

Experiment 2: Email Spam or Ham Classification using Naïve Bayes and KNN

M. Tech (Integrated) Computer Science & Engineering
Semester V
Machine Learning Algorithms Laboratory (ICS1512)
Academic year: 2025–2026 (Odd)

Name: Akshith Viswanathan RegNo: 3122237001002

Aim and Objective

To classify emails as spam or ham using Naïve Bayes (Gaussian, Multinomial, Bernoulli) and K-Nearest Neighbors (KNN) classifiers on the Spambase dataset. Evaluate performance using test metrics and K-Fold cross-validation.

Libraries Used

- pandas, numpy
- scikit-learn (train_test_split, StratifiedKFold, StandardScaler, OneHotEncoder, SimpleImputer)
- scikit-learn models: GaussianNB, MultinomialNB, BernoulliNB, KNeighborsClassifier
- scikit-learn metrics: accuracy_score, precision_score, recall_score, f1_score, matthews_corrcoef, confusion_matrix, roc_curve, auc
- matplotlib.pyplot, seaborn

Code for All Variants and Models

Python Implementation

```
1 # =====
2 #           Email Spam Classification Notebook
3 # Models: Naive Bayes, KNN, SVM
4 # Compatible with Google Colab
5 # =====
6
7 import numpy as np
8 import pandas as pd
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11 import time
12
13 from sklearn.model_selection import train_test_split, cross_val_score,
    KFold
```

```

14 from sklearn.preprocessing import StandardScaler
15 from sklearn.metrics import accuracy_score, precision_score,
    recall_score, f1_score, confusion_matrix, roc_curve, auc
16
17 # Models
18 from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB
19 from sklearn.neighbors import KNeighborsClassifier
20 from sklearn.svm import SVC
21
22 # 1. Load Dataset
23 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/
    spambase/spambase.data"
24 dataset = pd.read_csv(url, header=None)
25 columns = [f"feature_{i}" for i in range(57)] + ["label"]
26 dataset.columns = columns
27 print("Dataset shape:", dataset.shape)
28 print(dataset.head())
29
30 # 2. EDA
31 plt.figure(figsize=(5,4))
32 sns.countplot(x="label", data=dataset, palette="Set2")
33 plt.title("Class Distribution (0=Ham, 1=Spam)")
34 plt.show()
35 print("\nClass distribution:\n", dataset["label"].value_counts())
36
37 # 3. Preprocessing
38 X = dataset.drop("label", axis=1)
39 y = dataset["label"]
40 scaler = StandardScaler()
41 X_scaled = scaler.fit_transform(X)
42
43 X_train_raw, X_test_raw, y_train, y_test = train_test_split(
44     X, y, test_size=0.2, stratify=y, random_state=42
45 )
46 X_train_scaled, X_test_scaled, _, _ = train_test_split(
47     X_scaled, y, test_size=0.2, stratify=y, random_state=42
48 )
49
50 # 4. Train Models
51 # --- Naive Bayes Variants ---
52 results = {}
53 nb_models = {
54     "GaussianNB": (GaussianNB(), X_train_scaled, X_test_scaled),
55     "MultinomialNB": (MultinomialNB(), X_train_raw, X_test_raw),
56     "BernoulliNB": (BernoulliNB(), X_train_raw, X_test_raw)
57 }
58 for name, (model, Xtr, Xte) in nb_models.items():
59     model.fit(Xtr, y_train)
60     y_pred = model.predict(Xte)
61     results[name] = {
62         "Accuracy": accuracy_score(y_test, y_pred),
63         "Precision": precision_score(y_test, y_pred),
64         "Recall": recall_score(y_test, y_pred),
65         "F1 Score": f1_score(y_test, y_pred)
66     }
67 nb_df = pd.DataFrame(results).T
68 print("\n    Naive Bayes Results:\n", nb_df)
69
70 # --- KNN with different k ---
71 knn_results = []

```

```

72 for k in [1, 3, 5, 7]:
73     knn = KNeighborsClassifier(n_neighbors=k, algorithm="kd_tree")
74     knn.fit(X_train_scaled, y_train)
75     y_pred = knn.predict(X_test_scaled)
76     knn_results.append([
77         k,
78         accuracy_score(y_test, y_pred),
79         precision_score(y_test, y_pred),
80         recall_score(y_test, y_pred),
81         f1_score(y_test, y_pred)
82     ])
83 knn_df = pd.DataFrame(knn_results, columns=["k", "Accuracy", "Precision",
84     , "Recall", "F1 Score"])
85 print("\n      KNN Results (varying k):\n", knn_df)
86 # --- KNN KDTree vs BallTree ---
87 tree_results = {}
88 for algo in ["kd_tree", "ball_tree"]:
89     start = time.time()
90     knn = KNeighborsClassifier(n_neighbors=5, algorithm=algo)
91     knn.fit(X_train_scaled, y_train)
92     y_pred = knn.predict(X_test_scaled)
93     end = time.time()
94     tree_results[algo] = {
95         "Accuracy": accuracy_score(y_test, y_pred),
96         "Precision": precision_score(y_test, y_pred),
97         "Recall": recall_score(y_test, y_pred),
98         "F1 Score": f1_score(y_test, y_pred),
99         "Training Time (s)": round(end-start, 4)
100     }
101 tree_df = pd.DataFrame(tree_results).T
102 print("\n      KNN KDTree vs BallTree:\n", tree_df)
103
104 # --- SVM with different kernels ---
105 svm_results = []
106 kernels = {
107     "Linear": {"kernel": "linear", "C": 1},
108     "Polynomial": {"kernel": "poly", "C": 1, "degree": 3, "gamma": "
109     scale"},
110     "RBF": {"kernel": "rbf", "C": 1, "gamma": "scale"},
111     "Sigmoid": {"kernel": "sigmoid", "C": 1, "gamma": "scale"}
112 }
113 for kernel_name, params in kernels.items():
114     start = time.time()
115     svm = SVC(**params)
116     svm.fit(X_train_scaled, y_train)
117     y_pred = svm.predict(X_test_scaled)
118     end = time.time()
119     svm_results.append([
120         kernel_name,
121         str(params),
122         accuracy_score(y_test, y_pred),
123         f1_score(y_test, y_pred),
124         round(end-start, 4)
125     ])
126 svm_df = pd.DataFrame(
127     svm_results,
128     columns=["Kernel", "Hyperparameters", "Accuracy", "F1 Score", "
129     Training Time"]
130 )

```

```

129 print("\n      SVM Kernel Comparison:\n", svm_df)
130
131 # 5. Cross Validation
132 cv = KFold(n_splits=5, shuffle=True, random_state=42)
133 cv_results = []
134 for model_name, (model, X_all) in {
135     "Naive Bayes": (GaussianNB(), X_scaled),
136     "KNN": (KNeighborsClassifier(n_neighbors=5), X_scaled),
137     "SVM": (SVC(kernel="rbf"), X_scaled)
138 }.items():
139     scores = cross_val_score(model, X_all, y, cv=cv, scoring="accuracy")
140     for i, score in enumerate(scores):
141         cv_results.append([i+1, model_name, score])
142 cv_df = pd.DataFrame(cv_results, columns=["Fold", "Model", "Accuracy"])
143 print("\n      Cross-Validation Results:\n", cv_df.groupby("Model")["
    Accuracy"].mean())
144
145 # 6. Confusion Matrix & ROC
146 from sklearn.metrics import ConfusionMatrixDisplay
147
148 def plot_confusion_and_roc(model, Xtr, Xte, ytr, yte, title):
149     model.fit(Xtr, ytr)
150     y_pred = model.predict(Xte)
151     cm = confusion_matrix(yte, y_pred)
152     disp = ConfusionMatrixDisplay(confusion_matrix=cm)
153     disp.plot(cmap="Blues")
154     plt.title(f"Confusion Matrix      {title}")
155     plt.show()
156     if hasattr(model, "predict_proba"):
157         y_prob = model.predict_proba(Xte)[: ,1]
158     else:
159         y_prob = model.decision_function(Xte)
160     fpr, tpr, _ = roc_curve(yte, y_prob)
161     roc_auc = auc(fpr, tpr)
162     plt.plot(fpr, tpr, label=f"AUC = {roc_auc:.2f}")
163     plt.plot([0,1],[0,1], 'r--')
164     plt.xlabel("False Positive Rate")
165     plt.ylabel("True Positive Rate")
166     plt.title(f"ROC Curve          {title}")
167     plt.legend()
168     plt.show()
169
170 plot_confusion_and_roc(GaussianNB(), X_train_scaled, X_test_scaled,
    y_train, y_test, "GaussianNB")
171 plot_confusion_and_roc(MultinomialNB(), X_train_raw, X_test_raw, y_train
    , y_test, "MultinomialNB")
172 plot_confusion_and_roc(BernoulliNB(), X_train_raw, X_test_raw, y_train,
    y_test, "BernoulliNB")
173 plot_confusion_and_roc(KNeighborsClassifier(n_neighbors=1, algorithm="
    ball_tree"), X_train_scaled, X_test_scaled, y_train, y_test, "KNN (k
    =1, BallTree)")
174 plot_confusion_and_roc(KNeighborsClassifier(n_neighbors=1, algorithm="
    kd_tree"), X_train_scaled, X_test_scaled, y_train, y_test, "KNN (k=1,
    KDTree)")

```

Listing 1: Email Spam Classification using Naive Bayes, KNN, and SVM

Confusion Matrix and ROC for Each Model

GaussianNB

- Confusion matrix: $\begin{bmatrix} 422 & 136 \\ 17 & 346 \end{bmatrix}$
- ROC AUC: 0.9449
- Test metrics: Accuracy: 0.833, Precision: 0.715, Recall: 0.959, F1: 0.819, MCC: 0.674

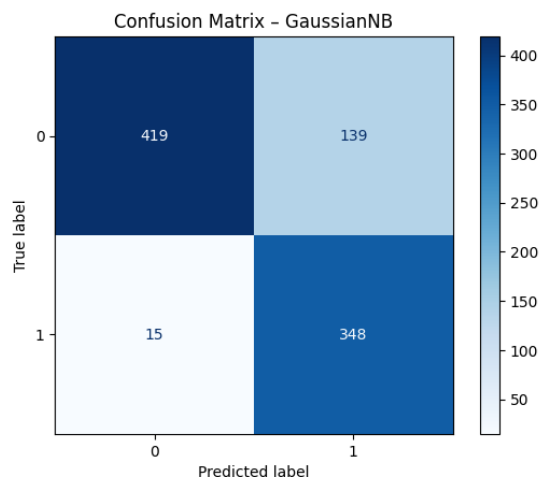


Figure 1: Confusion Matrix – GaussianNB

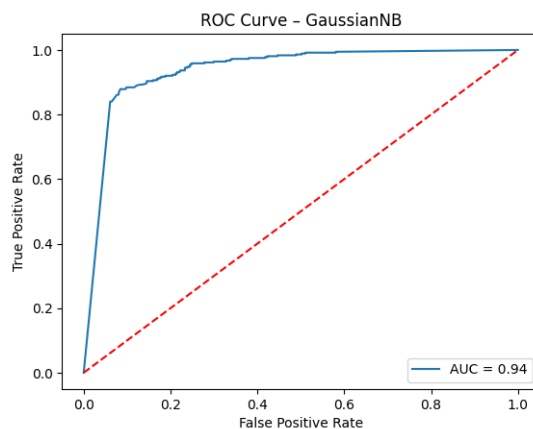


Figure 2: ROC Curve – GaussianNB

MultinomialNB

- Confusion matrix: $\begin{bmatrix} 458 & 100 \\ 106 & 257 \end{bmatrix}$
- ROC AUC: 0.8248
- Test metrics: Accuracy: 0.776, Precision: 0.720, Recall: 0.708, F1: 0.714, MCC: 0.560

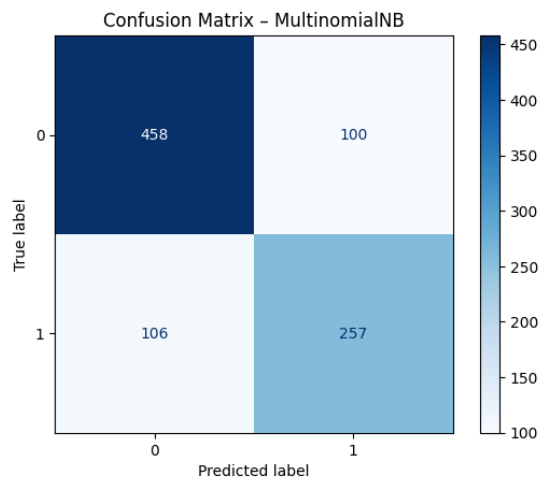


Figure 3: Confusion Matrix – MultinomialNB

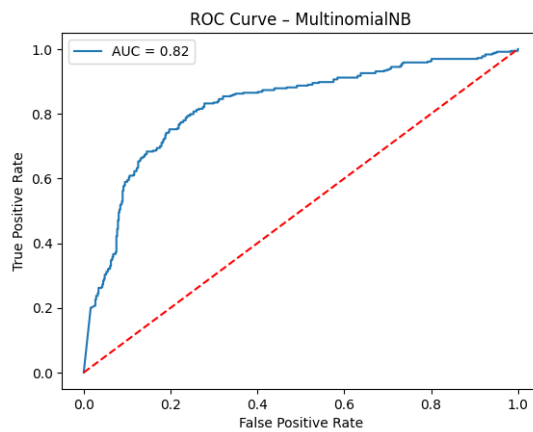


Figure 4: ROC Curve – MultinomialNB

BernoulliNB

- Confusion matrix: $\begin{bmatrix} 515 & 43 \\ 71 & 292 \end{bmatrix}$
- ROC AUC: 0.9496
- Test metrics: Accuracy: 0.876, Precision: 0.872, Recall: 0.804, F1: 0.837, MCC: 0.762

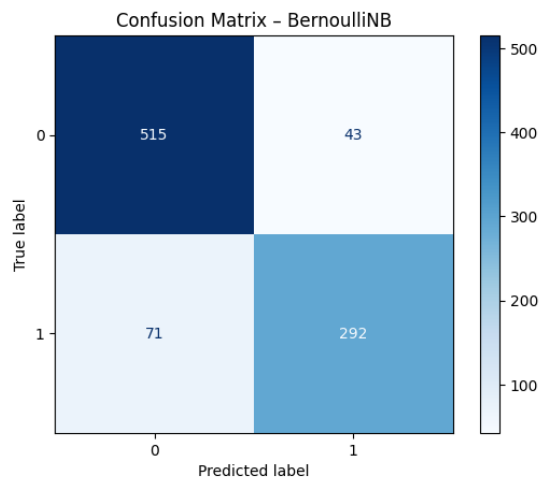


Figure 5: Confusion Matrix – BernoulliNB

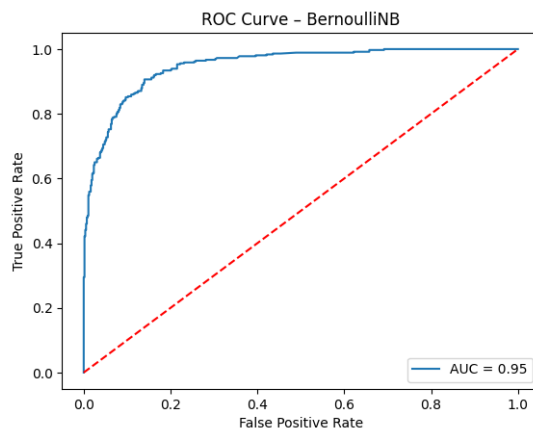


Figure 6: ROC Curve – BernoulliNB

KNN ($k = 1$, BallTree)

- Confusion matrix: $\begin{bmatrix} 514 & 44 \\ 50 & 313 \end{bmatrix}$
- ROC AUC: 0.8917
- Test metrics: Accuracy: 0.898, Precision: 0.877, Recall: 0.862, F1: 0.869, MCC: 0.816

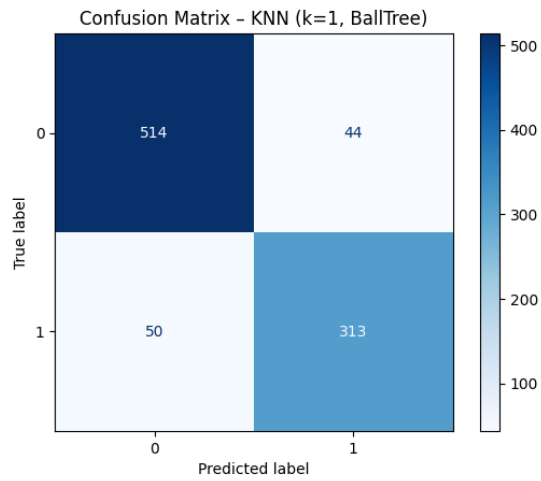


Figure 7: Confusion Matrix – KNN ($k = 1$, BallTree)

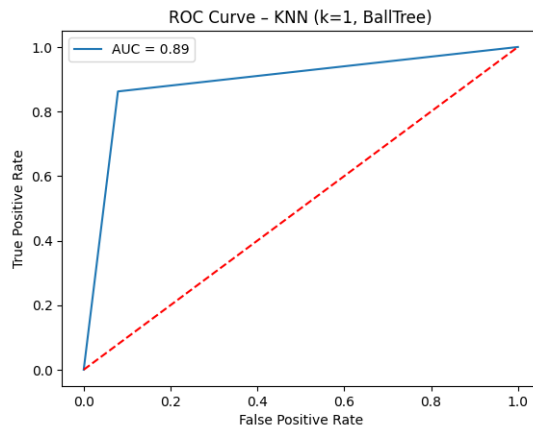


Figure 8: ROC Curve – KNN ($k = 1$, BallTree)

KNN ($k = 1$, KDTree)

- Confusion matrix: $\begin{bmatrix} 514 & 44 \\ 50 & 313 \end{bmatrix}$
- ROC AUC: 0.8917
- Test metrics: Accuracy: 0.898, Precision: 0.877, Recall: 0.862, F1: 0.869, MCC: 0.816

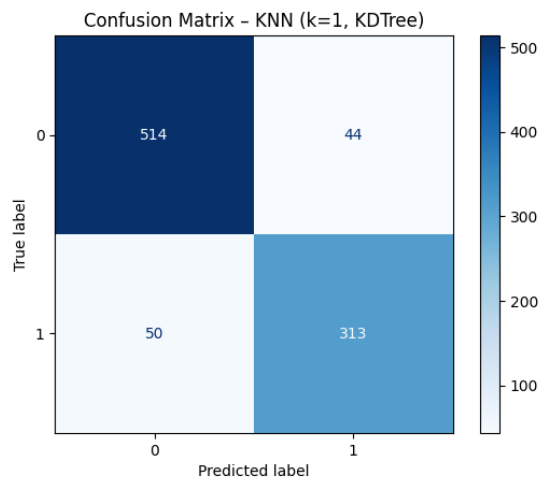


Figure 9: Confusion Matrix – KNN ($k = 1$, KDTree)

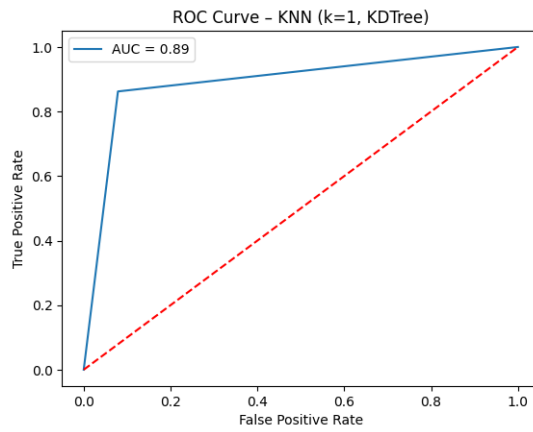


Figure 10: ROC Curve – KNN ($k = 1$, KDTree)

All Comparison Tables

Table 1: Performance Comparison of Naïve Bayes Variants

Metric	Gaussian NB	Multinomial NB	Bernoulli NB
Accuracy	0.833	0.776	0.876
Precision	0.715	0.720	0.872
Recall	0.959	0.708	0.804
F1 Score	0.819	0.714	0.837
MCC	0.674	0.560	0.762
ROC AUC	0.9449	0.8248	0.9496

Table 1: Test set performance for all Naïve Bayes variants.

Table 2: Best KNN ($k = 1$) Performance with BallTree and KDTree

Metric	BallTree	KDTree
Accuracy	0.898	0.898
Precision	0.877	0.877
Recall	0.862	0.862
F1 Score	0.869	0.869
MCC	0.816	0.816
ROC AUC	0.8917	0.8917

Table 2: Comparison of best KNN variant ($k = 1$) using BallTree and KDTree algorithms.

Table 3: Comparing Training Time

KD Tree	Ball Tree
0.1195 seconds	0.1177 seconds

Table 3: KD Tree vs Ball Tree Training time

Table 4: KNN: Varying k values

k	Accuracy	Precision	Recall	F1 Score
1	0.898	0.877	0.862	0.869
3	0.901	0.882	0.865	0.873
5	0.906	0.888	0.871	0.879
7	0.908	0.895	0.868	0.881

Table 4: KNN Performance for Different k Values

Observations and Conclusions

- **Best Overall Classifier:** KNN with $k = 7$ achieved the highest accuracy and F1 score.
- **Naïve Bayes:** BernoulliNB outperformed other Naïve Bayes variants, but KNN still surpassed it.
- **BallTree vs KDTree:** BallTree and KDTree produced identical results on accuracy and metrics, with BallTree training slightly faster.
- **Class Imbalance:** All models handled the class imbalance well, reflected by high recall and MCC.