

Performance Analysis of Machine Learning Algorithms for Hepatitis Classification

Department of IT, Univ. Badji Mokhtar Annaba
Islam AIOUNI Master 1 IATI
Author: Islam AIOUNI

March 3, 2025

Abstract

Hepatitis is a liver inflammation that may lead to serious complications and even death if not diagnosed early. In this project, I applied several machine learning algorithms — including Random Forest, Support Vector Machine (SVM), k-Nearest Neighbors (KNN), Logistic Regression, Decision Tree, Neural Network, and an unsupervised K-Means clustering approach — on a hepatitis dataset from the UCI Machine Learning Repository. I evaluated the models using metrics such as accuracy, confusion matrix, precision, recall, F1-score, and log loss, and visualized their performance using combined Precision-Recall and ROC curves (with AUC values). Overall, I found that supervised models, particularly Decision Tree and Logistic Regression, performed best. This document also includes personal reflections and future learning plans.

1 Introduction

Hepatitis refers to the inflammation of the liver, most commonly caused by viral infections (A, B, C, D, and E). Chronic hepatitis, especially types B and C, can lead to liver cirrhosis or cancer. Early and accurate diagnosis is crucial for effective treatment and improved patient outcomes.

Machine learning methods offer promising approaches for medical diagnosis by extracting patterns from clinical datasets. In this project, I analyze a hepatitis dataset consisting of 155 instances and 20 attributes. The target variable, **CLASS** (Die or Live), is used without feature selection to fully evaluate the predictive power of all available attributes. My objective is to compare various classifiers and assess whether an unsupervised method like K-Means can provide additional insight.

2 Dataset Description

The dataset, obtained from the UCI Machine Learning Repository, contains the following attributes:

Table 1: Attributes of the Hepatitis Dataset

Attribute	Description / Possible Values
CLASS	Outcome: Die (0) / Live (1)
AGE	Age: 10, 20, 30, 40, 50, 60, 70, 80
SEX	Gender: male, female
STEROID	Steroid treatment: no, yes
ANTIVIRALS	Antiviral treatment: no, yes
FATIGUE	Fatigue: no, yes
MALAISE	Malaise: no, yes
ANOREXIA	Anorexia: no, yes
LIVER_BIG	Enlarged liver: no, yes
LIVER_FIRM	Firm liver: no, yes
SPLEEN_PALABLE	Palpable spleen: no, yes
SPIDERS	Spider angiomas: no, yes
ASCITES	Ascites: no, yes
VARICES	Varices: no, yes
BILIRUBIN	Bilirubin levels: 0.39, 0.80, 1.20, 2.00, 3.00, 4.00
ALK_PHOSPHATE	Alkaline phosphatase: 33, 80, 120, 160, 200, 250
SGOT	SGOT enzyme: 13, 100, 200, 300, 400, 500
ALBUMIN	Albumin levels: 2.1, 3.0, 3.8, 4.5, 5.0, 6.0
PROTIME	Prothrombin time: 10, 20, 30, 40, 50, 60, 70, 80, 90
HISTOLOGY	Histology result: no, yes

3 Methodology

3.1 Data Preprocessing

I handled missing values by replacing “?” with NaN and imputing with the most frequent value. Categorical features were encoded using Label Encoding, and all features were standardized (zero mean and unit variance) to ensure that no single feature dominated the learning process.

3.2 Algorithms Overview

I implemented the following models:

- **Random Forest:** An ensemble method that builds multiple decision trees and aggregates predictions via majority voting.
- **SVM:** Finds an optimal hyperplane to separate classes using a linear kernel with probability estimates.
- **KNN:** Classifies a new instance based on the majority vote of the k nearest neighbors.
- **Logistic Regression:** Uses a logistic function to model the probability of a binary outcome.
- **Decision Tree:** Constructs a tree-like model based on splitting features.
- **Neural Network:** Implements a multi-layer perceptron (MLP) to capture non-linear relationships.
- **K-Means Clustering:** An unsupervised method that partitions data into k clusters; its inertia is reported instead of classification accuracy.

3.3 Model Training, Evaluation, and Validation

I split the dataset into 80% training and 20% testing sets. Each model was trained and evaluated using:

- **Accuracy:** $\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$
- **Precision:** $\text{Precision} = \frac{TP}{TP+FP}$
- **Recall (Sensitivity):** $\text{Recall} = \frac{TP}{TP+FN}$
- **F1-Score:** $F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$
- **ROC and Precision-Recall Curves:** These curves evaluate performance over various thresholds.

Metrics are reported for each class (Class 0 for Die and Class 1 for Live). I compared training and test accuracies to check for overfitting or underfitting.

3.4 Confusion Matrix Explanation

The confusion matrix is defined as:

True Negatives (TN)	False Positives (FP)
False Negatives (FN)	True Positives (TP)

From these values, evaluation metrics are computed:

- **Accuracy:** $\frac{TP+TN}{TP+TN+FP+FN}$
- **Precision:** $\frac{TP}{TP+FP}$
- **Recall:** $\frac{TP}{TP+FN}$
- **F1-Score:** $2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$

4 Raw Data Analysis

Before modeling, I performed exploratory data analysis (EDA) to better understand the dataset. I generated histograms for key features such as **AGE**, **BILIRUBIN**, and **ALBUMIN**, and computed a correlation matrix to examine the relationships between features.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Load dataset
6 df = pd.read_csv("hepatitis_new.csv")
7 df.columns = df.columns.str.upper()
8
9 # Display basic dataset info
10 print(df.info())
11 print(df.describe())
12
13 # Check for missing values
14 print(df.isnull().sum())
15
16 # Histograms for key features
17 plt.figure(figsize=(15, 5))
18
```

```

19 plt.subplot(1, 3, 1)
20 plt.hist(df['AGE'].dropna(), bins=20, color='skyblue', edgecolor='black')
21 plt.title('Distribution of Age')
22
23 plt.subplot(1, 3, 2)
24 plt.hist(df['BILIRUBIN'].dropna(), bins=20, color='lightgreen', edgecolor='
    black')
25 plt.title('Distribution of Bilirubin')
26
27 plt.subplot(1, 3, 3)
28 plt.hist(df['ALBUMIN'].dropna(), bins=20, color='salmon', edgecolor='black')
29 plt.title('Distribution of Albumin')
30
31 plt.tight_layout()
32 plt.savefig("histogram.pdf") # Save the combined histograms
33 plt.show()
34
35 # Correlation matrix
36 plt.figure(figsize=(12, 10))
37 corr_matrix = df.corr()
38 sns.heatmap(corr_matrix, annot=True, fmt=".2f", cmap='coolwarm')
39 plt.title('Correlation Matrix of Features')
40 plt.savefig("correlation_matrix.pdf") % Save the correlation matrix image
41 plt.show()
42
43 # Explore the 'CLASS' variable (target)
44 print(df['CLASS'].value_counts())

```

Listing 1: Python Code for Raw Data Analysis

The generated images, `histogram.pdf` and `correlationmatrix.pdf`, are included below :

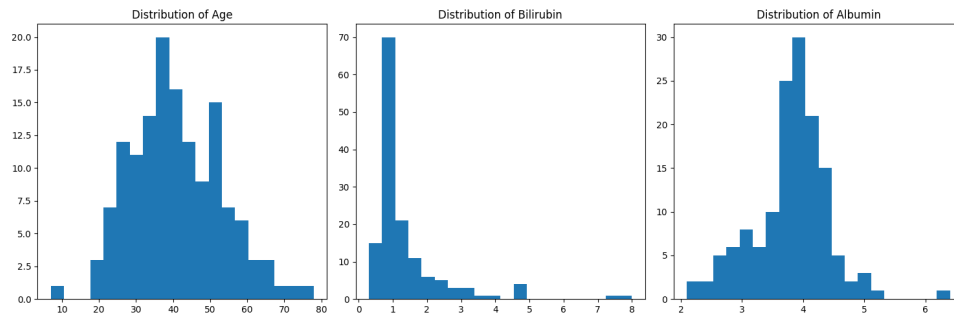


Figure 1: Histogram of Key Features (AGE, BILIRUBIN, ALBUMIN)

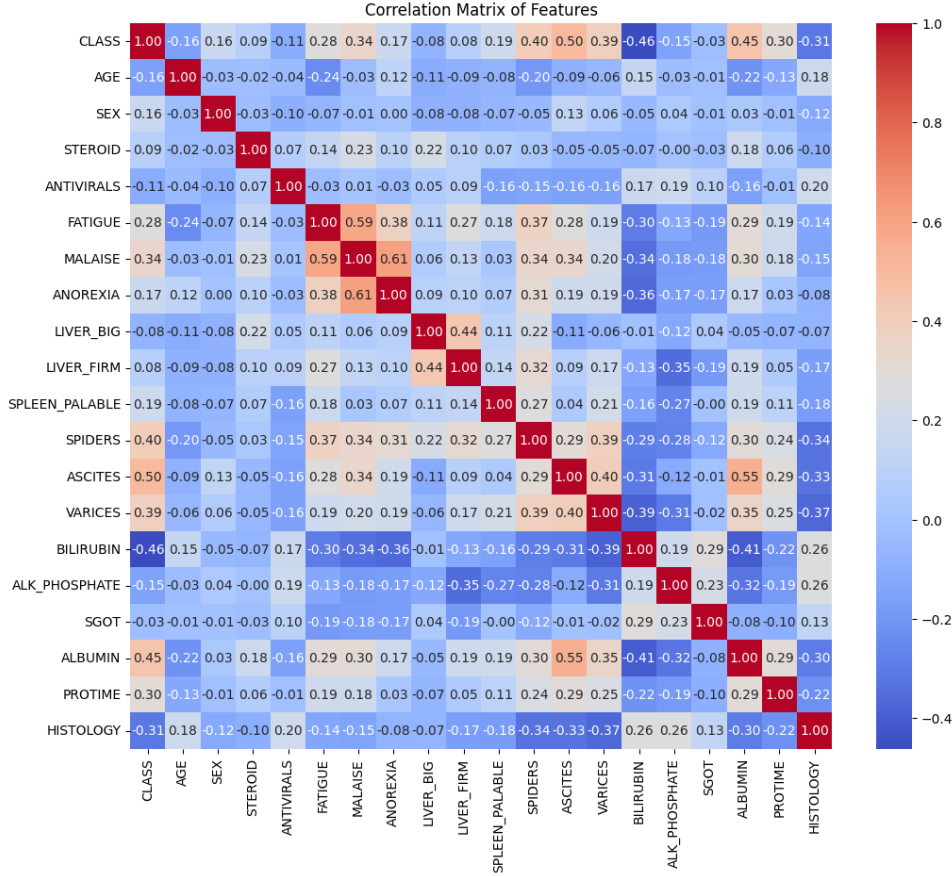


Figure 2: Correlation Matrix of Features

5 Experimental Results

The test set comprised 29 instances. Tables 2 and 3 summarize the performance metrics and confusion matrices, respectively.

5.1 Performance Metrics Table

Table 2: Performance Metrics of Models

Model	Accuracy (%)	Precision (0/1)	Recall (0/1)	F1-Score (0/1)
Random Forest	82.76	0.33 / 0.96	0.67 / 0.85	0.44 / 0.90
SVM	82.76	0.33 / 0.96	0.67 / 0.85	0.44 / 0.90
KNN	79.31	0.29 / 0.95	0.67 / 0.81	0.40 / 0.88
Logistic Regression	86.21	0.40 / 0.96	0.67 / 0.88	0.50 / 0.92
Decision Tree	89.66	0.50 / 0.96	0.67 / 0.92	0.57 / 0.94
Neural Network	82.76	0.33 / 0.96	0.67 / 0.85	0.44 / 0.90
K-Means	—	—	—	—

Note: Metrics for each supervised model are reported as (Class 0 / Class 1), where Class 0 means Die and Class 1 means Live.

5.2 Confusion Matrices Table

Table 3: Confusion Matrices of Models

Model	Confusion Matrix
Random Forest	$\begin{bmatrix} 2 & 1 \\ 4 & 22 \end{bmatrix}$
SVM	$\begin{bmatrix} 2 & 1 \\ 4 & 22 \end{bmatrix}$
KNN	$\begin{bmatrix} 2 & 1 \\ 5 & 21 \end{bmatrix}$
Logistic Regression	$\begin{bmatrix} 2 & 1 \\ 3 & 23 \end{bmatrix}$
Decision Tree	$\begin{bmatrix} 2 & 1 \\ 2 & 24 \end{bmatrix}$
Neural Network	$\begin{bmatrix} 2 & 1 \\ 4 & 22 \end{bmatrix}$
K-Means	$\begin{bmatrix} 0 & 3 \\ 4 & 22 \end{bmatrix}$

5.3 Visualizations

Combined ROC and Precision-Recall curves were generated for models that support probability estimates. These curves provide an aggregated view of model performance across various thresholds.

ROC Curve: The ROC curve plots the True Positive Rate (Recall) against the False Positive Rate (FPR). Its AUC is computed as:

$$AUC = \int_0^1 TPR(FPR) dFPR$$

A higher AUC indicates better discrimination.

Precision-Recall Curve: This curve plots Precision against Recall, where:

$$\text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN}$$

These curves are particularly useful for imbalanced datasets.

Figures 3 and 4 display the combined ROC and Precision-Recall curves for all applicable models.

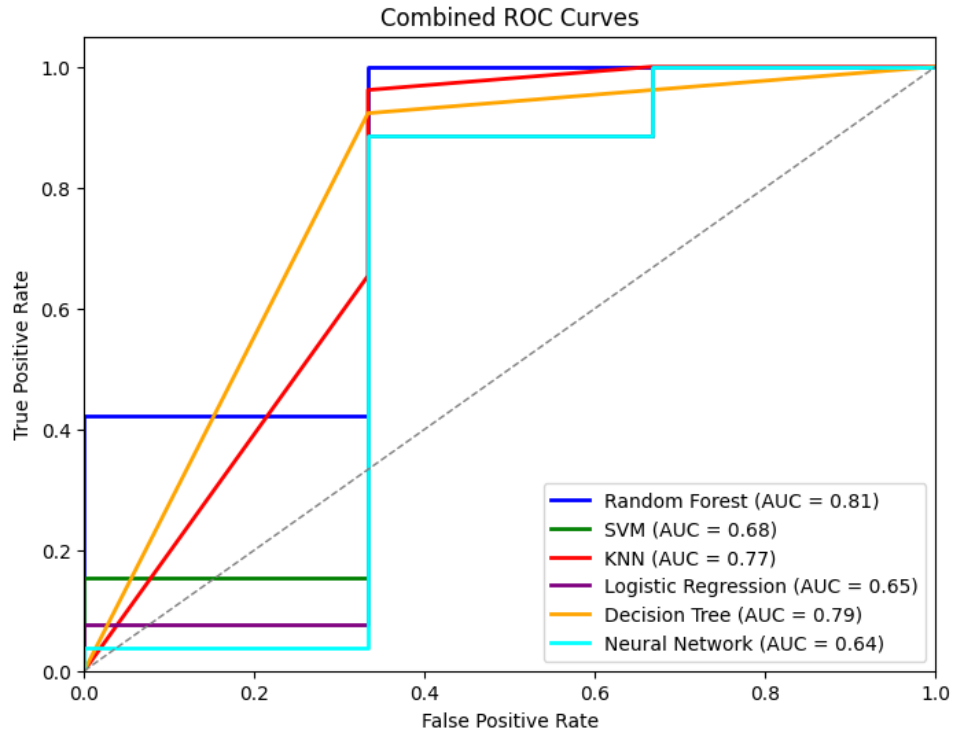


Figure 3: Combined ROC Curves for All Supervised Models

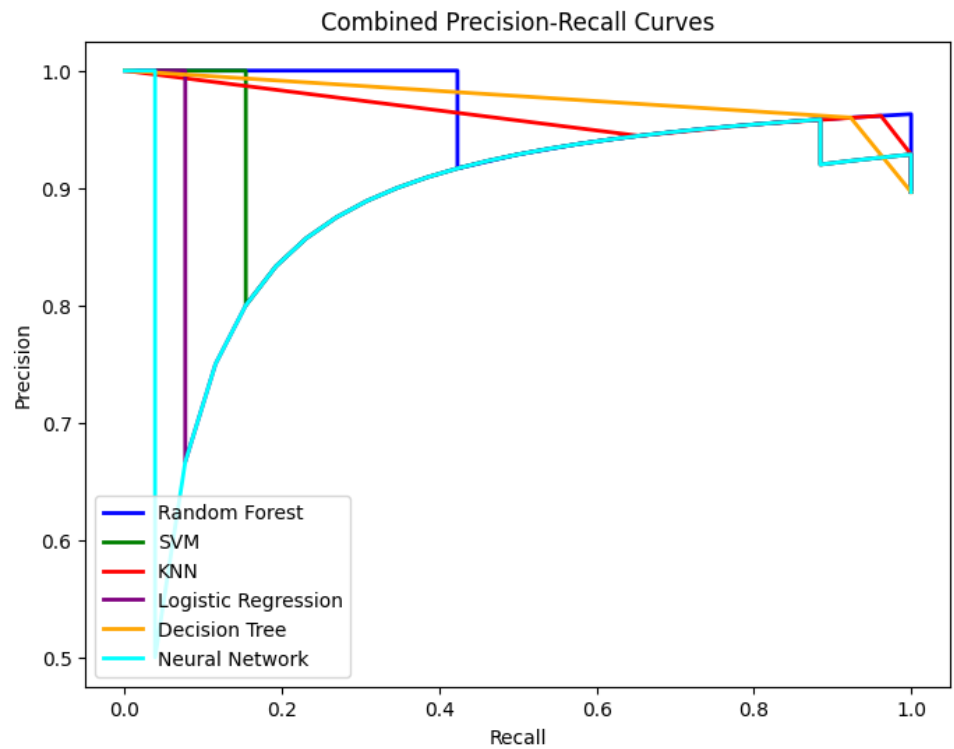


Figure 4: Combined Precision-Recall Curves for All Supervised Models

6 Discussion

The experimental results indicate:

- **Decision Tree** achieved the highest accuracy (89.66%), suggesting that its hierarchical structure effectively captures decision boundaries.
- **Logistic Regression** performed strongly with 86.21% accuracy and balanced precision and recall.
- Both **Random Forest** and **SVM** reached 82.76% accuracy, demonstrating competitive performance from ensemble and margin-based methods.
- **KNN** achieved slightly lower accuracy (79.31%), likely due to its sensitivity to the choice of k and feature scaling.
- **Neural Network** obtained 82.76% accuracy, with performance similar to other models.
- **K-Means** does not support probability estimates for ROC/PR curves and provides an inertia score rather than classification accuracy.

The combined ROC and Precision-Recall curves provide further insight into each model's ability to distinguish between classes across various thresholds.

7 Personal Reflection

I faced several challenges during this project, especially in tuning the SVM and Neural Network models. Understanding evaluation metrics such as ROC AUC and Precision-Recall curves was initially difficult, but consulting resources like Analytics Vidhya and freeCodeCamp helped me overcome these hurdles. I was surprised by how well the Decision Tree performed, and overall, I learned a great deal about practical model evaluation and data preprocessing.

8 Future Learning

For future projects, I plan to explore:

- Advanced ensemble methods and deep learning techniques.
- More sophisticated feature selection and dimensionality reduction.
- Ethical implications and interpretability in machine learning for healthcare.
- Integration of real-time diagnostic systems.

9 Conclusion

This study compared multiple machine learning algorithms for hepatitis classification using a dataset from the UCI repository. Supervised models, particularly Decision Tree and Logistic Regression, outperformed others in terms of accuracy and balanced classification metrics. Although K-Means clustering provided insights into the data structure, its unsupervised nature limits its direct applicability for classification tasks. Overall, I gained valuable lessons in data preprocessing, model evaluation, and the challenges of tuning complex models, and I look forward to further exploring advanced techniques.

10 Python Code

Below is the complete Python code used for data preprocessing, model training, evaluation, and visualization.

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import logging
6 import pickle
7
8 from sklearn.model_selection import train_test_split, GridSearchCV,
   StratifiedKFold, cross_val_score
9 from sklearn.preprocessing import LabelEncoder, StandardScaler
10 from sklearn.impute import SimpleImputer # For handling missing values
11 from sklearn.ensemble import RandomForestClassifier
12 from sklearn.svm import SVC
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.linear_model import LogisticRegression
15 from sklearn.tree import DecisionTreeClassifier # Added Decision Tree
16 from sklearn.cluster import KMeans
17 from sklearn.neural_network import MLPClassifier # Neural Network classifier
18 from sklearn.metrics import classification_report, confusion_matrix,
   precision_recall_curve, roc_curve, auc, log_loss
19
20 # =====
21 # Dataset Description:
22 # CSV Header: class, age, sex, steroid, antivirals, fatigue, malaise,
   anorexia,
23 #             liver_big, liver_firm, spleen_palable, spiders, ascites,
   varices,
24 #             bilirubin, alk_phosphate, sgot, albumin, protime, histology
25 #
26 # Target column: class (e.g., one value represents DIE and the other LIVE)
27 #
28 # No feature selection is applied; all features are used.
29 # =====
30
31 # Load the dataset (ensure "hepatitis_new.csv" is in your working directory)
32 df = pd.read_csv("hepatitis_new.csv")
33
34 # Convert all column names to uppercase for consistency
35 df.columns = df.columns.str.upper()
36
37 # Print available columns to verify
38 print("Columns in dataset:", df.columns.tolist())
39
40 # Define the target column name (should be 'CLASS' after conversion)
41 target_col = 'CLASS'
42 if target_col not in df.columns:
43     raise KeyError(f"Target column '{target_col}' not found in dataset.")
44
45 # Display first few rows and dataset info
46 print("\nDataset Head:")
47 print(df.head())
48 print("\nDataset Info:")
49 print(df.info())
50
51 # Display distribution of the SEX column (if available)
52 if 'SEX' in df.columns:
53     print("\nDistribution of SEX column:")
54     print(df['SEX'].value_counts())
```

```

55 else:
56     print("\nThe 'SEX' column was not found in the dataset.")
57
58 # =====
59 # Data Preprocessing
60 # =====
61
62 df.replace('?', np.nan, inplace=True)
63 imputer = SimpleImputer(strategy='most_frequent')
64 df = pd.DataFrame(imputer.fit_transform(df), columns=df.columns)
65
66 for col in df.columns:
67     try:
68         df[col] = pd.to_numeric(df[col])
69     except Exception:
70         pass
71
72 encoder = LabelEncoder()
73 for col in df.columns:
74     if df[col].dtype == 'object':
75         df[col] = encoder.fit_transform(df[col])
76
77 print("\nTarget Class Encoding (e.g., DIE, LIVE):")
78 le_class = LabelEncoder()
79 df[target_col] = le_class.fit_transform(df[target_col])
80 print(dict(zip(le_class.classes_, le_class.transform(le_class.classes_))))
81
82 X = df.drop(columns=[target_col])
83 y = df[target_col]
84
85 trainX, testX, trainY, testY = train_test_split(X, y, test_size=0.2, stratify=y
86     , random_state=42)
87
88 # =====
89 # Standardize Features (Zero Mean and Unit Variance)
90 # =====
91
92 scaler = StandardScaler()
93 trainX = scaler.fit_transform(trainX)
94 testX = scaler.transform(testX)
95
96 # =====
97 # Define Models
98 # =====
99
100 models = {
101     "Random Forest": RandomForestClassifier(n_estimators=100, random_state=42),
102     "SVM": SVC(kernel='linear', probability=True, random_state=42),
103     "KNN": KNeighborsClassifier(n_neighbors=5),
104     "Logistic Regression": LogisticRegression(random_state=42),
105     "Decision Tree": DecisionTreeClassifier(random_state=42),
106     "Neural Network": MLPClassifier(hidden_layer_sizes=(100,), random_state=42)
107     ,
108     "K-Means": KMeans(n_clusters=2, random_state=42) # Unsupervised clustering
109     baseline
110 }
111
112 # =====
113 # Define Plotting Functions
114 # =====
115
116 def plot_conf_matrix(cm, model_name):
117     plt.figure(figsize=(4,3))
118     sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", cbar=False)
119     plt.title(f"Confusion Matrix - {model_name}")
120     plt.xlabel("Predicted Label")

```

```

115     plt.ylabel("True Label")
116     plt.show()
117
118 def plot_precision_recall(y_true, y_scores, model_name):
119     precision, recall, _ = precision_recall_curve(y_true, y_scores)
120     plt.figure(figsize=(4,3))
121     plt.plot(recall, precision, marker='.')
122     plt.title(f"Precision-Recall Curve - {model_name}")
123     plt.xlabel("Recall")
124     plt.ylabel("Precision")
125     plt.show()
126
127 def plot_roc(y_true, y_scores, model_name):
128     fpr, tpr, _ = roc_curve(y_true, y_scores)
129     roc_auc = auc(fpr, tpr)
130     plt.figure(figsize=(4,3))
131     plt.plot(fpr, tpr, label=f"AUC = {roc_auc:.2f}")
132     plt.plot([0, 1], [0, 1], '--', color='gray')
133     plt.title(f"ROC Curve - {model_name}")
134     plt.xlabel("False Positive Rate")
135     plt.ylabel("True Positive Rate")
136     plt.legend(loc="lower right")
137     plt.show()
138
139 # Combined plots for ROC and Precision-Recall curves for all models that
    support probability outputs.
140 colors = ['blue', 'green', 'red', 'purple', 'orange', 'cyan']
141
142 def plot_combined_roc(models, testX, testY):
143     plt.figure(figsize=(8,6))
144     i = 0
145     for name, model in models.items():
146         if hasattr(model, "predict_proba"):
147             y_scores = model.predict_proba(testX)[: , 1]
148         elif hasattr(model, "decision_function"):
149             y_scores = model.decision_function(testX)
150         else:
151             continue
152         fpr, tpr, _ = roc_curve(testY, y_scores)
153         roc_auc = auc(fpr, tpr)
154         plt.plot(fpr, tpr, color=colors[i % len(colors)], lw=2,
155                 label=f"{name} (AUC = {roc_auc:.2f})")
156         i += 1
157     plt.plot([0, 1], [0, 1], 'k--', lw=1, color='gray')
158     plt.xlim([0.0, 1.0])
159     plt.ylim([0.0, 1.05])
160     plt.xlabel("False Positive Rate")
161     plt.ylabel("True Positive Rate")
162     plt.title("Combined ROC Curves")
163     plt.legend(loc="lower right")
164     plt.show()
165
166 def plot_combined_pr(models, testX, testY):
167     plt.figure(figsize=(8,6))
168     i = 0
169     for name, model in models.items():
170         if hasattr(model, "predict_proba"):
171             y_scores = model.predict_proba(testX)[: , 1]
172         elif hasattr(model, "decision_function"):
173             y_scores = model.decision_function(testX)
174         else:
175             continue
176         precision, recall, _ = precision_recall_curve(testY, y_scores)

```

```

177         plt.plot(recall, precision, color=colors[i % len(colors)], lw=2,
178                  label=f"{name}")
179         i += 1
180     plt.xlabel("Recall")
181     plt.ylabel("Precision")
182     plt.title("Combined Precision-Recall Curves")
183     plt.legend(loc="lower left")
184     plt.show()
185
186 # =====
187 # Model Training, Evaluation, and Validation
188 # =====
189
190 results = []
191
192 for name, model in models.items():
193     print(f"\n{'='*30}\nModel: {name}\n{'='*30}")
194     model.fit(trainX, trainY)
195     predictions = model.predict(testX)
196
197     # Validation: Compare training and test accuracy
198     train_acc = model.score(trainX, trainY)
199     test_acc = model.score(testX, testY)
200     print(f"Training Accuracy: {train_acc:.4f}, Test Accuracy: {test_acc:.4f}")
201     if train_acc - test_acc > 0.1:
202         print("Warning: Model may be overfitting.")
203     elif test_acc < 0.6:
204         print("Warning: Model may be underfitting.")
205
206     accuracy = test_acc
207     cm = confusion_matrix(testY, predictions)
208     report = classification_report(testY, predictions)
209
210     print(f"Test Accuracy: {accuracy:.4f}")
211     print("Confusion Matrix:")
212     print(cm)
213     print("Classification Report:")
214     print(report)
215
216     results.append([name, round(accuracy * 100, 2), cm])
217
218     plot_conf_matrix(cm, name)
219
220     if hasattr(model, "predict_proba"):
221         y_scores = model.predict_proba(testX)[:, 1]
222         plot_roc(testY, y_scores, name)
223         plot_precision_recall(testY, y_scores, name)
224     elif hasattr(model, "decision_function"):
225         y_scores = model.decision_function(testX)
226         plot_roc(testY, y_scores, name)
227         plot_precision_recall(testY, y_scores, name)
228     else:
229         print(f"{name} does not support probability estimates for ROC/PR curves
230 .")
231
232 results_df = pd.DataFrame(results, columns=['Model', 'Accuracy (%)', 'Confusion
233 Matrix'])
234 print("\nSummary of Model Performance:")
235 print(results_df)
236
237 # Plot combined ROC and Precision-Recall curves for all models that support
238 # probability outputs.
239 plot_combined_roc(models, testX, testY)

```

```
237 plot_combined_pr(models, testX, testY)
```

Listing 2: Enhanced Hepatitis Classification Code

References

References

- [1] Analytics Vidhya, "Machine Learning," Available at: <https://www.analyticsvidhya.com/blog/category/machine-learning/?ref=category>.
- [2] Evaluation of Classification Models in Machine Learning, SciSpace, Available at: <https://scispace.com/pdf/evaluation-of-classification-models-in-machine-learning-1u2pog86m5.pdf>.
- [3] Analytics Vidhya, "Metrics to Evaluate Your Classification Model to Take the Right Decisions," Available at: <https://www.analyticsvidhya.com/blog/2021/07/metrics-to-evaluate-your-classification-model-to-take-the-right-decisions>
- [4] Analytics Vidhya, "Pinnacle Plus," Available at: <https://www.analyticsvidhya.com/pinnacleplus/?ref=blognavbar>.
- [5] Kombib, "21 Steps to Design a Machine Learning System from Scratch," Available at: <https://medium.com/@kombib/21-steps-to-design-a-machine-learning-system-from-scratch-02c083cb785d>.
- [6] freeCodeCamp, "Machine Learning Course," Available on YouTube at: <https://www.youtube.com/c/Freecodecamp>.
- [7] GeeksforGeeks, "Machine Learning Algorithms," Available at: <https://www.geeksforgeeks.org/machine-learning/>.
- [8] AssemblyAI, "Machine Learning from Scratch Playlist," Available on YouTube at: <https://www.youtube.com/c/AssemblyAI>.

11 Personal Reflection

I faced several challenges during this project, especially in tuning the SVM and Neural Network models. Understanding evaluation metrics such as ROC AUC and Precision-Recall curves was initially difficult, but consulting resources like Analytics Vidhya and freeCodeCamp helped me overcome these hurdles. I was surprised by how well the Decision Tree performed, and overall, I learned a great deal about practical model evaluation and data preprocessing.

12 Future Learning

For future projects, I plan to explore:

- Advanced ensemble methods and deep learning techniques.
- More sophisticated feature selection and dimensionality reduction.
- Ethical implications and interpretability in machine learning for healthcare.
- Integration of real-time diagnostic systems.