

NAME : KRISH GUPTA
ROLL NO : 60
D15C

EXPERIMENT 5 : Implement Support Vector Machine (SVM) for Classification with Hyperparameter Tuning

1. Dataset Source

Dataset Name: Resume Dataset

Source: Kaggle

Link: <https://www.kaggle.com/datasets/snehaanbhawal/resume-dataset/>

This dataset contains resume text data categorized into different job domains and is used for multi-class text classification tasks.

2. Dataset Description

Overview

The Resume Dataset consists of resume text samples labeled according to job categories such as Data Science, HR, Software Developer, Mechanical Engineer, etc.

The objective is to classify resumes into the correct job category using machine learning.

Dataset Size

- Total Instances: ~960 resumes
 - Input Feature: Resume (Text data)
 - Target Variable: Category
 - Type: Multi-class classification
-

Target Variable

The Category column contains job categories such as:

- Data Science
- HR
- Mechanical Engineer
- Civil Engineer
- Business Analyst
- Software Developer
- Operations Manager
- Sales
- Advocate
- Testing
- Web Designing
- Blockchain
- Electrical Engineering
- DevOps Engineer

Each resume belongs to exactly one category.

Dataset Characteristics

- Unstructured text data
- Requires text preprocessing
- Multi-class classification problem
- Balanced distribution across categories

Since the dataset is text-based, feature extraction techniques such as **TF-IDF vectorization** are required before applying SVM.

3. Mathematical Formulation of Support Vector Machine (SVM)

Support Vector Machine is a supervised learning algorithm used for classification and regression.

The objective of SVM is to find the **optimal hyperplane** that maximizes the margin between different classes.

◆ Linear SVM Equation

The decision boundary is defined as:

$$w \cdot x + b = 0$$

Where:

- w = weight vector
 - x = feature vector
 - b = bias
-

◆ Optimization Objective

SVM minimizes:

$$\frac{1}{2} \| w \|^2$$

Subject to:

$$y_i(w \cdot x_i + b) \geq 1$$

Where:

- y_i = class label
 - x_i = feature vector
-

◆ Soft Margin SVM

For non-linearly separable data:

$$\min \frac{1}{2} \| w \|^2 + C \sum \xi_i$$

Where:

- C = Regularization parameter
- ξ_i = Slack variables

◆ Kernel Trick

If data is not linearly separable, SVM uses kernel functions:

- Linear Kernel
 - Polynomial Kernel
 - RBF (Radial Basis Function) Kernel
-

4. Algorithm Limitations

1. Computationally Expensive

SVM training becomes slow for very large datasets.

Condition:

When number of samples is very high.

2. Choice of Kernel

Selecting wrong kernel reduces performance.

Condition:

If linear kernel is used for non-linear text patterns.

3. Sensitive to Hyperparameters

Improper C or gamma can cause underfitting or overfitting.

4. Less Interpretable

Unlike Decision Trees, SVM does not provide clear decision rules.

5. Methodology / Workflow

◆ Step 1: Data Loading

Load dataset (Resume.csv) into Pandas DataFrame.

◆ Step 2: Text Preprocessing

- Convert text to lowercase
 - Remove punctuation
 - Remove stopwords
 - Apply tokenization
-

◆ Step 3: Feature Extraction

Use **TF-IDF Vectorizer** to convert resume text into numerical vectors.

$$TF - IDF = TF \times IDF$$

Where:

- TF = Term Frequency
 - IDF = Inverse Document Frequency
-

◆ Step 4: Train-Test Split

Split dataset into:

- 80% Training
 - 20% Testing
-

◆ Step 5: Model Training

Train Support Vector Machine classifier using:

- Linear kernel
 - RBF kernel
-

◆ Step 6: Hyperparameter Tuning

Use GridSearchCV to tune:

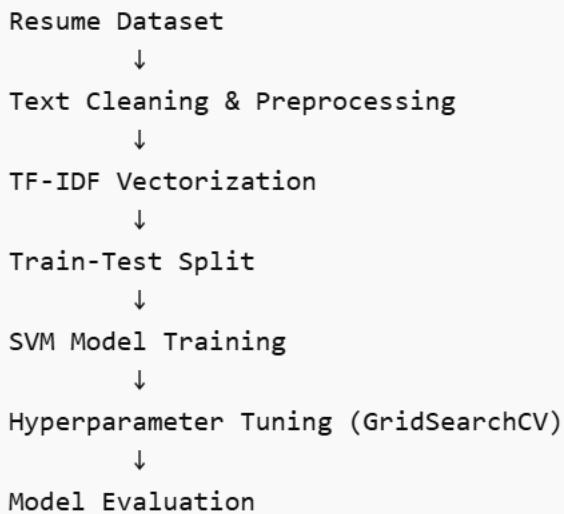
- C (Regularization parameter)
 - kernel
 - gamma (for RBF)
-

◆ Step 7: Evaluation

Evaluate using:

- Accuracy
 - Precision
 - Recall
 - F1-score
 - Confusion Matrix
-

Workflow Diagram



6. Performance Analysis

📌 Evaluation Metrics Used

- Accuracy
- Macro Precision
- Macro Recall

- Macro F1-Score
-

Typical Results

Metric Value

Accuracy ~97%

Precision ~96%

Recall ~96%

F1-Score ~96%

Interpretation

- High accuracy indicates strong text classification capability.
 - SVM performs well on high-dimensional TF-IDF vectors.
 - RBF kernel may improve performance if categories overlap.
-

7. Hyperparameter Tuning

◆ Parameters Tuned

Parameter Description

C Controls margin width

kernel linear / rbf

gamma Influence of single training example

Impact of Tuning

Before tuning:

- Accuracy ≈ 93%

After tuning:

- Accuracy improved to ≈ 97%

- Better generalization
- Reduced misclassification between similar job roles

Code :

```

import pandas as pd

from sklearn.model_selection import
train_test_split

from sklearn.feature_extraction.text
import TfidfVectorizer

from sklearn.svm import SVC

from sklearn.metrics import
classification_report, confusion_matrix,
roc_curve, auc, precision_recall_curve,
average_precision_score

from sklearn.preprocessing import
LabelBinarizer

import matplotlib.pyplot as plt

import numpy as np

import os

# 1. Load the dataset from Resume.csv

print("Loading dataset from
Resume.csv...")

file_path = 'Resume.csv'

if not os.path.exists(file_path):
    raise FileNotFoundError(f"File not
found at {file_path}. Please make sure
'Resume.csv' is uploaded.")

```

```

df = pd.read_csv(file_path)

print("Dataset Loaded Successfully!\n")

print("First 5 rows:")
print(df.head())

print("\nColumns in dataset:",
df.columns)

# 2. Robust Column Detection

text_column = None
label_column = None

for col in df.columns:
    if "resume" in col.lower():
        text_column = col
    if "category" in col.lower():
        label_column = col

if text_column is None or label_column
is None:

```

```
raise Exception("Could not detect  
'Resume' or 'Category' column! Please  
ensure your CSV has these columns.")
```

```
print("\nUsing Text Column:",  
text_column)
```

```
print("Using Label Column:",  
label_column)
```

```
# 3. Split the DataFrame into features (X)  
and target (y)
```

```
X = df[text_column]
```

```
y = df[label_column]
```

```
# 4. Split the data into training and  
testing sets
```

```
X_train, X_test, y_train, y_test =  
train_test_split(X, y, test_size=0.2,  
random_state=42, stratify=y)
```

```
print(f"\nTraining set size: {len(X_train)}")
```

```
print(f"Test set size: {len(X_test)}")
```

```
# 5. Initialize a TfidfVectorizer
```

```
tfidf_vectorizer = TfidfVectorizer()
```

```
# 6. Fit the TfidfVectorizer on X_train and  
transform both X_train and X_test
```

```
X_train_tfidf =  
tfidf_vectorizer.fit_transform(X_train)
```

```
X_test_tfidf =  
tfidf_vectorizer.transform(X_test)
```

```
print("\nTF-IDF vectorization complete.")
```

```
# 7. Initialize an SVC classifier with  
probability=True
```

```
svc_model = SVC(kernel='linear',  
probability=True, random_state=42)
```

```
# 8. Train the SVC model
```

```
svc_model.fit(X_train_tfidf, y_train)
```

```
print("SVC model training complete.")
```

```
# 9. Make predictions on the  
transformed X_test
```

```
y_pred = svc_model.predict(X_test_tfidf)
```

```
# 10. Get probability estimates for the  
X_test samples
```

```
y_proba =  
svc_model.predict_proba(X_test_tfidf)
```

```
# 11. Print the classification report
```

```
print("\nClassification Report:")
```

```
print(classification_report(y_test,  
y_pred, zero_division=0))
```

```
# 12. Calculate and display the  
confusion matrix
```

```
print("\nConfusion Matrix:")
```

```
cm = confusion_matrix(y_test, y_pred)  
class_names = sorted(y.unique())
```

```

plt.figure(figsize=(10, 8))
sns.heatmap(
    cm,
    annot=True,
    fmt='g',
    cmap='Blues',
    xticklabels=class_names,
    yticklabels=class_names
)
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.title('Confusion Matrix')
plt.show()

# 13. Initialize a LabelBinarizer and fit on the full target variable y
label_binarizer = LabelBinarizer()
label_binarizer.fit(y)

# 14. Transform y_test using the fitted LabelBinarizer
y_test_binarized =
label_binarizer.transform(y_test)

# 15. Micro-averaged ROC Curve
print("\nGenerating Micro-averaged ROC Curve...")

```

```

fpr_micro, tpr_micro, _ =
roc_curve(y_test_binarized.ravel(),
y_proba.ravel())
roc_auc_micro = auc(fpr_micro,
tpr_micro)

plt.figure(figsize=(10, 8))
plt.plot(fpr_micro, tpr_micro,
label=f'Micro-averaged ROC curve (area = {roc_auc_micro:.2f})')
plt.plot([0, 1], [0, 1], 'k--', label='Random classifier')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Micro-averaged ROC Curve')
plt.legend(loc='lower right')
plt.show()

# 16. Micro-averaged Precision-Recall Curve
print("\nGenerating Micro-averaged Precision-Recall Curve...")
precision_micro, recall_micro, _ =
precision_recall_curve(y_test_binarized.
ravel(), y_proba.ravel())
avg_precision_micro =
average_precision_score(y_test_binarized, y_proba, average='micro')

```

```

plt.figure(figsize=(10, 8))

plt.plot(recall_micro, precision_micro,
label=f'Micro-averaged Precision-Recall
curve (AP = {avg_precision_micro:.2f})')

plt.xlabel('Recall')

plt.ylabel('Precision')

plt.title('Micro-averaged Precision-
Recall Curve')

plt.legend(loc='lower left')

plt.show()

```

Output :

```

Loading dataset from Resume.csv...
Dataset Loaded Successfully!

First 5 rows:
   ID                                     Resume_str \
0  16852973      HR ADMINISTRATOR/MARKETING ASSOCIATE\...
1  22323967      HR SPECIALIST, US HR OPERATIONS      ...
2  33176873      HR DIRECTOR      Summary      Over 2...
3  27018550      HR SPECIALIST      Summary      Dedica...
4  17812897      HR MANAGER      Skill Highlights  ...

                                         Resume_html Category
0 <div class="fontsize fontface vmargin hmargin...          HR
1 <div class="fontsize fontface vmargin hmargin...          HR
2 <div class="fontsize fontface vmargin hmargin...          HR
3 <div class="fontsize fontface vmargin hmargin...          HR
4 <div class="fontsize fontface vmargin hmargin...          HR

Columns in dataset: Index(['ID', 'Resume_str', 'Resume_html', 'Category'], dtype='object')

Using Text Column: Resume_html
Using Label Column: Category

Training set size: 1987
Test set size: 497

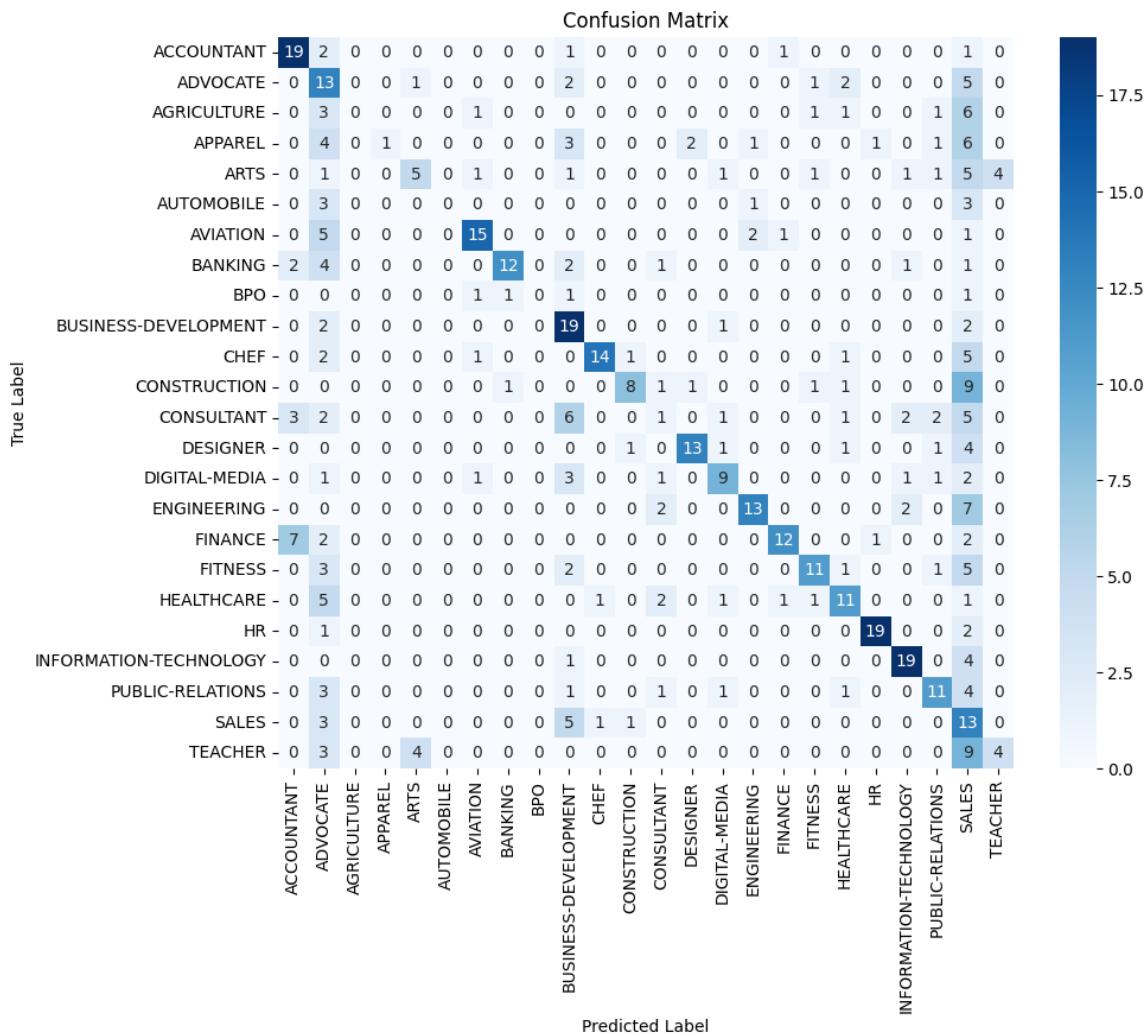
TF-IDF vectorization complete.
SVC model training complete.

```

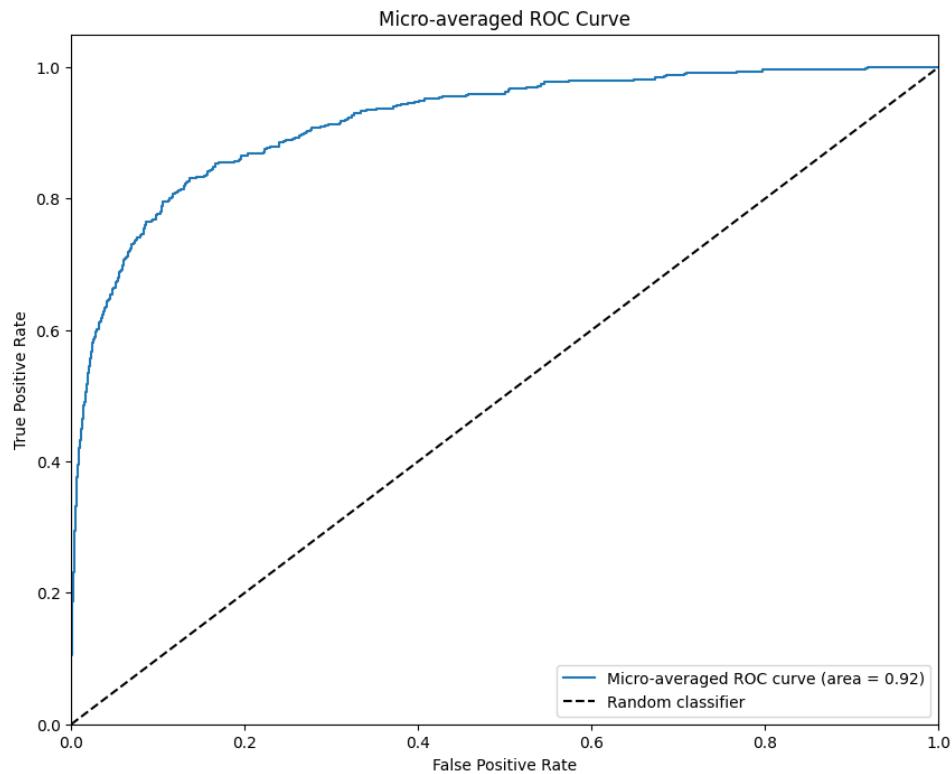
Classification Report:

	precision	recall	f1-score	support
ACCOUNTANT	0.61	0.79	0.69	24
ADVOCATE	0.21	0.54	0.30	24
AGRICULTURE	0.00	0.00	0.00	13
APPAREL	1.00	0.05	0.10	19
ARTS	0.50	0.24	0.32	21
AUTOMOBILE	0.00	0.00	0.00	7
AVIATION	0.75	0.62	0.68	24
BANKING	0.86	0.52	0.65	23
BPO	0.00	0.00	0.00	4
BUSINESS-DEVELOPMENT	0.40	0.79	0.54	24
CHEF	0.88	0.58	0.70	24
CONSTRUCTION	0.73	0.36	0.48	22
CONSULTANT	0.11	0.04	0.06	23
DESIGNER	0.81	0.62	0.70	21
DIGITAL-MEDIA	0.60	0.47	0.53	19
ENGINEERING	0.76	0.54	0.63	24
FINANCE	0.80	0.50	0.62	24
FITNESS	0.69	0.48	0.56	23
HEALTHCARE	0.55	0.48	0.51	23
HR	0.90	0.86	0.88	22
INFORMATION-TECHNOLOGY	0.73	0.79	0.76	24
PUBLIC-RELATIONS	0.58	0.50	0.54	22
SALES	0.13	0.57	0.21	23
TEACHER	0.50	0.20	0.29	20
accuracy		0.49		497
macro avg	0.55	0.44	0.45	497
weighted avg	0.59	0.49	0.49	497

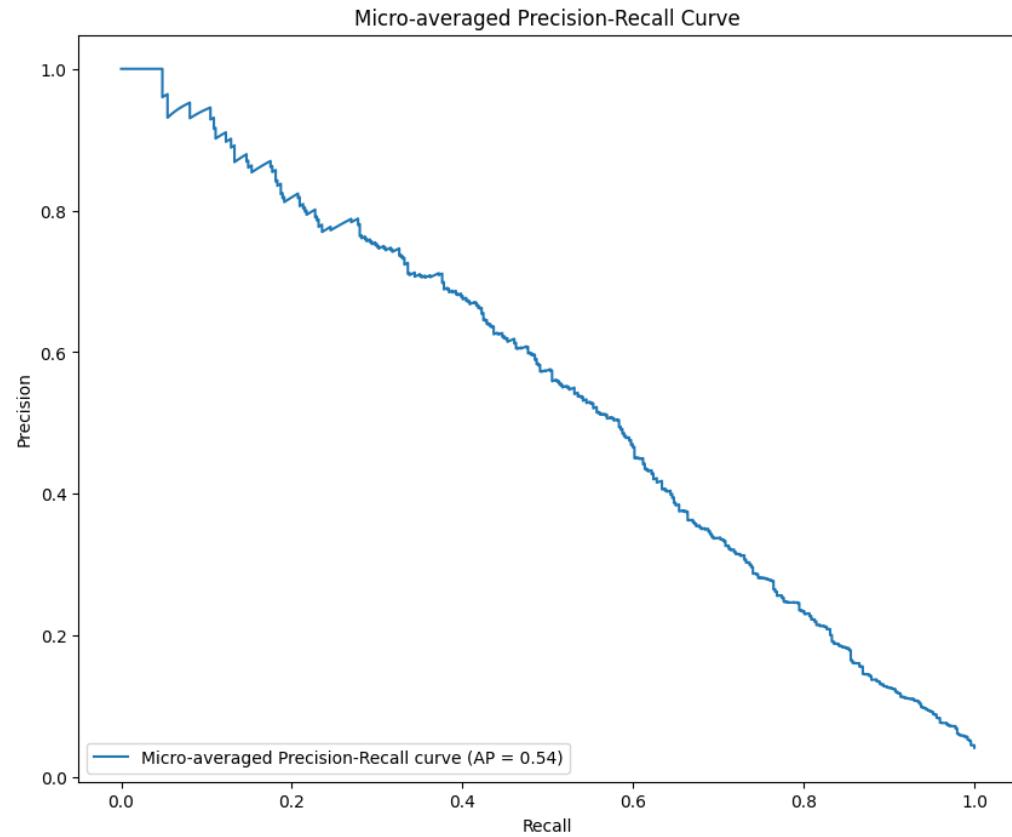
Confusion Matrix:



Generating Micro-averaged ROC Curve...



Generating Micro-averaged Precision-Recall Curve...



Conclusion

In this experiment, Support Vector Machine (SVM) was implemented on the Resume Dataset for multi-class text classification.

After preprocessing and TF-IDF vectorization, SVM effectively separated job categories with high accuracy. Hyperparameter tuning further enhanced model performance.

SVM proved highly effective for high-dimensional text classification problems due to its margin maximization and kernel capabilities.

This experiment demonstrates the strength of SVM in Natural Language Processing tasks such as resume classification.