['great time', 0.1385246599128933],
'concept 25': ['great time', 0.1385246599128933],
'concept 26': ['great time', 0.1385246599128933],
'concept 27': ['great time', 0.1385246599128933],
'concept 28': ['great time', 0.1385246599128933],
'concept 29': ['great time', 0.1385246599128933],
'concept 30': ['great time', 0.1385246599128933],
'concept 31': ['great time', 0.1385246599128933],
'concept 32': ['great time', 0.1385246599128933],
'concept 33': ['great time', 0.1385246599128933],
'concept 34': ['great time', 0.1385246599128933],
'concept 35': ['great time', 0.1385246599128933],
'concept 36': ['great time', 0.1385246599128933],
'concept 37': ['great time', 0.1385246599128933],
'concept 38': ['great time', 0.1385246599128933],
'concept 39': ['great time', 0.1385246599128933],
'concept 40': ['great time', 0.1385246599128933],
'concept 41': ['great time', 0.1385246599128933],
'concept 42': ['great time', 0.1385246599128933],
'concept 43': ['great time', 0.1385246599128933],
'concept 44': ['great time', 0.1385246599128933],
'concept 45': ['great time', 0.1385246599128933],
'concept 46': ['great time', 0.1385246599128933],
'concept 47': ['great time', 0.1385246599128933],
'concept 48': ['great time', 0.1385246599128933],
'concept 49': ['great time', 0.1385246599128933],
'concept 50': ['great time', 0.1385246599128933],
'concept 51': ['great time', 0.1385246599128933],
'concept 52': ['great time', 0.1385246599128933],
'concept 53': ['great time', 0.1385246599128933],
'concept 54': ['great time', 0.1385246599128933],
'concept 55': ['great time', 0.1385246599128933],
'concept 56': ['great time', 0.1385246599128933],
'concept 57': ['great time', 0.1385246599128933],
'concept 58': ['great time', 0.1385246599128933],
'concept 59': ['great time', 0.1385246599128933],
'concept 60': ['great time', 0.1385246599128933],
'concept 61': ['great time', 0.1385246599128933],
'concept 62': ['great time', 0.1385246599128933],
'concept 63': ['great time', 0.1385246599128933],
'concept 64': ['great time', 0.1385246599128933],
'concept 65': ['great time', 0.1385246599128933],
'concept 66': ['great time', 0.1385246599128933],
'concept 67': ['great time', 0.1385246599128933],
'concept 68': ['great time', 0.1385246599128933],
'concept 69': ['great time', 0.1385246599128933],
'concept 70': ['great time', 0.1385246599128933],
'concept 71': ['great time', 0.1385246599128933],
'concept 72': ['great time', 0.1385246599128933],
'concept 73': ['great time', 0.1385246599128933],
'concept 74': ['great time', 0.1385246599128933],
'concept 75': ['great time', 0.1385246599128933],
'concept 76': ['great time', 0.1385246599128933],
'concept 77': ['great time', 0.1385246599128933],
'concept 78': ['great time', 0.1385246599128933],
'concept 79': ['great time', 0.1385246599128933],
'concept 80': ['great time', 0.1385246599128933],
'concept 81': ['great time', 0.1385246599128933],
'concept 82': ['great time', 0.1385246599128933],
'concept 83': ['great time', 0.1385246599128933],
'concept 84': ['great time', 0.1385246599128933],
'concept 85': ['great time', 0.1385246599128933],
'concept 8
```


[98]: X=transformer.fit_transform(X).toarray()

Vocabulary

In [99]: vocab=vectorizer.get_feature_names()
pd.DataFrame(X.columns=vocab).head()

Out[99]:

	abl	abl	abl	abl	abl	abl	abl	abl	...	yummi	yummi	yummi	yummi	yuzu	yuzu	yuzu	yuzu
	accommod	eat	enjoy	find	finish	get	make	order	seat	see	...	grill	love	servic	sushi	sauce	sorbet
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0

1 rows x 12500 columns

Splitting Data

In [100]: train_X, valid_X, train_y, valid_y = train_test_split(X, y, test_size=0.2, random_state=1)

Fitting to Model

In [101]: classifier = LogisticRegression(max_iter=100)
classifier.fit(train_X,train_y)

/opt/anaconda3/lib/python3.8/site-packages/sklearn/linear_model/_logistic.py:762: ConvergenceWarning:
lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html>
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_1 = _check_optimize_result(
LogisticRegression()

Out[101]: LogisticRegression()

In [102]: y_pred = classifier.predict(valid_X)

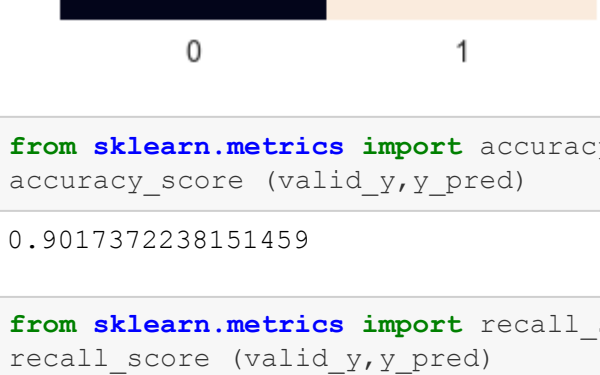
Model Evaluation

In [103]: cm = confusion_matrix(valid_y,y_pred)
cm

Out[103]: array([[5190, 2151],
[762, 21542]])

In [104]: sns.heatmap(cm, annot=True, fmt='g')

Out[104]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff44db1ceb0>



In [105]: from sklearn.metrics import accuracy_score
accuracy_score (valid_y,y_pred)

Out[105]: 0.9017372238151459

In [106]: from sklearn.metrics import recall_score
recall_score (valid_y,y_pred)

Out[106]: 0.9658357245337159

In [107]: from sklearn.metrics import f1_score
f1_score (valid_y,y_pred)

Out[107]: 0.9366697828119226

In [108]: from sklearn.metrics import precision_score
precision_score (valid_y,y_pred)

Out[108]: 0.909213691807707

In [109]: print(classification_report(valid_y, y_pred))

	precision	recall	f1-score	support
0	0.87	0.71	0.78	7341
1	0.91	0.97	0.94	22304
accuracy			0.90	29645
macro avg	0.89	0.84	0.86	29645
weighted avg	0.90	0.90	0.90	29645

ROC Curve

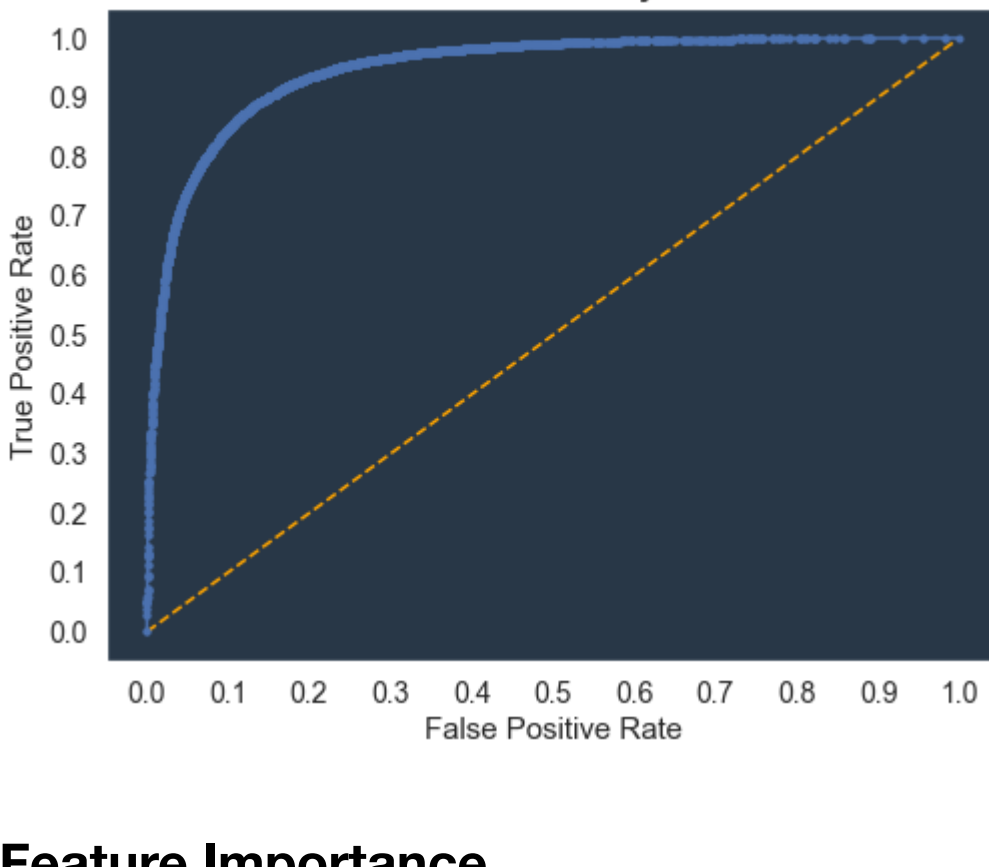
In [111]: lr_probs = classifier.predict_proba(valid_X)
lr_probs = lr_probs[:, 1]

In [112]: lr_auc = roc_auc_score(valid_y, lr_probs)

In [113]: print('Logistic: ROC AUC=%3f' % (lr_auc))
Logistic: ROC AUC=0.949

In [114]: lr_fpr, lr_tpr, _ = roc_curve(valid_y, lr_probs)

In [115]: fig = plt.figure(figsize=(8,6))
plt.plot([0,1], [0,1], color='orange', linestyle='--')
plt.plot(lr_fpr, lr_tpr, marker='.', label='Logistic')
plt.xticks(np.arange(0.0, 1.1, step=0.1))
plt.xlabel('False Positive Rate', fontsize=15)
plt.yticks(np.arange(0.0, 1.1, step=0.1))
plt.ylabel('True Positive Rate', fontsize=15)
plt.title('ROC Curve Analysis', fontweight='bold', fontsize=15)
plt.savefig('roc.png')

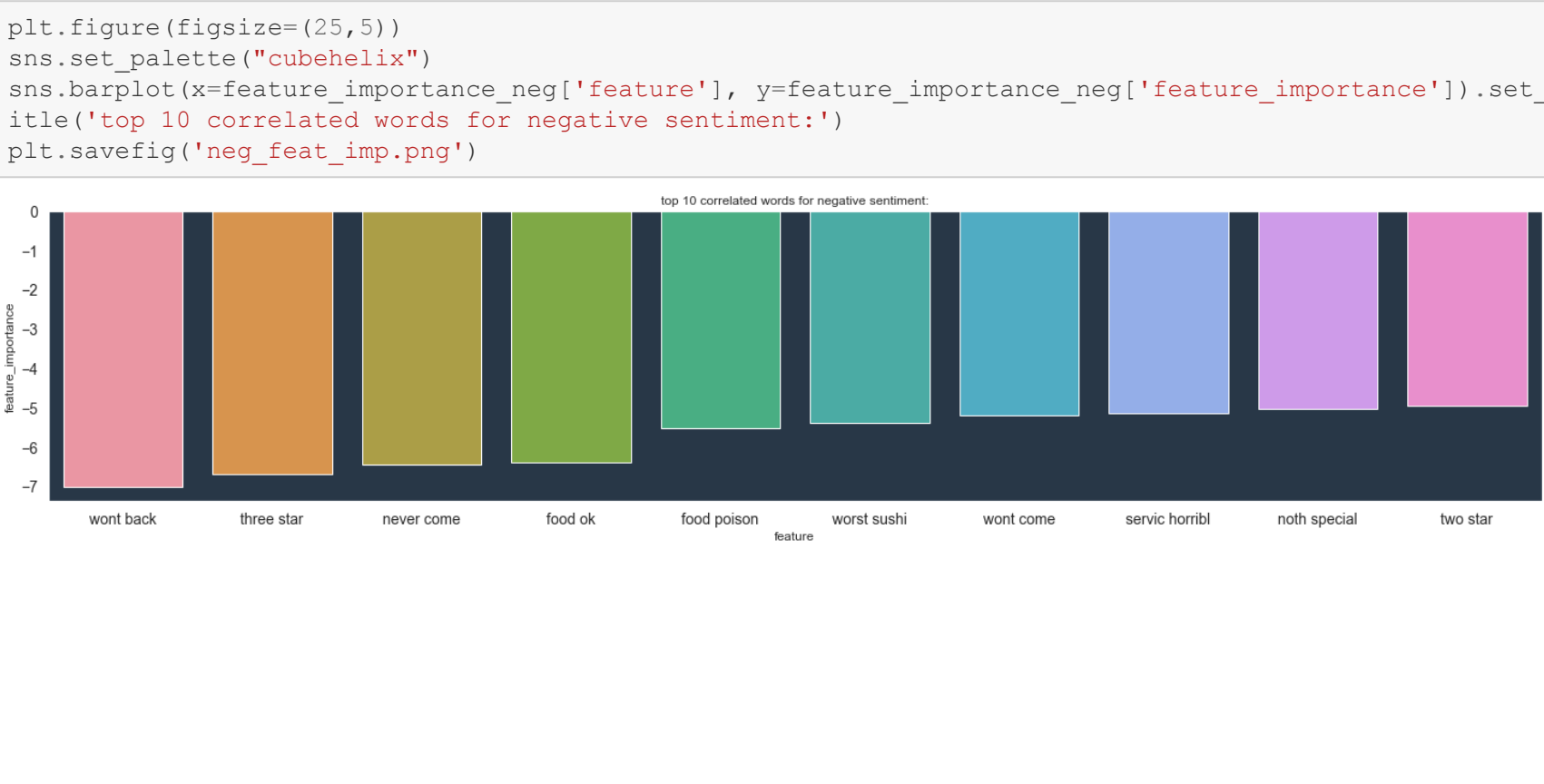


Feature Importance

In [116]: classifier.coef_
Out[116]: array([[0.39481493, -0.10126172, -0.04663601, ..., 1.00785253,
0.15980939, -2.5588048]])

In [117]: feature_importance=pd.DataFrame({'feature':vectorizer.get_feature_names(),'feature_importance':classif
ier.coef_[0]})
feature_importance_pos=feature_importance.sort_values('feature_importance',ascending=False).head(10)
feature_importance_neg = feature_importance.sort_values('feature_importance',ascending=True).head(10)

In [118]: plt.figure(figsize=(25,5))
sns.set_palette("cubehelix")
sns.barplot(x=feature_importance_pos['feature'], y=feature_importance_pos['feature_importance']).set_t
itle('top 10 correlated words for positive sentiment')
plt.savefig('pos_feat_imp.png')



In [119]: plt.figure(figsize=(25,5))
sns.set_palette("cubehelix")
sns.barplot(x=feature_importance_neg['feature'], y=feature_importance_neg['feature_importance']).set_t
itle('top 10 correlated words for negative sentiment')
plt.savefig('neg_feat_imp.png')

