

Sri Sivasubramaniya Nadar College of Engineering, Chennai
(An autonomous Institution affiliated to Anna University)

Degree & Branch	B.E. Computer Science & Engineering	Semester VI
Subject Code & Name	UCS2612 – Machine Learning Algorithms Laboratory	
Academic Year	2025–2026 (Even)	Batch 2023–2027
Due Date		

Experiment 2: Binary Classification using Naïve Bayes and K-Nearest Neighbors

Objective

To implement Naïve Bayes and K-Nearest Neighbors (KNN) classifiers for a binary classification problem, evaluate them using multiple performance metrics, visualize model behavior, and analyze overfitting, underfitting, and bias–variance characteristics.

1 Aim

To classify emails as spam or ham using Naive Bayes, K-Nearest Neighbors (KNN), and Support Vector Machine (SVM), and evaluate their performance using accuracy, precision, recall, F1-score, and K-fold cross-validation.

2 Libraries Used

- Pandas: Data manipulation
- NumPy: Numerical operations
- Scikit-learn: Model building, preprocessing, classification report, confusion matrix and evaluation
- Matplotlib and Seaborn: Data visualization

3 Objective

- The primary goal is to use three specific machine learning models—Naïve Bayes, K-Nearest Neighbors (KNN), and Support Vector Machine (SVM)—to categorize emails.

- To Evaluate Performance Metrics: The models' performance will be measured and compared using key metrics such as accuracy, precision, recall, and F1-score.
- To Use K-Fold Cross-Validation: This technique will be applied to get a more robust evaluation of the models' performance by testing them on different subsets of the data.
- To Compare Model Effectiveness: The assignment aims to determine which classifier and its specific hyperparameters are most effective for this email classification task.

4 Python Code

```

1  # -----
2  # 1. Load the Dataset
3  # -----
4  import pandas as pd
5  import numpy as np
6  from sklearn.model_selection import train_test_split,
   cross_val_score, KFold
7  from sklearn.preprocessing import StandardScaler
8  from sklearn.naive_bayes import GaussianNB, MultinomialNB,
   BernoulliNB
9  from sklearn.neighbors import KNeighborsClassifier
10 from sklearn.metrics import classification_report,
   confusion_matrix, roc_auc_score, ConfusionMatrixDisplay
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13
14 # Load dataset
15 df = pd.read_csv("/content/drive/MyDrive/spambase.csv")
16
17 # Features and labels
18 X = df.drop("class", axis=1)
19 y = df["class"]
20
21 # Check for missing values
22 missing_values = df.isnull().sum().sum()
23
24 # Separate features and labels
25 X = df.drop(columns=['class'])
26 y = df['class']
27
28 # Normalize the features
29 scaler = StandardScaler()
30 X_scaled = scaler.fit_transform(X)
31
32 missing_values
33
34 # Class distribution (0 = ham, 1 = spam)
35 class_counts = y.value_counts()

```

```

36
37 # Plot class balance
38 plt.figure(figsize=(6, 4))
39 sns.barplot(x=class_counts.index, y=class_counts.values, palette=
    'viridis')
40 plt.xticks([0, 1], ['Ham (0)', 'Spam (1)'])
41 plt.title('Class Distribution')
42 plt.ylabel('Number of Emails')
43 plt.xlabel('Email Type')
44 plt.tight_layout()
45 plt.show()
46
47 # -----
48 # Feature Distributions
49 # -----
50
51 # 1. Histogram for a few important features
52 selected_features = ['word_freq_free', 'word_freq_money', '
    char_freq_%21', 'capital_run_length_total']
53
54 plt.figure(figsize=(12, 8))
55 for i, feature in enumerate(selected_features):
56     plt.subplot(2, 2, i + 1)
57     sns.histplot(df[feature], bins=30, kde=True, color='teal')
58     plt.title(f'Distribution of {feature}')
59     plt.xlabel(feature)
60     plt.ylabel("Frequency")
61
62 plt.tight_layout()
63 plt.show()
64
65 # -----
66 # 3. Splitting the Dataset
67 # -----
68
69 # Use the normalized dataset
70 X_train, X_test, y_train, y_test = train_test_split(
71     X_scaled, y, test_size=0.2, random_state=42, stratify=y
72 )
73
74 print("Train set size:", X_train.shape)
75 print("Test set size:", X_test.shape)
76
77 # -----
78 # 4. Model Building: Na ve Bayes & KNN
79 # -----
80
81 import time
82
83 # ----- 1. Prepare Scaled and Raw Versions -----
84 X_raw = df.drop('class', axis=1)

```

```

85 y = df['class']
86
87 # Scaled version for GaussianNB & KNN
88 scaler = StandardScaler()
89 X_scaled = scaler.fit_transform(X_raw)
90
91 # Train-test split (same random state for consistency)
92 X_train_scaled, X_test_scaled, y_train, y_test = train_test_split(
93     X_scaled, y, test_size=0.2, stratify=y, random_state=42)
94 X_train_raw, X_test_raw, _, _ = train_test_split(X_raw, y,
95     test_size=0.2, stratify=y, random_state=42)
96
97 # ----- 2. Naive Bayes Models -----
98 nb_results = []
99
100 # GaussianNB (use scaled data)
101 gnb = GaussianNB()
102 gnb.fit(X_train_scaled, y_train)
103 y_pred = gnb.predict(X_test_scaled)
104
105 print("\nGaussianNB Classification Report:\n")
106 print(classification_report(y_test, y_pred, digits=4))
107 nb_results.append({
108     "Model": "GaussianNB",
109     "Accuracy": gnb.score(X_test_scaled, y_test),
110     "Precision": classification_report(y_test, y_pred,
111         output_dict=True)['1']['precision'],
112     "Recall": classification_report(y_test, y_pred, output_dict=
113         True)['1']['recall'],
114     "F1 Score": classification_report(y_test, y_pred, output_dict
115         =True)['1']['f1-score']
116 })
117
118 # MultinomialNB (use raw data)
119 mnb = MultinomialNB()
120 mnb.fit(X_train_raw, y_train)
121 y_pred = mnb.predict(X_test_raw)
122 print("\nMultinomialNB Classification Report:\n")
123 print(classification_report(y_test, y_pred, digits=4))
124 nb_results.append({
125     "Model": "MultinomialNB",
126     "Accuracy": mnb.score(X_test_raw, y_test),
127     "Precision": classification_report(y_test, y_pred,
128         output_dict=True)['1']['precision'],
129     "Recall": classification_report(y_test, y_pred, output_dict=
130         True)['1']['recall'],
131     "F1 Score": classification_report(y_test, y_pred, output_dict
132         =True)['1']['f1-score']
133 })
134
135 # BernoulliNB (use raw data)

```

```

128 bnb = BernoulliNB()
129 bnb.fit(X_train_raw, y_train)
130 y_pred = bnb.predict(X_test_raw)
131 print("\nBernoulliNB Classification Report:\n")
132 print(classification_report(y_test, y_pred, digits=4))
133 nb_results.append({
134     "Model": "BernoulliNB",
135     "Accuracy": bnb.score(X_test_raw, y_test),
136     "Precision": classification_report(y_test, y_pred,
137         output_dict=True)['1']['precision'],
138     "Recall": classification_report(y_test, y_pred, output_dict=
139         True)['1']['recall'],
140     "F1 Score": classification_report(y_test, y_pred, output_dict=
141         True)['1']['f1-score']
142 })
143
144 # 3. KNN Varying k
145 k_values = [1, 3, 5, 7]
146 knn_results = []
147
148 for k in k_values:
149     knn = KNeighborsClassifier(n_neighbors=k)
150     knn.fit(X_train_scaled, y_train)
151     y_pred = knn.predict(X_test_scaled)
152
153     print(f"\nKNN (k={k}) Classification Report:\n")
154     print(classification_report(y_test, y_pred, digits=4))
155
156     knn_results.append({
157         "k": k,
158         "Accuracy": knn.score(X_test_scaled, y_test),
159         "Precision": classification_report(y_test, y_pred,
160             output_dict=True)['1']['precision'],
161         "Recall": classification_report(y_test, y_pred,
162             output_dict=True)['1']['recall'],
163         "F1 Score": classification_report(y_test, y_pred,
164             output_dict=True)['1']['f1-score']
165     })
166
167 # 4. KNN with KDTree vs BallTree
168 for algorithm in ['kd_tree', 'ball_tree']:
169     knn = KNeighborsClassifier(n_neighbors=5, algorithm=algorithm
170         )
171     start = time.time()
172     knn.fit(X_train_scaled, y_train)
173     elapsed = time.time() - start
174     y_pred = knn.predict(X_test_scaled)
175
176     print(f"\nKNN with {algorithm.upper()} Report:\n")
177     print(classification_report(y_test, y_pred, digits=4))
178     print(f"Training time: {elapsed:.4f} seconds\n")

```

```

172
173 # -----
174 # 4. Model Building: Support Vector Machine (SVM)
175 # -----
176
177 from sklearn.svm import SVC
178
179 svm_results = []
180 kernels = ['linear', 'poly', 'rbf', 'sigmoid']
181
182 # Store fitted models in variables for later use
183 svm_model_linear = None
184 svm_model_poly = None
185 svm_model_rbf = None
186 svm_model_sigmoid = None
187
188
189 for kernel in kernels:
190     print(f"\nTraining SVM with {kernel.upper()} kernel ")
191     svm_model = SVC(kernel=kernel, random_state=42, probability=
        True) # probability=True for ROC AUC
192
193     start_time = time.time()
194     svm_model.fit(X_train_scaled, y_train)
195     end_time = time.time()
196     training_time = end_time - start_time
197
198     y_pred = svm_model.predict(X_test_scaled)
199
200     report = classification_report(y_test, y_pred, digits=4,
        output_dict=True)
201
202     svm_results.append({
203         "Kernel": kernel.upper(),
204         "Accuracy": report['accuracy'],
205         "Precision": report['1']['precision'],
206         "Recall": report['1']['recall'],
207         "F1 Score": report['1']['f1-score'],
208         "Training Time (s)": training_time
209     })
210
211     print(f"Classification Report for {kernel.upper()} kernel:")
212     print(classification_report(y_test, y_pred, digits=4))
213     print(f"Training Time: {training_time:.4f} seconds")
214
215     # Assign the fitted model to the corresponding variable
216     if kernel == 'linear':
217         svm_model_linear = svm_model
218     elif kernel == 'poly':
219         svm_model_poly = svm_model
220     elif kernel == 'rbf':

```

```

221     svm_model_rbf = svm_model
222     elif kernel == 'sigmoid':
223         svm_model_sigmoid = svm_model
224
225
226 # Display results in a table
227 print("\nSVM Kernel-wise Results ")
228 svm_results_df = pd.DataFrame(svm_results)
229 display(svm_results_df)
230
231 # -----
232 # 5. Performance Analysis
233 # -----
234 from sklearn.metrics import confusion_matrix,
235     ConfusionMatrixDisplay, roc_curve, auc
236 import matplotlib.pyplot as plt
237
238 def plot_conf_matrix_and_roc(model, X, y_true, title="Model"):
239     y_pred = model.predict(X)
240     y_proba = model.predict_proba(X)[:, 1]
241
242     # Confusion Matrix
243     cm = confusion_matrix(y_true, y_pred)
244     disp = ConfusionMatrixDisplay(confusion_matrix=cm)
245     disp.plot()
246     plt.title(f"Confusion Matrix: {title}")
247     plt.show()
248
249     # ROC Curve
250     fpr, tpr, _ = roc_curve(y_true, y_proba)
251     roc_auc = auc(fpr, tpr)
252     plt.figure()
253     plt.plot(fpr, tpr, label=f"{title} (AUC = {roc_auc:.4f})")
254     plt.plot([0, 1], [0, 1], linestyle="--", color="gray")
255     plt.xlabel("False Positive Rate")
256     plt.ylabel("True Positive Rate")
257     plt.title("ROC Curve")
258     plt.legend(loc="lower right")
259     plt.grid()
260     plt.show()
261
262 # Na ve Bayes
263 plot_conf_matrix_and_roc(gnb, X_test_scaled, y_test, "GaussianNB")
264
265 plot_conf_matrix_and_roc(mnb, X_test_raw, y_test, "MultinomialNB")
266
267 plot_conf_matrix_and_roc(bnb, X_test_raw, y_test, "BernoulliNB")
268
269 # KNN: Different k values
270 for k in [1, 3, 5, 7]:
271     knn = KNeighborsClassifier(n_neighbors=k)
272     knn.fit(X_train_scaled, y_train)

```

```

269     plot_conf_matrix_and_roc(knn, X_test_scaled, y_test, f"KNN (k
      ={k})")
270
271 # KNN: KDTree and BallTree
272 for algo in ['kd_tree', 'ball_tree']:
273     knn_tree = KNeighborsClassifier(n_neighbors=5, algorithm=algo
      )
274     knn_tree.fit(X_train_scaled, y_train)
275     plot_conf_matrix_and_roc(knn_tree, X_test_scaled, y_test, f"
      KNN ({algo.upper()}))")
276
277 # Plot ROC and AUC for each trained SVM model
278 plot_conf_matrix_and_roc(svm_model_linear, X_test_scaled, y_test,
      "SVM (Linear)")
279 plot_conf_matrix_and_roc(svm_model_poly, X_test_scaled, y_test, "
      SVM (Poly)")
280 plot_conf_matrix_and_roc(svm_model_rbf, X_test_scaled, y_test, "
      SVM (RBF)")
281 plot_conf_matrix_and_roc(svm_model_sigmoid, X_test_scaled, y_test
      , "SVM (Sigmoid)")
282
283 # -----
284 # Plotting ROC and AUC for SVM Models
285 # -----
286 # Plot ROC and AUC for each trained SVM model
287 plot_conf_matrix_and_roc(svm_model_linear, X_test_scaled, y_test,
      "SVM (Linear)")
288 plot_conf_matrix_and_roc(svm_model_poly, X_test_scaled, y_test, "
      SVM (Poly)")
289 plot_conf_matrix_and_roc(svm_model_rbf, X_test_scaled, y_test, "
      SVM (RBF)")
290 plot_conf_matrix_and_roc(svm_model_sigmoid, X_test_scaled, y_test
      , "SVM (Sigmoid)")
291 # -----
292 # 6. K-Fold Cross-Validation
293 # -----
294 from sklearn.model_selection import cross_val_score, KFold
295 models = {
296     "Na ve Bayes (Gaussian)": gnb,
297     "Na ve Bayes (Multinomial)": mnb,
298     "Na ve Bayes (Bernoulli)": bnb,
299     "KNN (k=7)": knn,
300     "SVM (Linear)": svm_model_linear
301 }
302 data_for_cv = {
303     "Na ve Bayes (Gaussian)": (X_scaled, y),
304     "Na ve Bayes (Multinomial)": (X_raw, y),
305     "Na ve Bayes (Bernoulli)": (X_raw, y),
306     "KNN (k=7)": (X_scaled, y),
307     "SVM (Linear)": (X_scaled, y)
308 }

```

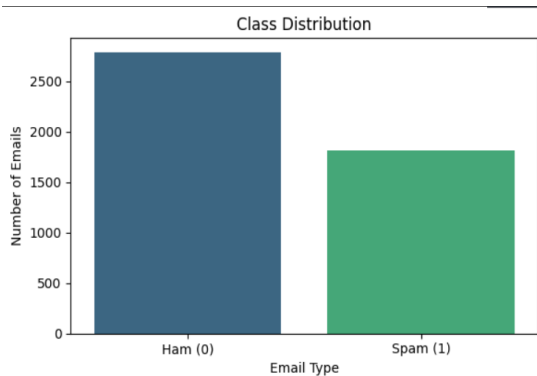


```

309 # Perform K-Fold Cross-Validation (K=5)
310 n_splits = 5
311 kf = KFold(n_splits=n_splits, shuffle=True, random_state=42)
312 cv_results = {}
313 for name, model in models.items():
314     X_cv, y_cv = data_for_cv[name]
315     print(f"\nPerforming {n_splits}-Fold Cross-Validation for {
316           name}...")
317     scores = cross_val_score(model, X_cv, y_cv, cv=kf, scoring='
318           accuracy')
319     cv_results[name] = scores
320     print(f"Scores: {scores}")
321     print(f"Average Accuracy: {scores.mean():.4f}")
322 # Display results in a table
323 print("\n--- K-Fold Cross-Validation Results (K=5) ---")
324 cv_results_df = pd.DataFrame(cv_results)
325 cv_results_df.index = [f"Fold {i+1}" for i in range(n_splits)]
326 cv_results_df.loc["Average"] = cv_results_df.mean()
327 display(cv_results_df)

```

5 Output Screenshots



(a) Class Distribution of Dataset

GaussianNB Classification Report:				
	precision	recall	f1-score	support
0	0.9654	0.7509	0.8448	558
1	0.7146	0.9587	0.8188	363
accuracy			0.8328	921
macro avg	0.8400	0.8548	0.8318	921
weighted avg	0.8666	0.8328	0.8345	921
MultinomialNB Classification Report:				
	precision	recall	f1-score	support
0	0.8121	0.8208	0.8164	558
1	0.7199	0.7080	0.7139	363
accuracy			0.7763	921
macro avg	0.7660	0.7644	0.7651	921
weighted avg	0.7757	0.7763	0.7760	921
BernoulliNB Classification Report:				
	precision	recall	f1-score	support
0	0.8788	0.9229	0.9003	558
1	0.8716	0.8044	0.8367	363
accuracy			0.8762	921
macro avg	0.8752	0.8637	0.8685	921
weighted avg	0.8760	0.8762	0.8753	921

(b) Naive Bayes Output

Figure 1: Comparison of Class Distribution and Naive Bayes Output

KNN (k=1) Classification Report:

	precision	recall	f1-score	support
0	0.9113	0.9211	0.9162	558
1	0.8768	0.8623	0.8694	363
accuracy			0.8929	921
macro avg	0.8940	0.8917	0.8928	921
weighted avg	0.8977	0.8979	0.8978	921

KNN (k=3) Classification Report:

	precision	recall	f1-score	support
0	0.9133	0.9247	0.9190	558
1	0.8820	0.8650	0.8734	363
accuracy			0.9012	921
macro avg	0.8976	0.8949	0.8962	921
weighted avg	0.9016	0.9012	0.9010	921

KNN (k=5) Classification Report:

	precision	recall	f1-score	support
0	0.9168	0.9283	0.9225	558
1	0.8876	0.8705	0.8790	363
accuracy			0.9055	921
macro avg	0.9022	0.8994	0.9008	921
weighted avg	0.9053	0.9055	0.9054	921

KNN (k=7) Classification Report:

	precision	recall	f1-score	support
0	0.9156	0.9337	0.9246	558
1	0.8949	0.8678	0.8811	363
accuracy			0.9077	921
macro avg	0.9053	0.9007	0.9028	921
weighted avg	0.9079	0.9077	0.9074	921

(a) KNN Output

KNN with KD_TREE Report:

	precision	recall	f1-score	support
0	0.9168	0.9283	0.9225	558
1	0.8876	0.8705	0.8790	363
accuracy			0.9055	921
macro avg	0.9022	0.8994	0.9008	921
weighted avg	0.9053	0.9055	0.9054	921

Training time: 0.0188 seconds

KNN with BALL_TREE Report:

	precision	recall	f1-score	support
0	0.9168	0.9283	0.9225	558
1	0.8876	0.8705	0.8790	363
accuracy			0.9055	921
macro avg	0.9022	0.8994	0.9008	921
weighted avg	0.9053	0.9055	0.9054	921

Training time: 0.0123 seconds

(b) KDTREE Vs BALLTREE Output

Figure 2: Comparison of KNN and Tree-based Outputs

Training SVM with LINEAR kernel

Classification Report for LINEAR kernel:

	precision	recall	f1-score	support
0	0.9349	0.9516	0.9432	558
1	0.9235	0.8981	0.9106	363
accuracy			0.9305	921
macro avg	0.9292	0.9248	0.9269	921
weighted avg	0.9304	0.9305	0.9303	921

Training Time: 3.6872 seconds

Training SVM with POLY kernel

Classification Report for POLY kernel:

	precision	recall	f1-score	support
0	0.7376	0.9875	0.8444	558
1	0.9598	0.4601	0.6220	363
accuracy			0.7796	921
macro avg	0.8487	0.7238	0.7332	921
weighted avg	0.8252	0.7796	0.7568	921

Training Time: 2.8629 seconds

(a) SVM Output

Training SVM with RBF kernel

Classification Report for RBF kernel:

	precision	recall	f1-score	support
0	0.9270	0.9552	0.9409	558
1	0.9277	0.8843	0.9055	363
accuracy			0.9273	921
macro avg	0.9274	0.9197	0.9232	921
weighted avg	0.9273	0.9273	0.9269	921

Training Time: 1.7237 seconds

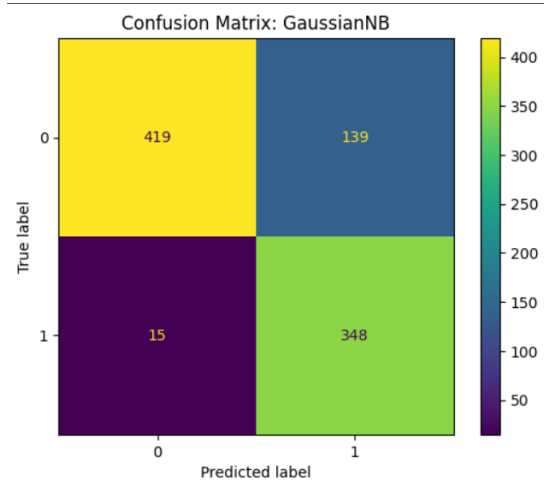
Training SVM with SIGMOID kernel

Classification Report for SIGMOID kernel:

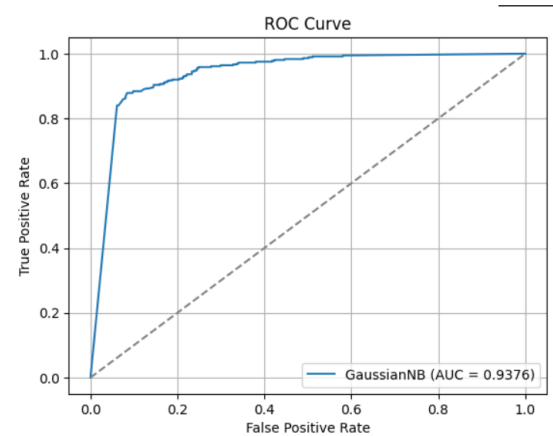
	precision	recall	f1-score	support
0	0.9034	0.9050	0.9042	558
1	0.8536	0.8512	0.8524	363
accuracy			0.8838	921
macro avg	0.8785	0.8781	0.8783	921
weighted avg	0.8838	0.8838	0.8838	921

Training Time: 1.8486 seconds

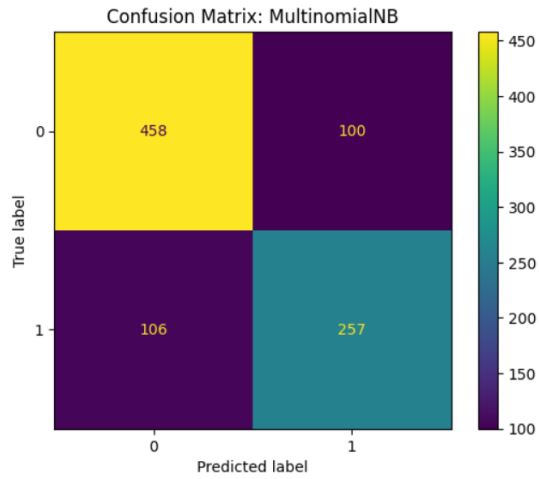
(b) SVM Output



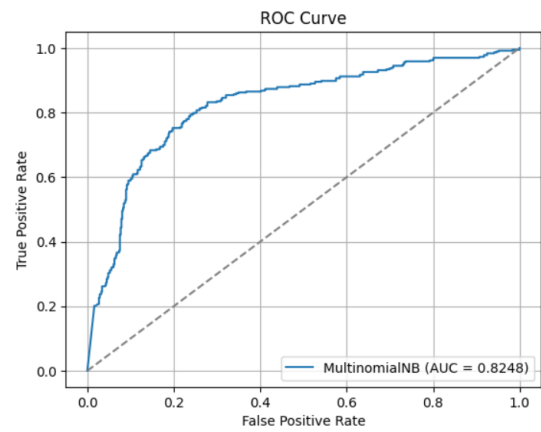
(a) GaussianNB confusion matrix



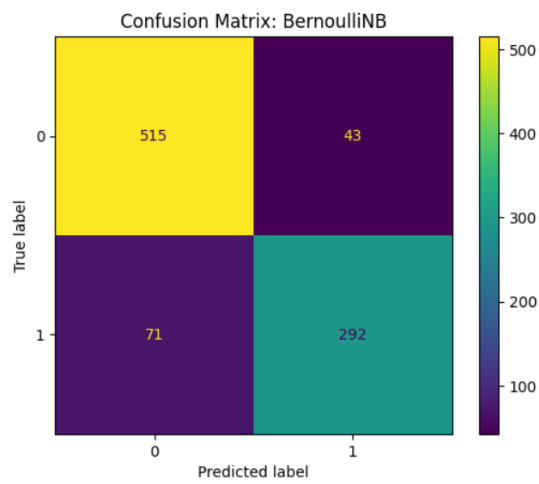
(b) GaussianNB ROC AUC



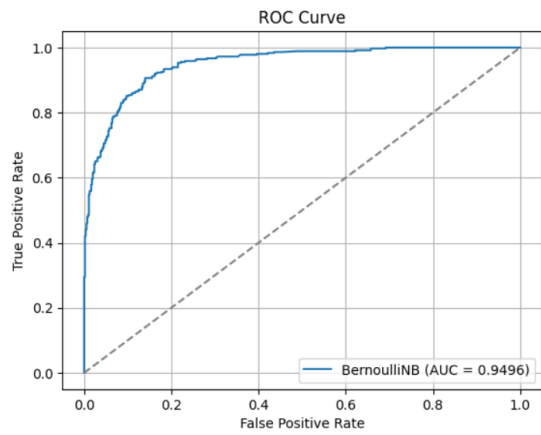
(a) MultinomialNB confusion matrix



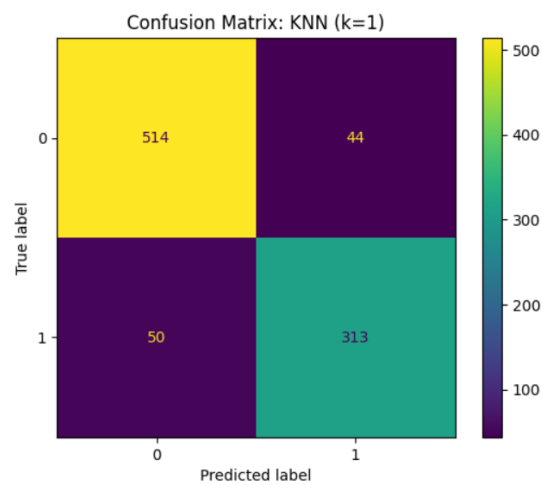
(b) MultinomialNB ROC AUC



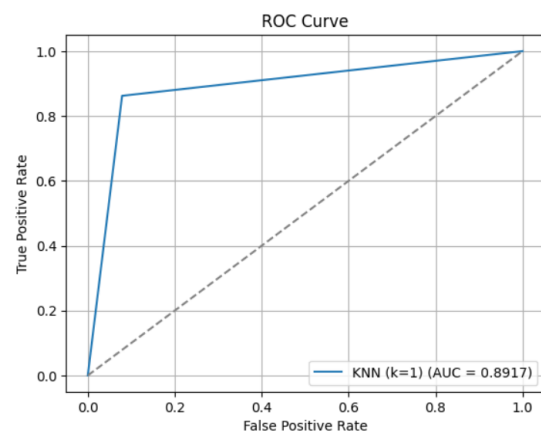
(a) BernoulliNB confusion matrix



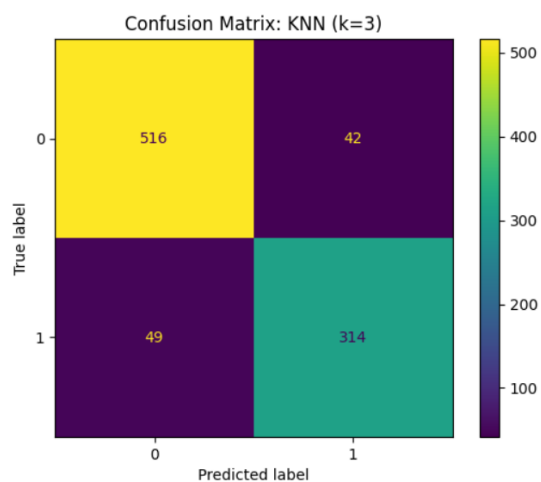
(b) BernoulliNB ROC AUC



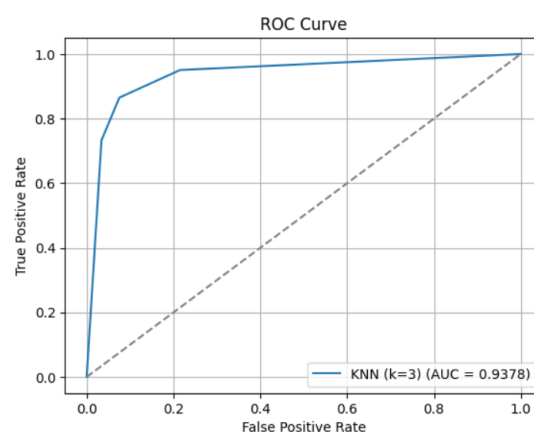
(a) KNN (k=1) confusion matrix



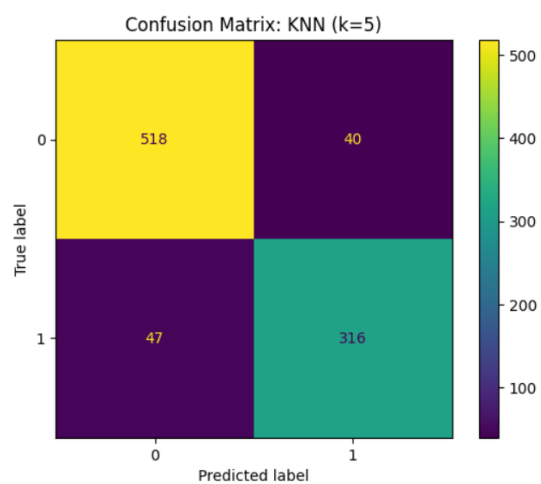
(b) KNN (k=1) ROC AUC



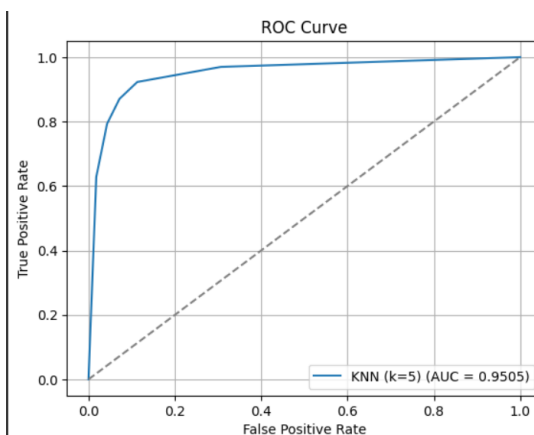
(a) KNN (k=3) confusion matrix



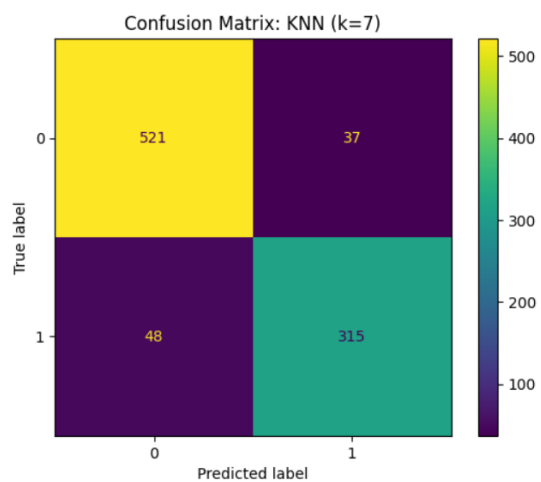
(b) KNN (k=3) ROC AUC



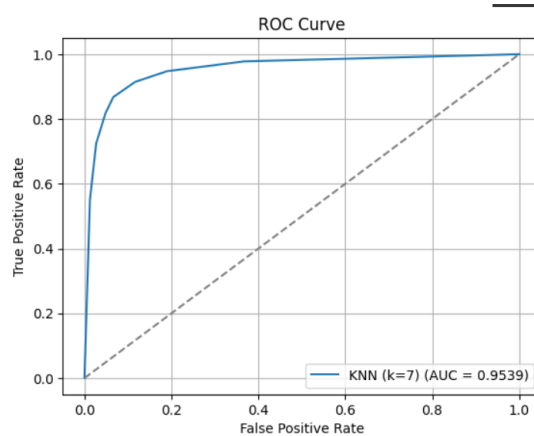
(a) KNN (k=5) confusion matrix



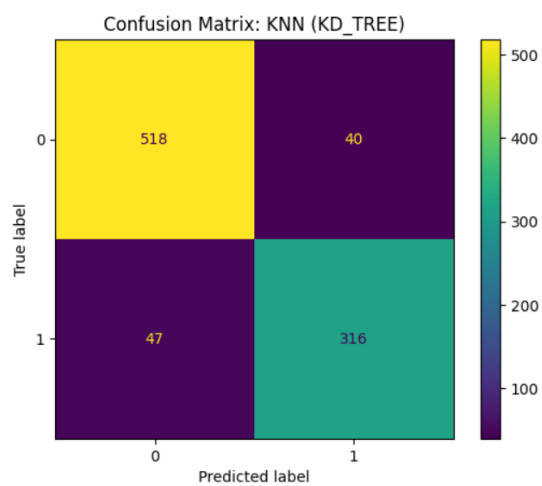
(b) KNN (k=5) ROC AUC



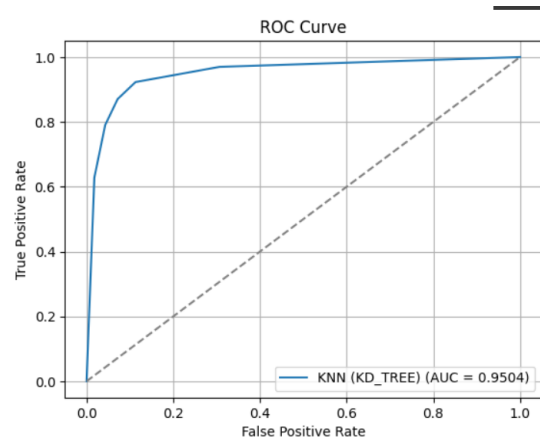
(a) KNN (k=7) confusion matrix



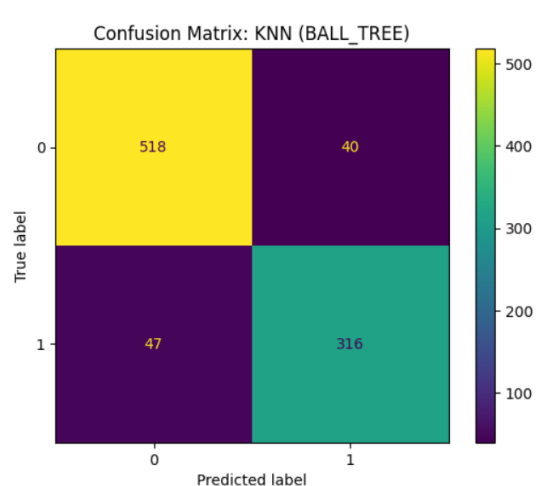
(b) KNN (k=7) ROC AUC



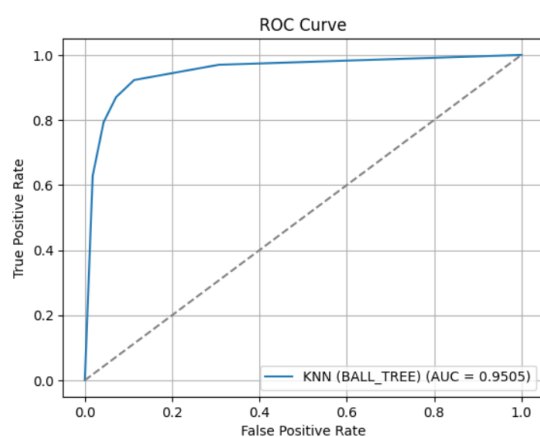
(a) KNN KDTREE confusion matrix



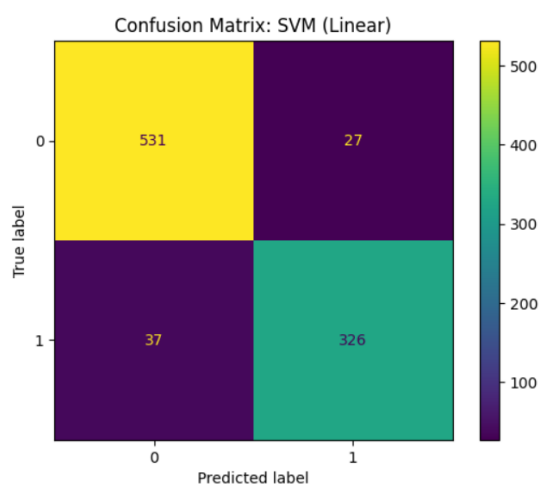
(b) KNN KDTREE ROC AUC



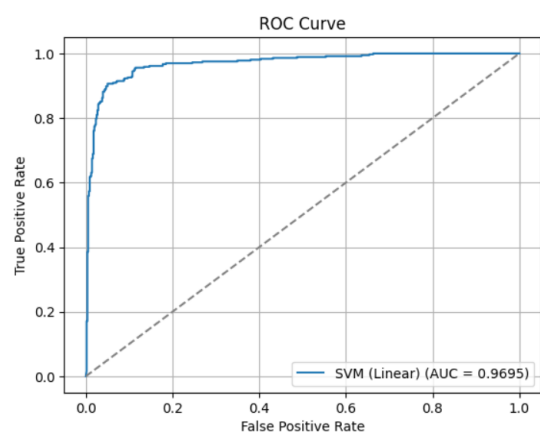
(a) KNN BALLTREE confusion matrix



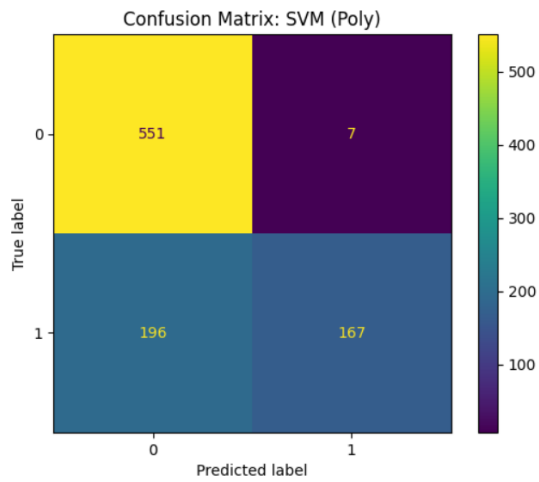
(b) KNN BALLTREE ROC AUC



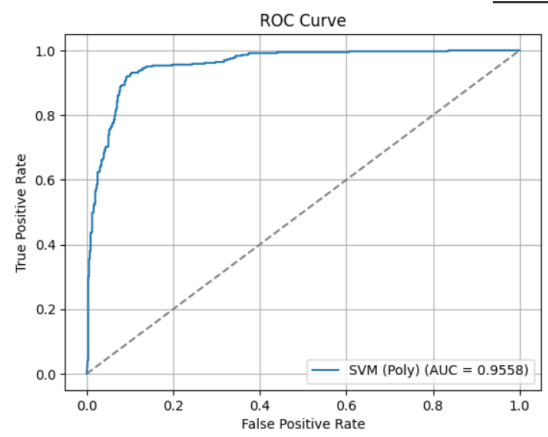
(a) SVM Linear confusion matrix



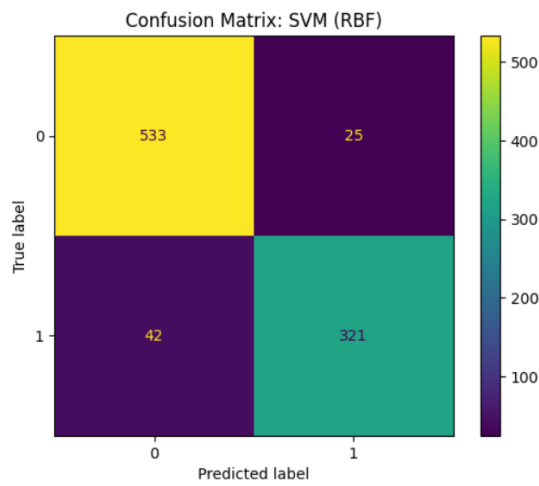
(b) SVM Linear ROC AUC



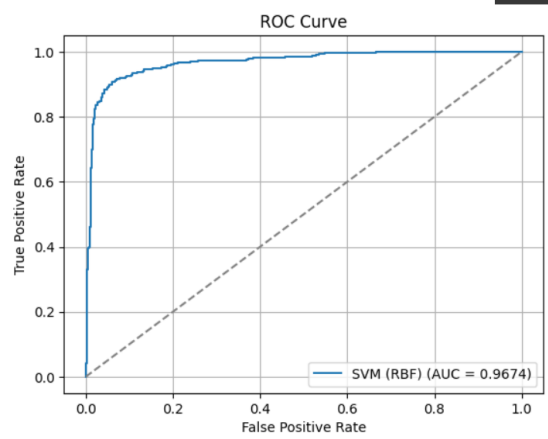
(a) SVM Poly confusion matrix



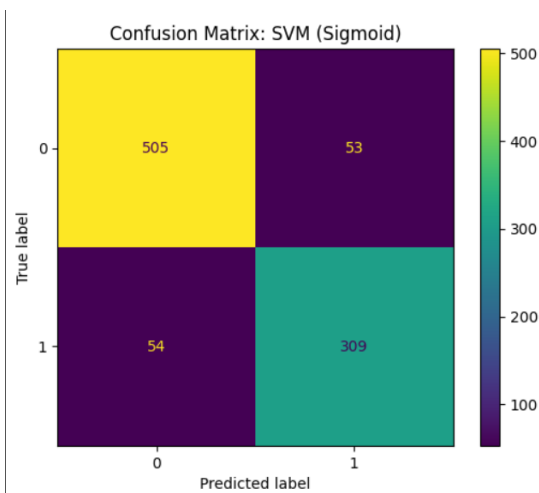
(b) SVM Poly ROC AUC



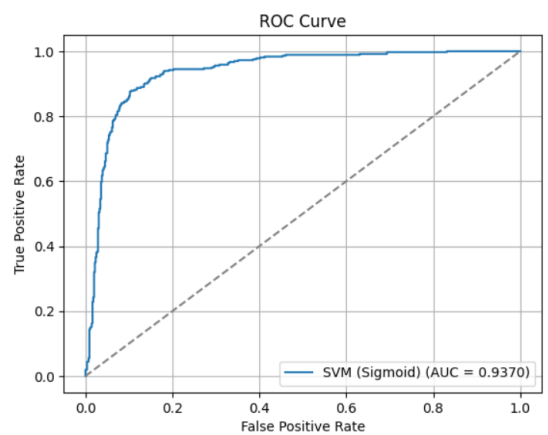
(a) SVM RBF confusion matrix



(b) SVM RBF ROC AUC



(a) SVM sigmoid confusion matrix



(b) SVM sigmoid ROC AUC

	Kernel	Accuracy	Precision	Recall	F1 Score	Training Time (s)
0	LINEAR	0.930510	0.923513	0.898072	0.910615	3.687203
1	POLY	0.779587	0.959770	0.460055	0.621974	2.862876
2	RBF	0.927253	0.927746	0.884298	0.905501	1.723683
3	SIGMOID	0.883822	0.853591	0.851240	0.852414	1.848590

(a) SVM Kernel-wise Results

	Naive Bayes (Gaussian)	Naive Bayes (Multinomial)	Naive Bayes (Bernoulli)	KNN (k=7)	SVM (Linear)
Fold 1	0.821933	0.766102	0.880565	0.895765	0.925081
Fold 2	0.803261	0.776087	0.890217	0.913043	0.928261
Fold 3	0.794565	0.780435	0.883696	0.913043	0.916304
Fold 4	0.822826	0.814130	0.886957	0.900000	0.936957
Fold 5	0.833696	0.801087	0.890217	0.911957	0.930435
Average	0.815256	0.791568	0.886330	0.906762	0.927408

(b) 5-Fold Cross-Validation

6 Observation

1) Which classifier had the best average accuracy?

- Based on the 5-fold cross-validation results, the **SVM (Linear)** classifier achieved the best average accuracy of **0.9274**. Other classifiers and their average accuracies were: **KNN (k=7)** at 0.9068, **Naive Bayes (Bernoulli)** at 0.8863, **Naive Bayes (Gaussian)** at 0.8153, and **Naive Bayes (Multinomial)** at 0.7916.

2) Which Naive Bayes variant worked best?

- The **Bernoulli Naive Bayes** variant performed the best among the three Naive Bayes models. It had an average accuracy of **0.8863** from 5-fold cross-validation, a precision of 0.8536, a recall of 0.8512, and an F1 score of 0.8524 on the test set. This was significantly higher than Gaussian Naive Bayes (average accuracy 0.8153) and Multinomial Naive Bayes (average accuracy 0.7916). The Bernoulli model also had the highest AUC score of **0.9496** among the Naive Bayes variants.

3) How did KNN accuracy vary with k and tree type?

- The accuracy of the KNN classifier generally **improved as the value of k increased**. On the test set, the accuracy was 0.9055 for k=5 and 0.9054 for k=7. The Area Under the Curve (AUC) also increased with k, from 0.8917 at k=1, to 0.9378 at k=3, 0.9505 at k=5, and **0.9539 at k=7**.
- The choice of tree type (**KDTREE** vs. **BALLTREE**) for KNN with k=5 had **no impact on accuracy or performance metrics**; both models achieved an accuracy of 0.9055 and a near-identical AUC (0.9504 and 0.9505, respectively). The only difference was that the Ball Tree algorithm had a slightly faster training time (0.0123 seconds) compared to the KD Tree algorithm (0.0188 seconds).

4) Which SVM kernel was most effective?

- The **linear kernel** was the most effective for the SVM classifier. It achieved the highest accuracy of **0.9305** on the test set, followed by the RBF kernel at 0.9273. The linear kernel also had the best overall performance metrics, including an F1 score of 0.9106, while also having the highest AUC of **0.9695**.
- While the RBF kernel had a higher recall for ham emails (class 0) at 0.9552 compared to the linear kernel's 0.9516, the linear kernel's performance for spam emails (class 1) was stronger, with a precision of 0.9235 compared to RBF's 0.9277, and recall of 0.8981 compared to RBF's 0.8843.

5) How did hyperparameters influence performance?

- Hyperparameters significantly influenced the performance of the models.
 - **KNN**: The hyperparameter **k** (the number of neighbors) directly affected the model's accuracy, precision, recall, and AUC. A larger k value (up to 7) generally resulted in higher accuracy and AUC, suggesting that considering more neighbors led to a more robust classification. The algorithm hyperparameter had a negligible effect on performance but did influence training time.
 - **SVM**: The choice of **kernel** was the most crucial hyperparameter for SVM. The linear and RBF kernels performed best, with the **linear kernel** having the highest accuracy and AUC. In contrast, the poly (polynomial) kernel performed poorly, with a low accuracy of 0.7796 and a low F1 score of 0.6220 for spam classification. The sigmoid kernel also performed worse than the linear and RBF kernels, with a lower accuracy of 0.8838. This demonstrates that some kernels are better suited for this specific dataset and classification problem.

7 Learning Outcomes

From this assignment,

- We gained practical experience in applying three fundamental machine learning models—Naïve Bayes, K-Nearest Neighbors (KNN), and Support Vector Machine (SVM)—to a real-world classification problem.
- We learned the importance of preparing data for machine learning, including handling missing values, standardizing or normalizing features, and splitting the dataset for training and testing.
- We explored how model performance is influenced by different hyperparameters, such as the k value in KNN and the kernel type in SVM.
- We learned to evaluate model effectiveness using a variety of metrics like accuracy, precision, recall, and F1-score. We also learned to visualize performance using confusion matrices and ROC curves to gain a deeper understanding of model behavior.
- We understood the significance of using K-fold cross-validation to get a more reliable estimate of a model's performance and avoid overfitting to a single train-test split.
- We were able to compare the strengths and weaknesses of different algorithms and their variants to determine which one is most suitable for a given dataset.