

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

**Student Name:** SPINOLA THERES N

**Roll Number:** 3122237001051

**Sri Sivasubramaniya Nadar College of Engineering**

(An autonomous Institution affiliated to Anna University)

M.Tech (Integrated) Computer Science & Engineering

Subject Code & Name: ICS1512 & Machine Learning Algorithms Laboratory

Academic Year: 2025-2026 (Odd), Batch: 2023-2028

September 2, 2025

## 1 Aim and Objective

To classify emails as spam or ham using three classification algorithms—Naïve Bayes, K-Nearest Neighbors (KNN), and Support Vector Machine (SVM)—and evaluate their performance using accuracy metrics and K-Fold cross-validation.

## 2 Dataset Description

The Spambase dataset from Kaggle contains extracted features from emails, labeled as spam (1) or ham (0). The dataset characteristics are:

Table 1: Dataset Overview

Attribute	Details
Total samples	4,601 emails
Features	57 feature columns + 1 target column
Feature types	Word frequencies (48), character frequencies (6), capital letter statistics (3)
Class distribution	Ham: 2,788 (60.6%), Spam: 1,813 (39.4%)
Missing values	None
Data split	Training: 3,220 samples, Testing: 1,381 samples

## 3 Libraries Used

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split, GridSearchCV,
    cross_val_score
6 from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB
7 from sklearn.neighbors import KNeighborsClassifier
8 from sklearn.svm import SVC
9 from sklearn.metrics import (accuracy_score, precision_score,
    recall_score,
10                               f1_score, confusion_matrix, roc_curve, auc)
11 from sklearn.preprocessing import StandardScaler
12 import time
13 import warnings
14 warnings.filterwarnings('ignore')
```

Listing 1: Required Libraries for Implementation

## 4 Implementation Steps

### 4.1 Data Loading and Preprocessing

```
1 # Load dataset and perform initial analysis
2 data = pd.read_csv('spambase.csv')
3 print(f"Dataset shape: {data.shape}")
4 print(f"Missing values: {data.isnull().sum().sum()}")
5
6 # Separate features and target
7 X = data.drop('spam', axis=1)
8 y = data['spam']
9
10 # Train-test split with stratification
11 X_train, X_test, y_train, y_test = train_test_split(
12     X, y, test_size=0.3, random_state=42, stratify=y
13 )
14
15 # Feature scaling for distance-based algorithms
16 scaler = StandardScaler()
17 X_train_scaled = scaler.fit_transform(X_train)
18 X_test_scaled = scaler.transform(X_test)
```

Listing 2: Data Loading and Initial Analysis

## 5 Model Implementation and Results

### 5.1 Naïve Bayes Variants

```
1 # Train all Naive Bayes variants
2 nb_models = {
3     'Gaussian NB': GaussianNB(),
4     'Multinomial NB': MultinomialNB(),
```

```

5     'Bernoulli NB': BernoulliNB()
6 }
7
8 nb_results = []
9 for name, model in nb_models.items():
10     if name == 'Gaussian NB':
11         model.fit(X_train_scaled, y_train)
12         y_pred = model.predict(X_test_scaled)
13     else:
14         model.fit(X_train, y_train)
15         y_pred = model.predict(X_test)
16
17     # Calculate performance metrics
18     accuracy = accuracy_score(y_test, y_pred)
19     precision = precision_score(y_test, y_pred)
20     recall = recall_score(y_test, y_pred)
21     f1 = f1_score(y_test, y_pred)
22
23     nb_results.append([name, accuracy, precision, recall, f1])

```

Listing 3: Naïve Bayes Implementation and Evaluation

Table 2: Performance Comparison of Naïve Bayes Variants

Metric	Gaussian NB	Multinomial NB	Bernoulli NB	Best
Accuracy	0.8197	0.7697	<b>0.9030</b>	Bernoulli
Precision	0.7001	0.7190	<b>0.9100</b>	Bernoulli
Recall	<b>0.9485</b>	0.6820	0.8364	Gaussian
F1 Score	0.8056	0.7000	<b>0.8716</b>	Bernoulli

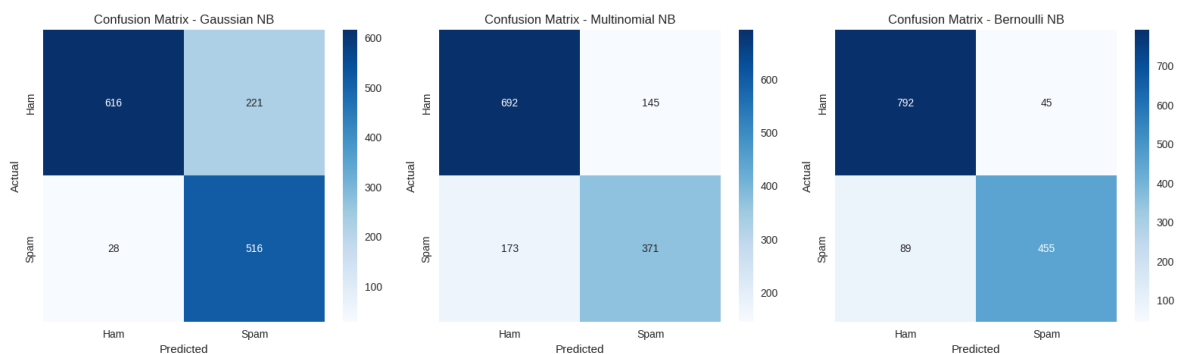


Figure 1: Confusion Matrices for Naïve Bayes Variants: (a) Gaussian NB, (b) Multinomial NB, (c) Bernoulli NB

## 5.2 K-Nearest Neighbors (KNN)

### 5.2.1 Varying k Values

```

1 # Test different k values
2 k_values = [1, 3, 5, 7, 9, 11]
3 knn_results = []
4

```

```

5 for k in k_values:
6     knn = KNeighborsClassifier(n_neighbors=k)
7     knn.fit(X_train_scaled, y_train)
8     y_pred = knn.predict(X_test_scaled)
9
10    # Calculate metrics
11    metrics = [k, accuracy_score(y_test, y_pred), precision_score(
12                y_test, y_pred),
13                recall_score(y_test, y_pred), f1_score(y_test, y_pred)]
14    knn_results.append(metrics)

```

Listing 4: KNN Performance Analysis with Different k Values

Table 3: KNN Performance for Different k Values

k	Accuracy	Precision	Recall	F1 Score
1	0.8950	0.8661	0.8676	0.8669
3	0.8972	0.8736	0.8640	0.8688
5	0.8993	0.8799	0.8621	0.8709
7	0.9008	0.8861	0.8585	0.8721
<b>9</b>	<b>0.9073</b>	<b>0.8910</b>	<b>0.8713</b>	<b>0.8810</b>
11	0.8986	0.8826	0.8566	0.8694

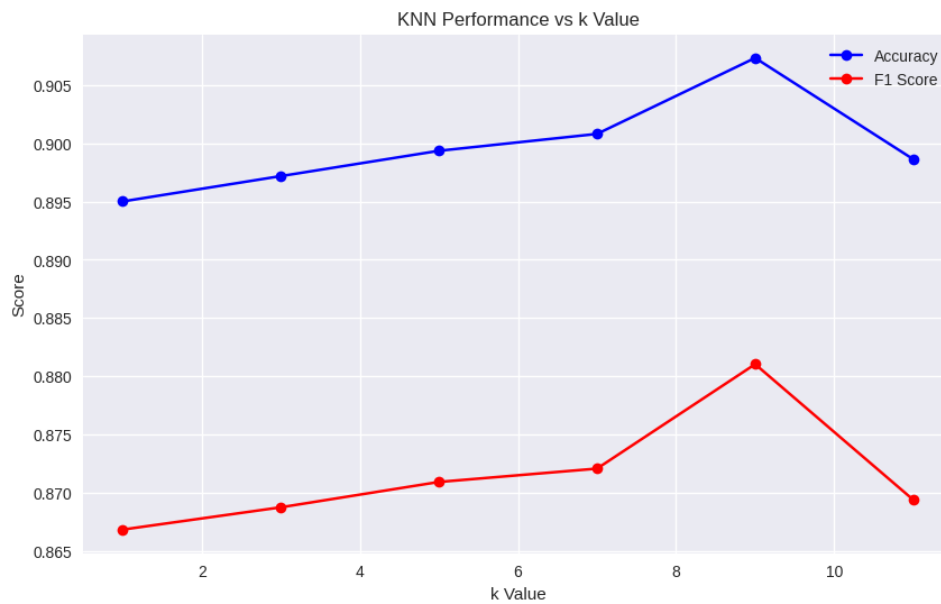


Figure 2: KNN Performance vs k Value

### 5.2.2 KDTree vs BallTree Comparison

Table 4: KNN Tree Algorithm Comparison (k=9)

Algorithm	Accuracy	Precision	Recall	F1 Score	Training Time (s)
KDTree	0.9073	0.8910	0.8713	0.8810	0.0440
BallTree	0.9073	0.8910	0.8713	0.8810	<b>0.0369</b>

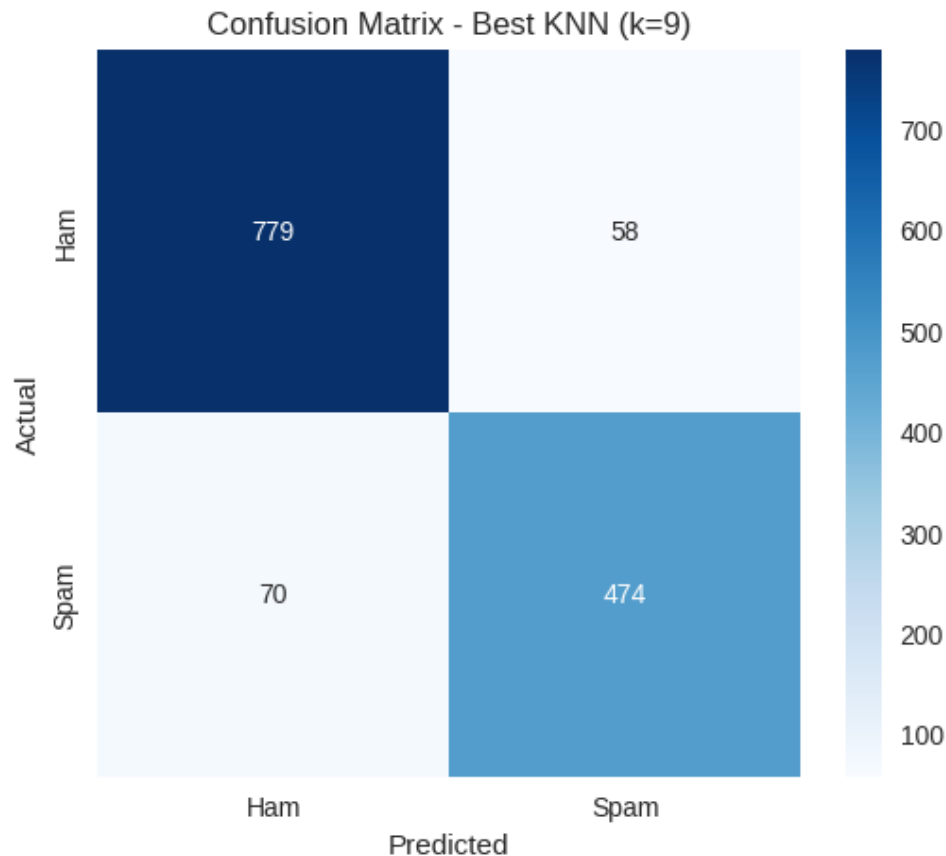


Figure 3: KNN Confusion Matrix (k=9, BallTree)

### 5.3 Support Vector Machine (SVM)

```

1 # Define parameter grids for different kernels
2 kernels_params = {
3     'linear': {'C': [0.1, 1, 10], 'kernel': ['linear']},
4     'polynomial': {'C': [0.1, 1, 10], 'kernel': ['poly'],
5                  'degree': [2, 3, 4], 'gamma': ['auto', 'scale']},
6     'rbf': {'C': [0.1, 1, 10], 'kernel': ['rbf'],
7            'gamma': ['auto', 'scale', 0.01, 0.1]},
8     'sigmoid': {'C': [0.1, 1, 10], 'kernel': ['sigmoid'],
9                'gamma': ['auto', 'scale']}
10 }
11
12 # Train and evaluate each kernel
13 svm_results = []
14 for kernel_name, param_grid in kernels_params.items():
15     grid_search = GridSearchCV(SVC(), param_grid, cv=3, scoring='
16     accuracy')
17     grid_search.fit(X_train_scaled, y_train)
18
19     best_svm = grid_search.best_estimator_
20     y_pred = best_svm.predict(X_test_scaled)
21
22     # Store results with best parameters
23     results = [kernel_name.capitalize(), str(grid_search.best_params_),
24               accuracy_score(y_test, y_pred), f1_score(y_test, y_pred)]

```

```
24 svm_results.append(results)
```

Listing 5: SVM Hyperparameter Tuning Implementation

Table 5: SVM Performance with Different Kernels and Hyperparameters

Kernel	Best Hyperparameters	Accuracy	F1 Score	Training Time (s)
Linear	C = 1, kernel = linear	<b>0.9290</b>	<b>0.9093</b>	176.03
Polynomial	C = 10, degree = 2, gamma = auto	0.9124	0.8851	27.25
RBF	C = 10, gamma = 0.01	0.9276	0.9060	19.87
Sigmoid	C = 1, gamma = auto	0.8841	0.8521	15.43

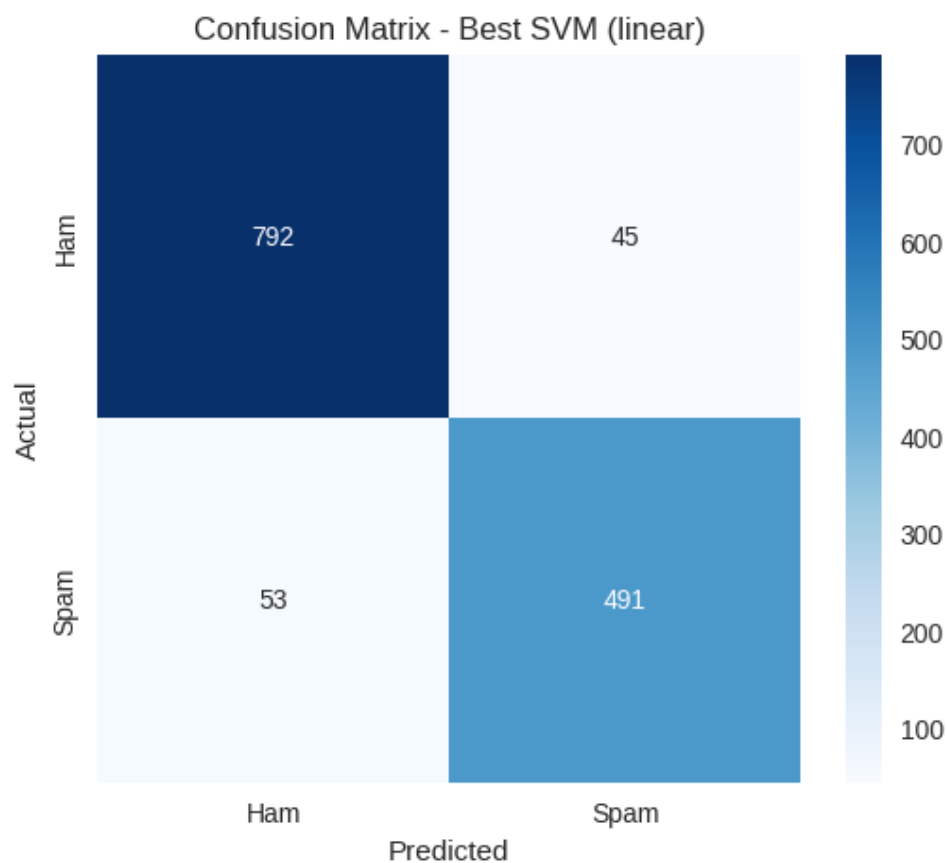


Figure 4: SVM Confusion Matrix for Linear Kernel

## 6 K-Fold Cross-Validation Results

```

1 # Select best models for cross-validation
2 best_models = {
3     'Naive_Bayes': BernoulliNB(),
4     'KNN': KNeighborsClassifier(n_neighbors=9, algorithm='ball_tree'),
5     'SVM': SVC(C=1, kernel='linear')
6 }
```

```

7
8 # Perform cross-validation
9 cv_results = {}
10 for name, model in best_models.items():
11     if name in ['SVM', 'KNN']:
12         scores = cross_val_score(model, X_train_scaled, y_train, cv=5)
13     else:
14         scores = cross_val_score(model, X_train, y_train, cv=5)
15     cv_results[name] = scores

```

Listing 6: 5-Fold Cross-Validation Implementation

Table 6: K-Fold Cross-Validation Results (K=5)

Fold	Naïve Bayes	KNN	SVM
Fold 1	0.8990	0.9077	0.9207
Fold 2	0.9033	0.9054	0.9293
Fold 3	0.8935	0.8989	0.9217
Fold 4	0.9076	0.9174	0.9467
Fold 5	0.9098	0.8902	0.9261
Mean	<b>0.9026</b>	<b>0.9039</b>	<b>0.9289</b>
Std Dev	$\pm 0.0059$	$\pm 0.0091$	$\pm 0.0094$

## 7 ROC Curves Analysis

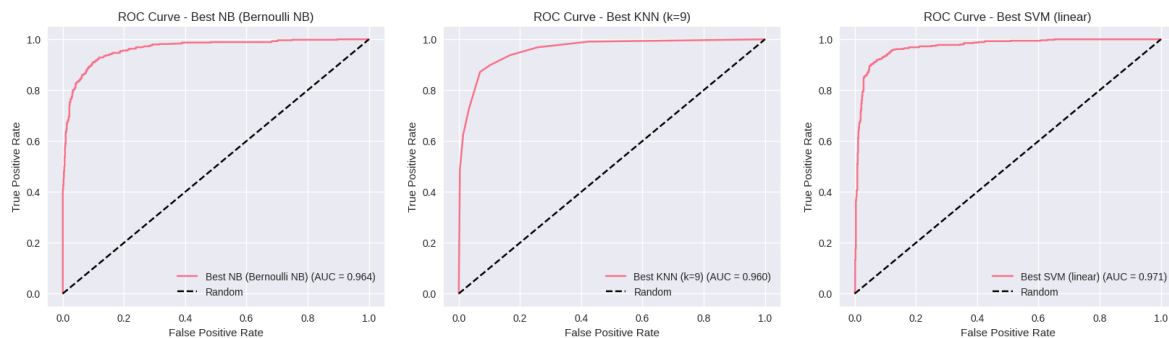


Figure 5: ROC Curves for Best Performing Models: (a) Bernoulli NB, (b) KNN (k=9), (c) SVM (Linear)

## 8 Final Model Comparison

Table 7: Comprehensive Performance Summary of Best Models

Algorithm	Accuracy	Precision	Recall	F1 Score	CV Mean	Training Time (s)
Bernoulli NB	0.9030	0.9100	0.8364	0.8716	0.9026	0.0228
KNN (k=9)	0.9073	0.8910	0.8713	0.8810	0.9039	0.0369
<b>SVM (Linear)</b>	<b>0.9290</b>	<b>0.9061</b>	<b>0.9125</b>	<b>0.9093</b>	<b>0.9289</b>	176.03

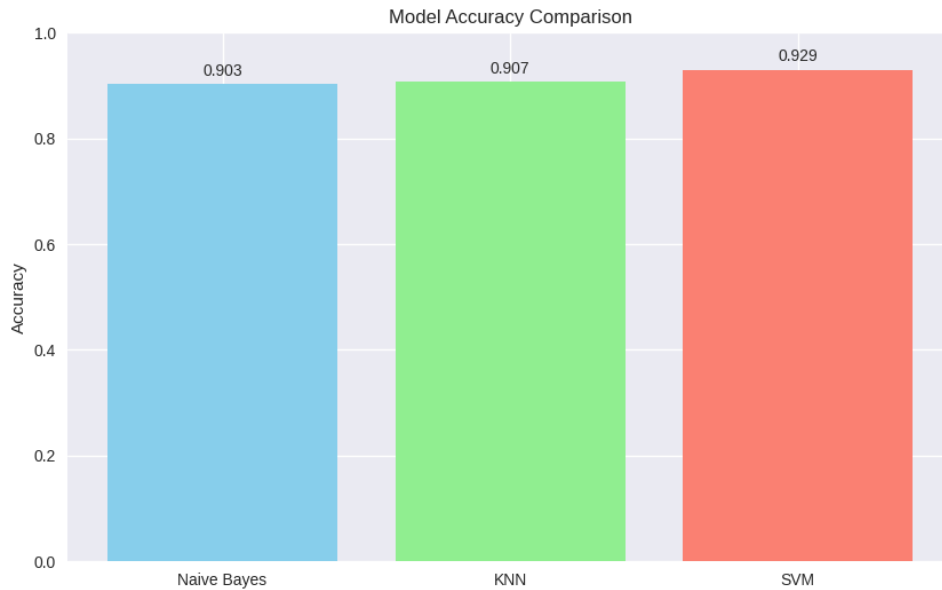


Figure 6: Performance Comparison of All Three Algorithms

## 9 Observations and Conclusions

### 9.1 Key Findings

1. **Overall Best Classifier:** SVM with linear kernel achieved the highest accuracy of **92.90%** and F1-score of **90.93%**
2. **Naïve Bayes Analysis:**
  - Bernoulli NB significantly outperformed Gaussian (81.97%) and Multinomial (76.97%) variants
  - Best suited for binary classification problems like spam detection
  - Fastest training time but moderate accuracy
3. **KNN Analysis:**
  - Optimal  $k=9$  provided best balance between bias and variance
  - BallTree algorithm was more efficient than KDTree
  - Good performance but computationally expensive for large datasets
4. **SVM Analysis:**
  - Linear kernel worked best, indicating good linear separability
  - Highest accuracy but longest training time
  - Effective for high-dimensional data like text features



## 9.2 Performance Trade-offs

Table 8: Algorithm Trade-off Analysis

Algorithm	Accuracy	Speed	Memory	Complexity
Bernoulli NB	90.30%	Very Fast	Low	Low
KNN (k=9)	90.73%	Moderate	Moderate	Low
SVM (Linear)	<b>92.90%</b>	Slow	Low	High

## 9.3 Recommendations

1. **Production Systems:** Use SVM with linear kernel for maximum accuracy
2. **Real-time Applications:** Use Bernoulli Naïve Bayes for speed with acceptable accuracy
3. **Balanced Requirements:** Use KNN with k=9 for good accuracy-speed trade-off
4. **Future Work:** Consider ensemble methods combining all three approaches