

Name: Akhil Sreedhara
NJIT UCID: as3638
Email: as3638@njit.edu
Date: Nov 13, 2024
Professor: Dr. Yasser Abdulllah

GITHUB REPOSITORY: <https://github.com/svakhil00/DataMiningFinal>

*I had an existing github account so I just used that

Introduction

My code uses 3 different machine learning models to classify different flowers in the iris dataset. I first train Random Forest, Support Vector Machine, and LSTM on the dataset. I then make the models predict all the classes in the test set. Using the predictions and actual labels, I calculate various metrics(full list can be found below) and print the results. I repeat this 10 times to perform 10-fold cross-validation and print the overall performances of the models.

Installation: Please check the requirements file to see instructions on how to setup the project.

Algorithms Used:

1. Random Forest
2. Support Vector Machine
3. LSTM

Dataset: <https://archive.ics.uci.edu/dataset/53/iris>

I used the Iris dataset to categorize flowers. There are normally 3 classes of flowers, but I trimmed the data so that there are only 2. In my code, instead of loading the dataset from a csv, i pull the data through an api call. I still include the dataset in my project.

Metrics Calculated:

1. **True Positive** -
model correctly predicts positive class
2. **True Negative** -
model correctly predicts negative class
3. **False Positive** -
model incorrectly predicts positive class
4. **False Negative** -
model incorrectly predicts negative class
5. **True Positive Rate** -
proportion of true positives that were correctly identified
6. **True Negative Rate** -
proportion of negatives predicted negative
7. **False Positive Rate** -
proportion of negatives predicted as positive
8. **False Negative Rate** -
proportion of positives predicted as negative
9. **Recall** -
same as tpr
10. **Precision** -
quality of the positive prediction
11. **F1-Score** -
harmonic recall of precision and recall
12. **Accuracy** -
proportion of correctly predicted cases(positive and negative)
13. **Error Rate** -
proportion of incorrectly predicted cases
14. **Specificity** -
same as tnr
15. **Negative Predictive Value** -
proportion of negative predictions that were correct
16. **False Discovery Rate** -
proportion of positive predictions that were incorrect

17. Balanced Accuracy -

more balanced when dataset is imbalanced

18. True Skill Statistics -

measures how well the model predicts compared to randomly guessing

19. Heidke Skill Score -

model's skill compared to random chance

20. Brier Score -

(mean squared error) calculates the average error

21. Brier Skill Score -

measure of how well a model performed compared to the baseline -1: worse 0: same 1: better

Code Walkthrough:

All required dependencies used by program. Installation instructions can be found in requirements file.

sklearn - used for 3 models, k-fold, getting tp, fp, fn, tp

numpy - math

pandas - data manipulation

matplotlib - graph

tensorflow - lstm

```
1 # imports
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from sklearn.svm import LinearSVC
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.preprocessing import MinMaxScaler # MinMaxScaler
8 from sklearn.model_selection import KFold
9 from sklearn.metrics import confusion_matrix
10 from tensorflow.keras.models import Sequential
11 from tensorflow.keras.layers import Dense, LSTM, Input
```

Function that takes in the labels for the test set along with the predicted labels from each model. Calculates and returns all metrics explained above.

```
1 # Calculates all Metrics
2 def evaluatePerformance(yTest, yPredictions, yProbabilities=None):
3     # Confusion Matrix
4     tn, fp, fn, tp = confusion_matrix(yTest, yPredictions).ravel()
5     # Positive
6     p = tp + fn
7     # Negative
8     n = tn + fp
9     # True Positive Rate
10    tpr = tp / p
11    # True Negative rate
12    tnr = tn / n
13    # False Positive Rate
14    fpr = fp / n
15    # False Negative Rate
16    fnr = fn / p
17    # Recall
18    r = tp / p
19    # Precision
20    precision = tp / (tp + fp)
21    # F1 Measure
22    f1 = 2 * (precision * r) / (precision + r)
23    # Accuracy
24    acc = (tp + tn) / (p + n)
25    # Error Rate
26    e = (fp + fn) / (p + n)
27    # Specificity
28    spc = tn / (fp + tn)
29    # Negative Predictive Value
30    npv = tn / (tn + fn)
31    # False Discovery Rate
32    fdr = fp / (fp + tp)
33    # Balanced Accuracy
```

```

34     bacc = .5 * (tp / (tp + fn) + tn / (tn + fp))
35     # True Skill Statistics
36     tss = tp / (tp + fn) - fp / (fp + tn)
37     # Heidke Skill Score
38     hss = 2 * (tp * tn - fp * fn) / ((tp + fn) * (fn + tn) + (tp + fp) * (fp + tn))
39     # Brier Score
40     bs = -123456789
41     if yProbabilities is not None:
42         bs = np.mean((yTest - yProbabilities) ** 2)
43     else:
44         bs = np.mean((yTest - yPredictions) ** 2)
45     # Brier Skill Score (The formula in the notes forgot to do "1 -" in the beginning)
46     bss = 1 - (bs / np.mean((yTest - np.mean(yTest)) ** 2))
47
48     return {
49         'True Positive': tp,
50         'True Negative': tn,
51         'False Positive': fp,
52         'False Negative': fn,
53         'True Positive Rate': tpr,
54         'True Negative Rate': tnr,
55         'False Positive Rate': fpr,
56         'False Negative Rate': fnr,
57         'Recall': r,
58         'Precision': precision,
59         'F1-Score': f1,
60         'Accuracy': acc,
61         'Error Rate': e,
62         'Specificity': spc,
63         'Negative Predictive Value': npv,
64         'False Discovery Rate': fdr,
65         'Balanced Accuracy': bacc,
66         'True Skill Statistics': tss,
67         'Heidke Skill Score': hss,
68         'Brier Score': bs,
69         'Brier Skill Score': bss
70     }

```

Function takes the metrics of each model in performance and the model name as a string. It then prints the data in a tabular format making it easier to read.

```

1 # Prints Metrics in Tabular Formant
2 def printMetrics(performance, modelName):
3     df = pd.DataFrame(performance, index=[0])
4     print(f'\n{modelName}')
5     print(df.T.to_string(header=False))

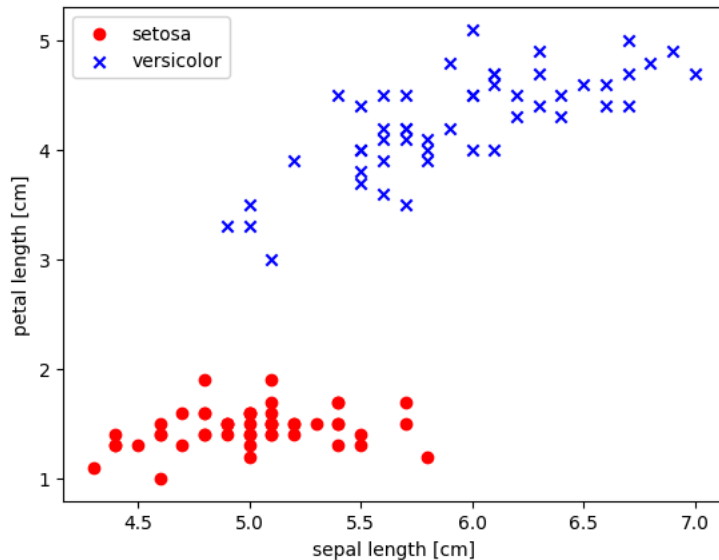
```

Calls an api to get the dataset and plots the data using 2 of the features. I called an api instead of reading a csv because it was easier. The actual csv is also included in the project files.

```

1 df = pd.read_csv('sreedhara_akhil_iris.csv', header=None, encoding='utf-8')
2
3 y = df.iloc[0:100, 4].values
4 y = np.where(y == 'Iris-setosa', 0, 1)
5 X = df.iloc[:100, [0, 1, 2, 3]].values
6
7 # Plot data
8 plt.scatter(X[:50, 0], X[:50, 2],
9             color='red', marker='o', label='setosa')
10 plt.scatter(X[50:100, 0], X[50:100, 2],
11            color='blue', marker='x', label='versicolor')
12 plt.xlabel('sepal length [cm]')
13 plt.ylabel('petal length [cm]')
14 plt.legend(loc='upper left')
15 plt.show()

```



Creating 10 folds of training and testing. Initializes arrays to store all metric information.

```
1 kf = KFold(n_splits=10, shuffle=True, random_state=42)
2
3 rfMetrics: list[dict] = []
4 svcMetrics: list[dict] = []
5 lstmMetrics: list[dict] = []
```

This loop iterates through the 10 folds. On each iteration it trains each of the models using the training set, predicts the classes for the test set, and calculates and prints the performance metrics for each model.

*** For a clearer view of the output please check outputs.txt

```
1 # K-Fold 10 times
2 for i, (train_index, test_index) in enumerate(kf.split(X), start=1):
3     # Splitting the data
4     X_train, X_test = X[train_index], X[test_index]
5     y_train, y_test = y[train_index], y[test_index]
6
7     # Normalize the data
8     scaler = MinMaxScaler()
9     X_train_scaled = scaler.fit_transform(X_train)
10    X_test_scaled = scaler.transform(X_test)
11
12    # Reshape data for LSTM (samples, timesteps, features)
13    T = 1 # Set timesteps to 1 for basic LSTM
14    X_train_lstm = X_train_scaled.reshape((X_train_scaled.shape[0], T, X_train_scaled.shape[1]))
15    X_test_lstm = X_test_scaled.reshape((X_test_scaled.shape[0], T, X_test_scaled.shape[1]))
16
17    # Build the LSTM model
18    lstm = Sequential()
19    lstm.add(Input(shape=(T, X_train_scaled.shape[1])))
20    lstm.add(LSTM(units=50, return_sequences=False))
21    lstm.add(Dense(units=1, activation='sigmoid'))
22
23    # Compile the model
24    lstm.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
25
26    # Train the model
27    lstm.fit(X_train_lstm, y_train, epochs=50, batch_size=8, validation_split=0.1, verbose=0)
28
29    # Predict with LSTM
30    lstmProbabilities = lstm.predict(X_test_lstm)
31    lstmPredictions = (lstmProbabilities > 0.5).astype(int)
32
33    # Random Forest Classifier
34    rf = RandomForestClassifier()
35    rf.fit(X_train, y_train)
36    rfPredictions = rf.predict(X_test)
37
```

```

38 # Support Vector Machine Classifier
39 svc = LinearSVC(max_iter=10000)
40 svc.fit(X_train, y_train)
41 svcPredictions = svc.predict(X_test)
42
43 # Calculate Performances
44 performance = evaluatePerformance(y_test, rfPredictions)
45 printMetrics(performance, 'Random Forest Metrics:')
46 rfMetrics.append(performance)
47
48 performance = evaluatePerformance(y_test, svcPredictions)
49 printMetrics(performance, 'Support Vector Machine Metrics:')
50 svcMetrics.append(performance)
51
52 performance = evaluatePerformance(y_test, lstmPredictions, lstmProbabilities)
53 printMetrics(performance, 'Long Short Term Memory Metrics:')
54 lstmMetrics.append(performance)
55

```

```

... Precision          1.0
  F1-Score            1.0
  Accuracy            1.0
  Error Rate          0.0
  Specificity          1.0
  Negative Predictive Value 1.0
  False Discovery Rate 0.0
  Balanced Accuracy    1.0
  True Skill Statistics 1.0
  Heidke Skill Score   1.0
  Brier Score          0.0
  Brier Skill Score    1.0

```

Support Vector Machine Metrics:

```

True Positive          4.0
True Negative          6.0
False Positive          0.0
False Negative          0.0
True Positive Rate      1.0
True Negative Rate      1.0
False Positive Rate      0.0
False Negative Rate      0.0
Recall                  1.0
Precision                1.0
F1-Score                1.0
Accuracy                1.0
Error Rate              0.0
Specificity              1.0
Negative Predictive Value 1.0
False Discovery Rate      0.0
Balanced Accuracy        1.0
True Skill Statistics     1.0
Heidke Skill Score       1.0
Brier Score              0.0
Brier Skill Score        1.0

```

Long Short Term Memory Metrics:

```

True Positive          4.000000
True Negative          6.000000
False Positive          0.000000
False Negative          0.000000
True Positive Rate      1.000000
True Negative Rate      1.000000
False Positive Rate      0.000000
False Negative Rate      0.000000
Recall                  1.000000
Precision                1.000000
F1-Score                1.000000
Accuracy                1.000000
Error Rate              0.000000
Specificity              1.000000
Negative Predictive Value 1.000000
False Discovery Rate      0.000000
Balanced Accuracy        1.000000
True Skill Statistics     1.000000
Heidke Skill Score       1.000000
Brier Score              0.467046
Brier Skill Score       -0.946025

```

Once all 10 of the folds have been trained and tested, we find the overall performance of the models.

```
1 # Overall Performances
2 printMetrics(pd.DataFrame(rfMetrics).mean().to_dict(), 'Overall Random Forest Metrics:')
3 printMetrics(pd.DataFrame(svcMetrics).mean().to_dict(), 'Overall Support Vector Machine Metrics:')
4 printMetrics(pd.DataFrame(lstmMetrics).mean().to_dict(), 'Overall Long Short Term Memory Metrics:')
```

Analysis:

Overall, all 3 models performed perfectly on the test sets. They were able to correctly predict the class for every single piece of data in the dataset. I believe this to be because my dataset was pretty small easy to classify. Even though they performed perfectly, the lstm had a lower Brier Score and Brier Skill Score than the other two models. This is because lstm calculates probabilities instead of just a hard label. Using the probability we can tell how confident the model is in a prediction. Since the other 2 models are predicting with 100% confidence always, they give off the illusion of a greater brier score. In larger data sets, they would have a few wrong predictions and their true brier scores would show. The brier skill score for lstm came out to be negative which indicates that it performs worse than the baseline. I am curious to see what would happen on larger and more complicated datasets.

*** For a clearer view of the output please check outputs.txt