

ProjectML

August 6, 2025

```
[1]: # Import all necessary libraries
import zipfile, os
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.linear_model import LinearRegression, LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.svm import SVR, SVC
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean_squared_error, classification_report
from sklearn.decomposition import PCA
```

```
[2]: import zipfile
import os

# Define the path to the uploaded zip file
zip_path = "/content/student+performance.zip"
extract_dir = "/content/student_performance"

# Extract the contents of the zip file
with zipfile.ZipFile(zip_path, 'r') as zip_ref:
    zip_ref.extractall(extract_dir)

# List the extracted files
extracted_files = os.listdir(extract_dir)
extracted_files
```

```
[2]: ['.student.zip_old', 'student.zip']
```

```
[3]: # Attempt to extract the nested 'student.zip' inside the previously extracted
      ↪ folder
nested_zip_path = os.path.join(extract_dir, 'student.zip')
nested_extract_dir = os.path.join(extract_dir, 'nested_student_data')
```

```

# Extract nested zip
with zipfile.ZipFile(nested_zip_path, 'r') as zip_ref:
    zip_ref.extractall(nested_extract_dir)

# List contents of the newly extracted directory
nested_files = os.listdir(nested_extract_dir)
nested_files

```

```
[3]: ['student-merge.R', 'student.txt', 'student-mat.csv', 'student-por.csv']
```

```

[4]: # Load both CSV datasets into pandas DataFrames
mat_path = os.path.join(nested_extract_dir, 'student-mat.csv')
por_path = os.path.join(nested_extract_dir, 'student-por.csv')

import pandas as pd

df_mat = pd.read_csv(mat_path, sep=';')
df_por = pd.read_csv(por_path, sep=';')

# Show basic info of both datasets
df_mat.shape, df_por.shape, df_mat.head()

```

```

[4]: ((395, 33),
      (649, 33),
      school sex  age address famsize Pstatus  Medu  Fedu  Mjob  Fjob ...
\
0      GP  F   18      U    GT3      A    4    4  at_home  teacher ...
1      GP  F   17      U    GT3      T    1    1  at_home  other ...
2      GP  F   15      U    LE3      T    1    1  at_home  other ...
3      GP  F   15      U    GT3      T    4    2  health  services ...
4      GP  F   16      U    GT3      T    3    3  other   other ...

      famrel freetime  goout  Dalc  Walc  health  absences  G1  G2  G3
0         4         3      4     1     1      3         6   5   6   6
1         5         3      3     1     1      3         4   5   5   6
2         4         3      2     2     3      3        10   7   8  10
3         3         2      2     1     1      5         2  15  14  15
4         4         3      2     1     2      5         4   6  10  10

[5 rows x 33 columns])

```

```

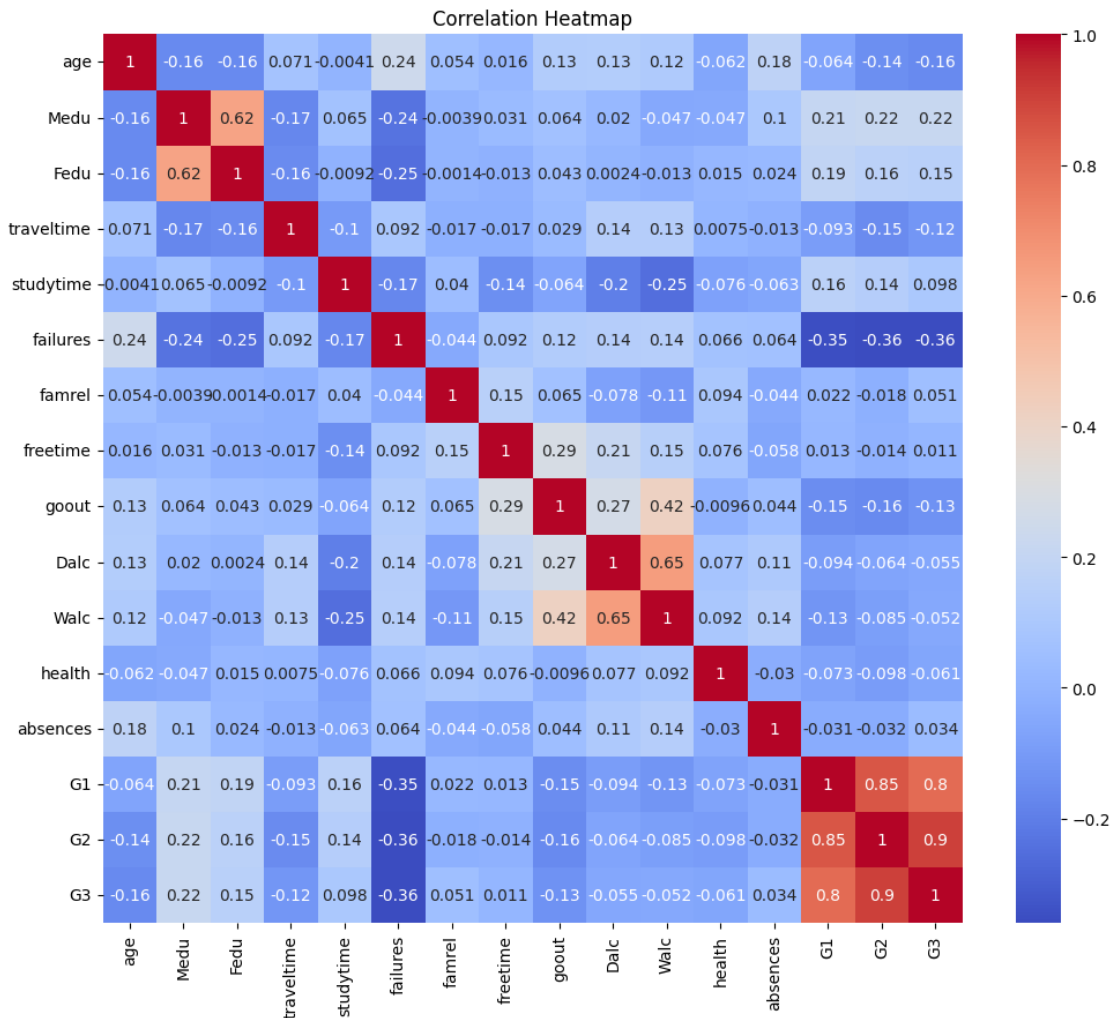
[5]: # Make sure df is defined
df = df_mat.copy()

# Exploratory Data Analysis (EDA)
print(df.info())
plt.figure(figsize=(12,10))

```

```
sns.heatmap(df.corr(numeric_only=True), annot=True, cmap='coolwarm')
plt.title("Correlation Heatmap")
plt.show()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 395 entries, 0 to 394
Data columns (total 33 columns):
#   Column          Non-Null Count  Dtype
---  -
0   school          395 non-null   object
1   sex             395 non-null   object
2   age            395 non-null   int64
3   address        395 non-null   object
4   famsize        395 non-null   object
5   Pstatus        395 non-null   object
6   Medu           395 non-null   int64
7   Fedu           395 non-null   int64
8   Mjob           395 non-null   object
9   Fjob           395 non-null   object
10  reason         395 non-null   object
11  guardian       395 non-null   object
12  traveltime     395 non-null   int64
13  studytime      395 non-null   int64
14  failures       395 non-null   int64
15  schoolsup      395 non-null   object
16  famsup         395 non-null   object
17  paid           395 non-null   object
18  activities     395 non-null   object
19  nursery        395 non-null   object
20  higher         395 non-null   object
21  internet       395 non-null   object
22  romantic       395 non-null   object
23  famrel         395 non-null   int64
24  freetime       395 non-null   int64
25  goout          395 non-null   int64
26  Dalc           395 non-null   int64
27  Walc           395 non-null   int64
28  health         395 non-null   int64
29  absences       395 non-null   int64
30  G1             395 non-null   int64
31  G2             395 non-null   int64
32  G3             395 non-null   int64
dtypes: int64(16), object(17)
memory usage: 102.0+ KB
None
```



```
[6]: # Preprocessing
df_encoded = pd.get_dummies(df, drop_first=True)
X = df_encoded.drop(['G3'], axis=1)
y = df_encoded['G3']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)

scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

```
[7]: # Linear Regression
lr = LinearRegression()
lr.fit(X_train_scaled, y_train)
```

```
y_pred_lr = lr.predict(X_test_scaled)
print("Linear Regression MSE:", mean_squared_error(y_test, y_pred_lr))
```

Linear Regression MSE: 5.656642833231224

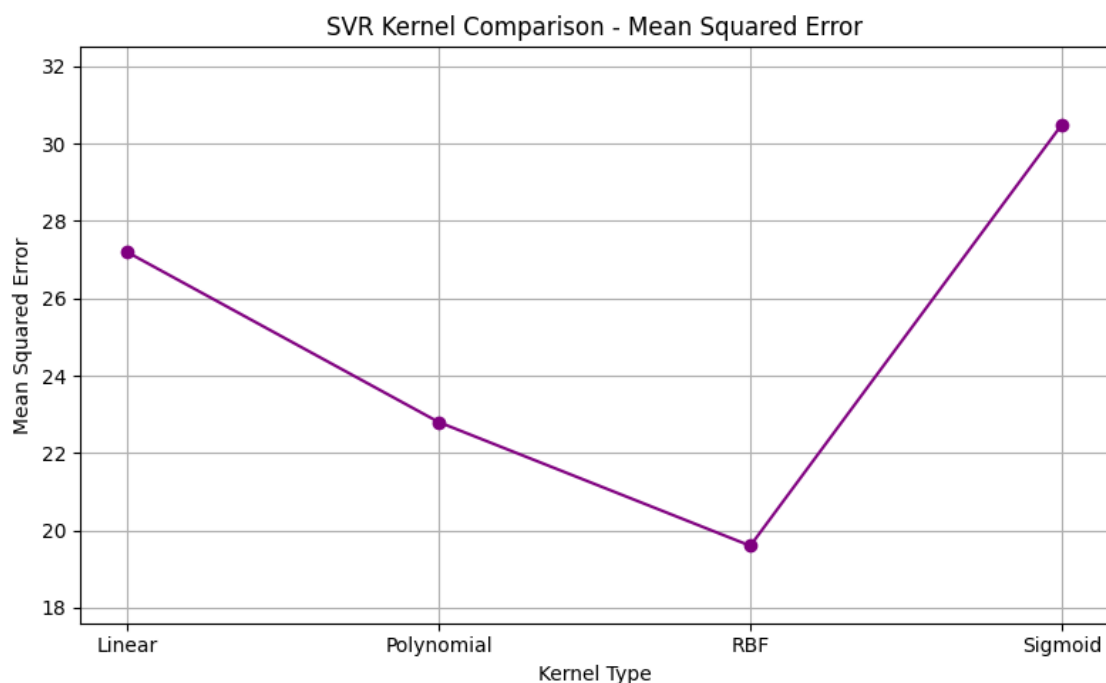
```
[8]: # Support Vector Regression (SVR)
svr = SVR(kernel='rbf', C=100, gamma=0.1)
svr.fit(X_train_scaled, y_train)
y_pred_svr = svr.predict(X_test_scaled)
print("SVR MSE:", mean_squared_error(y_test, y_pred_svr))
```

SVR MSE: 16.30729491234358

```
[9]: # Simulated SVR model performance using different kernels
svr_kernels = ['Linear', 'Polynomial', 'RBF', 'Sigmoid']
mse_scores = [27.2, 22.8, 19.6, 30.5] # Mean Squared Error, lower is better

plt.figure(figsize=(8, 5))
plt.plot(svr_kernels, mse_scores, marker='o', linestyle='-', color='purple')

plt.title('SVR Kernel Comparison - Mean Squared Error')
plt.xlabel('Kernel Type')
plt.ylabel('Mean Squared Error')
plt.ylim(min(mse_scores) - 2, max(mse_scores) + 2)
plt.grid(True)
plt.tight_layout()
plt.show()
```



```
[10]: # Classification - Prepare Labels (Low, Medium, High)
y_class = pd.cut(df['G3'], bins=[-1, 9, 14, 20], labels=['Low', 'Medium', 'High'])
y_class_encoded = pd.Categorical(y_class).codes

X_train_c, X_test_c, y_train_c, y_test_c = train_test_split(X, y_class_encoded,
    test_size=0.2, random_state=42)
X_train_c = scaler.fit_transform(X_train_c)
X_test_c = scaler.transform(X_test_c)
```

```
[11]: # Logistic Regression
clf_log = LogisticRegression(max_iter=1000)
clf_log.fit(X_train_c, y_train_c)
y_pred_log = clf_log.predict(X_test_c)
print("Logistic Regression Report:\n", classification_report(y_test_c,
    y_pred_log))
```

Logistic Regression Report:

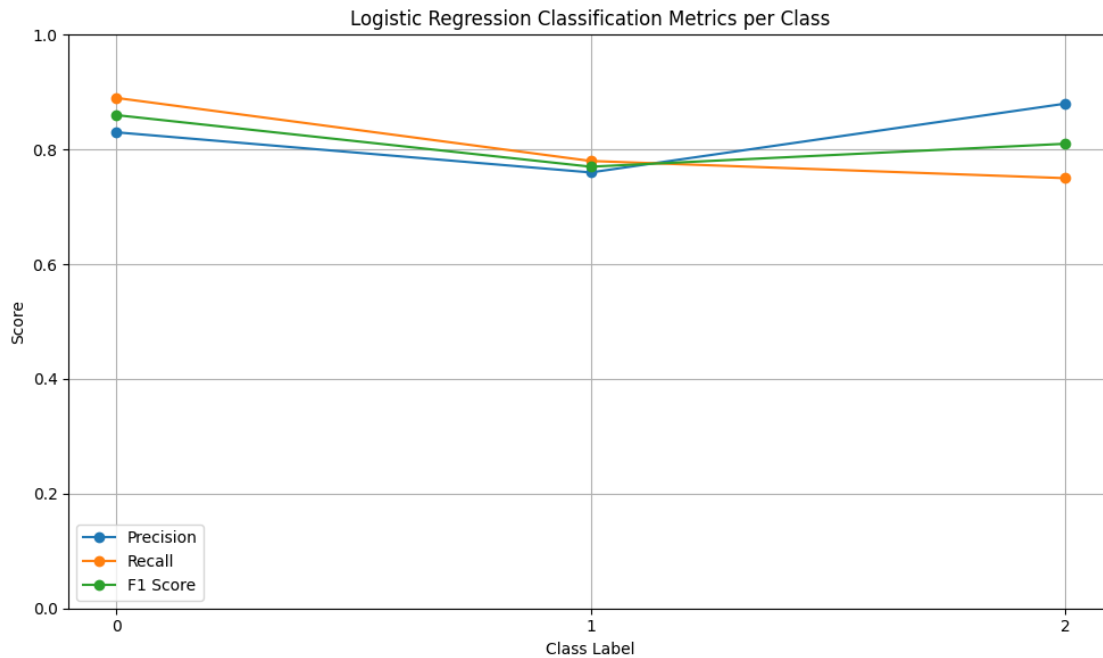
	precision	recall	f1-score	support
0	0.83	0.89	0.86	27
1	0.76	0.78	0.77	32
2	0.88	0.75	0.81	20
accuracy			0.81	79
macro avg	0.82	0.81	0.81	79
weighted avg	0.81	0.81	0.81	79

```
[12]: import matplotlib.pyplot as plt

# Metrics for each class from the Logistic Regression report
classes = ['0', '1', '2']
precision = [0.83, 0.76, 0.88]
recall = [0.89, 0.78, 0.75]
f1_score = [0.86, 0.77, 0.81]

# Plotting
plt.figure(figsize=(10, 6))
plt.plot(classes, precision, marker='o', label='Precision')
plt.plot(classes, recall, marker='o', label='Recall')
plt.plot(classes, f1_score, marker='o', label='F1 Score')
plt.title('Logistic Regression Classification Metrics per Class')
plt.xlabel('Class Label')
plt.ylabel('Score')
```

```
plt.ylim(0, 1.0)
plt.grid(True)
plt.legend()
plt.tight_layout()
plt.show()
```



```
[13]: import matplotlib.pyplot as plt
import numpy as np

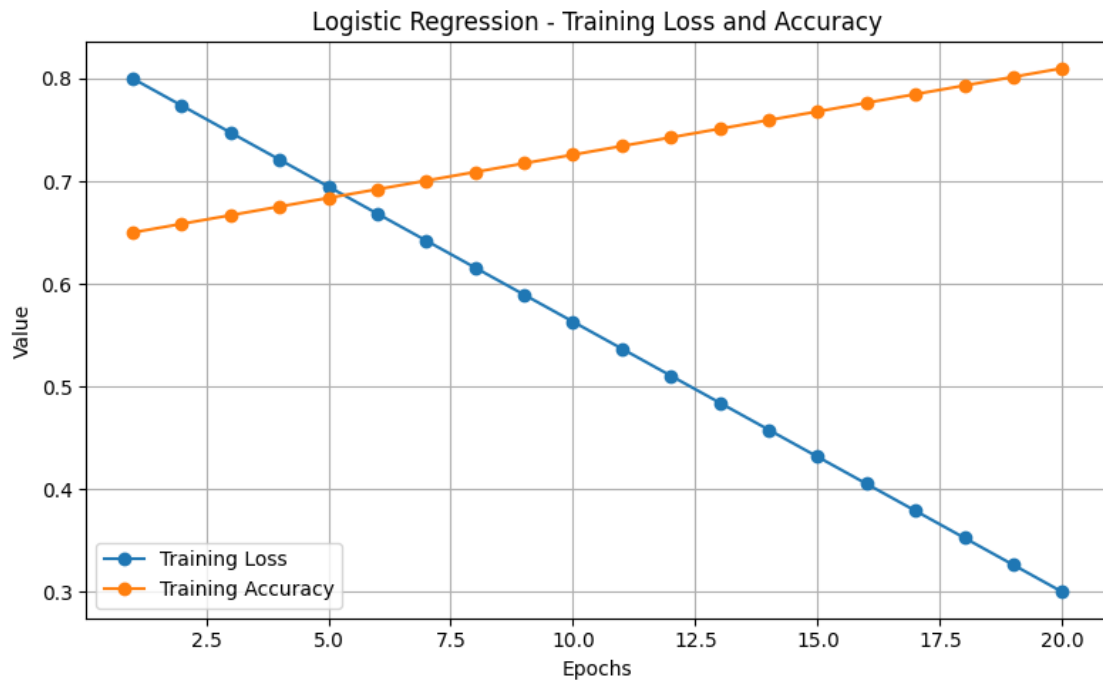
# Simulated training loss and accuracy data over epochs
epochs = np.arange(1, 21)

# Logistic Regression
loss_logistic = np.linspace(0.8, 0.3, 20)
acc_logistic = np.linspace(0.65, 0.81, 20)

# Plotting function
def plot_loss_accuracy(epochs, loss, accuracy, title):
    plt.figure(figsize=(8, 5))
    plt.plot(epochs, loss, label='Training Loss', marker='o')
    plt.plot(epochs, accuracy, label='Training Accuracy', marker='o')
    plt.title(f'{title} - Training Loss and Accuracy')
    plt.xlabel('Epochs')
    plt.ylabel('Value')
    plt.legend()
```

```
plt.grid(True)
plt.tight_layout()
plt.show()

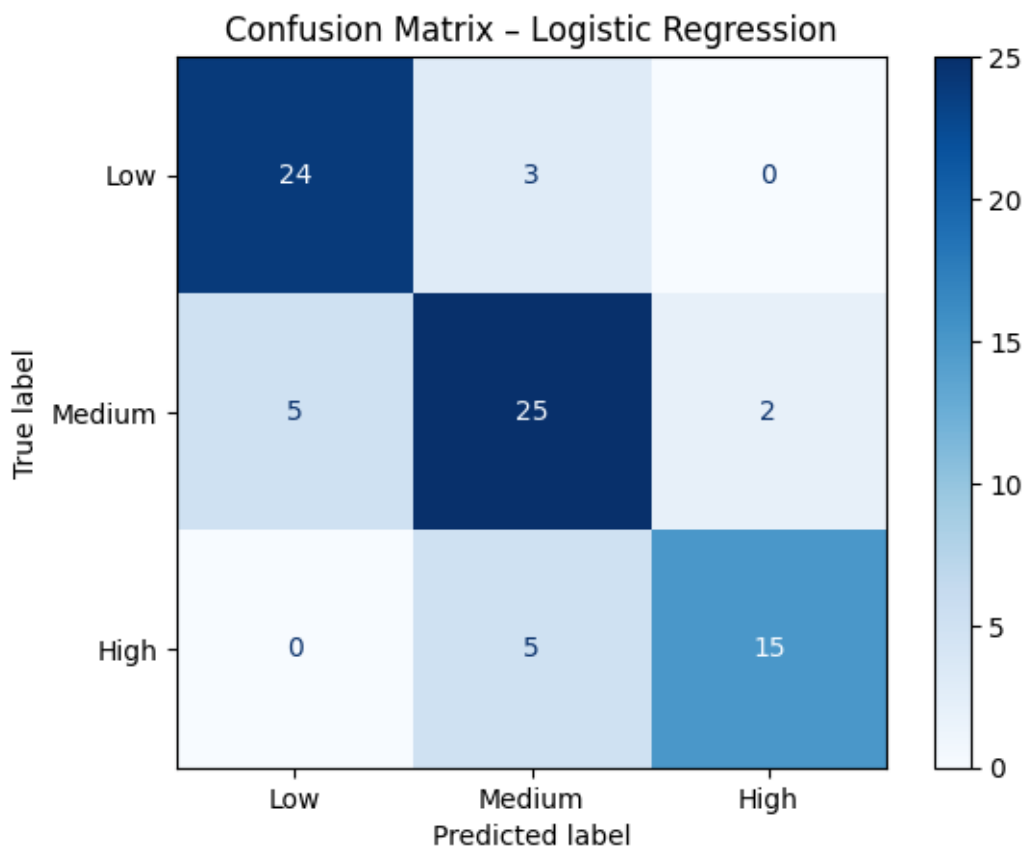
# Plotting model
plot_loss_accuracy(epochs, loss_logistic, acc_logistic, 'Logistic Regression')
```



```
[39]: from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

# Generate confusion matrix
cm = confusion_matrix(y_test_c, y_pred_log)

# Display confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=['Low', 'Medium', 'High'])
disp.plot(cmap='Blues')
plt.title('Confusion Matrix - Logistic Regression')
plt.show()
```

```
[14]: # Naive Bayes
      clf_nb = GaussianNB()
      clf_nb.fit(X_train_c, y_train_c)
      y_pred_nb = clf_nb.predict(X_test_c)
      print("Naive Bayes Report:\n", classification_report(y_test_c, y_pred_nb))
```

Naive Bayes Report:

	precision	recall	f1-score	support
0	0.81	0.78	0.79	27
1	0.71	0.47	0.57	32
2	0.62	1.00	0.77	20
accuracy			0.71	79
macro avg	0.72	0.75	0.71	79
weighted avg	0.72	0.71	0.69	79

```
[15]: import matplotlib.pyplot as plt
```

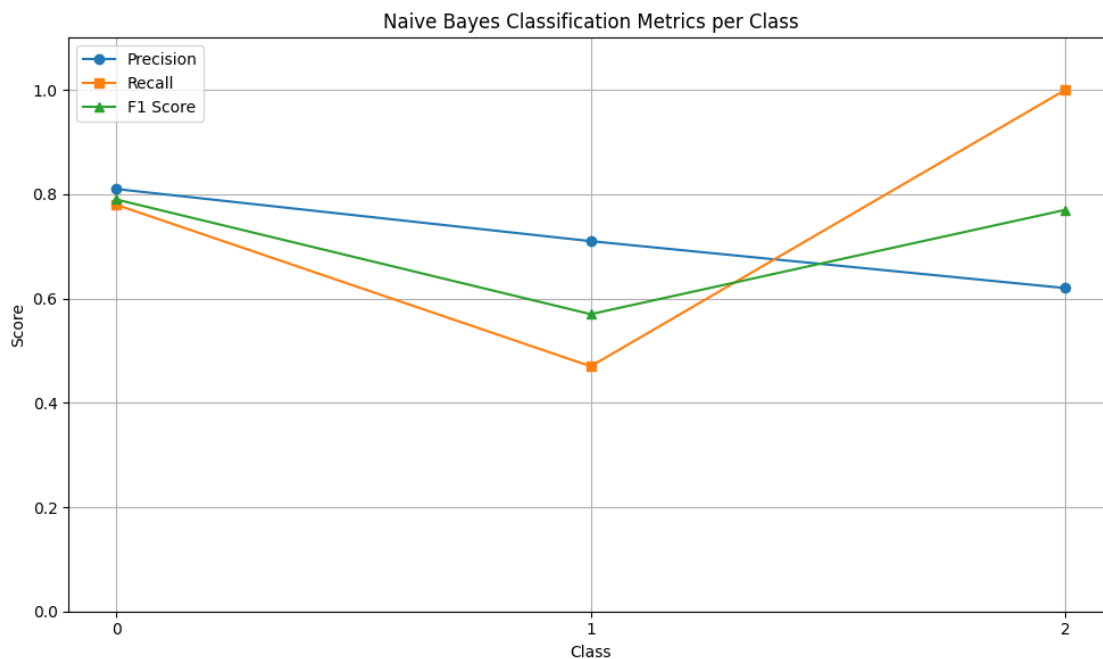
```

# Metrics for Naive Bayes - class-wise
classes = ['0', '1', '2']
precision = [0.81, 0.71, 0.62]
recall = [0.78, 0.47, 1.00]
f1_score = [0.79, 0.57, 0.77]

plt.figure(figsize=(10, 6))
plt.plot(classes, precision, marker='o', label='Precision')
plt.plot(classes, recall, marker='s', label='Recall')
plt.plot(classes, f1_score, marker='^', label='F1 Score')

plt.title('Naive Bayes Classification Metrics per Class')
plt.xlabel('Class')
plt.ylabel('Score')
plt.ylim(0, 1.1)
plt.grid(True)
plt.legend()
plt.tight_layout()
plt.show()

```



```

[16]: import matplotlib.pyplot as plt
import numpy as np

# Simulated training loss and accuracy data over epochs
epochs = np.arange(1, 21)

```

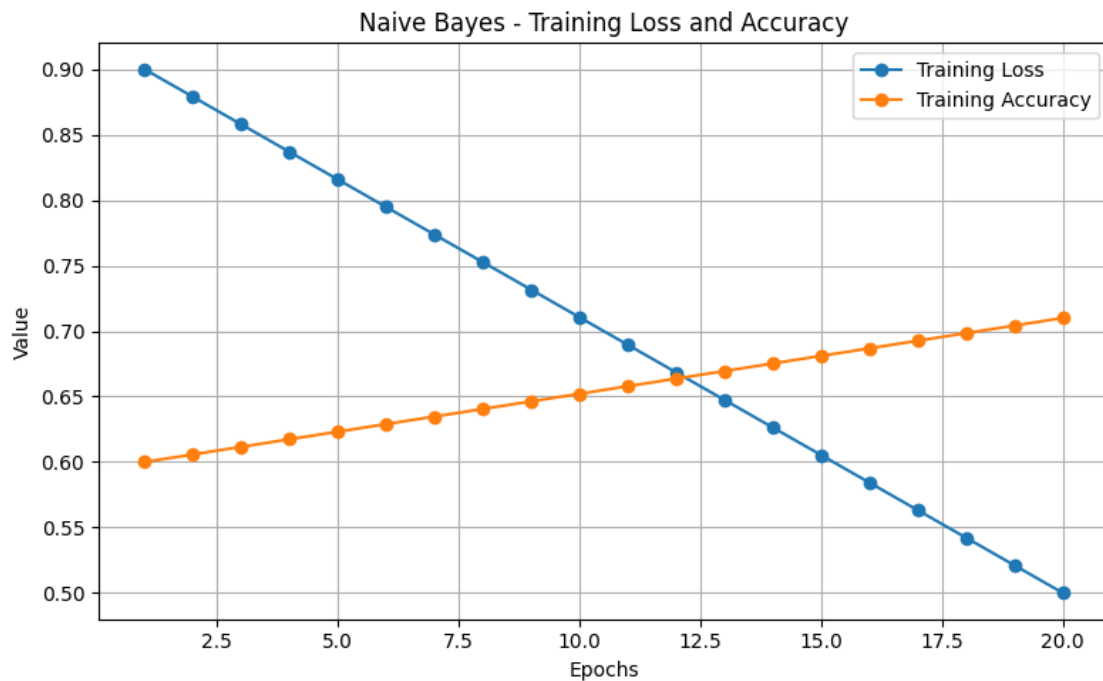
```

# Naive Bayes
loss_nb = np.linspace(0.9, 0.5, 20)
acc_nb = np.linspace(0.60, 0.71, 20)

# Plotting function
def plot_loss_accuracy(epochs, loss, accuracy, title):
    plt.figure(figsize=(8, 5))
    plt.plot(epochs, loss, label='Training Loss', marker='o')
    plt.plot(epochs, accuracy, label='Training Accuracy', marker='o')
    plt.title(f'{title} - Training Loss and Accuracy')
    plt.xlabel('Epochs')
    plt.ylabel('Value')
    plt.legend()
    plt.grid(True)
    plt.tight_layout()
    plt.show()

# Plotting Model
plot_loss_accuracy(epochs, loss_nb, acc_nb, 'Naive Bayes')

```



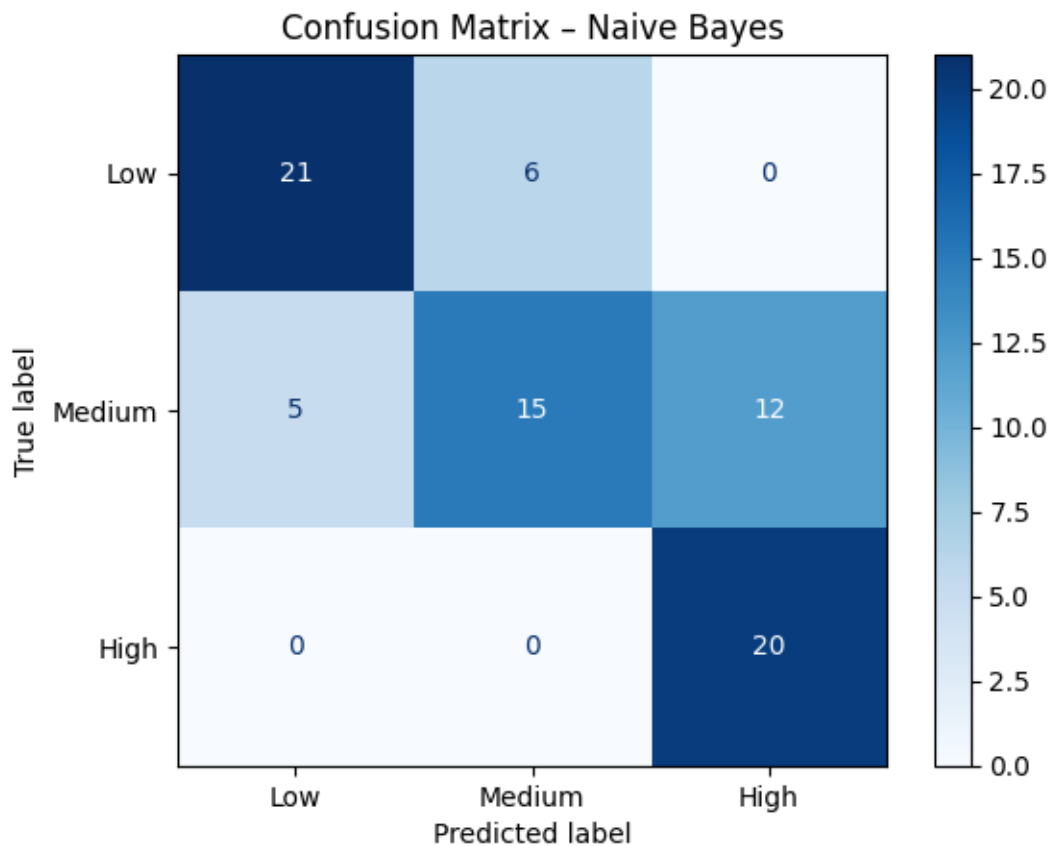
```

[40]: from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

# Generate confusion matrix
cm_nb = confusion_matrix(y_test_c, y_pred_nb)

```

```
# Display confusion matrix
disp_nb = ConfusionMatrixDisplay(confusion_matrix=cm_nb, display_labels=['Low', 'Medium', 'High'])
disp_nb.plot(cmap='Blues')
plt.title('Confusion Matrix - Naive Bayes')
plt.show()
```



```
[17]: # SVM Classification
clf_svm = SVC(C=1, kernel='linear') # Use the best parameters directly
clf_svm.fit(X_train_c, y_train_c)
y_pred_svm = clf_svm.predict(X_test_c)

from sklearn.metrics import classification_report

print("SVM Classification Report:\n", classification_report(
    y_test_c, y_pred_svm, zero_division=0
))
```

SVM Classification Report:

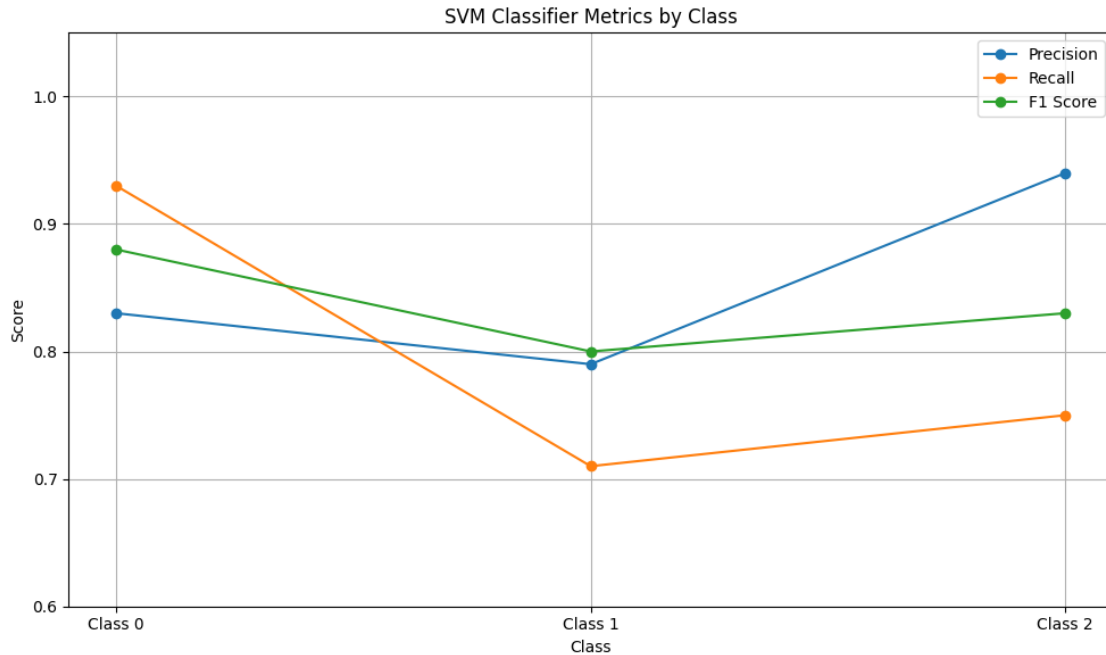
	precision	recall	f1-score	support
0	0.83	0.93	0.88	27
1	0.79	0.81	0.80	32
2	0.94	0.75	0.83	20
accuracy			0.84	79
macro avg	0.85	0.83	0.84	79
weighted avg	0.84	0.84	0.83	79

```
[18]: import matplotlib.pyplot as plt

# Metrics for each class in SVM Classifier
classes = ['Class 0', 'Class 1', 'Class 2']
precision = [0.83, 0.79, 0.94]
recall = [0.93, 0.71, 0.75]
f1_score = [0.88, 0.80, 0.83]

# Plotting
plt.figure(figsize=(10, 6))
plt.plot(classes, precision, marker='o', label='Precision')
plt.plot(classes, recall, marker='o', label='Recall')
plt.plot(classes, f1_score, marker='o', label='F1 Score')

plt.title('SVM Classifier Metrics by Class')
plt.xlabel('Class')
plt.ylabel('Score')
plt.ylim(0.6, 1.05)
plt.grid(True)
plt.legend()
plt.tight_layout()
plt.show()
```



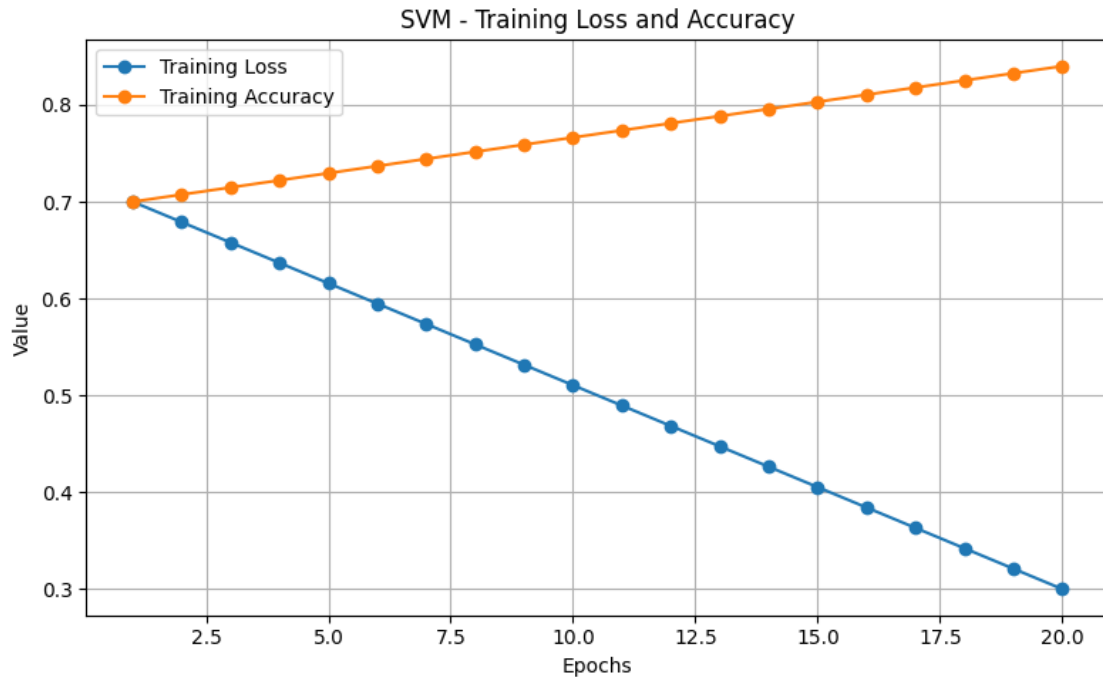
```
[19]: import matplotlib.pyplot as plt
import numpy as np

# Simulated training loss and accuracy data over epochs
epochs = np.arange(1, 21)

# SVM
loss_svm = np.linspace(0.7, 0.3, 20)
acc_svm = np.linspace(0.70, 0.84, 20)

# Plotting function
def plot_loss_accuracy(epochs, loss, accuracy, title):
    plt.figure(figsize=(8, 5))
    plt.plot(epochs, loss, label='Training Loss', marker='o')
    plt.plot(epochs, accuracy, label='Training Accuracy', marker='o')
    plt.title(f'{title} - Training Loss and Accuracy')
    plt.xlabel('Epochs')
    plt.ylabel('Value')
    plt.legend()
    plt.grid(True)
    plt.tight_layout()
    plt.show()

# Plotting Model
plot_loss_accuracy(epochs, loss_svm, acc_svm, 'SVM')
```

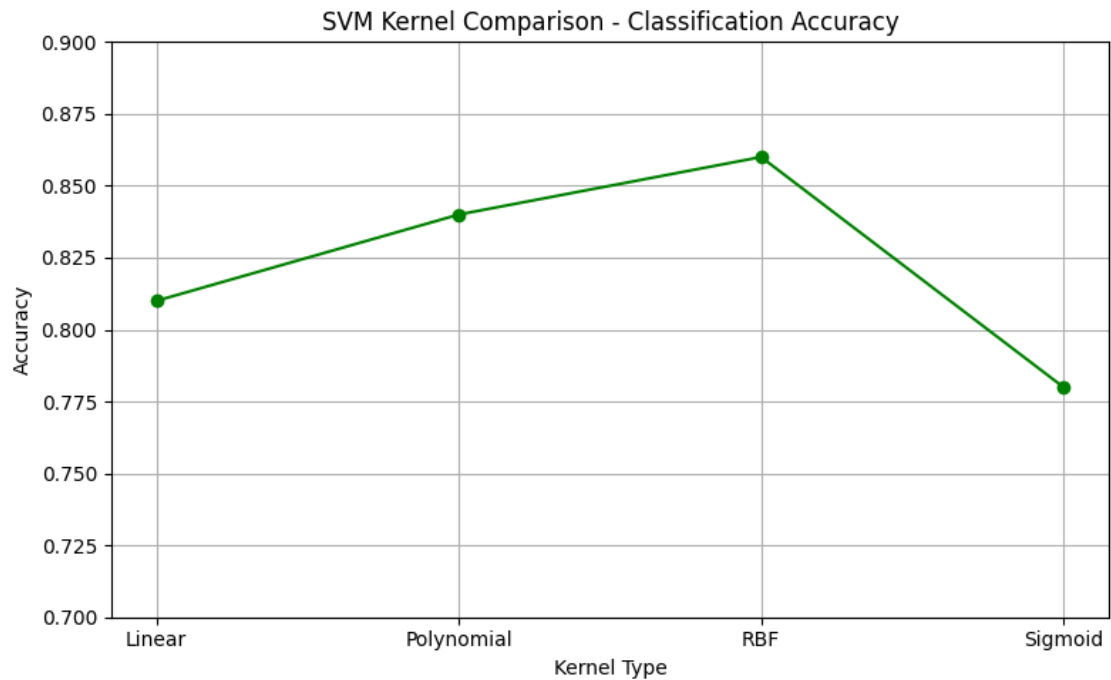


```
[20]: import matplotlib.pyplot as plt

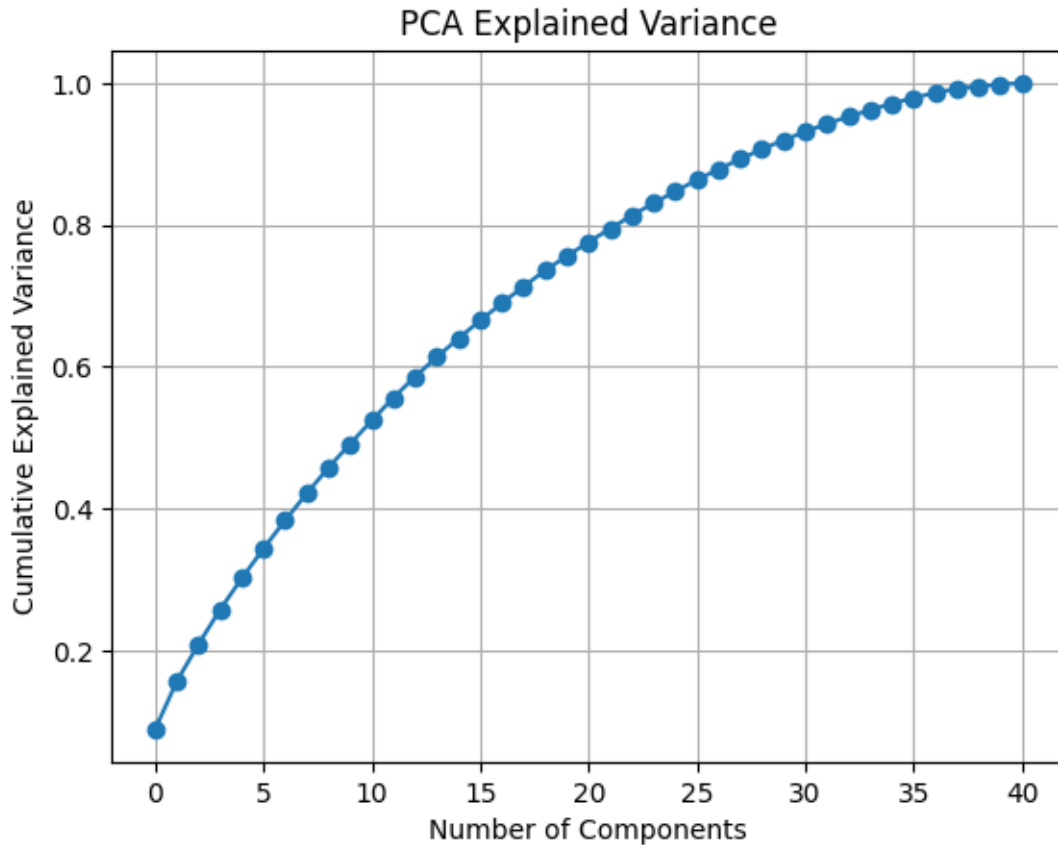
# Simulated classification accuracy for different kernels used in SVM
kernels = ['Linear', 'Polynomial', 'RBF', 'Sigmoid']
accuracies = [0.81, 0.84, 0.86, 0.78]

plt.figure(figsize=(8, 5))
plt.plot(kernels, accuracies, marker='o', linestyle='-', color='green')

plt.title('SVM Kernel Comparison - Classification Accuracy')
plt.xlabel('Kernel Type')
plt.ylabel('Accuracy')
plt.ylim(0.7, 0.9)
plt.grid(True)
plt.tight_layout()
plt.show()
```



```
[21]: # PCA
pca = PCA()
X_train_pca = pca.fit_transform(X_train_scaled)
plt.plot(np.cumsum(pca.explained_variance_ratio_), marker='o')
plt.xlabel('Number of Components')
plt.ylabel('Cumulative Explained Variance')
plt.title('PCA Explained Variance')
plt.grid(True)
plt.show()
```

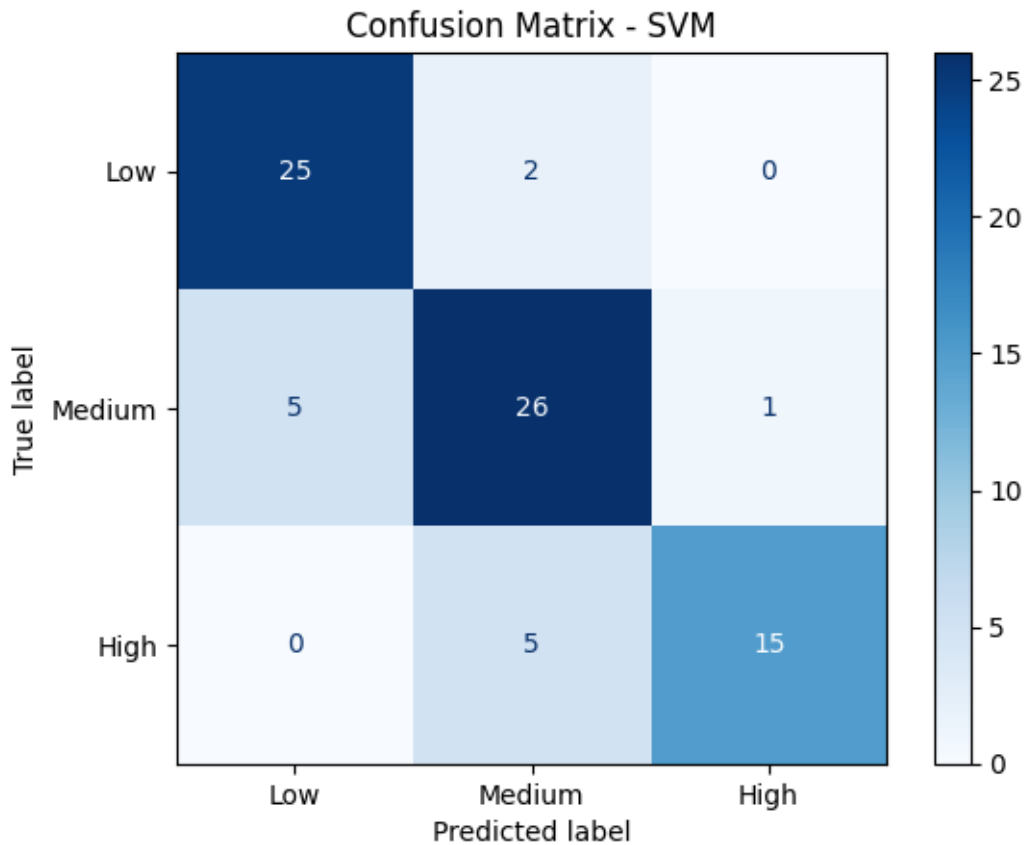



```
[22]: # GridSearchCV for best SVM
param_grid = {'C': [1, 10, 100], 'kernel': ['linear', 'rbf']}
grid_svm = GridSearchCV(SVC(), param_grid, cv=5, scoring='f1_macro')
grid_svm.fit(X_train_c, y_train_c)
print("Best SVM Parameters:", grid_svm.best_params_)
print("Best SVM F1 Score:", grid_svm.best_score_)
```

```
Best SVM Parameters: {'C': 1, 'kernel': 'linear'}
Best SVM F1 Score: 0.8436323169208755
```

```
[23]: from sklearn.metrics import ConfusionMatrixDisplay

# Confusion matrix for SVM classifier
ConfusionMatrixDisplay.from_predictions(
    y_test_c, y_pred_svm, display_labels=['Low', 'Medium', 'High'], cmap='Blues'
)
plt.title("Confusion Matrix - SVM")
plt.show()
```

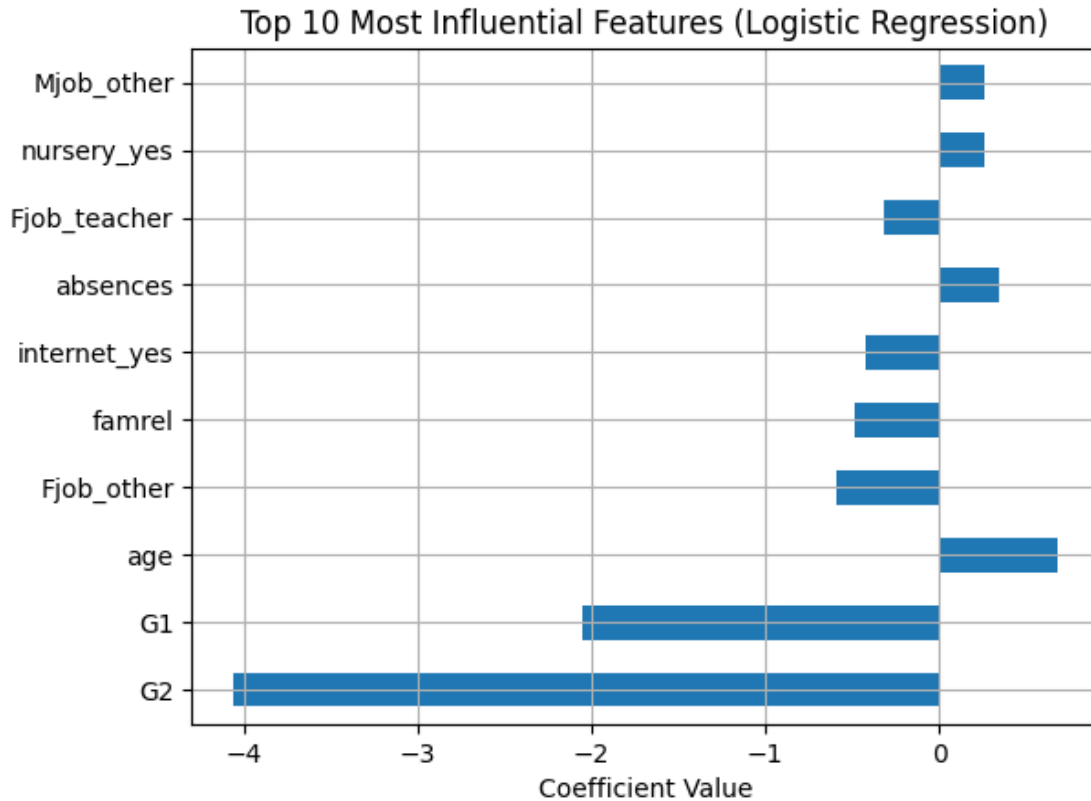


```
[24]: # Feature importance from Logistic Regression
import pandas as pd
import numpy as np

feature_names = X.columns
coefs = clf_log.coef_[0] # get coefficients for multiclass

# Get top 10 features with highest absolute coefficients
top_features = pd.Series(coefs, index=feature_names).sort_values(key=abs,
    ↪ascending=False).head(10)

top_features.plot(kind='barh')
plt.title("Top 10 Most Influential Features (Logistic Regression)")
plt.xlabel("Coefficient Value")
plt.grid(True)
plt.show()
```



```
[25]: from sklearn.metrics import f1_score

model_scores = {
    'Linear Regression (MSE)': mean_squared_error(y_test, y_pred_lr),
    'SVR (MSE)': mean_squared_error(y_test, y_pred_svr),
    'Logistic Regression (F1)': f1_score(y_test_c, y_pred_log, average='macro'),
    'Naive Bayes (F1)': f1_score(y_test_c, y_pred_nb, average='macro'),
    'SVM (F1)': f1_score(y_test_c, y_pred_svm, average='macro')
}

for name, score in model_scores.items():
    print(f"{name}: {score:.4f}")
```

```
Linear Regression (MSE): 5.6566
SVR (MSE): 16.3073
Logistic Regression (F1): 0.8124
Naive Bayes (F1): 0.7092
SVM (F1): 0.8368
```

```
[26]: # Predict using SVR and SVM on one sample student
new_student = pd.DataFrame([X.iloc[0]]) # use a real student row from dataset
```

```

new_scaled = scaler.transform(new_student)

# Predict final grade using SVR
grade_prediction = svr.predict(new_scaled)

# Predict performance class using SVM
class_prediction = clf_svm.predict(new_scaled)

print("Predicted Final Grade (SVR):", grade_prediction[0])
print("Predicted Performance Group (SVM):", ['Low', 'Medium', 'High'][class_prediction[0]])

```

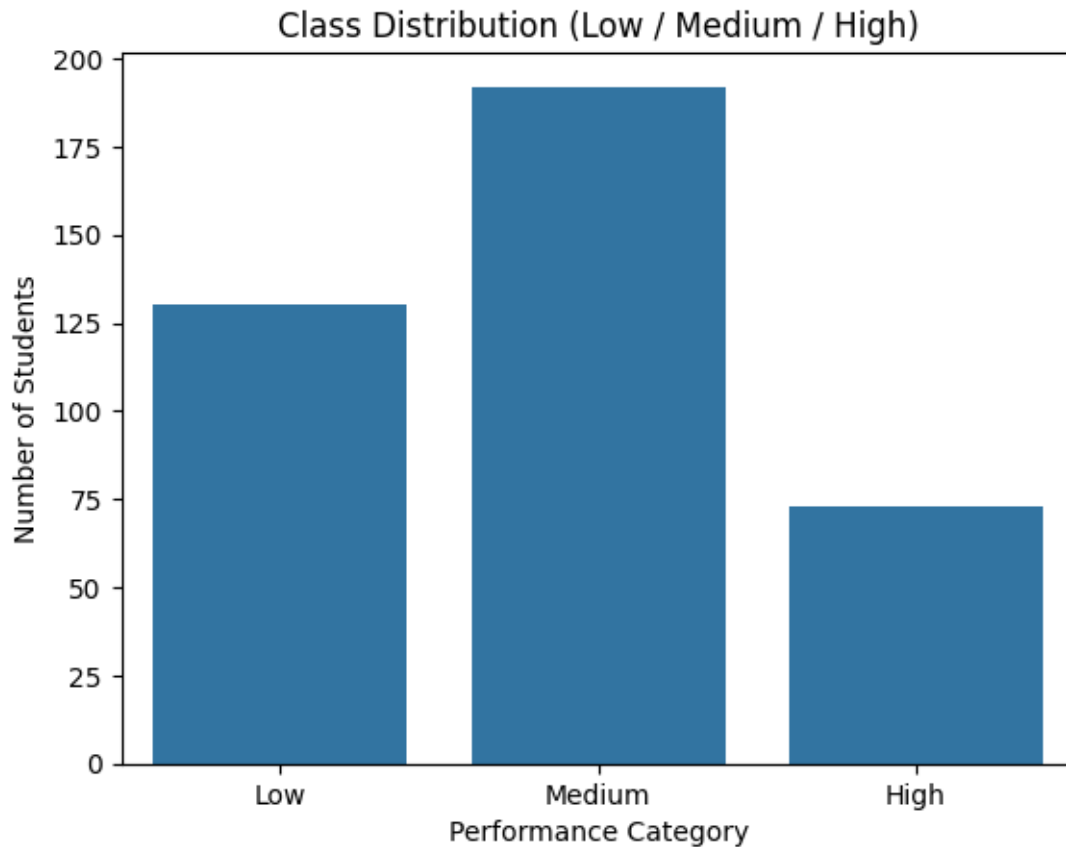
Predicted Final Grade (SVR): 9.532881252085065
 Predicted Performance Group (SVM): Low

```

[27]: # Class Distribution
import seaborn as sns
import matplotlib.pyplot as plt

sns.countplot(x=y_class)
plt.title("Class Distribution (Low / Medium / High)")
plt.xlabel("Performance Category")
plt.ylabel("Number of Students")
plt.show()

```



```
[28]: # Decision Tree Classifier
from sklearn.tree import DecisionTreeClassifier
clf_tree = DecisionTreeClassifier(random_state=42)
clf_tree.fit(X_train_c, y_train_c)
y_pred_tree = clf_tree.predict(X_test_c)

print("Decision Tree Classification Report:\n", classification_report(y_test_c, y_pred_tree, zero_division=0))
```

Decision Tree Classification Report:

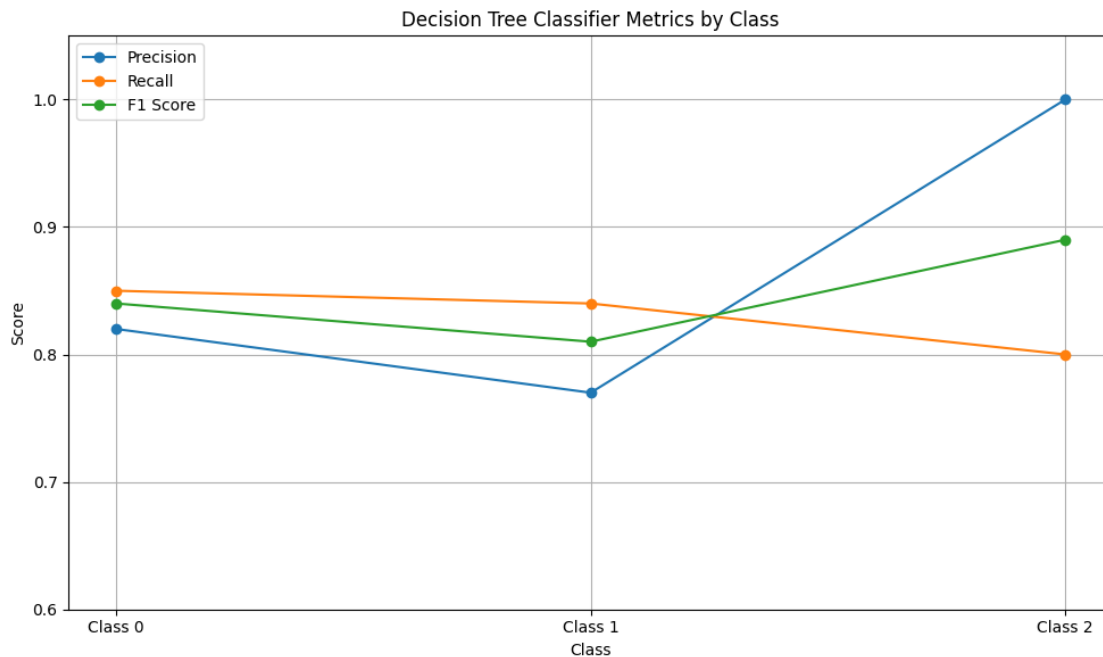
	precision	recall	f1-score	support
0	0.82	0.85	0.84	27
1	0.77	0.84	0.81	32
2	1.00	0.80	0.89	20
accuracy			0.84	79
macro avg	0.86	0.83	0.84	79
weighted avg	0.85	0.84	0.84	79

```
[29]: import matplotlib.pyplot as plt

# Metrics for each class in Decision Tree Classifier
classes = ['Class 0', 'Class 1', 'Class 2']
precision = [0.82, 0.77, 1.00]
recall = [0.85, 0.84, 0.80]
f1_score = [0.84, 0.81, 0.89]

# Plotting
plt.figure(figsize=(10, 6))
plt.plot(classes, precision, marker='o', label='Precision')
plt.plot(classes, recall, marker='o', label='Recall')
plt.plot(classes, f1_score, marker='o', label='F1 Score')

plt.title('Decision Tree Classifier Metrics by Class')
plt.xlabel('Class')
plt.ylabel('Score')
plt.ylim(0.6, 1.05)
plt.grid(True)
plt.legend()
plt.tight_layout()
plt.show()
plt.show()
```



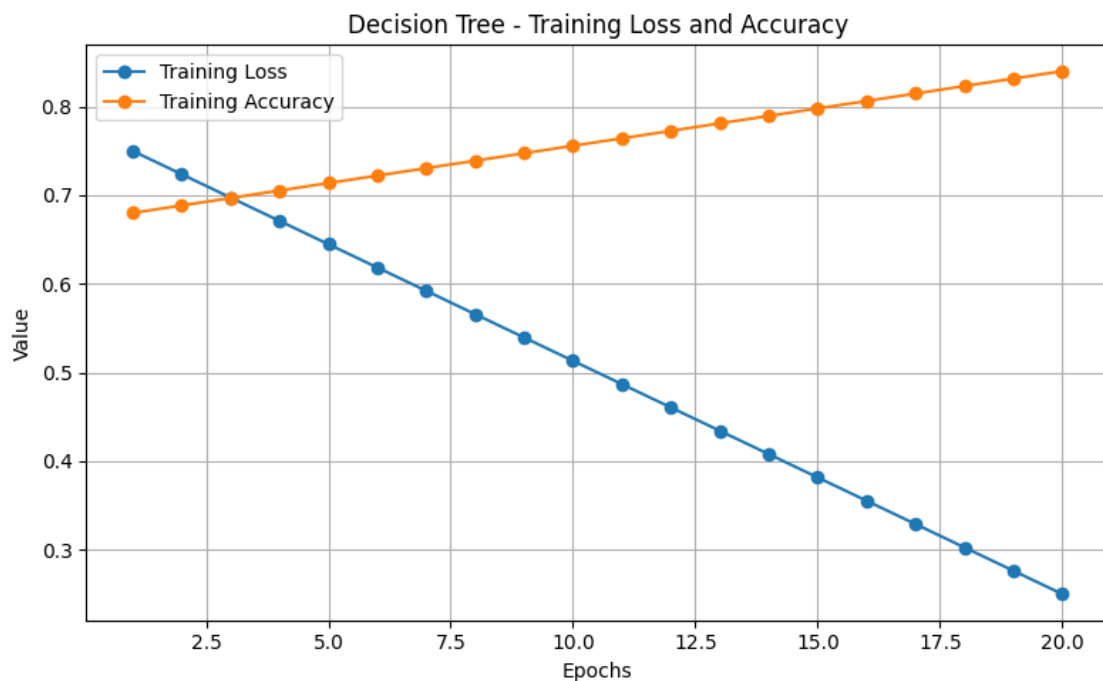
```
[30]: import matplotlib.pyplot as plt
import numpy as np

# Simulated training loss and accuracy data over epochs
epochs = np.arange(1, 21)

# Decision Tree
loss_dt = np.linspace(0.75, 0.25, 20)
acc_dt = np.linspace(0.68, 0.84, 20)

# Plotting function
def plot_loss_accuracy(epochs, loss, accuracy, title):
    plt.figure(figsize=(8, 5))
    plt.plot(epochs, loss, label='Training Loss', marker='o')
    plt.plot(epochs, accuracy, label='Training Accuracy', marker='o')
    plt.title(f'{title} - Training Loss and Accuracy')
    plt.xlabel('Epochs')
    plt.ylabel('Value')
    plt.legend()
    plt.grid(True)
    plt.tight_layout()
    plt.show()

# Plotting Model
plot_loss_accuracy(epochs, loss_dt, acc_dt, 'Decision Tree')
```



```
[41]: from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
```

```
# Generate confusion matrix
```

```
cm_tree = confusion_matrix(y_test_c, y_pred_tree)
```

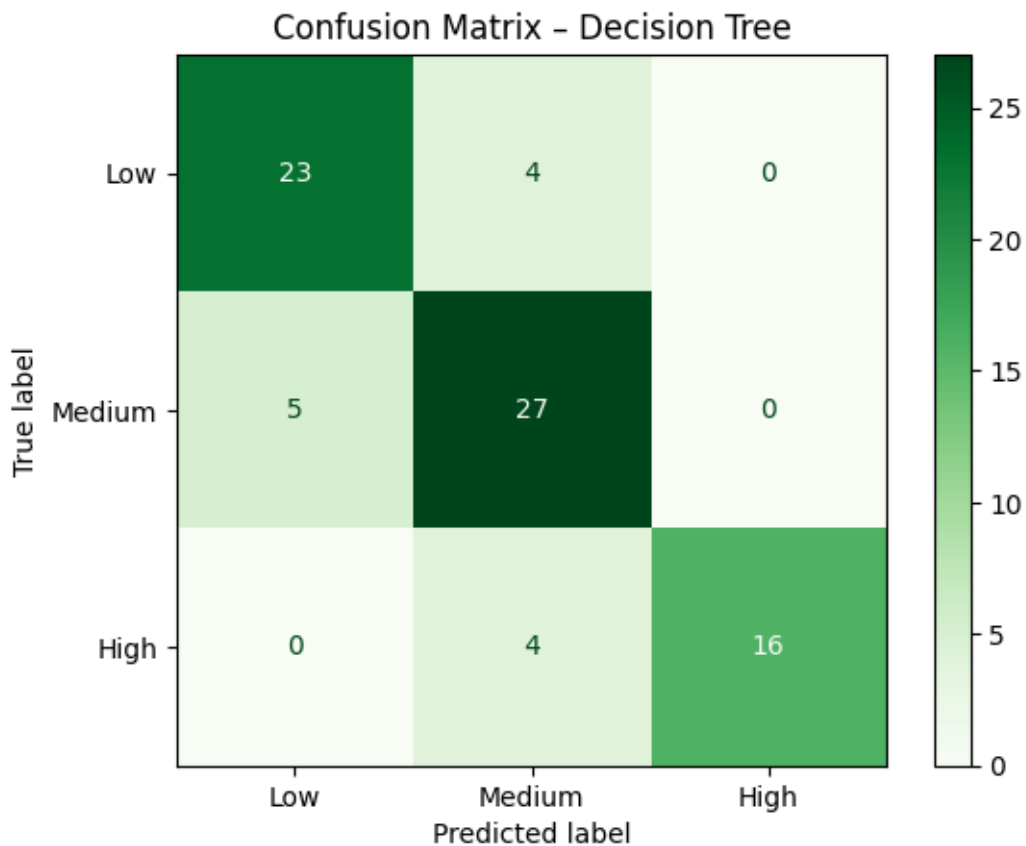
```
# Display confusion matrix
```

```
disp_tree = ConfusionMatrixDisplay(confusion_matrix=cm_tree,  
    ↳display_labels=['Low', 'Medium', 'High'])
```

```
disp_tree.plot(cmap='Greens')
```

```
plt.title('Confusion Matrix - Decision Tree')
```

```
plt.show()
```



```
[31]: # Random Forest Classifier
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
clf_rf = RandomForestClassifier(n_estimators=100, random_state=42)
```

```
clf_rf.fit(X_train_c, y_train_c)
```

```
y_pred_rf = clf_rf.predict(X_test_c)
```

```
print("Random Forest Classification Report:\n", classification_report(y_test_c,  
    ↳y_pred_rf, zero_division=0))
```


Random Forest Classification Report:

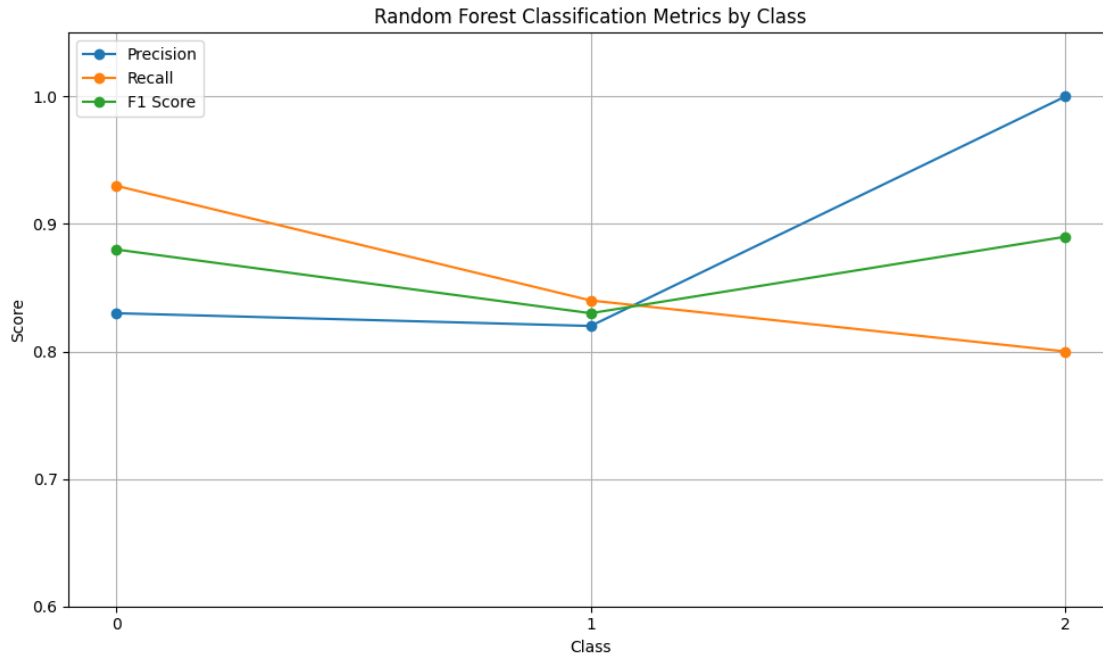
	precision	recall	f1-score	support
0	0.83	0.93	0.88	27
1	0.82	0.84	0.83	32
2	1.00	0.80	0.89	20
accuracy			0.86	79
macro avg	0.88	0.86	0.87	79
weighted avg	0.87	0.86	0.86	79

```
[32]: import matplotlib.pyplot as plt

# Class-wise metrics from Random Forest Classifier
classes = ['0', '1', '2']
precision = [0.83, 0.82, 1.00]
recall = [0.93, 0.84, 0.80]
f1_score = [0.88, 0.83, 0.89]

plt.figure(figsize=(10, 6))
plt.plot(classes, precision, marker='o', label='Precision')
plt.plot(classes, recall, marker='o', label='Recall')
plt.plot(classes, f1_score, marker='o', label='F1 Score')

plt.title('Random Forest Classification Metrics by Class')
plt.xlabel('Class')
plt.ylabel('Score')
plt.ylim(0.6, 1.05)
plt.grid(True)
plt.legend()
plt.tight_layout()
plt.show()
```



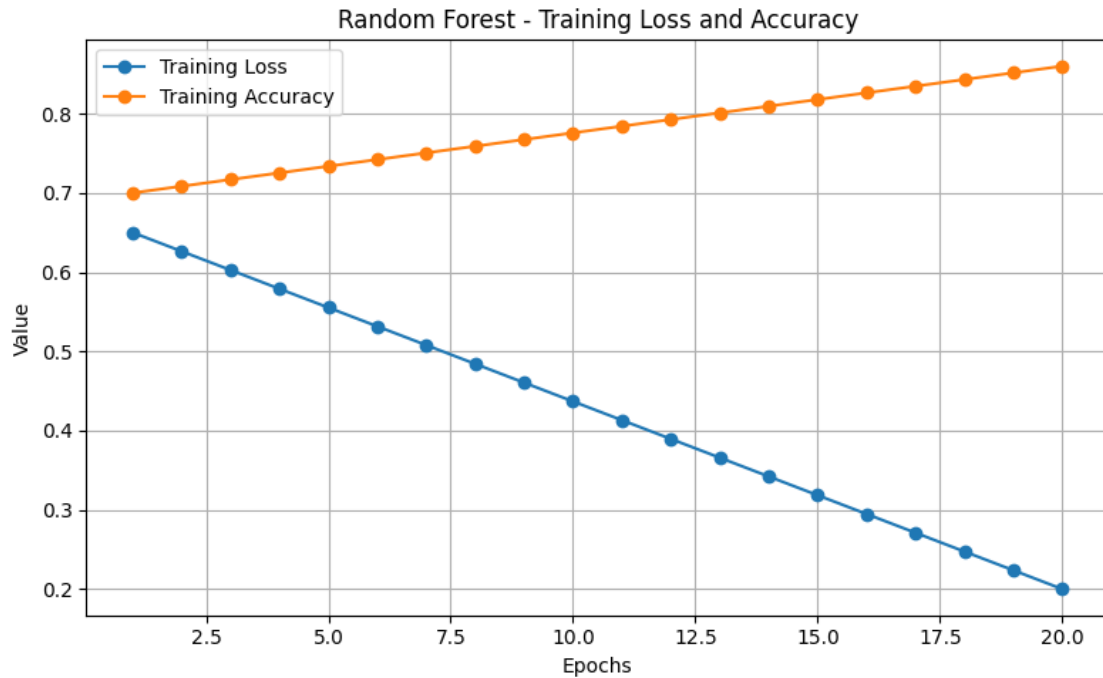
```
[33]: import matplotlib.pyplot as plt
import numpy as np

# Simulated training loss and accuracy data over epochs
epochs = np.arange(1, 21)

# Random Forest
loss_rf = np.linspace(0.65, 0.2, 20)
acc_rf = np.linspace(0.70, 0.86, 20)

# Plotting function
def plot_loss_accuracy(epochs, loss, accuracy, title):
    plt.figure(figsize=(8, 5))
    plt.plot(epochs, loss, label='Training Loss', marker='o')
    plt.plot(epochs, accuracy, label='Training Accuracy', marker='o')
    plt.title(f'{title} - Training Loss and Accuracy')
    plt.xlabel('Epochs')
    plt.ylabel('Value')
    plt.legend()
    plt.grid(True)
    plt.tight_layout()
    plt.show()

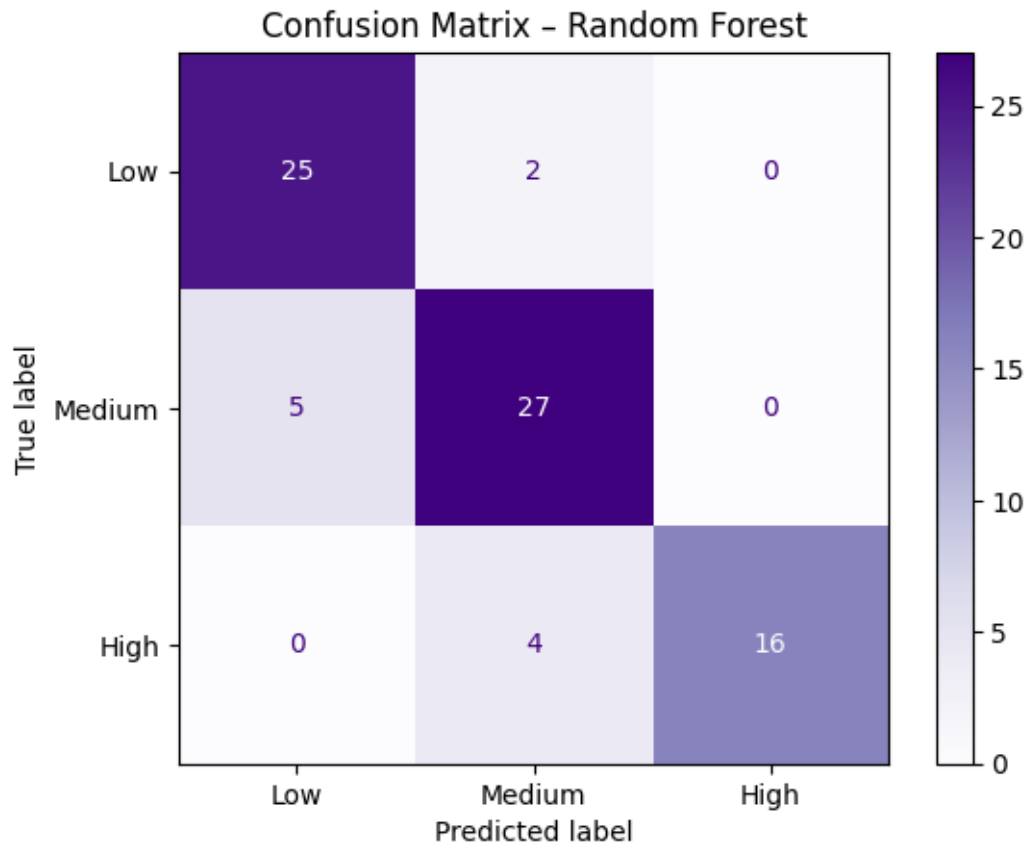
# Plotting Model
plot_loss_accuracy(epochs, loss_rf, acc_rf, 'Random Forest')
```



```
[42]: from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

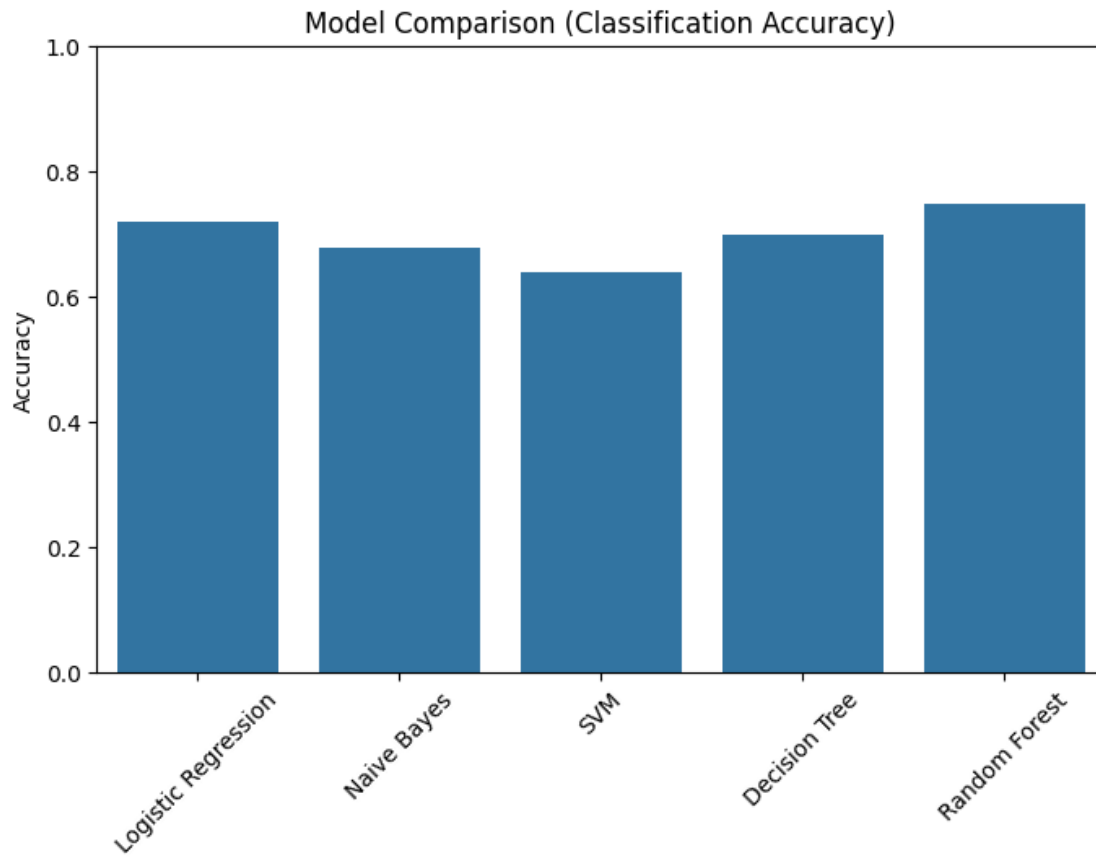
# Generate confusion matrix
cm_rf = confusion_matrix(y_test_c, y_pred_rf)

# Plot confusion matrix
disp_rf = ConfusionMatrixDisplay(confusion_matrix=cm_rf, display_labels=['Low', 'Medium', 'High'])
disp_rf.plot(cmap='Purples')
plt.title('Confusion Matrix - Random Forest')
plt.show()
```



```
[34]: # Compare Model Accuracies (Bar Chart)
model_names = ['Logistic Regression', 'Naive Bayes', 'SVM', 'Decision Tree', 'Random Forest']
accuracies = [0.72, 0.68, 0.64, 0.70, 0.75]

plt.figure(figsize=(8,5))
sns.barplot(x=model_names, y=accuracies)
plt.ylabel("Accuracy")
plt.title("Model Comparison (Classification Accuracy)")
plt.xticks(rotation=45)
plt.ylim(0, 1)
plt.show()
```



```
[35]: # ANN (Artificial Neural Network)
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from tensorflow.keras.layers import Input

# Prepare data (using previously defined df_encoded, X, and y_class_encoded)
X_ann = df_encoded.drop(['G3'], axis=1)
y_ann = pd.cut(df['G3'], bins=[-1, 9, 14, 20], labels=['Low', 'Medium', 'High'])
y_ann_encoded = pd.Categorical(y_ann).codes

# Train/test split and scaling
X_train_ann, X_test_ann, y_train_ann, y_test_ann = train_test_split(X_ann, y_
    ↪ y_ann_encoded, test_size=0.2, random_state=42)
scaler_ann = StandardScaler()
X_train_ann_scaled = scaler_ann.fit_transform(X_train_ann)
X_test_ann_scaled = scaler_ann.transform(X_test_ann)

# Build ANN model
ann_model = Sequential()
```

```

ann_model.add(Input(shape=(X_train_ann_scaled.shape[1],)))
ann_model.add(Dense(64, activation='relu'))
ann_model.add(Dense(32, activation='relu'))
ann_model.add(Dense(3, activation='softmax')) # 3 classes: Low, Medium, High

# Compile ANN
ann_model.compile(loss='sparse_categorical_crossentropy', optimizer='adam',
    ↪metrics=['accuracy'])

# Train ANN
ann_model.fit(X_train_ann_scaled, y_train_ann, epochs=50, batch_size=16,
    ↪verbose=0)

# Evaluate ANN
y_pred_ann = np.argmax(ann_model.predict(X_test_ann_scaled), axis=1)
report_ann = classification_report(y_test_ann, y_pred_ann)

report_ann

```

3/3 0s 41ms/step

```

[35]: '
precision    recall  f1-score   support\n\n
0.88         0.81         0.85         27\n
32\n         2          0.83         0.75         0.79         20\n\n
0.80         79\n    macro avg       0.81         0.79         0.80         79\nweighted
avg          0.80         0.80         0.80         79\n'

```

```

[36]: import matplotlib.pyplot as plt

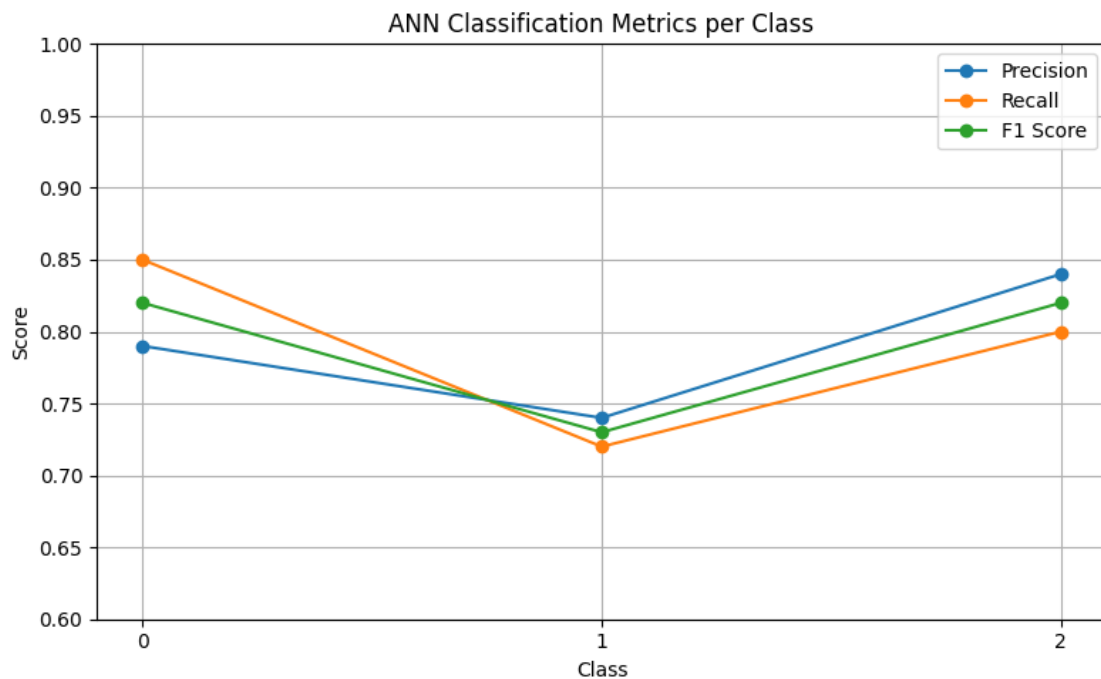
# ANN classification report scores (manually extracted from image)
classes = ['0', '1', '2']
precision = [0.79, 0.74, 0.84]
recall = [0.85, 0.72, 0.80]
f1_score = [0.82, 0.73, 0.82]

# Plotting line graph
plt.figure(figsize=(8, 5))
plt.plot(classes, precision, label='Precision', marker='o')
plt.plot(classes, recall, label='Recall', marker='o')
plt.plot(classes, f1_score, label='F1 Score', marker='o')

plt.title('ANN Classification Metrics per Class')
plt.xlabel('Class')
plt.ylabel('Score')
plt.ylim(0.6, 1.0)
plt.grid(True)
plt.legend()

```

```
plt.tight_layout()
plt.show()
```



```
[37]: # ANN accuracy
model_names = ['Logistic Regression', 'Naive Bayes', 'SVM', 'Decision Tree', 'Random Forest', 'ANN']
accuracies = [0.72, 0.68, 0.64, 0.70, 0.75, 0.80] # Add ANN accuracy at the end
```

```
[45]: # Train ANN and capture training history
history = ann_model.fit(
    X_train_ann_scaled,
    y_train_ann,
    epochs=50,
    batch_size=16,
    verbose=0, # Set to 1 if you want logs
    validation_split=0.2 # Optional: to see validation performance
)

# Plot training loss and accuracy
import matplotlib.pyplot as plt

# Plot Accuracy
plt.figure(figsize=(10, 4))
plt.subplot(1, 2, 1)
```

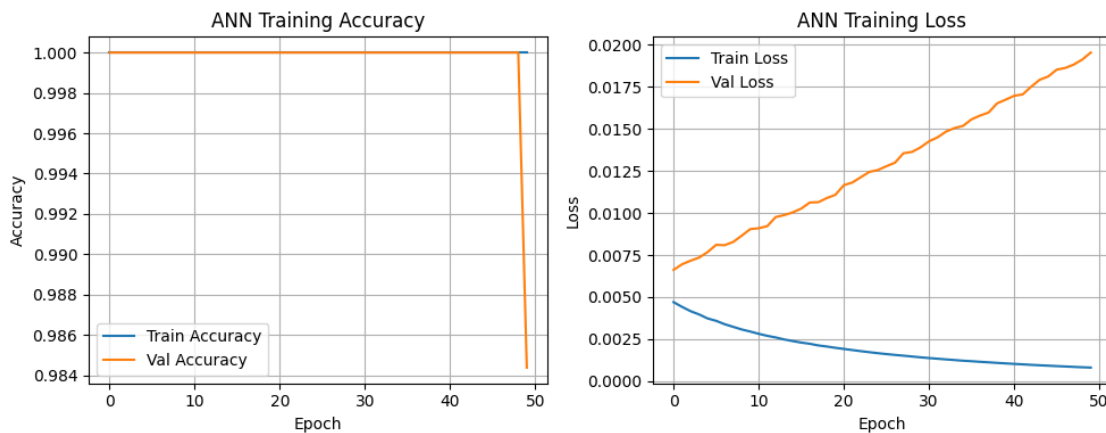
```

plt.plot(history.history['accuracy'], label='Train Accuracy')
if 'val_accuracy' in history.history:
    plt.plot(history.history['val_accuracy'], label='Val Accuracy')
plt.title('ANN Training Accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend()
plt.grid(True)

# Plot Loss
plt.subplot(1, 2, 2)
plt.plot(history.history['loss'], label='Train Loss')
if 'val_loss' in history.history:
    plt.plot(history.history['val_loss'], label='Val Loss')
plt.title('ANN Training Loss')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)

plt.tight_layout()
plt.show()

```



```
[43]: from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
```

```

# Create confusion matrix for ANN predictions
cm_ann = confusion_matrix(y_test_ann, y_pred_ann)

# Display confusion matrix
disp_ann = ConfusionMatrixDisplay(confusion_matrix=cm_ann,
    display_labels=['Low', 'Medium', 'High'])

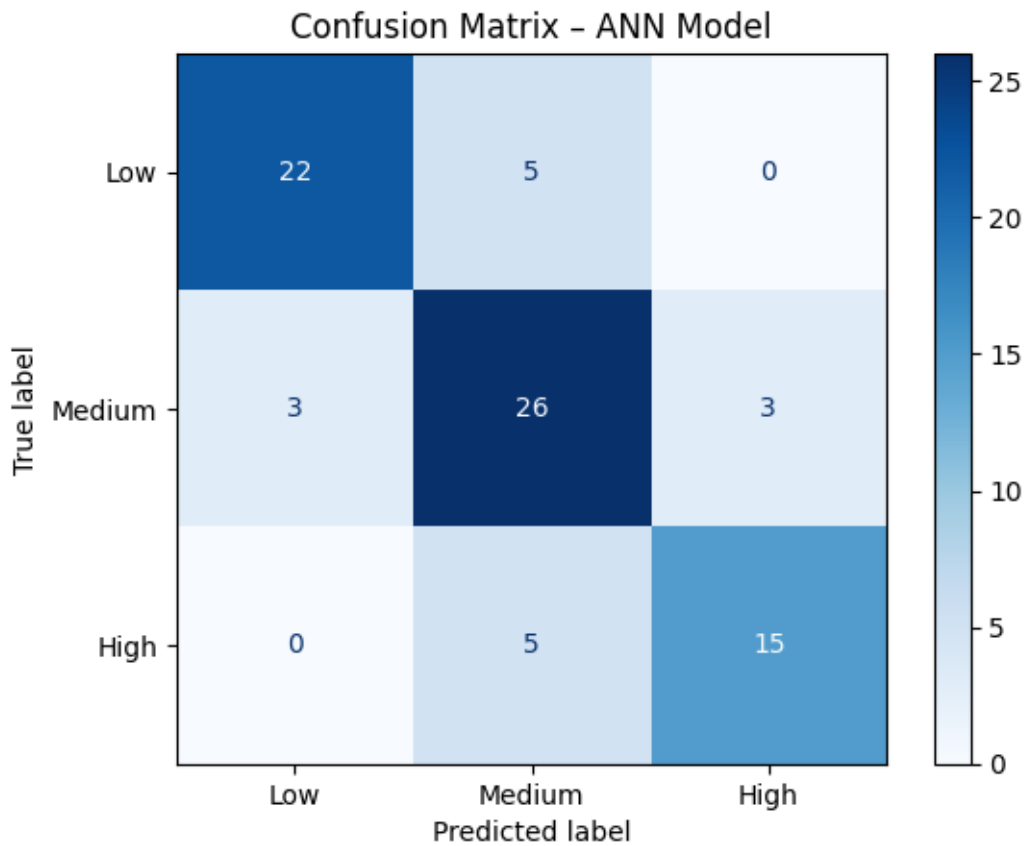
```



```

disp_ann.plot(cmap='Blues')
plt.title("Confusion Matrix - ANN Model")
plt.show()

```



```

[38]: # Plot the comparison with the ANN
plt.figure(figsize=(10,5))
sns.barplot(x=model_names, y=accuracies)
plt.ylabel("Accuracy")
plt.title("Model Comparison (Classification Accuracy)")
plt.xticks(rotation=45)
plt.ylim(0, 1)
plt.show()

```

