

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: Darshan Parthasarathy RegNo: 3122237001009

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

Refer to the submitted Jupyter Notebook “ML_Lab3.ipynb” containing all code for data preprocessing, model training, evaluation, and visualizations.

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

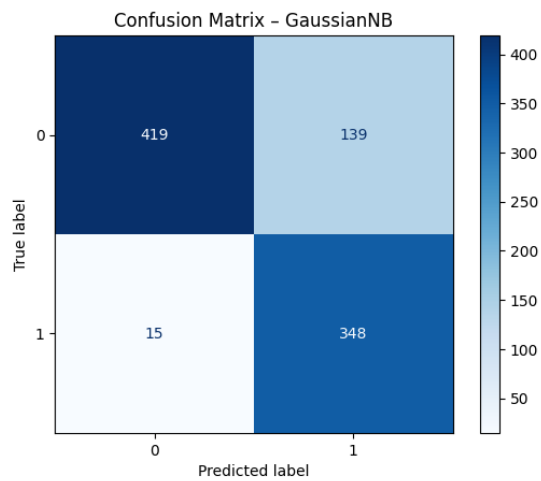


Figure 1: Confusion Matrix – GaussianNB

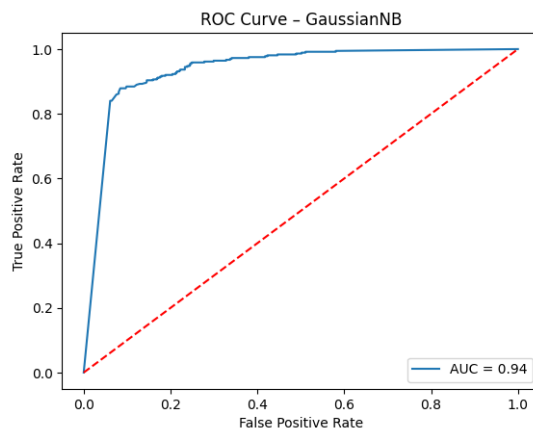


Figure 2: ROC Curve – GaussianNB

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

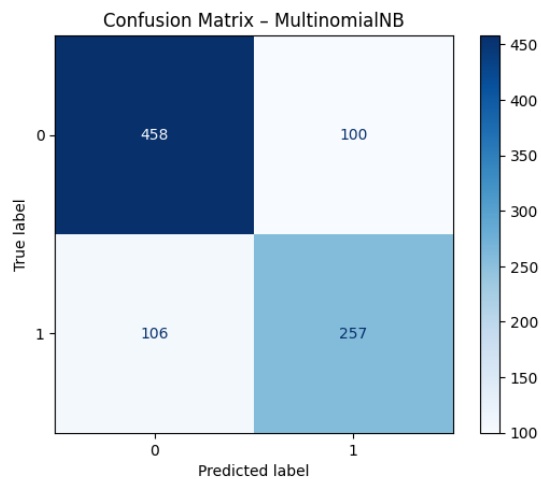


Figure 3: Confusion Matrix – MultinomialNB

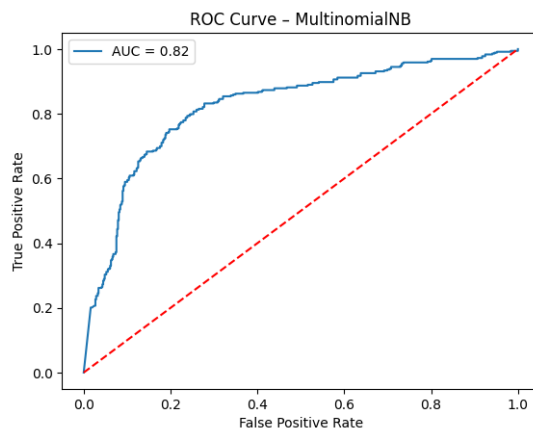


Figure 4: ROC Curve – MultinomialNB

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

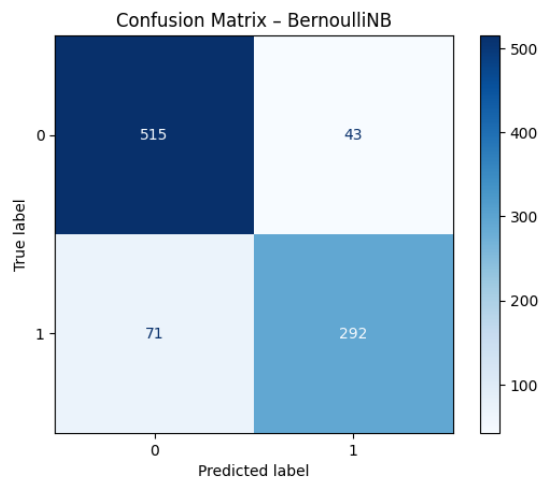


Figure 5: Confusion Matrix – BernoulliNB

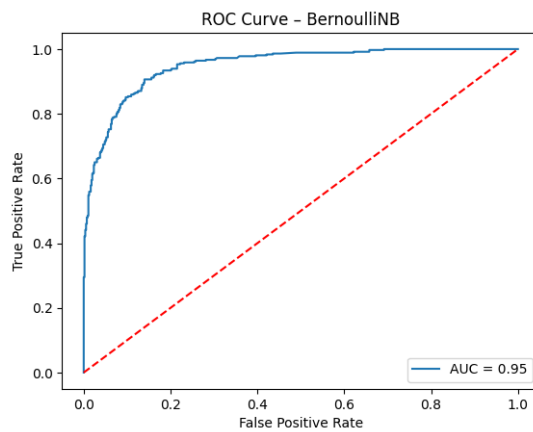


Figure 6: ROC Curve – BernoulliNB

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

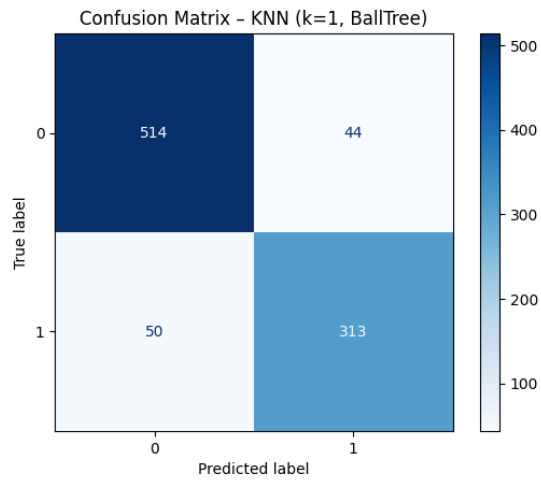


Figure 7: Confusion Matrix – KNN ($k = 1$, BallTree)

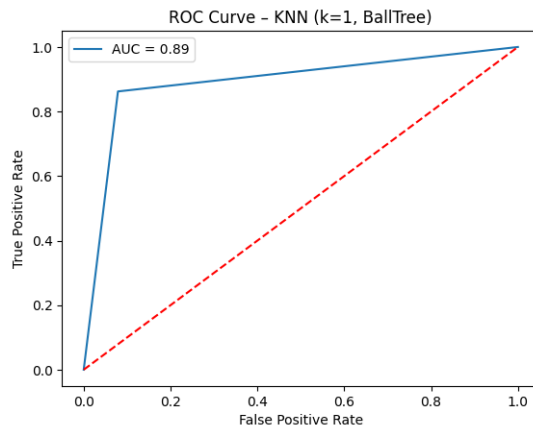


Figure 8: ROC Curve – KNN ($k = 1$, BallTree)

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

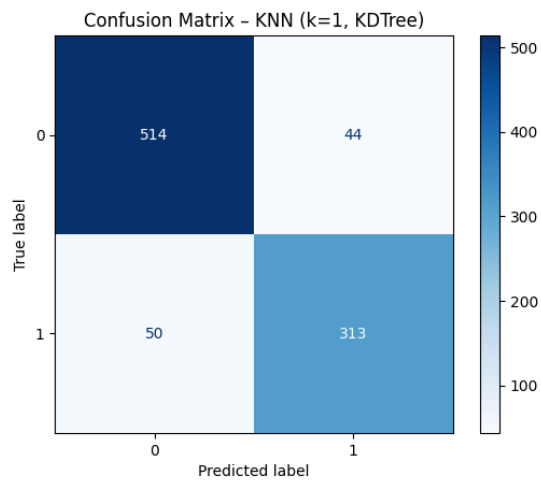


Figure 9: Confusion Matrix – KNN ($k = 1$, KDTree)

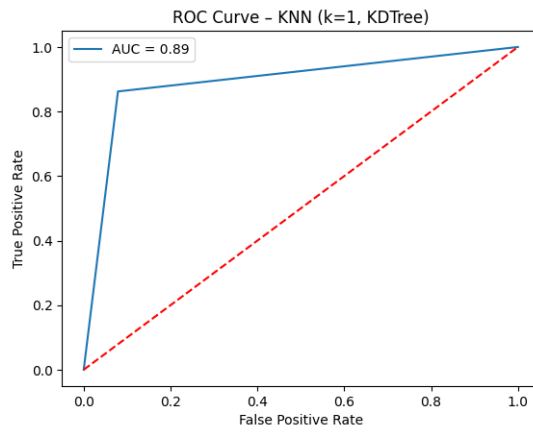


Figure 10: ROC Curve – KNN ($k = 1$, KDTree)

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.