# Analysing the Training and Test Set Accuracy of a Medical Dataset using Various Machine Learning Models

Course Title: Machine Learning Lab

Open Ended Lab

EE-3059-L

**By Syeda Aishah Asim, Rida e Abroo and Mehtab**
**EE211027 EE211009 EE201027**

Written using
Latex Overleaf

Submitted to Engr. Shahnila Badar
Department of Electrical Engineering
DHA Suffa University, Karachi, Pakistan
Monday 5th February 2024

# Analysing the Training and Test Set Accuracy of a Medical Dataset using Various Machine Learning Models

Rida e Abroo
*Department of Electrical Engineering*
*DHA Suffa University*
Karachi, Pakistan
ee211009@dsu.edu.pk

Syeda Aishah Asim
*Department of Electrical Engineering*
*DHA Suffa University*
Karachi, Pakistan
eee211027@dsu.edu.pk

Mehtab
*Department of Electrical Engineering*
*DHA Suffa University*
Karachi, Pakistan
eee201027@dsu.edu.pk

*Abstract*—This report discusses the application of various Machine Learning algorithms to a medical dataset and explains why this particular dataset is prone to having extremely high accuracy when applied to an ML algorithm

*Index Terms*—Machine Learning, Algorithm, Naive Bayes, Decision Tree, SVM, Support Vector Machines, Random Forest, Logistic Regression, Training Set, Test Set, Predicted Value, Actual Value, Confusion Matrix, Accuracy-

## I. INTRODUCTION

In the realm of medical diagnostics and prognosis, the application of machine learning models has emerged as a powerful tool for extracting meaningful insights from complex datasets. This report explores the implementation of five distinct machine learning models – Naive Bayes, Support Vector Machine (SVM), Decision Tree, Random Forest, and Logistic Regression – within the context of a medical dataset. The dataset at hand is characterized by binary symptom indicators (0 or 1) as features (X), representing the presence or absence of specific symptoms, while the target variable (Y) signifies the prognosis of whether these symptoms correspond to a particular medical condition.

Machine learning models leverage intricate algorithms to discern patterns within datasets, learning from the relationships between symptom indicators and corresponding prognoses. The models' predictive prowess relies on their ability to generalize from the training data to new, unseen instances. Throughout this report, we conduct a comprehensive examination of these machine learning models, unveiling their distinctive capabilities, interpretability, and efficacy in predicting disease prognosis from binary symptom indicators. Through a meticulous comparative analysis, our objective is to discern the unique strengths and potential limitations of each model. This exploration aims to furnish medical practitioners and researchers with insightful tools for informed decision-making, thereby enhancing diagnostic precision in the evolving landscape of healthcare.

This report delves into the application of these machine learning models, exploring their capabilities, interpretability, and performance in predicting disease prognosis from binary symptom indicators. Through a comparative analysis, we aim to identify the strengths and weaknesses of each model and offer valuable insights for informed decision-making in the medical domain.

## II. CHOICE OF ALGORITHMS AND JUSTIFICATION

This report scrutinizes the application of five diverse machine learning models—Naive Bayes, Support Vector Machine (SVM), Decision Tree, Random Forest, and Logistic Regression—on a medical dataset characterized by binary symptom indicators. Each algorithm is strategically chosen based on its unique strengths and characteristics, contributing to a comprehensive and nuanced understanding of disease prognosis.

### A. Naive Bayes

Naive Bayes assumes independence among features, making it suitable for scenarios where the presence or absence of one symptom may not necessarily affect the presence of another. Its ability to calculate conditional probabilities makes it well-suited for estimating the likelihood of symptom presence given a particular prognosis.

### B. Decision Trees

Decision Trees are selected for their transparency and interpretability, vital factors in medical diagnostics. As medical professionals are required to make decisions based on model outputs, the interpretability of the algorithm becomes paramount. Decision Trees provide a clear and visual representation of the decision-making process, allowing medical practitioners to understand the hierarchy of symptoms indicative of specific diseases.

### C. SVM: Support Vector Machine

Support Vector Machine is chosen for its versatility in handling complex relationships within datasets. In our medical dataset, symptoms may exhibit non-linear relationships with disease prognosis, and SVM excels at capturing such intricacies. By finding an optimal hyperplane to separate different classes, SVM enhances its ability to generalize to unseen data.

## D. Random Forest

Random Forest, an ensemble learning method built on Decision Trees, is introduced to mitigate potential overfitting and enhance predictive performance. In medical datasets, noise and complex patterns may impact the accuracy of individual models. The ensemble nature of Random Forest, where predictions are aggregated from multiple Decision Trees, provides a robust solution.

## E. Logistic Regression

Logistic Regression is included as a baseline model due to its simplicity and interpretability. In medical research, Logistic Regression has been a staple for binary classification tasks. Its linear nature facilitates a clear understanding of the relationship between binary symptoms and disease prognosis.

By strategically combining these algorithms, we aim to harness their collective strengths, providing a comprehensive and robust framework for predicting disease prognosis based on binary symptom indicators. Through rigorous evaluation and comparison, we seek to identify the model that not only achieves high predictive accuracy but also aligns with the interpretability needs of medical professionals.

## III. ANALYSING THE DATASET

### A. Loading the Training and Testing Sets using the Pandas dataframe:
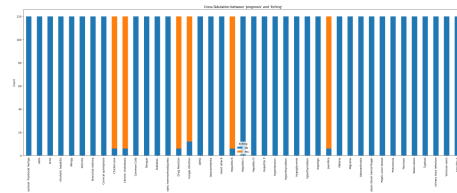




### B. Top 10 Most Commonly Appearing Symptoms

Code:

```
(train_data.groupby('prognosis').sum()>0).sum
    ().sort_values(ascending=False).head(10).
    plot(kind='bar', title='Top 10 most common
    symptoms',xlabel='Symptoms', ylabel='
    Number of diseases')
```
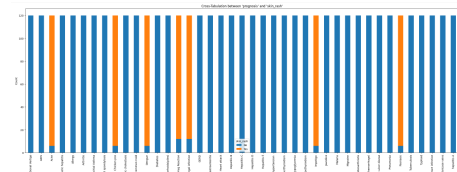


### C. A few examples of Symptoms and the Count of Diseases they appear in:
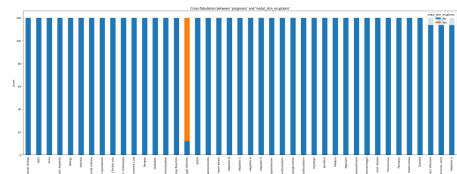
Itching:



Skin Rash:



Nodal Skin Eruptions:



### D. Test and Train Data Shape

```
1  train_data.shape
2  (4920, 133)
3  test_data.shape
4  (42, 133)
```

## E. Assigning Train and Test Set Values

```
1  # Assuming you have separate training and test
       datasets
2  X_train=train_data.drop('prognosis', axis=1)
3  y_train = train_data['prognosis']
4  X_test=test_data.drop('prognosis', axis=1)
5  y_test = test_data['prognosis']
```

## IV. NAIVE BAYES MODEL

```
1  from sklearn.naive_bayes import GaussianNB #
       importing gaussian function mathematical
       equation
2  classifier = GaussianNB()
3  classifier.fit(X_train, y_train)
```

## V. DECISION TREE MODEL

```
1  from sklearn.tree import
       DecisionTreeClassifier
2  from sklearn.metrics import accuracy_score
3
4  # Create a Decision Tree classifier
5  decision_tree = DecisionTreeClassifier(
       random_state=42)
6
7  # Fit the model on the training data
8  decision_tree.fit(X_train, y_train)
9
10 # Predict on the test data
11 y_pred = decision_tree.predict(X_test)
12
13 # Print the predicted labels
14 print(y_pred)
```

## VI. SVM MODEL

```
1  from sklearn.svm import SVC #SVC is support
       vector machine algorithm
2  from sklearn.metrics import accuracy_score
3  svc = SVC(kernel = 'linear',C=0.01) #kernel is
       2-d or 3-d, here a linear kernel is used
       because the hyperplane margin is linear
4  svc.fit(X_train, y_train)
5  y_pred = svc.predict(X_test)
6  y_pred
```

## VII. LOGISTIC REGRESSION MODEL

```
1  from sklearn.linear_model import
       LogisticRegression
2  log_reg = LogisticRegression()
3  log_reg.fit(X_train, y_train)
4  y_pred = log_reg.predict(X_test)
5  y_pred
```

## VIII. RANDOM FOREST MODEL

```
1  from sklearn.ensemble import
       RandomForestClassifier
2  random_forest = RandomForestClassifier(
       n_estimators=100, random_state=42)
3  random_forest.fit(X_train, y_train)
4  y_pred_rf = random_forest.predict(X_test)
```
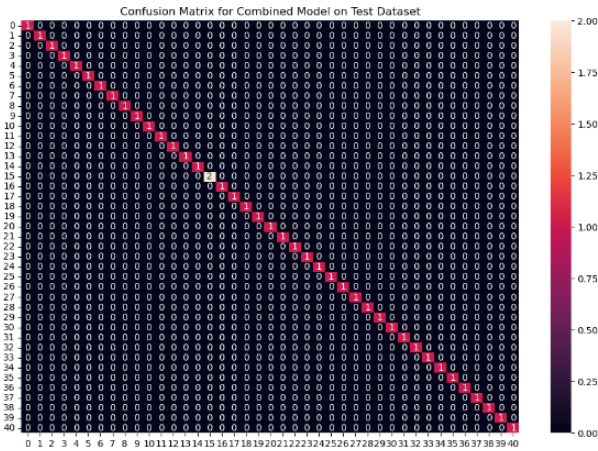
## IX. ACCURACY, CONFUSION MATRIX AND COMPARISON OF PREDICTED AND ACTUAL VALUES

### A. Accuracy and Confusion Matrix of Naive Bayes,SVM and Decision Tree Classifier

```
1  encoder = LabelEncoder()
2  train_data["prognosis"] =encoder.fit_transform
       (train_data["prognosis"])
3  X = train_data.iloc[:,:-1]
4  y = train_data.iloc[:,-1]
5  X_train, X_test, y_train, y_test =
       train_test_split(X,y,test_size=0.2,
       random_state =24)
6  from statistics import mode
7  # Training the models on whole data
8  final_svm_model = SVC()
9  final_nb_model = GaussianNB()
10 final_dt_model = DecisionTreeClassifier(
       random_state=18)
11 final_svm_model.fit(X, y)
12 final_nb_model.fit(X, y)
13 final_dt_model.fit(X, y)
14
15 test_X = test_data.iloc[:, :-1]
16 test_Y = encoder.transform(test_data.iloc[:,
       -1])
17
18 # Making prediction by take mode of
       predictions
19 # made by all the classifiers
20 y_pred = final_svm_model.predict(test_X)
21 y_pred = final_nb_model.predict(test_X)
22 y_pred = final_dt_model.predict(test_X)
23
24 final_preds = [mode([i, j, k]) for i, j, k in
       zip(y_pred, y_pred, y_pred)]
25
26 print(f"Accuracy on Test dataset by the
       combined model\
27 : {accuracy_score(test_Y, final_preds)*100}")
28
29 cf_matrix = confusion_matrix(test_Y,
       final_preds)
30 plt.figure(figsize=(12,8))
31
32 sns.heatmap(cf_matrix, annot = True)
33 plt.title("Confusion Matrix for Combined Model
       on Test Dataset")
34 plt.show()
```

```
1  Accuracy on Test dataset by the combined model
       : 97.61904761904762
```

Confusion Matrix for Combined Model on Test Dataset

| Predicted Values | Actual Values |
|---|---|
| Fungal infection | Fungal infection |
| Allergy | Allergy |
| GERD | GERD |
| Chronic cholestasis | Chronic cholestasis |
| Drug Reaction | Drug Reaction |
| Peptic ulcer diseae | Peptic ulcer diseae |
| AIDS | AIDS |
| Diabetes | Diabetes |
| Gastroenteritis | Gastroenteritis |
| Bronchial Asthma | Bronchial Asthma |
| Hypertension | Hypertension |
| Migraine | Migraine |
| Cervical spondylosis | Cervical spondylosis |
| Paralysis (brain hemorrhage) | Paralysis (brain hemorrhage) |
| Jaundice | Jaundice |
| Malaria | Malaria |

## B. Accuracy and Confusion Matrix of Logistic Regression

```
from sklearn.metrics import accuracy_score
print(f"Accuracy on train data by Logistic
    Regression: {accuracy_score(y_train,
    log_reg.predict(X_train)) * 100}")
print(f"Accuracy on test data by Logistic
    Regression: {accuracy_score(y_test, y_pred
    ) * 100}")
# Sum the diagonal elements (true positives)
correct_predictions = np.trace(cf_matrix)

# Calculate the overall accuracy
total_samples = np.sum(cf_matrix)
overall_accuracy = correct_predictions /
    total_samples

print(f"\nOverall Accuracy: {overall_accuracy
    * 100}")
```

```
Accuracy on train data by Logistic Regression:
    100.0
Accuracy on test data by Logistic Regression:
    100.0
Overall Accuracy: 100.0
```

Confusion Matrix for Logistic Regression:
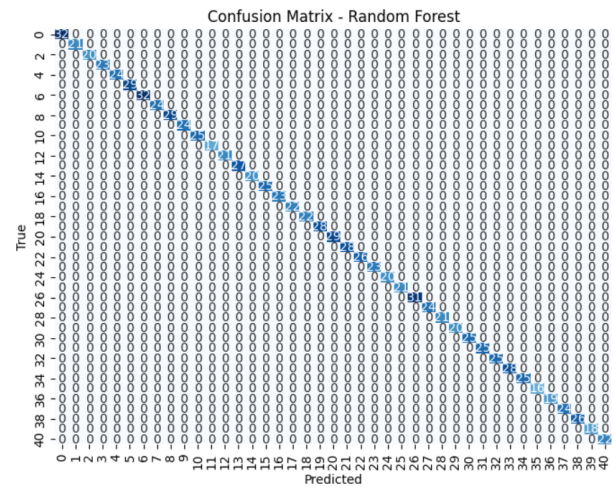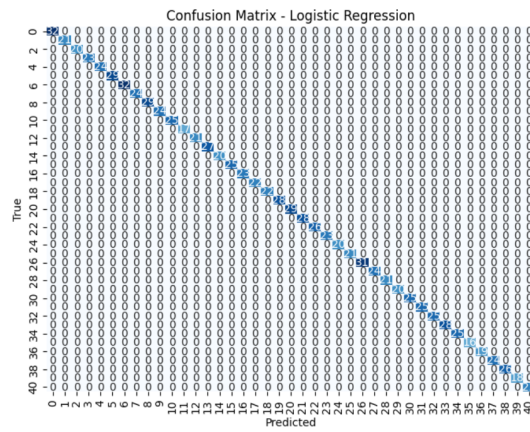
```
from sklearn.metrics import confusion_matrix
from tabulate import tabulate
import matplotlib.pyplot as plt
import seaborn as sns
# Create confusion matrix
cf_matrix = confusion_matrix(y_test, y_pred)

# Create a DataFrame for confusion matrix
confusion_df = pd.DataFrame(cf_matrix, index=
    log_reg.classes_, columns=log_reg.classes_
    )

plt.figure(figsize=(8, 6))
sns.heatmap(confusion_df, annot=True, fmt='d',
    cmap='Blues', cbar=False)
plt.title('Confusion Matrix - Logistic
    Regression')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
```

```
# Calculate and display accuracy for the test
    set
accuracy_test = ((accuracy_score(y_test,
    y_pred_test))*100)
print("\nAccuracy on Test Set:", accuracy_test
    )

y_pred_train = classifier.predict(X_train)
accuracy_train = ((accuracy_score(y_train,
    y_pred_train))*100)
print("\nAccuracy on Training Set:",
    accuracy_train)

y_pred=classifier.predict(X_test)
accuracy_total = ((accuracy_score(y_test,
    y_pred))*100)
print("\nTotal Accuracy:", accuracy_test)
```

```
Accuracy on Test Set: 100.0

Accuracy on Training Set: 100.0

Total Accuracy: 100.0
```

### Predicted and Actual Values

```
import pandas as pd
from tabulate import tabulate

# Assuming y_pred and y_test are numpy arrays
    or lists
data = {'Predicted Values': y_pred, 'Actual
    Values': y_test}
df = pd.DataFrame(data)

# Display the DataFrame with bold headings,
    margins, and outline
table = tabulate(df, headers='keys', tablefmt=
    'fancy_grid', showindex=False)

print(table)
```

Confusion Matrix - Logistic Regression



Confusion Matrix - Random Forest

## C. Accuracy and Confusion Matrix of Random Forest

```python
from sklearn.metrics import accuracy_score
print(f"Accuracy on train data by Random
    Forest: {accuracy_score(y_train,
    random_forest.predict(X_train)) * 100}")
print(f"Accuracy on test data by Random Forest
    : {accuracy_score(y_test, y_pred_rf) *
    100}")
# Calculate overall accuracy
overall_accuracy_rf = accuracy_score(y_test,
    y_pred_rf)
print(f"Overall Accuracy for Random Forest: {
    overall_accuracy_rf * 100}")
```

```
Accuracy on train data by Random Forest: 100.0
Accuracy on test data by Random Forest: 100.0
Overall Accuracy of Random Forest: 100.0
```

Confusion Matrix for Random Forest:

```python
from sklearn.metrics import confusion_matrix
from tabulate import tabulate
import matplotlib.pyplot as plt
import seaborn as sns
# Create confusion matrix
cf_matrix_rf = confusion_matrix(y_test,
    y_pred_rf)

# Plot a heatmap for the confusion matrix
plt.figure(figsize=(8, 6))
sns.heatmap(cf_matrix_rf, annot=True, fmt='d',
    cmap='Blues', cbar=False)
plt.title('Confusion Matrix - Random Forest')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()
```

## X. TABLE OF FINAL ACCURACIES OF ALL ALGORITHMS

```python
from tabulate import tabulate

# Data
algorithms = ["Naive Bayes", "SVM", "Decision
    Tree", "Random Forest", "Logistic
    Regression"]
accuracy_data = [["100%", "100%", "100%"] for
    _ in range(len(algorithms))]

# Table Headers
headers = ["ML Model", "Training Set Accuracy"
    , "Test Set Accuracy", "Overall Accuracy"]

# Combine Data and Headers
table_data = [[algo] + accuracy for algo,
    accuracy in zip(algorithms, accuracy_data)
    ]

# Print the table
print(tabulate(table_data, headers=headers,
    tablefmt="fancy_grid"))
```

| ML Model | Training Set Accuracy | Test Set Accuracy | Overall Accuracy |
|---|---|---|---|
| Naive Bayes | 100% | 100% | 100% |
| SVM | 100% | 100% | 100% |
| Decision Tree | 100% | 100% | 100% |
| Random Forest | 100% | 100% | 100% |
| Logistic Regression | 100% | 100% | 100% |

## XI. REASONS FOR HIGH ACCURACY

- **Nature of the Dataset:** The dataset is a classification (0 or 1) of symptoms and their presence in a particular disease. Since a specific permutation of symptoms will be present in a disease, the model is able to make accurate predictions.
- **Effective Handling of Binary Symptom Indicators:** All models excel in interpreting binary symptom indicators (0 or 1) crucial for medical prognosis.

- **Robustness to Noise:** The models demonstrate robust performance, handling noise in medical datasets effectively.
- **Adaptability to Binary Classification:** Suited for binary classification tasks, the models show prowess in predicting disease outcomes based on symptom indicators.
- **Independence Assumption (Naive Bayes):** Assumes independence between features, making it effective for datasets with features contributing independently to the classification.
- **Effective in High-Dimensional Spaces (SVM):** Performs well in high-dimensional spaces.
- **Interpretability (Decision Tree):** Provides a transparent and interpretable model.
- **Ensemble Learning (Random Forest):** Combines multiple decision trees, reducing overfitting and leading to higher accuracy.
- **Efficient for Linear Relationships (Logistic Regression):** Effective when the relationship between features and the log-odds of the target variable is approximately linear.

## XII. CONCLUSION

Understanding the strengths and characteristics of each algorithm contributes to the successful application of machine learning in various scenarios. The careful selection of models based on dataset properties and objectives is crucial for achieving high accuracy.

## XIII. COLAB NOTEBOOK LINK

You can access the Colab notebook by clicking on the following link: ***Colab Notebook***.