

Analysis of Restaurant Reviews

```
1  # Importing Libraries
2  import numpy as np
3  import pandas as pd
4
5  # Import dataset
6  dataset = pd.read_csv('/content/Restaurant_Reviews.tsv', delimiter = '\t')
```

```
1  # library to clean data
2  import re
3
4  # Natural Language Tool Kit
5  import nltk
6
7  nltk.download('stopwords')
8
9  # to remove stopword
10 from nltk.corpus import stopwords
11
12 # for Stemming propose
13 from nltk.stem.porter import PorterStemmer
14
15 # Initialize empty array
16 # to append clean text
17 corpus = []
18
19 # 1000 (reviews) rows to clean
20 for i in range(0, 1000):
21
22     # column : "Review", row ith
23     review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][i])
24
25     # convert all cases to lower cases
26     review = review.lower()
27
28     # split to array(default delimiter is " ")
29     review = review.split()
30
31     # creating PorterStemmer object to
32     # take main stem of each word
33     ps = PorterStemmer()
34
35     # loop for stemming each word
36     # in string array at ith row
37     review = [ps.stem(word) for word in review
38               if not word in set(stopwords.words('english'))]
39
40     # rejoin all string array elements
41     # to create back into a string
42     review = ' '.join(review)
43
44     # append each string to create
45     # array of clean text
46     corpus.append(review)
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
```

```
1  # Creating the Bag of Words model
2  from sklearn.feature_extraction.text import CountVectorizer
3
4  # To extract max 1500 feature.
5  # "max_features" is attribute to
6  # experiment with to get better results
7  cv = CountVectorizer(max_features = 1500)
8
9  # X contains corpus (dependent variable)
10 X = cv.fit_transform(corpus).toarray()
11
12 # y contains answers if review
13 # is positive or negative
14 y = dataset.iloc[:, 1].values
```

```
1  # Splitting the dataset into
2  # the Training set and Test set
3  from sklearn.model_selection import train_test_split
4
5  # experiment with "test_size"
6  # to get better results
7  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25)
```

```
1  # Fitting Random Forest Classification
2  # to the Training set
3  from sklearn.ensemble import RandomForestClassifier
4
5  # n_estimators can be said as number of
6  # trees, experiment with n_estimators
7  # to get better results
8  model = RandomForestClassifier(n_estimators = 501,
9                                criterion = 'entropy')
10
11 model.fit(X_train, y_train)
```

```
RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
```

```
criterion='entropy', max_depth=None, max_features='auto',
max_leaf_nodes=None, max_samples=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=501,
n_jobs=None, oob_score=False, random_state=None,
verbose=0, warm_start=False)
```

```
1 # Predicting the Test set results
2 y_pred = model.predict(X_test)
3
4 y_pred

array([0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 1, 0,
       1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1,
       0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 0,
       0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0, 0,
       1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0,
       0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0,
       0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0,
       0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 0,
       0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1,
       0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0,
       0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0,
       0, 0, 0, 0, 1, 0, 0, 0])

1 # Making the Confusion Matrix
2 from sklearn.metrics import confusion_matrix, classification_report
3
4 cm = confusion_matrix(y_test, y_pred)
5
6 print("Confusion Matrix:\n",cm)
```

```
Confusion Matrix:
[[122  15]
 [ 42  71]]
```

```
1 report=classification_report(y_test, y_pred)
2 print("Classification Report:\n",report)

Classification Report:
              precision    recall  f1-score   support

         0       0.74      0.89      0.81        137
         1       0.83      0.63      0.71        113

 accuracy          0.77          250
 macro avg       0.78      0.76      0.76          250
weighted avg       0.78      0.77      0.77          250
```

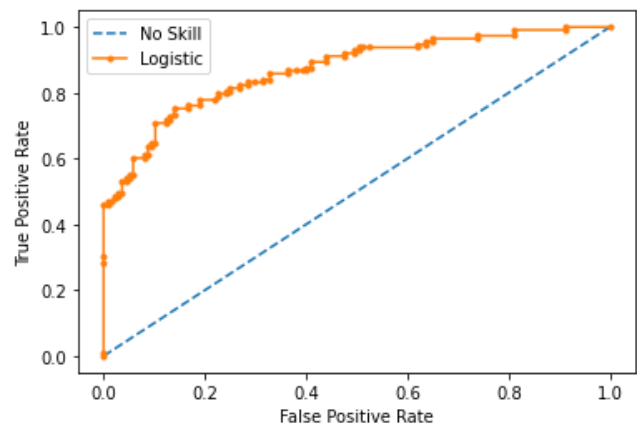
```
1 from sklearn.linear_model import LogisticRegression
2 # fit a model
3 lr = LogisticRegression(penalty='l2',C=0.5)
4 lr.fit(X_train, y_train)
5 y_pred = lr.predict(X_test)
```

Running the example prints the ROC AUC for the logistic regression model and the no skill classifier that only predicts 0 for all examples.

```
1 from sklearn.metrics import roc_auc_score
2 # predict probabilities
3 lr_probs = lr.predict_proba(X_test)
4 # keep probabilities for the positive outcome only
5 lr_probs = lr_probs[:, 1]
6 # calculate scores
7 ns_auc = roc_auc_score(y_test, ns_probs)
8 lr_auc = roc_auc_score(y_test, lr_probs)
9 # summarize scores
10 print('No Skill: ROC AUC=%.3f' % (ns_auc))
11 print('Logistic: ROC AUC=%.3f' % (lr_auc))

No Skill: ROC AUC=0.500
Logistic: ROC AUC=0.871
```

```
1 # calculate roc curves
2 ns_fpr, ns_tpr, _ = roc_curve(y_test, ns_probs)
3 lr_fpr, lr_tpr, _ = roc_curve(y_test, lr_probs)
4 # plot the roc curve for the model
5 plt.plot(ns_fpr, ns_tpr, linestyle='--', label='No Skill')
6 plt.plot(lr_fpr, lr_tpr, marker='.', label='Logistic')
7 # axis labels
8 plt.xlabel('False Positive Rate')
9 plt.ylabel('True Positive Rate')
10 # show the legend
11 plt.legend()
12 # show the plot
13 plt.show()
```



✓ 0s completed at 7:36 PM

