

Lo Kok Fu

DSIF 2



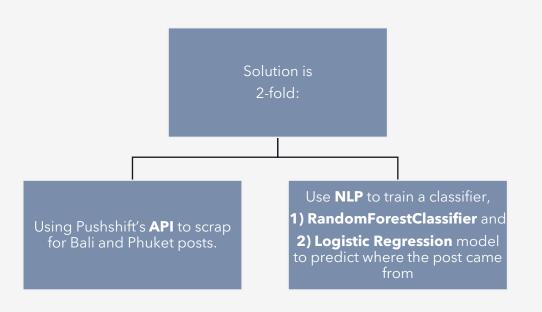
Agenda

- Problem Statement
- Introduction
- Obtaining Data
- Data Exploration
- Data Processing / Engineering
- Modelling
- Results
- Recommendations



Problem Statement

Utilize NLP to train a classifier on a binary classification problem to predict which subreddit post it comes from. Gather our data through Gather from subreddit using API.





Introduction

- Topic classifier: Tourist beaches.
- Bali & Phuket
- Similar point of Interest
- Familiarity
- COVID-19. Let's see what we scrape





```
1 # create function for web scrapping submissions:
 3 def scrapping(url, loops, subreddit):
       start time = time.time() # get the time in seconds since epoch
        params = {'subreddit': subreddit,
        'size': 100,
        'before': round(start time)
 9
10
        for i in range(loops):
12
           current time = time.time()
13
            #requesting data
14
           res = requests.get(url, params)
15
           print(f'res status {i+1}: ', res.status code)
16
17
            data = res.json()
18
            posts = data['data']
19
            post df = pd.DataFrame(posts)
20
            df.append(post df)
21
            #get oldest post time and use as before parameter in next request
22
            old = post df['created utc'].min()
23
            params['before'] = old
24
           time.sleep(1)
25
            reddit posts = pd.concat(df)
26
27
            filename = subreddit + ' submission.csv'
28
        return reddit posts.to csv('./datasets/' + filename, index=False)
```

```
# scrape for bali submission data set from Reddit.
url = 'https://api.pushshift.io/reddit/submission/search'
loops = 8 # no. of Loops to scrap
subreddit = 'bali' # subreddit topic

scrapping(url, loops, subreddit)

res status 1: 200
res status 2: 200
res status 3: 200
res status 4: 200
res status 4: 200
res status 5: 200
res status 6: 200
res status 7: 200
res status 7: 200
res status 8: 200

1 bali_submission = pd.read_csv('./datasets/bali_submission.csv')
```

Obtaining Data

- Subreddit limited scrapping no.
- Created a function to loop scrapping
- Looping from oldest post scrapped
- Concat
- File save
- Data Cleaning
 - Stop words
 - Punctuations
 - Lower case
 - Removed duplicates
 - Removed int
 - Filled in null values
 - Merged 'self_text' & 'title
 - Concat comments to df
- 1600 for each topic, Total 3200
- Balance: 3075
- Data imbalance
 - 1 0.505691
 - 0 0.494309





Data Exploration



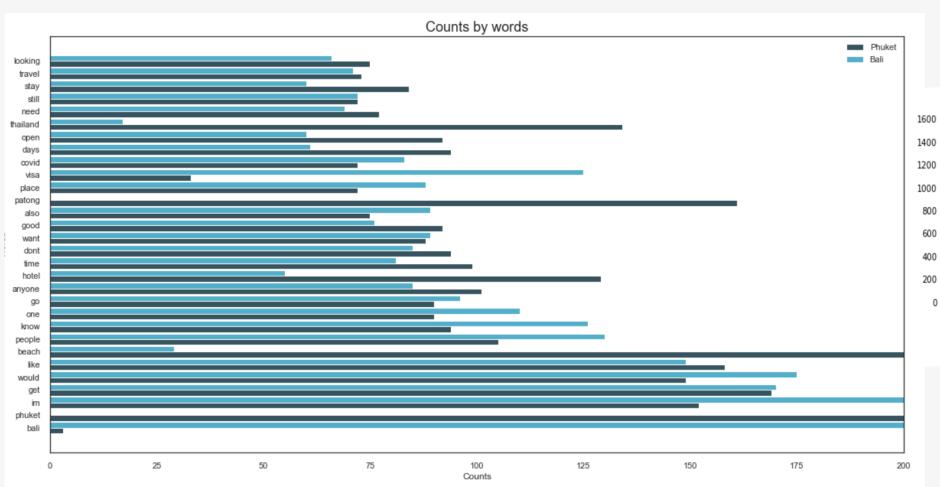
- 'Patong' is a name of a place/beach in phuket
- 'hotel', Recommendations?
- 'Bangkok' the capital of thailand.
- 'sandbox' linked to Covid-19
- 'open'. Opening of Borders?

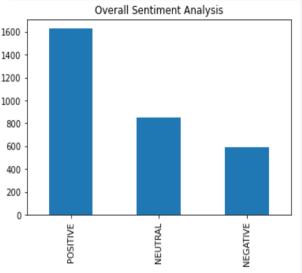
- 'Quarantine'& 'Jakarta'. This may be due to essential visitors are to serve quarantine in Jakarta before visiting bali.
- 'Villa'. Holiday villas are spruced all over the island. thus the frequent mentioned of it.
- 'Visa'. Foreign Visitors to Bali will require visas.





Data Exploration







Data Processing / Engineering

- Word Count
- Sentiments values (Vader)
- Function transformer
 - Standard Scaler
 - Count vectorizer
- Pipeline to combine features and estimator

```
أعس
                    ine to combine features and estimator
                _eline([
             _ures', FeatureUnion([('numeric_features', Pipeline([('selector', get_numeric_a.
                ('text_features', Pipeline([('selector', get_text_data),('cvec', CountVectorize.
        ('rf', RandomForestClassifier())
   params = { 'rf_n_estimators': [80,90,110,120],
              'rf__max_depth':[5,10,15,20],
                'features text_features cvec stop_words' : [None, stop]}
10
11
12
   gs_rf = GridSearchCV(pipe_rf, params, cv=5, verbose = 1, n_jobs = -1)
14
   gs_rf.fit(X_train, y_train)
```





Modelling

- Random Forest Classifier
 - Gridsearch CV with Hyperparameter Tuning:
 - n estimators = 110
 - max_depth = 20
- Logistic Regression
 - Gridsearch CV with Hyperparameter Tuning:
 - Penalty = I2
 - Solver = 'sag'
 - C value = 10

Modelling 1: RandomForestClassifier

Modelling 2: Logistic Regression

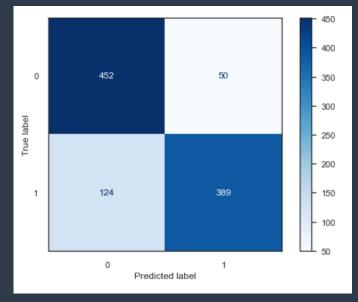
```
pipe = Pipeline([
        ('features', FeatureUnion([
                ('numeric_features', Pipeline([
                    ('selector', get numeric data),
                    ('ss', StandardScaler())
                ('text_features', Pipeline([
                    ('selector', get_text_data),
                    ('cvec', CountVectorizer())
               ]))
12
        ('logreg', LogisticRegression())
13 ])
14
15 params = {
                'features text features cvec stop words' : [None, stop],
17
                'logreg_penalty' : ['l1', 'l2'],
18
                'logreg_solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'],
19
                'logreg_C' : [100, 10, 1.0, 0.1, 0.01]
20
21
    gs = GridSearchCV(pipe, params, cv=5, verbose = 1, n_jobs = -1)
   gs.fit(X_train, y_train)
```

Fitting 5 folds for each of 100 candidates, totalling 500 fits

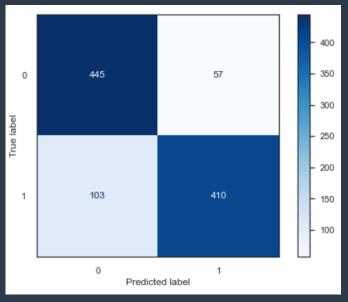


Results

Model	Train Score	Test Score	Specificity	F1 Score	AUC
Random Forest Classifier	0.92718	0.82856	0.90039	0.81722	0.94435
Logistic Regression	0.94611	0.84236	0.8864	0.83673	0.93579

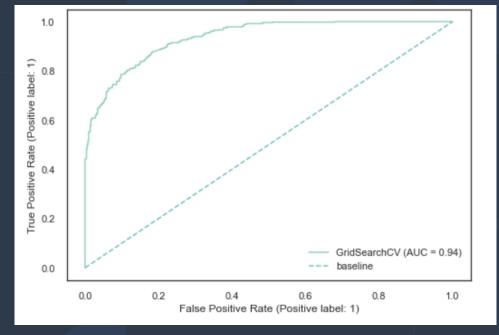


Random Forest Classifier Confusion Matrix

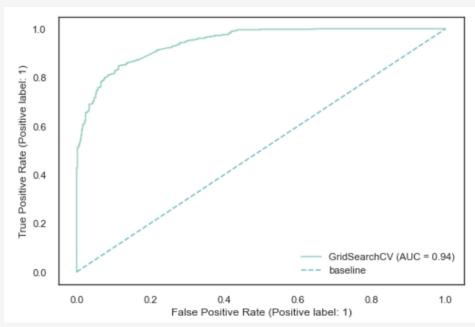


Logistic Regression Confusion Matrix





Logistic Regression ROC Curve



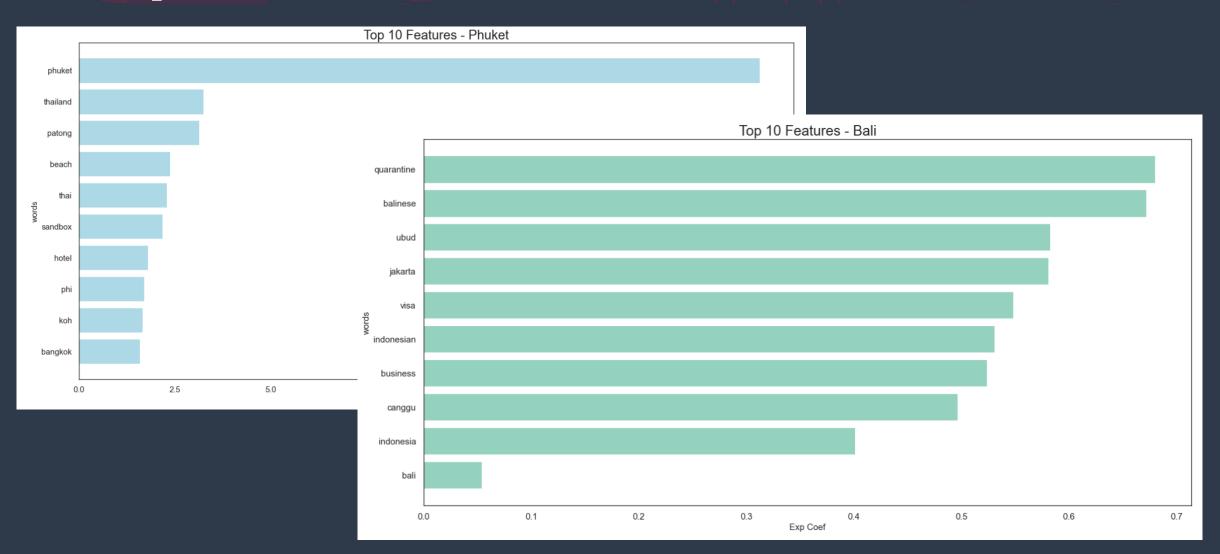
Random Forest ROC Curve

Results ROC Curve

Model	Train Score	Test Score	Specificity	F1 Score	AUC
Random Forest Classifier	0.92718	0.82856	0.90039	0.81722	0.94435
Logistic Regression	0.94611	0.84236	0.8864	0.83673	0.93579

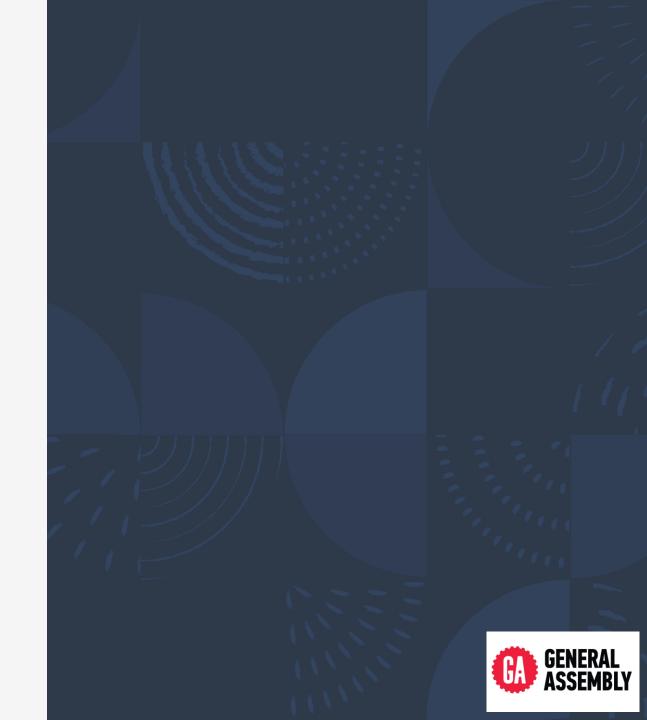


Top 10 Features



Recommendations

- Logistic Regression
 - Best F1 scoring
 - Accuracy Score
 - Slight difference in scores
- More time, Faster CPU
 - Explore hyperparameter Tuning
 - Feature engineering
 - TFIDF
 - Other models



Thank you

