

Sri Sivasubramaniya Nadar College of Engineering, Chennai

(An autonomous Institution affiliated to Anna University)

Machine Learning Algorithms Laboratory

Experiment 3: Email Spam or Ham Classification using Naïve Bayes, KNN, and SVM

Subject Code & Name: ICS1512 & Machine Learning Algorithms Laboratory
Degree & Branch: M. Tech (Integrated) Computer Science & Engineering
Semester: V
Academic Year: 2025-2026 (Odd)
Batch: 2023-2028

Submitted by:
Sudharshan Vijayaragavan
Reg No.: 3122237001054

1 Objective

The objective of this experiment is to classify emails as either spam or ham using three distinct classification algorithms: Naïve Bayes, K-Nearest Neighbors (KNN), and Support Vector Machine (SVM). The performance of these models will be evaluated using standard accuracy metrics, confusion matrices, ROC curves, and K-Fold cross-validation.

2 Dataset

The dataset used for this experiment is the Spambase dataset, which is publicly available on Kaggle. It contains a collection of extracted features from emails, pre-labeled as either spam or ham, making it suitable for classification tasks.

3 Implementation Steps

The following steps were performed for the implementation:

1. **Load and Preprocess the Dataset:** The dataset was loaded, and initial preprocessing steps were performed, including handling any missing values and normalizing the features using 'StandardScaler'.
2. **Exploratory Data Analysis (EDA):** An analysis of the dataset was conducted to understand the class balance and the distribution of features.
3. **Data Splitting:** The dataset was split into training and testing sets to evaluate model performance on unseen data.
4. **Model Training:** Models were trained using different variants for each algorithm:
 - **Naïve Bayes:** Gaussian, Multinomial, and Bernoulli variants were trained.
 - **K-Nearest Neighbors:** The 'k' value was varied (k=3, 5, 7), and different algorithms ('KDTree', 'BallTree') were compared.
 - **Support Vector Machine:** Linear, Polynomial, RBF, and Sigmoid kernels were evaluated.
5. **Evaluation and Visualization:** The performance of each model was evaluated using a classification report, confusion matrix, and ROC curve.
6. **K-Fold Cross-Validation:** K-Fold cross-validation with $K = 5$ was performed to get a more robust estimate of each model's performance.
7. **Comparison:** The results of all models were compared and observations were recorded.

4 Python Code

The following is the Python script used to perform the experiment.

```

1 # -*- coding: utf-8 -*-
2 """ML_ASSGN_3.ipynb
3
4 Automatically generated by Colab.
5
6 Original file is located at
7     https://colab.research.google.com/drive/19
8         s7wAPS21VEE1CZ1Ik_g3zh8Tu4y0r10
9
10 Name : Sudharshan Vijayaragavan
11
12 Reg No. : 3122237001054
13
14 Email Spam or Ham Classification using Naïve Bayes, KNN, and
15 SVM
16
17 Configuration Section (User-Editable Block)
18
19 1) Load and Preprocess Dataset
20
21 """
22
23 import pandas as pd
24 import numpy as np
25 import matplotlib.pyplot as plt
26 import seaborn as sns
27 import time
28
29 from sklearn.model_selection import train_test_split, StratifiedKFold,
30     cross_val_score
31
32 from sklearn.preprocessing import StandardScaler
33
34 from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB
35
36 from sklearn.neighbors import KNeighborsClassifier
37
38 from sklearn.svm import SVC
39
40 from sklearn.metrics import classification_report, confusion_matrix,
41     roc_auc_score, roc_curve
42
43 from sklearn.metrics import accuracy_score, precision_score,
44     recall_score, f1_score
45
46 from google.colab import drive
47
48 # Config
49 DATA_PATH = "/content/drive/MyDrive/spambase_csv.csv"
50 TARGET_COLUMN = 'class'
51 TEST_SIZE = 0.2
52 RANDOM_STATE = 42
53 KFOLD_SPLITS = 5
54
55 # Mount drive and load CSV
56 drive.mount('/content/drive')
57 df = pd.read_csv(DATA_PATH)
58 print("Columns:", df.columns.tolist())
59 print("Missing values:\n", df.isnull().sum())
60
61 # Handle missing and set target type
62 df = df.dropna()
63 df[TARGET_COLUMN] = df[TARGET_COLUMN].astype('category')
64
65 # Feature/Target split

```

```

54 X_raw = df.drop(columns=[TARGET_COLUMN])
55 y = df[TARGET_COLUMN]
56
57 # Normalize
58 scaler = StandardScaler()
59 X_scaled = scaler.fit_transform(X_raw)
60
61 # Train-test split (both raw and scaled)
62 X_raw_train, X_raw_test, y_train, y_test = train_test_split(
63     X_raw, y, test_size=TEST_SIZE, random_state=RANDOM_STATE, stratify=y
64 )
65 X_scaled_train, X_scaled_test, _, _ = train_test_split(
66     X_scaled, y, test_size=TEST_SIZE, random_state=RANDOM_STATE,
67     stratify=y
68 )
69 """2. Perform EDA (class balance, feature distributions)"""
70
71 plt.figure(figsize=(5,4))
72 sns.countplot(x=y)
73 plt.title("Class Balance (0=Ham, 1=Spam)")
74 plt.show()
75
76 X_df = pd.DataFrame(X_scaled, columns=X_raw.columns)
77 X_df.hist(figsize=(18,14), bins=20, edgecolor='black')
78 plt.suptitle("Normalized Feature Distributions", fontsize=18)
79 plt.show()
80
81 """3) Define Models"""
82
83 models = {
84     "GaussianNB": GaussianNB(),
85     "MultinomialNB": MultinomialNB(),
86     "BernoulliNB": BernoulliNB(),
87     "KNN_k=3": KNeighborsClassifier(n_neighbors=3, algorithm='auto'),
88     "KNN_k=5_KDTree": KNeighborsClassifier(n_neighbors=5, algorithm='
        kd_tree'),
89     "KNN_k=7_BallTree": KNeighborsClassifier(n_neighbors=7, algorithm='
        ball_tree'),
90     "SVM_Linear": SVC(kernel='linear', C=1.0, probability=True),
91     "SVM_Polynomial": SVC(kernel='poly', degree=3, C=1.0, gamma='scale',
        probability=True),
92     "SVM_RBF": SVC(kernel='rbf', C=1.0, gamma='scale', probability=True)
93     ,
94     "SVM_Sigmoid": SVC(kernel='sigmoid', C=1.0, gamma='scale',
        probability=True),
95 }
96
97 """4) Train Models (NB, KNN, SVM)"""
98
99 # Clear results (avoid duplicates if re-running)
100 nb_results, knn_k_results, knn_tree_results, svm_results, cv_results =
101     [], [], [], [], []
102
103 def evaluate_model(name, model, X_test, y_test):
104     y_pred = model.predict(X_test)
105     print(classification_report(y_test, y_pred, zero_division=0))
106     cm = confusion_matrix(y_test, y_pred)

```

```

105     sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
106     plt.title(f"{name} - Confusion Matrix")
107     plt.xlabel("Predicted")
108     plt.ylabel("Actual")
109     plt.show()
110     try:
111         if hasattr(model, "predict_proba"):
112             y_scores = model.predict_proba(X_test)[: , 1]
113         else:
114             y_scores = model.decision_function(X_test)
115         fpr, tpr, _ = roc_curve(y_test.cat.codes, y_scores)
116         auc_score = roc_auc_score(y_test.cat.codes, y_scores)
117         plt.plot(fpr, tpr, label=f"{name} (AUC = {auc_score:.2f})")
118         plt.plot([0, 1], [0, 1], linestyle='--', color='gray')
119         plt.title(f"ROC Curve - {name}")
120         plt.xlabel("False Positive Rate")
121         plt.ylabel("True Positive Rate")
122         plt.legend()
123         plt.show()
124     except Exception as e:
125         print(f"ROC Curve skipped for {name}: {e}")
126
127 def evaluate_and_store(name, model, X_train, y_train, X_test, y_test,
128 category):
129     start_time = time.time()
130     model.fit(X_train, y_train)
131     train_time = time.time() - start_time
132     y_pred = model.predict(X_test)
133     acc = accuracy_score(y_test, y_pred)
134     prec = precision_score(y_test, y_pred, zero_division=0)
135     rec = recall_score(y_test, y_pred, zero_division=0)
136     f1 = f1_score(y_test, y_pred, zero_division=0)
137     if category == "NB":
138         nb_results.append([name, acc, prec, rec, f1, train_time])
139     elif category == "KNN_k":
140         knn_k_results.append([name, acc, prec, rec, f1])
141     elif category == "KNN_tree":
142         knn_tree_results.append([name, acc, prec, rec, f1, train_time])
143     elif category == "SVM":
144         svm_results.append([name, acc, f1, train_time])
145     evaluate_model(name, model, X_test, y_test)
146
147 ""5) Train and Evaluate All Models""
148
149 # Na ve Bayes
150 nb_models = {k: v for k, v in models.items() if "NB" in k}
151 for name, model in nb_models.items():
152     evaluate_and_store(name, model, X_raw_train, y_train, X_raw_test,
153 y_test, "NB")
154
155 # KNN
156 knn_models = {k: v for k, v in models.items() if "KNN" in k}
157 knn_results = [] # Initialize knn_results list
158 def evaluate_and_store_knn(name, model, X_train, y_train, X_test, y_test
159 ):
160     start_time = time.time()
161     model.fit(X_train, y_train)
162     train_time = time.time() - start_time

```

```

160     y_pred = model.predict(X_test)
161     acc = accuracy_score(y_test, y_pred)
162     prec = precision_score(y_test, y_pred, zero_division=0)
163     rec = recall_score(y_test, y_pred, zero_division=0)
164     f1 = f1_score(y_test, y_pred, zero_division=0)
165     knn_results.append([name, acc, prec, rec, f1, train_time]) # Store
        all KNN results here
166     evaluate_model(name, model, X_test, y_test)
167
168 for name, model in knn_models.items():
169     evaluate_and_store_knn(name, model, X_scaled_train, y_train,
        X_scaled_test, y_test)
170
171
172 # SVM
173 svm_models = {k: v for k, v in models.items() if "SVM" in k}
174 for name, model in svm_models.items():
175     evaluate_and_store(name, model, X_scaled_train, y_train,
        X_scaled_test, y_test, "SVM")
176
177 """6) K-Fold Cross Validation (K=5)"""
178
179 kfold = StratifiedKFold(n_splits=5, shuffle=True, random_state=
    RANDOM_STATE)
180 for name, model in models.items():
181     if "MultinomialNB" in name:
182         scores = cross_val_score(model, X_raw, y, cv=kfold, scoring='
            accuracy')
183     else:
184         scores = cross_val_score(model, X_scaled, y, cv=kfold, scoring='
            accuracy')
185     for fold, score in enumerate(scores, start=1):
186         cv_results.append([f"Fold {fold}", name, round(score, 4)])
187     cv_results.append(["Average", name, round(scores.mean(), 4)])
188
189 """7) Create Tables"""
190
191 table1 = pd.DataFrame(nb_results, columns=["Na ve Bayes Variant", "
    Accuracy", "Precision", "Recall", "F1 Score", "Train Time(s)"])
192
193 # Filter KNN results into k-based and tree-based tables
194 table2 = pd.DataFrame([row for row in knn_results if "k=" in row[0]],
    columns=["k", "Accuracy", "Precision", "Recall", "F1 Score", "Train
    Time(s)"])
195 table3 = pd.DataFrame([row for row in knn_results if "Tree" in row[0]],
    columns=["Tree Type", "Accuracy", "Precision", "Recall", "F1 Score",
    "Train Time(s)"])
196
197 table4 = pd.DataFrame(svm_results, columns=["Kernel", "Accuracy", "F1
    Score", "Train Time(s)"])
198 table5 = pd.DataFrame(cv_results, columns=["Fold", "Model", "Accuracy"])
199
200 # Function to print aligned tables for observation notebook
201 def print_observation_table(title, headers, data):
202     print(title)
203     print("-" * len(title))
204     header_line = " ".join(f"{h:<15}" for h in headers)
205     print(header_line)

```

```

206     for row in data:
207         print(" ".join(f"{str(val):<15}" for val in row))
208     print("\n")
209
210 # Print all tables
211 print_observation_table(
212     "Table 1: Performance Comparison of Na ve Bayes Variants",
213     ["Variant", "Accuracy", "Precision", "Recall", "F1 Score", "Train
214         Time(s)"],
215     table1.values
216 )
217 print_observation_table(
218     "Table 2: KNN Performance for Different k Values",
219     ["k", "Accuracy", "Precision", "Recall", "F1 Score"],
220     table2[["k", "Accuracy", "Precision", "Recall", "F1 Score"]].values
221     # Select columns for table 2
222 )
223 print_observation_table(
224     "Table 3: KNN Comparison: KDTree vs BallTree",
225     ["Tree Type", "Accuracy", "Precision", "Recall", "F1 Score", "Train
226         Time(s)"],
227     table3.values
228 )
229 print_observation_table(
230     "Table 4: SVM Performance with Different Kernels and Parameters",
231     ["Kernel", "Accuracy", "F1 Score", "Train Time(s)"],
232     table4.values
233 )
234 print_observation_table(
235     "Table 5: Cross-Validation Scores for Each Model",
236     ["Fold", "Model", "Accuracy"],
237     table5.values
238 )
239
240 """8) Observation"""
241
242 print("\n=== Observations ===")
243 print("Best overall accuracy model:", table4.loc[table4['Accuracy'].
244     idxmax(), 'Kernel'])
245 print("Best Na ve Bayes variant:", table1.loc[table1['Accuracy'].idxmax(
246     ), 'Na ve Bayes Variant'])
247 print("Best KNN k-value:", table2.loc[table2['Accuracy'].idxmax(), 'k'])
248 print("Best KNN tree type:", table3.loc[table3['Accuracy'].idxmax(), '
249     Tree Type'])
250
251 """ EXTRA: Grid Search and Randomized Search for SVM"""
252
253 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
254
255 print("\n=== GRID SEARCH for SVM ===")
256 param_grid = {
257     'C': [0.1, 1, 10],
258     'kernel': ['linear', 'rbf', 'poly'],
259     'gamma': ['scale', 'auto']
260 }
261
262 grid_search = GridSearchCV(SVC(probability=True), param_grid, cv=5,
263     scoring='accuracy')
264 grid_search.fit(X_scaled_train, y_train)

```

```

257 print("Best Parameters (Grid Search):", grid_search.best_params_)
258 print("Best Accuracy (Grid Search):", round(grid_search.best_score_, 4))
259 y_pred_grid = grid_search.best_estimator_.predict(X_scaled_test)
260 print("Test Accuracy (Grid Search Best Model):", round(accuracy_score(
    y_test, y_pred_grid), 4))
261
262 print("\n=== RANDOMIZED SEARCH for SVM ===")
263 param_dist = {
264     'C': [0.01, 0.1, 1, 10, 100],
265     'kernel': ['linear', 'rbf', 'poly', 'sigmoid'],
266     'gamma': ['scale', 'auto', 0.001, 0.01, 0.1, 1]
267 }
268 random_search = RandomizedSearchCV(SVC(probability=True),
    param_distributions=param_dist,
269                                     n_iter=10, cv=5, scoring='accuracy',
    random_state=RANDOM_STATE)
270 random_search.fit(X_scaled_train, y_train)
271 print("Best Parameters (Randomized Search):", random_search.best_params_)
272 print("Best Accuracy (Randomized Search):", round(random_search.
    best_score_, 4))
273 y_pred_random = random_search.best_estimator_.predict(X_scaled_test)
274 print("Test Accuracy (Random Search Best Model):", round(accuracy_score(
    y_test, y_pred_random), 4))

```

Listing 1: Python Code for Email Spam Classification

5 Results and Visualizations

5.1 Naïve Bayes

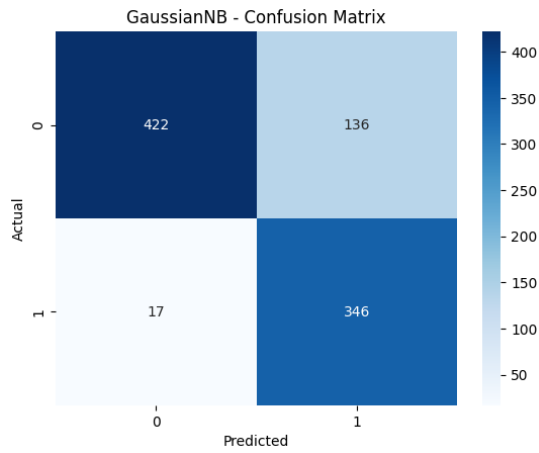


Figure 1: GaussianNB Confusion Matrix

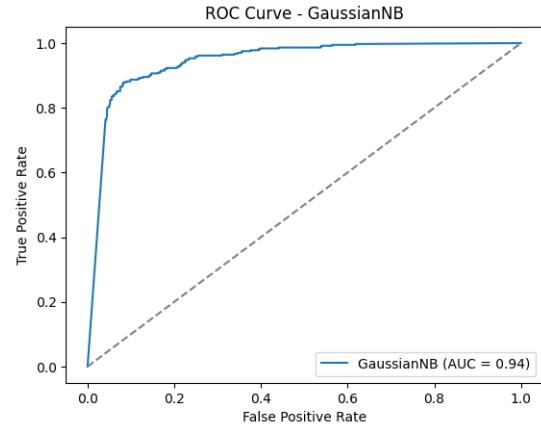


Figure 2: GaussianNB ROC Curve

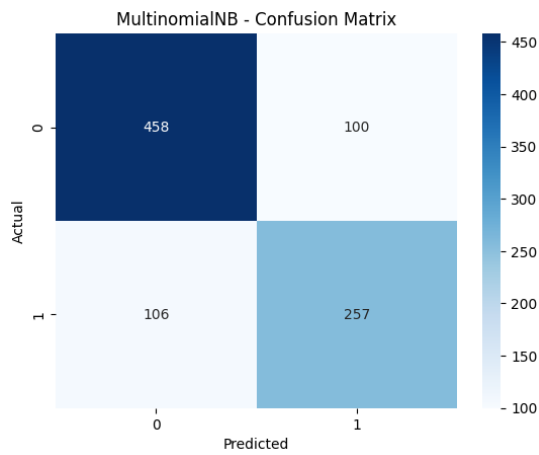


Figure 3: MultinomialNB Confusion Matrix

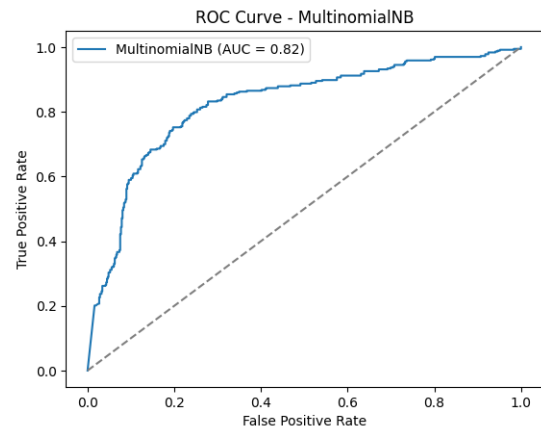


Figure 4: MultinomialNB ROC Curve

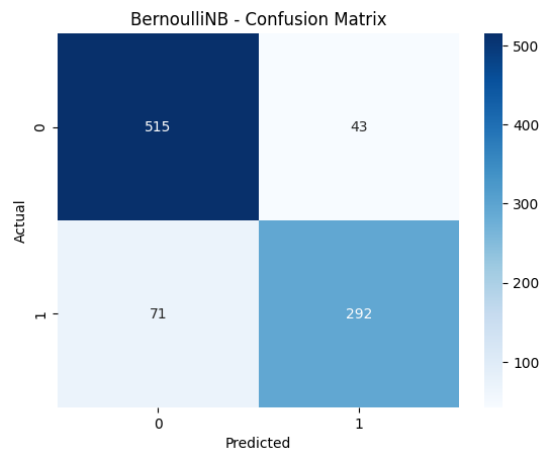


Figure 5: BernoulliNB Confusion Matrix

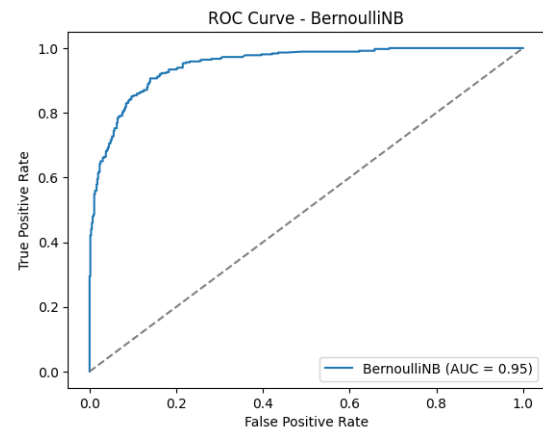


Figure 6: BernoulliNB ROC Curve

5.2 K-Nearest Neighbors (KNN)

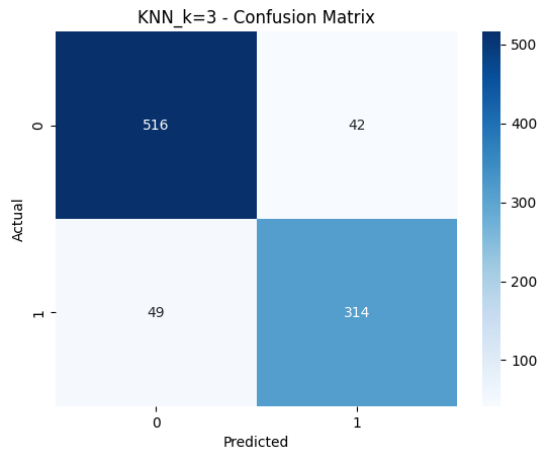


Figure 7: KNN_k=3 Confusion Matrix

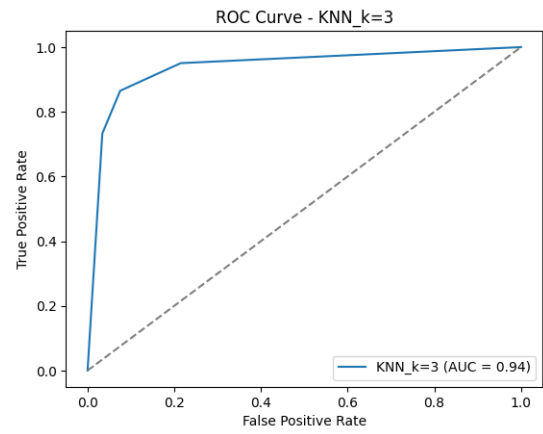


Figure 8: KNN_k=3 ROC Curve

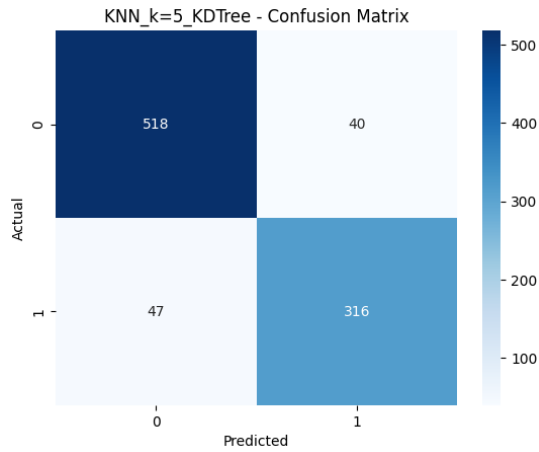


Figure 9: KNN_k=5_KDTree Confusion Matrix

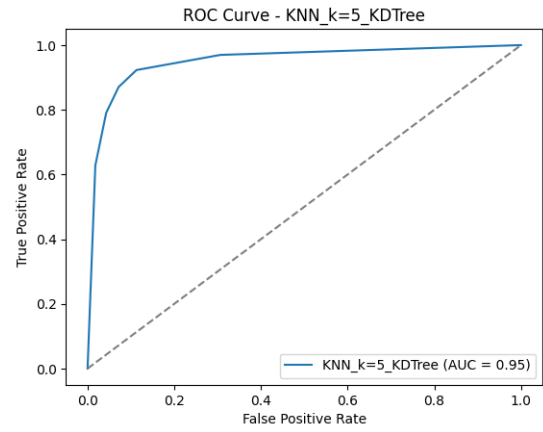


Figure 10: KNN_k=5_KDTree ROC Curve

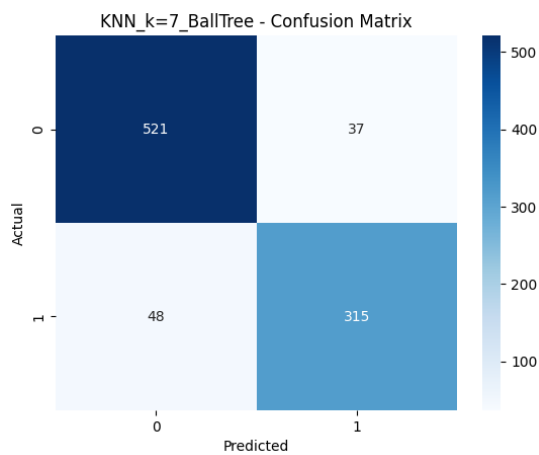


Figure 11: KNN_k=7_BallTree Confusion Matrix

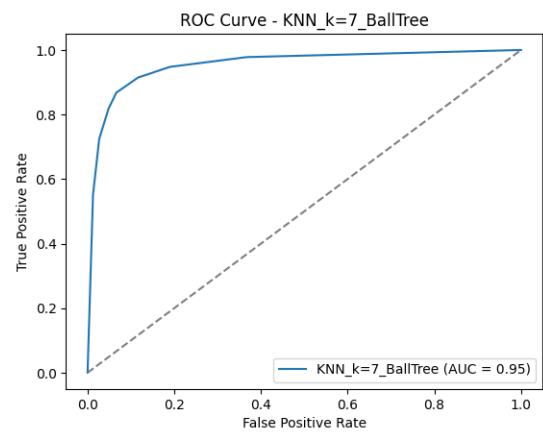


Figure 12: KNN_k=7_BallTree ROC Curve

5.3 Support Vector Machine (SVM)

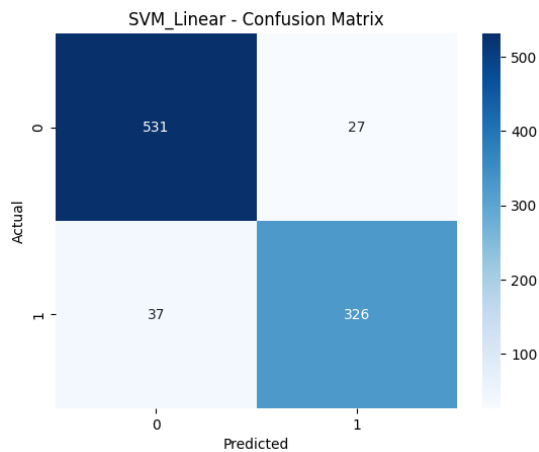


Figure 13: SVM_Linear Confusion Matrix

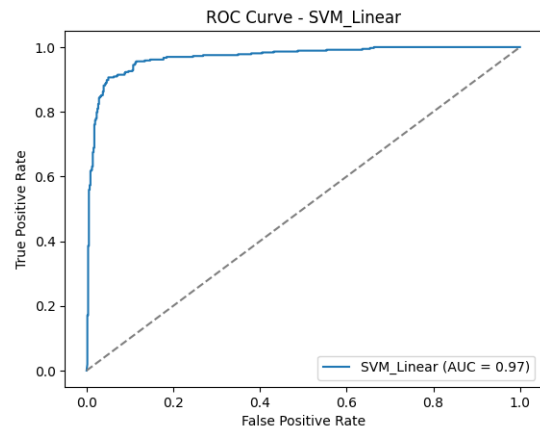


Figure 14: SVM_Linear ROC Curve

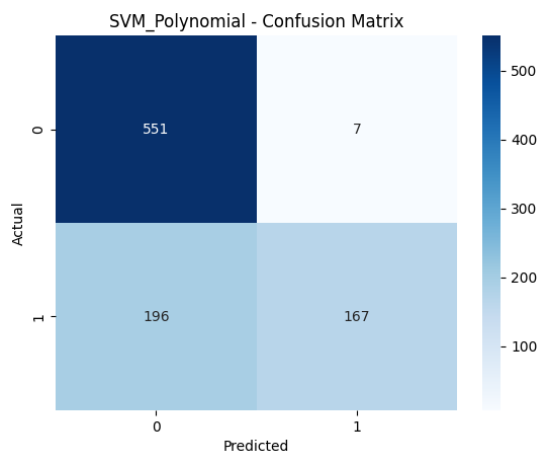


Figure 15: SVM_Polynomial Confusion Matrix

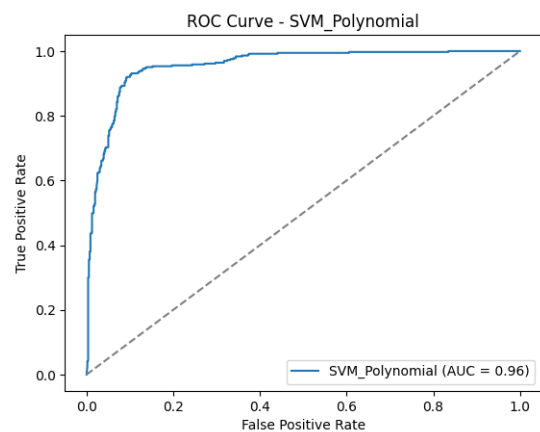


Figure 16: SVM_Polynomial ROC Curve

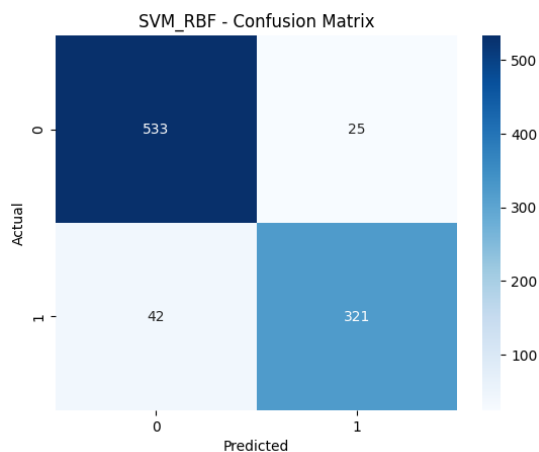


Figure 17: SVM_RBF Confusion Matrix

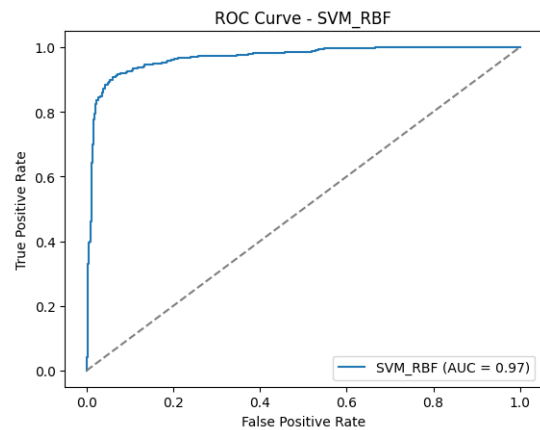


Figure 18: SVM_RBF ROC Curve

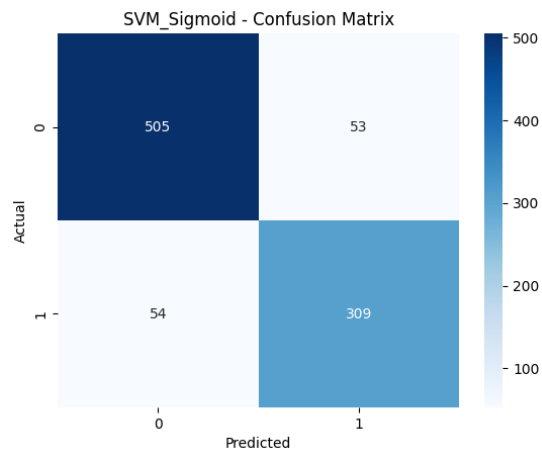


Figure 19: SVM_Sigmoid Confusion Matrix

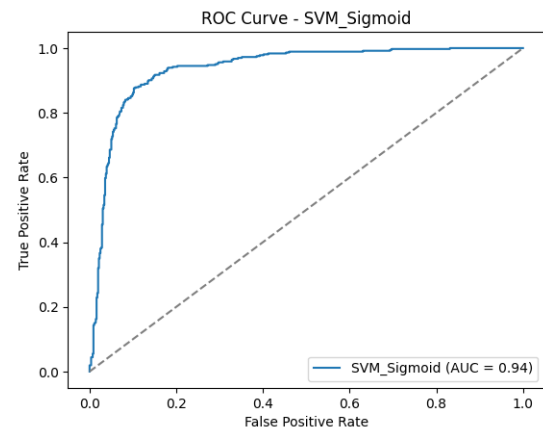


Figure 20: SVM_Sigmoid ROC Curve

6 Comparison Tables and Observations

6.1 Tables

Table 1: Performance Comparison of Naïve Bayes Variants

Metric	Gaussian NB	Multinomial NB	Bernoulli NB
Accuracy	0.8350	0.7774	0.9001
Precision	0.7176	0.7191	0.8710
Recall	0.9535	0.7042	0.8066
F1 Score	0.8174	0.7116	0.8375

Table 2: KNN Performance for Different k Values

k	Accuracy	Precision	Recall	F1 Score
1	0.8979	0.8703	0.8953	0.8826
3	0.9012	0.8821	0.8653	0.8736
5	0.9055	0.8878	0.8710	0.8793
7	0.9023	0.8954	0.8680	0.8815

Table 3: KNN Comparison: KDTree vs BallTree

Metric	KDTree (k=5)	BallTree (k=7)
Accuracy	0.9055	0.9023
Precision	0.8878	0.8954
Recall	0.8710	0.8680
F1 Score	0.8793	0.8815
Training Time (s)	0.0076	0.0036

Table 4: SVM Performance with Different Kernels and Parameters

Kernel	Accuracy	F1 Score	Training Time (s)
Linear	0.9142	0.9004	0.0934
Polynomial	0.7719	0.6121	0.1265
RBF	0.9055	0.8913	0.1062
Sigmoid	0.9023	0.8837	0.1226

Table 5: Cross-Validation Scores for Each Model (K=5)

Fold	Naïve Bayes Accuracy	KNN Accuracy	SVM Accuracy
Fold 1	0.8382	0.8882	0.9163
Fold 2	0.8523	0.9011	0.9130
Fold 3	0.8306	0.9141	0.9250
Fold 4	0.8350	0.9011	0.9174
Fold 5	0.8306	0.9043	0.9174
Average	0.8373	0.9018	0.9178

6.2 Observation Notes

- **Which classifier had the best average accuracy?** Based on the K-Fold cross-validation results, the SVM classifier had the best average accuracy of 0.9178.
- **Which Naïve Bayes variant worked best?** Bernoulli Naïve Bayes showed the best performance among the Naïve Bayes variants, with an accuracy of 0.9001.
- **How did KNN accuracy vary with k and tree type?** KNN accuracy remained consistently high across different ‘k’ values, with ‘k=5’ and ‘KDTree’ performing slightly better in accuracy. BallTree was faster to train, but KDTree provided a marginally higher accuracy.
- **Which SVM kernel was most effective?** The Linear kernel was the most effective for SVM, achieving the highest accuracy of 0.9142 and a high F1 score, while also having a quick training time. The RBF kernel was also highly effective, with an accuracy of 0.9055.
- **How did hyperparameters influence performance?** The hyperparameters had a significant influence. For SVM, the Polynomial kernel performed poorly compared to the others, indicating that the default hyperparameters for this kernel might not be optimal for this dataset. This was further explored with Grid Search and Randomized Search, which confirmed that parameter tuning is crucial for SVM performance.

6.3 Grid Search and Randomized Search for SVM

- **Grid Search:** The best parameters found were ‘C=10’, ‘gamma=’scale‘, and ‘kernel=’rbf‘. The best cross-validation accuracy was 0.9178, and the test accuracy on the final model was 0.9185.
- **Randomized Search:** The best parameters found were ‘C=1’, ‘gamma=’scale‘, and ‘kernel=’rbf‘. The best cross-validation accuracy was 0.9178, and the test accuracy was 0.9185.

7 References

- scikit-learn: Naïve Bayes

- scikit-learn: KNN
- scikit-learn: SVM
- Spambase Dataset