**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 |

**Keerthana R CSE-B 3122235001066**

**Experiment 4: Binary Classification using Linear and Kernel-Based Models**

# Objective

To classify emails as spam or ham (non-spam) using Logistic Regression and Support Vector Machine (SVM) classifiers, perform hyperparameter tuning, and analyze the effect of different kernels and regularization techniques on classification performance.

**Specific Goals:**

1. Implement Logistic Regression with L1 and L2 regularization
2. Implement SVM with multiple kernels (Linear, Polynomial, RBF, Sigmoid)
3. Perform hyperparameter tuning using Grid Search
4. Compare probabilistic vs margin-based classifiers
5. Analyze bias-variance trade-off
6. Evaluate models using multiple performance metrics
7. Perform k-fold cross-validation for robust evaluation

# Dataset Description

## Spambase Dataset

**Source:** UCI Machine Learning Repository / Kaggle
**Purpose:** Email spam detection
**Size:** 4,601 emails
**Features:** 57 numerical features + 1 binary target
**Classes:** 0 = Ham (non-spam), 1 = Spam

## Feature Categories

1. **Word Frequency (48 features):** Percentage of specific words (e.g., "free", "money", "business")
2. **Character Frequency (6 features):** Percentage of characters: ;, (, [, !, $, #
3. **Capital Letter Statistics (3 features):**
   - Average length of capital letter sequences
   - Longest capital letter sequence
   - Total number of capital letters

### Class Distribution

- Ham (Class 0): 2,788 emails (60.6%)
- Spam (Class 1): 1,813 emails (39.4%)
- Moderately imbalanced dataset

# Mathematical and Theoretical Foundation

### 1. Logistic Regression

Logistic Regression is a probabilistic linear classifier that models the probability of class membership using the sigmoid function.

### 1.1 Sigmoid Function

$$\sigma(z) = \frac{1}{1 + e^{-z}} \tag{1}$$

where $z = \mathbf{w}^T\mathbf{x} + b$ is the linear combination of features.

### 1.2 Probability Model

The probability that a sample $\mathbf{x}$ belongs to class 1:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T\mathbf{x}+b)}} \tag{2}$$

### 1.3 Decision Rule

$$\hat{y} = \begin{cases} 1 & \text{if } P(y = 1|\mathbf{x}) \geq 0.5 \\ 0 & \text{otherwise} \end{cases} \tag{3}$$

### 1.4 Loss Function

Logistic Regression minimizes the log-loss (binary cross-entropy):

$$\text{Loss} = -\frac{1}{m}\sum_{i=1}^{m}[y_i \log(\hat{p}_i) + (1 - y_i)\log(1 - \hat{p}_i)] \tag{4}$$

where $\hat{p}_i = P(y = 1|\mathbf{x}_i)$ is the predicted probability.

### 1.5 Regularization

**L1 Regularization (Lasso):**

$$\text{Loss}_{\text{L1}} = -\frac{1}{m}\sum_{i=1}^{m}[y_i \log(\hat{p}_i) + (1 - y_i)\log(1 - \hat{p}_i)] + \frac{1}{C}\|\mathbf{w}\|_1 \tag{5}$$

- Encourages sparsity (some coefficients $\to 0$)
- Performs automatic feature selection
- Useful when many features are irrelevant

**L2 Regularization (Ridge):**

$$\text{Loss}_{\text{L2}} = -\frac{1}{m}\sum_{i=1}^{m}[y_i\log(\hat{p}_i) + (1-y_i)\log(1-\hat{p}_i)] + \frac{1}{2C}\|\mathbf{w}\|^2 \tag{6}$$

- Shrinks all coefficients proportionally
- Does not perform feature selection
- Better for correlated features

### 1.6 Hyperparameter C

$C$ is the inverse of regularization strength:

- **Small C** (e.g., 0.01): Strong regularization, simpler model, higher bias
- **Large C** (e.g., 100): Weak regularization, complex model, higher variance

## 2. Support Vector Machine (SVM)

SVM is a margin-based classifier that finds the optimal hyperplane separating classes by maximizing the margin.

### 2.1 Linear SVM

For linearly separable data, SVM finds the hyperplane:

$$\mathbf{w}^T\mathbf{x} + b = 0 \tag{7}$$

that maximizes the margin:

$$\text{Margin} = \frac{2}{\|\mathbf{w}\|} \tag{8}$$

### 2.2 Optimization Problem

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{m}\xi_i \tag{9}$$

subject to:

$$y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0 \tag{10}$$

where $\xi_i$ are slack variables allowing misclassification.

### 2.3 SVM Kernels

For non-linearly separable data, SVM uses the kernel trick to map data to higher dimensions.

**Linear Kernel:**

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T\mathbf{x}_j \tag{11}$$

**Polynomial Kernel:**

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\gamma\mathbf{x}_i^T\mathbf{x}_j + r)^d \tag{12}$$

**RBF (Radial Basis Function) Kernel:**

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma\|\mathbf{x}_i - \mathbf{x}_j\|^2) \tag{13}$$

**Sigmoid Kernel:**

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\gamma\mathbf{x}_i^T\mathbf{x}_j + r) \tag{14}$$

**2.4 SVM Hyperparameters**

**C (Regularization Parameter):**

- **Small C**: Wide margin, more misclassifications, higher bias
- **Large C**: Narrow margin, fewer misclassifications, higher variance

**$\gamma$ (Kernel Coefficient):**

- **Small $\gamma$**: Smooth decision boundary, higher bias
- **Large $\gamma$**: Complex decision boundary, higher variance
- **'scale'**: $\gamma = \frac{1}{n_{\text{features}} \times \text{Var}(X)}$
- **'auto'**: $\gamma = \frac{1}{n_{\text{features}}}$

## 3. Performance Metrics

For binary classification with TP, TN, FP, FN:

**Accuracy:**

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \tag{15}$$

**Precision:**

$$\text{Precision} = \frac{TP}{TP + FP} \tag{16}$$

**Recall:**

$$\text{Recall} = \frac{TP}{TP + FN} \tag{17}$$

**F1 Score:**

$$F1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \tag{18}$$

# Preprocessing Steps

## 1. Data Loading and Exploration

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, GridSearchCV,
    cross_val_score
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.metrics import (accuracy_score, precision_score, recall_score,
                             f1_score, confusion_matrix, classification_report)
import time
import warnings
warnings.filterwarnings('ignore')

# Load dataset
df = pd.read_csv('spambase.csv')

# Display basic information
```

```
19 print("Dataset Shape:", df.shape)
20 print("\nFirst few rows:")
21 print(df.head())
22 print("\nDataset Info:")
23 print(df.info())
24 print("\nMissing Values:")
25 print(df.isnull().sum().sum())
26 print("\nClass Distribution:")
27 print(df['spam'].value_counts())
```

Listing 1: Loading the Spambase Dataset

## 2. Feature Standardization

```
1  # Separate features and target
2  X = df.drop('spam', axis=1)
3  y = df['spam']
4
5  # Split data (80-20)
6  X_train, X_test, y_train, y_test = train_test_split(
7      X, y, test_size=0.2, random_state=42, stratify=y
8  )
9
10 # Standardize features (critical for SVM)
11 scaler = StandardScaler()
12 X_train_scaled = scaler.fit_transform(X_train)
13 X_test_scaled = scaler.transform(X_test)
14
15 print(f"Training set: {X_train.shape}")
16 print(f"Test set: {X_test.shape}")
17 print(f"Class distribution in training: {y_train.value_counts()}")
```

Listing 2: Feature Scaling

## 3. Exploratory Data Analysis

```
1  # Class distribution
2  plt.figure(figsize=(8, 5))
3  y.value_counts().plot(kind='bar', color=['green', 'red'], alpha=0.7)
4  plt.title('Class Distribution: Ham vs Spam')
5  plt.xlabel('Class (0=Ham, 1=Spam)')
6  plt.ylabel('Count')
7  plt.xticks(rotation=0)
8  plt.show()
9
10 # Feature correlation with target
11 correlations = X.corrwith(y).sort_values(ascending=False)
12 print("Top 10 features correlated with spam:")
13 print(correlations.head(10))
14
15 # Correlation heatmap (top features)
16 plt.figure(figsize=(10, 8))
17 top_features = correlations.abs().nlargest(15).index
18 corr_matrix = df[list(top_features) + ['spam']].corr()
19 sns.heatmap(corr_matrix, annot=True, fmt='.2f', cmap='coolwarm')
20 plt.title('Correlation Heatmap of Top Features')
```

```
21 plt.tight_layout ()
22 plt.show ()
```
Listing 3: EDA Visualizations

# Implementation Details

## 1. Baseline Logistic Regression

```
1  # Train baseline Logistic Regression (L2, default C=1)
2  print("="*60)
3  print("Baseline Logistic Regression")
4  print("="*60)
5
6  start_time = time.time ()
7  lr_baseline = LogisticRegression(random_state=42, max_iter=1000)
8  lr_baseline.fit(X_train_scaled, y_train)
9  lr_baseline_time = time.time () - start_time
10
11 # Predictions
12 y_pred_lr_baseline = lr_baseline.predict(X_test_scaled)
13
14 # Evaluation
15 print(f"\nBaseline Logistic Regression Results:")
16 print(f"Training Time: {lr_baseline_time:.4f}s")
17 print(f"Accuracy: {accuracy_score(y_test, y_pred_lr_baseline):.4f}")
18 print(f"Precision: {precision_score(y_test, y_pred_lr_baseline):.4f}")
19 print(f"Recall: {recall_score(y_test, y_pred_lr_baseline):.4f}")
20 print(f"F1 Score: {f1_score(y_test, y_pred_lr_baseline):.4f}")
21 print("\nConfusion Matrix:")
22 print(confusion_matrix(y_test, y_pred_lr_baseline))
```
Listing 4: Baseline Logistic Regression

## 2. Logistic Regression with Hyperparameter Tuning

```
1  # Define parameter grid
2  param_grid_lr = {
3      'penalty': ['l1', 'l2'],
4      'C': [0.01, 0.1, 1, 10, 100],
5      'solver': ['liblinear', 'saga']
6  }
7
8  # Grid Search with 5-fold CV
9  print("\n" + "="*60)
10 print("Logistic Regression Hyperparameter Tuning")
11 print("="*60)
12
13 grid_lr = GridSearchCV(
14     LogisticRegression(random_state=42, max_iter=1000),
15     param_grid_lr,
16     cv=5,
17     scoring='accuracy',
18     n_jobs=-1,
19     verbose=1
```

```
20 )
21
22 start_time = time.time()
23 grid_lr.fit(X_train_scaled, y_train)
24 lr_tuning_time = time.time() - start_time
25
26 print(f"\nBest Parameters: {grid_lr.best_params_}")
27 print(f"Best CV Accuracy: {grid_lr.best_score_:.4f}")
28 print(f"Tuning Time: {lr_tuning_time:.4f}s")
29
30 # Best model
31 lr_best = grid_lr.best_estimator_
32 y_pred_lr_best = lr_best.predict(X_test_scaled)
33
34 # Evaluation
35 print(f"\nOptimized Logistic Regression Results:")
36 print(f"Test Accuracy: {accuracy_score(y_test, y_pred_lr_best):.4f}")
37 print(f"Precision: {precision_score(y_test, y_pred_lr_best):.4f}")
38 print(f"Recall: {recall_score(y_test, y_pred_lr_best):.4f}")
39 print(f"F1 Score: {f1_score(y_test, y_pred_lr_best):.4f}")
```

Listing 5: Tuning Logistic Regression

## 3. SVM with Different Kernels

```
1 # Test different kernels
2 kernels = ['linear', 'poly', 'rbf', 'sigmoid']
3 svm_results = {}
4
5 print("\n" + "="*60)
6 print("SVM Kernel Comparison")
7 print("="*60)
8
9 for kernel in kernels:
10     print(f"\nTraining SVM with {kernel} kernel...")
11
12     start_time = time.time()
13     if kernel == 'poly':
14         svm = SVC(kernel=kernel, degree=3, random_state=42)
15     else:
16         svm = SVC(kernel=kernel, random_state=42)
17
18     svm.fit(X_train_scaled, y_train)
19     train_time = time.time() - start_time
20
21     # Predictions
22     y_pred = svm.predict(X_test_scaled)
23
24     # Store results
25     svm_results[kernel] = {
26         'model': svm,
27         'accuracy': accuracy_score(y_test, y_pred),
28         'precision': precision_score(y_test, y_pred),
29         'recall': recall_score(y_test, y_pred),
30         'f1': f1_score(y_test, y_pred),
31         'time': train_time
32     }
33
```

```
34    print(f"Accuracy: {svm_results[kernel]['accuracy']:.4f}")
35    print(f"F1 Score: {svm_results[kernel]['f1']:.4f}")
36    print(f"Training Time: {train_time:.4f}s")
```

Listing 6: SVM Kernel Comparison

## 4. SVM Hyperparameter Tuning

```
1  # Focus on RBF kernel (typically best for spam detection)
2  param_grid_svm = {
3      'C': [0.1, 1, 10, 100],
4      'gamma': ['scale', 'auto'],
5      'kernel': ['rbf']
6  }
7
8  print("\n" + "="*60)
9  print("SVM (RBF) Hyperparameter Tuning")
10 print("="*60)
11
12 grid_svm = GridSearchCV(
13     SVC(random_state=42),
14     param_grid_svm,
15     cv=5,
16     scoring='accuracy',
17     n_jobs=-1,
18     verbose=1
19 )
20
21 start_time = time.time()
22 grid_svm.fit(X_train_scaled, y_train)
23 svm_tuning_time = time.time() - start_time
24
25 print(f"\nBest Parameters: {grid_svm.best_params_}")
26 print(f"Best CV Accuracy: {grid_svm.best_score_:.4f}")
27 print(f"Tuning Time: {svm_tuning_time:.4f}s")
28
29 # Best model
30 svm_best = grid_svm.best_estimator_
31 y_pred_svm_best = svm_best.predict(X_test_scaled)
32
33 # Evaluation
34 print(f"\nOptimized SVM Results:")
35 print(f"Test Accuracy: {accuracy_score(y_test, y_pred_svm_best):.4f}")
36 print(f"Precision: {precision_score(y_test, y_pred_svm_best):.4f}")
37 print(f"Recall: {recall_score(y_test, y_pred_svm_best):.4f}")
38 print(f"F1 Score: {f1_score(y_test, y_pred_svm_best):.4f}")
```

Listing 7: Tuning SVM (RBF Kernel)

## 5. K-Fold Cross-Validation

```
1  from sklearn.model_selection import cross_validate
2
3  # Define models to evaluate
4  models = {
5      'Logistic Regression': lr_best,
```

```
6      'SVM (RBF)': svm_best
7  }
8
9  # Perform 5-fold CV
10 print("\n" + "="*60)
11 print("5-Fold Cross-Validation")
12 print("="*60)
13
14 cv_results = {}
15 for name, model in models.items():
16     print(f"\nEvaluating {name}...")
17
18     cv_scores = cross_validate(
19         model, X_train_scaled, y_train,
20         cv=5,
21         scoring=['accuracy', 'precision', 'recall', 'f1'],
22         return_train_score=True
23     )
24
25     cv_results[name] = {
26         'test_accuracy': cv_scores['test_accuracy'],
27         'test_precision': cv_scores['test_precision'],
28         'test_recall': cv_scores['test_recall'],
29         'test_f1': cv_scores['test_f1']
30     }
31
32     print(f"Accuracy: {cv_scores['test_accuracy'].mean():.4f} (+/- {cv_scores['
       test_accuracy'].std():.4f})")
33     print(f"Precision: {cv_scores['test_precision'].mean():.4f} (+/- {cv_scores['
       test_precision'].std():.4f})")
34     print(f"Recall: {cv_scores['test_recall'].mean():.4f} (+/- {cv_scores['
       test_recall'].std():.4f})")
35     print(f"F1 Score: {cv_scores['test_f1'].mean():.4f} (+/- {cv_scores['test_f1
       '].std():.4f})")
```

Listing 8: 5-Fold Cross-Validation

# Visualizations

## 1. Confusion Matrices

```
1  from sklearn.metrics import ConfusionMatrixDisplay
2
3  fig, axes = plt.subplots(1, 2, figsize=(12, 5))
4
5  # Logistic Regression
6  ConfusionMatrixDisplay.from_predictions(
7      y_test, y_pred_lr_best, ax=axes[0], cmap='Blues'
8  )
9  axes[0].set_title(f'Logistic Regression\nAccuracy: {accuracy_score(y_test,
       y_pred_lr_best):.4f}')
10 axes[0].grid(False)
11
12 # SVM
13 ConfusionMatrixDisplay.from_predictions(
14     y_test, y_pred_svm_best, ax=axes[1], cmap='Greens'
15 )
```

```
16  axes[1].set_title(f'SVM (RBF)\nAccuracy: {accuracy_score(y_test, y_pred_svm_best)
        :.4f}')
17  axes[1].grid(False)
18
19  plt.tight_layout()
20  plt.show()
```

Listing 9: Confusion Matrix Visualization

## 2. ROC Curves

```
1   from sklearn.metrics import roc_curve, auc
2
3   # Get probability predictions
4   y_prob_lr = lr_best.decision_function(X_test_scaled)
5   y_prob_svm = svm_best.decision_function(X_test_scaled)
6
7   # Compute ROC curves
8   fpr_lr, tpr_lr, _ = roc_curve(y_test, y_prob_lr)
9   fpr_svm, tpr_svm, _ = roc_curve(y_test, y_prob_svm)
10
11  roc_auc_lr = auc(fpr_lr, tpr_lr)
12  roc_auc_svm = auc(fpr_svm, tpr_svm)
13
14  # Plot
15  plt.figure(figsize=(10, 6))
16  plt.plot(fpr_lr, tpr_lr, linewidth=2,
17           label=f'Logistic Regression (AUC = {roc_auc_lr:.3f})')
18  plt.plot(fpr_svm, tpr_svm, linewidth=2,
19           label=f'SVM RBF (AUC = {roc_auc_svm:.3f})')
20  plt.plot([0, 1], [0, 1], 'k--', linewidth=2, label='Random Classifier')
21
22  plt.xlabel('False Positive Rate', fontsize=12)
23  plt.ylabel('True Positive Rate', fontsize=12)
24  plt.title('ROC Curves Comparison', fontsize=14)
25  plt.legend(fontsize=11)
26  plt.grid(True, alpha=0.3)
27  plt.tight_layout()
28  plt.show()
```

Listing 10: ROC Curve Comparison

## 3. Performance Comparison

```
1   # Comparison data
2   models_names = ['Logistic Regression', 'SVM (RBF)']
3   metrics = ['Accuracy', 'Precision', 'Recall', 'F1 Score']
4
5   lr_scores = [
6       accuracy_score(y_test, y_pred_lr_best),
7       precision_score(y_test, y_pred_lr_best),
8       recall_score(y_test, y_pred_lr_best),
9       f1_score(y_test, y_pred_lr_best)
10  ]
11
12  svm_scores = [
```

```
13      accuracy_score(y_test, y_pred_svm_best),
14      precision_score(y_test, y_pred_svm_best),
15      recall_score(y_test, y_pred_svm_best),
16      f1_score(y_test, y_pred_svm_best)
17  ]
18
19  # Plot
20  x = np.arange(len(metrics))
21  width = 0.35
22
23  fig, ax = plt.subplots(figsize=(10, 6))
24  ax.bar(x - width/2, lr_scores, width, label='Logistic Regression', alpha=0.8)
25  ax.bar(x + width/2, svm_scores, width, label='SVM (RBF)', alpha=0.8)
26
27  ax.set_xlabel('Metrics', fontsize=12)
28  ax.set_ylabel('Score', fontsize=12)
29  ax.set_title('Performance Metrics Comparison', fontsize=14)
30  ax.set_xticks(x)
31  ax.set_xticklabels(metrics)
32  ax.legend(fontsize=11)
33  ax.grid(True, alpha=0.3, axis='y')
34  ax.set_ylim([0.85, 1.0])
35
36  plt.tight_layout()
37  plt.show()
```

Listing 11: Model Comparison Bar Chart

# Performance Results

## Hyperparameter Tuning Results

Table 1: Hyperparameter Tuning Summary

| Model | Search | Best Parameters | Best CV Acc |
|---|---|---|---|
| Logistic Regression | Grid Search | penalty='l2', C=10, solver='saga' | 0.9293 |
| SVM | Grid Search | kernel='rbf', C=10, gamma='scale' | 0.9315 |

**Analysis:**

- Logistic Regression optimal with L2 regularization and C=10
- SVM RBF kernel optimal with C=10 and gamma='scale'
- Both models achieved ¿93% CV accuracy
- Grid Search explored 20 combinations for LR, 8 for SVM