Cameron Lee ML Milestone 3

Dataset: Drug Review

Problem Statement

The problem statement that I am proposing is classifying the effectiveness of a drug. Accurately classifying the effectiveness of drugs is crucial for medical professionals, researchers, and patients to make informed decisions regarding treatment options. By classifying the effectiveness of these drugs, we can provide valuable insights into their efficacy, enabling healthcare providers to prescribe the most suitable medications, researchers to identify trends in drug performance, and patients to make well-informed decisions about their healthcare.

Solution

The problem statement that I am proposing is classifying the <u>effectiveness</u> of a drug. Therefore, effectiveness will be used as the target. In particular, I created a binary classification with the target focused on whether the drug is considered "effective" or "ineffective". These are the given features I will be focusing on to drive the study

urlDrugName	effectiveness	sideEffects	condition	rating	benefits Review	sideEffects Review	comments Review

In addition to those features, I implemented text preprocessing using the NLTK library for the reviews and performed sentiment analysis using the TextBlob library. These sentiment analysis fields show negative values for negative sentiment analysis and positive values for positive sentiment analysis. These columns were added to the test data with the following names:

sideEffectsSentiment	commentsSentiment	benefitsSentiment	avgSentiment
----------------------	-------------------	-------------------	--------------

I took the average score of sideEffectsSentiment, commentsSentiment, and benefitsSentiment and rounded the values to create a column called "avgSentiment" which corresponds to the average sentiment from all three types of reviews present in each datapoint.

I made the design decision for an "effective" classification, the data would have to meet the target criteria of the following:

If this criteria is not met, the classification of that datapoint would be considered "ineffective"

After a long process of data preprocessing and feature engineering that included text preprocessing, sentiment analysis implementation, handling missing values, encoding categorical columns, dropping unneeded columns, and splitting the datasets into their respective sets (assigning target values, test and train datasets were given), I was ready to begin classifiers to help me solve the problem.

In order to address the classification problem I used the three following classifiers:

Random Forest

I started by developing a Random Forest model as I have found that it works well by combining the power of multiple decision trees, leading to robust and accurate decisions. I began by using Randomized Search Cross validation with a broad set of hyperparameters in order to narrow down the range of which hyperparameters were the most relevant. This yielded the following:

```
Fitting 3 folds for each of 100 candidates, totalling 300 fits {'n_estimators': 400, 'min_samples_split': 10, 'min_samples_leaf': 1, 'max_features': 'sqrt', 'max_depth': 20, 'bootstrap': True} Random Forest Classifier Metrics:
Testing Accuracy: 0.6003861003861004
```

From there, I narrowed down the values and switched to GridSearch cross validation.

```
param_grid = {
    'bootstrap': [True],
    'max_depth': [30, 40, 50, 100],
    'max_features': ['sqrt'],
    'min_samples_leaf': [1, 2, 3],
    'min_samples_split': [8, 10, 12],
    'n_estimators': [350, 375, 400, 425, 450]
}
```

Then I ran the entire algorithm on the updated y_pred and tested it against the Y_test (target value I defined). This yielded the following:

At this point in time, I was testing the effectiveness of sentiment analysis. By including the average sentiment analysis on top of the values for each type of review already present, I saw a 5 percent increase in accuracy. This was the final report for Random Forest.

```
Random Forest Classifier Metrics:
Accuracy: 0.8455598455598455
Classification Report:
                      precision
                                           recall
                                                       f1-score
                                                                         support
                                                                                318
718
                  0
1
                              0.73
                                              0.80
                                                              0.76
                                              0.87
                                                              0.89
                              0.91
                                                                               1036
      accuracy
 macro avģ
veighted avg
                              0.82
0.85
                                                              0.82
                                                                               1036
                                              0.85
                                                              0.85
                                                                               1036
```

Gaussian Naive Bayes

I chose to use Naive Bayes due to the ability to create a strong distinction between classifications, specifically an "effective" review and an "ineffective". As the dataset has a lot of text data, it made sense for me to include those metrics such as sentiment analysis and look for positive and negative results.

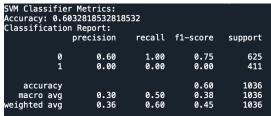
With this model, I had prepared the train and test datasets with the set targets. I then fit the data, predicted the values and this gave me the following:

```
GNB Classifier Metrics:
Accuracy: 0.7268339768339769
Classification Report:
                       precision
                                             recall
                                                       f1-score
                                                                           support
                               0.86
0.61
                                               0.66
0.83
                                                                0.74
0.71
                                                                                  625
411
                                                                0.73
0.73
0.73
                                                                                1036
1036
1036
       accuracy
                               0.74
0.76
                                               0.75
0.73
 macro avg
weighted avg
```

Through my research, I found there was not much tuning I could do with this model due to its simplistic nature. So I applied 'smoothing' to account for any zero probabilities and also included the average sentiment since I found success with it in Random Forest and saw about a 10 percent increase in accuracy. Here is the final report for Guassian Naive Bayes:

Support Vector Machine

Finally for this model, I had prepared the train and test datasets with the set targets. I fit the data, predicted the values and this gave me the following:



This was the lowest accuracy I had produced and since this was the last model I created, I felt that I had exhausted the possibility to improve the dataset in the preprocessing stage. I tuned the C and Gamma values of the SVC model, selected the best parameters and narrowed down the hyperparameters. This produced the final report for SVM with an increase of 9 percent accuracy.

SVM Classifier Metrics: Accuracy: 0.693050193050193 Classification Report:						
rt						
18						
18						
36						
36						
36						
,						

Training and Testing Dataset

The training and testing datasets were provided in the database. The following gives insight into the size of each respective dataset. After further analysis, this is about a 1:3 ratio of test data to train data. As mentioned above, the targets were set by a predefined "effectiveness" classification.

	X_train	Y_train	X_test	Y_test
# Entries	3107	3107	1036	1036

Results

By using accuracy as the main metric for the accuracy of classification of each of the models, here is the final result.

	RF	GNB	SVM
Accuracy	0.85	0.82	0.69

Random Forest was able to produce the best accuracy rate in classifying the effectiveness of the drugs. GNB was also effective in doing so with the implementation of smoothing. I believe that GNB gave great insight into the classification of this data and perhaps with additional feature engineering could prove to be very effective due to the simplicity and computation speed. Random Forest required extensive tuning compared to the rest of the models and computation time was the longest. SVM produced the worst accuracy. However with additional feature engineering along with tuning the weights of that additional data, could yield a greater accuracy.

If I were to take this investigation further, I would look to create additional sources of analysis to find strong correlations between certain features and perform additional feature engineering. I believe the models can be improved with additional data with careful attention to avoid overfitting or underfitting.

In conclusion, I was able to classify the effectiveness of drugs in this dataset by defining a target that the models would be looking to meet, implementing Random Forest, Gaussian Naive Bayes and Support Vector Machine, and analyzing the classification report.

Code: (Not intended to be part of the report with a 3 page limit, just wanted to include at the end for additional insight)

```
import pandas as pd
 2 import numpy as no
 4 ### Data Preprocessing ###
5 from sklearn.impute import SimpleImputer
6 from sklearn.preprocessing import LabelEncoder
  7 #sentiment analysis
 8 from textblob import TextBlob
 9 #text processing
10 import re
11 import nltk
13 ### Hyperparameter tuning ###
14 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
17 from sklearn import svm
18 from sklearn.ensemble import RandomForestClassifier
19 from sklearn.naive_bayes import GaussianNB
22 from sklearn.metrics import (
          accuracy_score,
classification_report
25 )
27 ###### READING IN DATA & PREPROCESSING ######
28 dtypes = {
           /pes = {
    "Unnamed: 0": "numerical",
    "urlDrugName": "categorical",
    "rating": "numerical",
    "effectiveness": "categorical",
    "sideEffects": "categorical",
    "condition": "categorical",
    "benefitsReview": "text",
    "sideEffectsReview": "text",
    "commentsReview": "text",
    "commentsReview": "text",
}
33
34
35
36
37
38 }
40 #Load training data
41 train_data = pd.read_csv('drugLib_raw/drugLibTrain_raw.tsv', sep='\t')
42 #Load test data
43 test_data = pd.read_csv('drugLib_raw/drugLibTest_raw.tsv', sep='\t')
46 ### Review Preprocessing (for text values) ###
47 def utils_preprocess_text(text, flg_stemm=False, flg_lemm=True, lst_stopwords=None):
48 ## clean (convert to lowercase and remove punctuations and characters and then strip)
             ## clean (convert to lowercase and remove punctuations at
text = re.sub(r'[^\w\s]', '', str(text).lower().strip())
49
50
51
             ## Tokenize (convert from string to list)
             ## TOKETIZE (CONVERT FYOM STRING to LIST)

ISt_text = text.split()

## remove Stopwords

if lst_stopwords is not None:

lst_text = [word for word in lst_text if word not in

lst_stopwords]
52
53
54
55
56
57
58
59
60
61
            ## Stemming (remove -ing, -ly, ...)
if flg_stemm == True:
    ps = nltk.stem.porter.PorterStemmer()
    lst_text = [ps.stem(word) for word in lst_text]
62
64
65
66
67
            ## Lemmatisation (convert the word into root word)
if flg_lemm == True:
    lem = nltk.stem.wordnet.WordNetlemmatizer()
    lst_text = [lem.lemmatize(word) for word in lst_text]
             ## back to string from list
text = " ".join(lst_text)
return text
68
69
73 train_data["sideEffectsReview"] = train_data["sideEffectsReview"].apply(lambda x: utils_preprocess_text(x, flg_stemm=False, flg_lemm=True))
74 train_data["commentsReview"] = train_data["commentsReview"].apply(lambda x: utils_preprocess_text(x, flg_stemm=False, flg_lemm=True))
75 train_data["benefitsReview"] = train_data["benefitsReview"].apply(lambda x: utils_preprocess_text(x, flg_stemm=False, flg_lemm=True))
78 train_data['sideEffectsReview'] = train_data['sideEffectsReview'].astype(str)
79 train_data["sideEffectsSentiment"] = train_data["sideEffectsReview"].apply(lambda x: TextBlob(x).sentiment.polarity)
81 train_data['commentsReview'] = train_data['commentsReview'].astype(str)
82 train_data["commentsSentiment"] = train_data["commentsReview"].apply(lambda x: TextBlob(x).sentiment.polarity)
85 train_data['benefitsReview'] = train_data['benefitsReview'].astype(str)
86 train_data["benefitsSentiment"] = train_data["benefitsReview"].apply(lambda x: TextBlob(x).sentiment.polarity)
88 ######################
90 test_data["sideEffectsReview"] = test_data["sideEffectsReview"].apply(lambda x: utils_preprocess_text(x, flg_stemm=False, flg_lemm=True))
91 test_data["commentsReview"] = test_data["commentsReview"].apply(lambda x: utils_preprocess_text(x, flg_stemm=False, flg_lemm=True))
92 test_data["benefitsReview"] = test_data["benefitsReview"].apply(lambda x: utils_preprocess_text(x, flg_stemm=False, flg_lemm=True))
.
95 test_data['sideEffectsReview'] = test_data['sideEffectsReview'].astype(str)
96 test_data["sideEffectsSentiment"] = test_data["sideEffectsReview"].apply(lambda x: TextBlob(x).sentiment.polarity)
    test_data['commentsReview'] = test_data['commentsReview'].astype(str)
test_data["commentsSentiment"] = test_data["commentsReview"].apply(lambda x: TextBlob(x).sentiment.polarity)
.02 test_data['benefitsReview'] = test_data['benefitsReview'].astype(str)
.03 test_data["benefitsSentiment"] = test_data["benefitsReview"].apply(lambda x: TextBlob(x).sentiment.polarity)
```

```
.
106 #overall sentiment classificaiton
107 train_data["avgSentiment"] = (train_data["sideEffectsSentiment"]+train_data["sideEffectsSentiment"]+train_data["sideEffectsSentiment"]) /
108 train_data["avgSentiment"] = round(train_data["avgSentiment"])
109
110 test_data["avgSentiment"] = (test_data["sideEffectsSentiment"]+test_data["sideEffectsSentiment"]) / 3
111 test_data["avgSentiment"] = round(test_data["avgSentiment"])
113 ## Handling the missing values and assigning old column names
114 train_imp, test_imp = [
115 pd.DataFrame(
              SimpleImputer(strategy="most_frequent").fit_transform(df), columns=df.columns
117
          for df in (train_data, test_data)
119 ]
120
21 ## Encoding the categorical columns
122 for i in ["urlDrugName", "condition", "sideEffects"]:
123    train_imp[i] = LabelEncoder().fit_transform(train_imp[i])
124    test_imp[i] = LabelEncoder().fit_transform(test_imp[i])
 .26 ### drop text columns
127 train_imp.drop("benefitsReview", axis=1, inplace=True)
128 train_imp.drop("sideEffectsReview", axis=1, inplace=True)
129 train_imp.drop("commentsReview", axis=1, inplace=True)
L30
131 test_imp.drop("benefitsReview", axis=1, inplace=True)
132 test_imp.drop("sideEffectsReview", axis=1, inplace=True)
133 test_imp.drop("commentsReview", axis=1, inplace=True)
L34
135 train_imp.columns = train_imp.columns.astype(str)
136 test_imp.columns = test_imp.columns.astype(str)
 .37
 .38 #determine target based off of effectiveness and sentiment
139 #Design decision: if highly effective or considerably effective and has
140 #positive sentiment or neutral sentiment present -> effective classification
L41 #else -> ineffective classification
 42
ls0
151 train_imp["target"] = (targetTrain).astype(int)
152 test_imp["target"] = (targetTest).astype(int)
153
L55 ## Splitting the train and test datasets into feature variables
```

```
60 ###### Models (Fitting, Tuning, Evaluation) ######
162 def print_metrics(y_true, y_pred):
        # Accuracy
accuracy = accuracy_score(y_true, y_pred)
print("Accuracy:", accuracy)
163
164
165
166
167
         # Classification Report
        class_report = classification_report(y_true, y_pred)
print("Classification Report:")
.68
69
         print(class_report)
73 ###### SVM ######
74 svm_clf = svm.SVC()
 76 param_grid = {'C': [0.1, 1, 2],
77 'gamma': [1, 0.1, 0.01],
78 'kernel': ['rbf']}
81 grid = GridSearchCV(svm_clf, param_grid, refit = True, verbose = 3, n_jobs=-1)
   grid.fit(X_train,Y_train)
85 best_model = arid.best_estimator_
 87 svm_y_pred = best_model.predict(X_test)
 print("SVM Classifier Metrics:")
 print_metrics(Y_test, svm_y_pred)
 2 ###### GNB ######
 93 gnb = GaussianNB()
 cv=3,
verbose=1,
scoring='accuracy')
 gs_NB.fit(X_train,Y_train)
   best_params = gs_NB.best_params_
 07 best_model = gs_NB.best_estimator_
  gnb_y_pred = best_model.predict(X_test)
 10 print("GNB Classifier Metrics:")
11 print_metrics(Y_test, gnb_y_pred)
13 ###### Random Forest ######
13 ###### Random Forest ######

14 param_grid = {

15 'bootstrap': [True],

16 'max_depth': [30, 40, 50, 100],

17 'max_features': ['sapt'],

18 'min_samples_leaf': [1, 2, 3],

19 'min_samples_split': [8, 10, 12],

20 'n_estimators': [350, 375, 400, 425, 450]
220
223 rf = RandomForestClassifier()
.27
_225 #Tune Hyperparameters
226 rand_rf = GridSearchCV(estimator = rf, param_grid = param_grid, cv = 3, verbose=2, n_jobs = -1)
228 rand_rf.fit(X_train, Y_train)
229 best_params = rand_rf.best_params_
230 best_model = rand_rf.best_estimator_
32 rf_y_pred = best_model.predict(X_test)
```