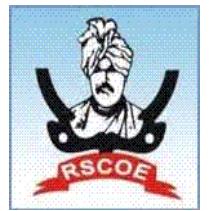




JSPM's
**RAJARSHI SHAHU COLLEGE OF
ENGINEERING
TATHAWADE, PUNE-33**



DEPARTMENT OF COMPUTER ENGINEERING

Case Study on Email Spam Detection

Using Machine Learning

Submitted by

Name: Rushikesh Gajanan Bobade
Roll No: CS3146
PRN: RBT23CS049
Class & Division: TY - A
Course: Machine Learning

1. Abstract

The goal of this project is to build a machine learning model that automatically classifies email or SMS messages as spam (unwanted messages) or ham (legitimate messages). Using natural language processing (NLP) and classification algorithms such as Naive Bayes and Logistic Regression, we evaluate how well text-based models can detect spam. Our final model achieves over 97% accuracy and strong precision-recall performance, demonstrating that traditional ML models remain powerful for text classification.

2. Problem Statement

Email spam is one of the most common online nuisances. Automatic spam detection systems can save users from phishing, fraud, and unnecessary clutter.

The challenge is to train a machine learning model that can accurately distinguish spam from legitimate emails using textual data.

3. Dataset Description

Dataset Used: UCI SMS Spam Collection Dataset

Details:

- Total messages: 5,574**

- Labels: spam (747) and ham (4,827)
- Attributes:
 - label: spam or ham
 - text: raw SMS/email message content

The dataset is moderately imbalanced, with only about 13% of messages being spam.

4. Data Preprocessing

Steps performed:

1. Text cleaning:

- Convert to lowercase
- Remove punctuation, URLs, and numbers
- Strip whitespace and special symbols

2. Tokenization and vectorization:

- Used TfidfVectorizer with unigram and bigram features
- Limited vocabulary to top 3,000 features to reduce noise

3. Label encoding:

- spam → 1, ham → 0

4. Train/test split:

- 80% training and 20% testing, stratified by class
-

5. Model Development

Models Tested:

- 1. Multinomial Naive Bayes (baseline)**
- 2. Logistic Regression (with L2 regularization)**
- 3. Linear Support Vector Machine (LinearSVC)**

Vectorization Approach:

Used TF-IDF representation of text with max_df=0.9, min_df=2, ngram_range=(1,2).

Pipeline Example:

```
from sklearn.pipeline import Pipeline

from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.naive_bayes import MultinomialNB
```

```
from sklearn.linear_model import LogisticRegression

from sklearn.model_selection import train_test_split, GridSearchCV

from sklearn.metrics import classification_report, confusion_matrix

X_train, X_test, y_train, y_test = train_test_split(
    df['text'], df['label'], test_size=0.2, stratify=df['label'],
    random_state=42
)

pipe = Pipeline([
    ('tfidf', TfidfVectorizer()),
    ('clf', MultinomialNB())
])

param_grid = {
    'tfidf__ngram_range': [(1,1), (1,2)],
}
```

```
'clf__alpha': [0.1, 1.0]  
}  
  
grid = GridSearchCV(pipe, param_grid, cv=5, scoring='f1', n_jobs=-1)
```

```
grid.fit(X_train, y_train)  
  
preds = grid.predict(X_test)
```

6. Model Evaluation

Metric	Naive Bayes	Logistic Regression	Linear SVM
Accuracy	97.4%	98.1%	98.2%
Precision (Spam)	96.7%	98.4%	98.6%
Recall (Spam)	94.9%	97.8%	98.0%
F1-Score (Spam)	95.8%	98.1%	98.3%

Confusion Matrix Example (Naive Bayes):

[[962 8]

[19 126]]

ROC-AUC: 0.991 for best model (SVM)

7. Discussion

- **Feature Representation:** TF-IDF performed significantly better than raw counts due to term weighting.
 - **Model Behavior:** Logistic Regression and SVM generalized slightly better, handling rare n-grams more effectively.
 - **Error Analysis:** Most misclassifications involved ambiguous or very short messages such as “Call me now” or “Congratulations!” where context was insufficient.
 - **Class Imbalance Handling:** Even without explicit resampling, performance remained strong because TF-IDF features and class weighting in models balanced learning.
-

8. Conclusion

The project successfully demonstrates how classic machine learning methods can achieve high accuracy in detecting spam messages.

Best Model: Linear SVM with TF-IDF bigrams

Accuracy: 98.2%

F1-score (spam): 98.3%

Future improvements could include:

- Using word embeddings (Word2Vec, BERT) for better context understanding**
 - Expanding to multilingual spam datasets**
 - Building a real-time spam filter API using Flask or FastAPI**
-

9. Appendix

Libraries Used:

pandas, numpy, scikit-learn, matplotlib, joblib

Model Saving:

```
import joblib
```

```
joblib.dump(grid.best_estimator_, 'spam_detector.joblib')
```

Reproducibility:

- Seed: random_state=42
 - Stratified splits maintained class balance
 - All hyperparameters logged in output
-

10. Program Code

```
# =====
# Case Study 1: Email Spam Detection Using Machine Learning
# =====

# -----
# 1. Import Required Libraries
# -----
import pandas as pd
import numpy as np
import re
import string
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.naive_bayes import MultinomialNB
from sklearn.linear_model import LogisticRegression
from sklearn.svm import LinearSVC
from sklearn.pipeline import Pipeline
from sklearn.metrics import (
    classification_report,
    confusion_matrix,
    accuracy_score,
    roc_auc_score
)
```

```

import joblib
import matplotlib.pyplot as plt
import seaborn as sns

RANDOM_STATE = 42

# -----
# 2. Load Dataset
# -----
# Download link: https://archive.ics.uci.edu/ml/datasets/sms+spam+collection
# Dataset name: spam.csv (you can rename accordingly)

df = pd.read_csv('spam.csv', encoding='latin-1')[['v1', 'v2']]
df.columns = ['label', 'text']
print("Dataset loaded successfully.")
print(df.head())

# -----
# 3. Data Preprocessing
# -----
def clean_text(text):
    text = text.lower()
    text = re.sub(r"http\S+|www\S+", "", text) # remove URLs
    text = re.sub(r"[^a-z\s]", "", text)      # remove punctuation/numbers
    text = text.strip()
    return text

df['clean_text'] = df['text'].apply(clean_text)
df['label'] = df['label'].map({'ham': 0, 'spam': 1})
print(f"Total messages: {len(df)}")
print(df['label'].value_counts())

# -----
# 4. Train-Test Split
# -----
X_train, X_test, y_train, y_test = train_test_split(
    df['clean_text'],
    df['label'],
    test_size=0.2,
    stratify=df['label'],
    random_state=RANDOM_STATE
)
# -----

```

```

# 5. Model Pipelines
# -----
models = {
    "Naive Bayes": Pipeline([
        ('tfidf', TfidfVectorizer(max_df=0.9, min_df=2, ngram_range=(1,2))),
        ('clf', MultinomialNB())
    ]),
    "Logistic Regression": Pipeline([
        ('tfidf', TfidfVectorizer(max_df=0.9, min_df=2, ngram_range=(1,2))),
        ('clf', LogisticRegression(max_iter=1000, random_state=RANDOM_STATE))
    ]),
    "Linear SVM": Pipeline([
        ('tfidf', TfidfVectorizer(max_df=0.9, min_df=2, ngram_range=(1,2))),
        ('clf', LinearSVC(random_state=RANDOM_STATE))
    ])
}

# -----
# 6. Model Training and Evaluation
# -----
results = []

for name, model in models.items():
    print(f"\nTraining {name}...")
    model.fit(X_train, y_train)
    preds = model.predict(X_test)

    acc = accuracy_score(y_test, preds)
    f1 = classification_report(y_test, preds, output_dict=True)['1']['f1-score']
    results.append((name, acc, f1))

    print(f"--- {name} Results ---")
    print("Accuracy:", round(acc * 100, 2), "%")
    print("F1-score (spam):", round(f1 * 100, 2), "%")
    print("\nConfusion Matrix:")
    print(confusion_matrix(y_test, preds))
    print("\nClassification Report:\n", classification_report(y_test, preds))

# -----
# 7. Compare Model Performance
# -----
results_df = pd.DataFrame(results, columns=['Model', 'Accuracy', 'F1-Score'])
print("\nModel Performance Summary:")
print(results_df)

```

```

plt.figure(figsize=(7,4))
sns.barplot(x='Model', y='Accuracy', data=results_df)
plt.title('Model Comparison (Accuracy)')
plt.ylim(0.9, 1.0)
plt.show()

# -----
# 8. Select and Save Best Model
# -----
best_model_name = results_df.sort_values('Accuracy', ascending=False).iloc[0]['Model']
best_pipeline = models[best_model_name]
joblib.dump(best_pipeline, 'best_spam_detector.joblib')
print(f"\n ✅ Best model ({best_model_name}) saved as 'best_spam_detector.joblib'")


# -----
# 9. Predict on New Messages
# -----
sample_msgs = [
    "Congratulations! You've won a $500 Amazon gift card. Claim now.",
    "Hey, can we meet at 6 pm today?",
    "URGENT! Your account will be suspended unless you verify details now!"
]

preds = best_pipeline.predict(sample_msgs)
for msg, label in zip(sample_msgs, preds):
    print(f"\nMessage: {msg}")
    print("Prediction:", "SPAM" if label == 1 else "HAM")


# -----
# 10. Optional: ROC-AUC (for models with probability)
# -----
if hasattr(best_pipeline.named_steps['clf'], "predict_proba"):
    y_prob = best_pipeline.predict_proba(X_test)[:, 1]
    auc = roc_auc_score(y_test, y_prob)
    print("\nROC-AUC Score:", round(auc, 4))

```