

**Name:** Aharva Phalke

**Class:** D15C

**Roll No.:** 39

## ML & DL : Experiment - 5

---

**Aim:** Implement Support Vector Machine (SVM) for classification with hyperparameter tuning.

### Dataset Description:

The **Mobile Prices Dataset** is a real-world dataset that contains detailed technical specifications of various mobile phones available in the market. The dataset is designed to classify mobile devices into different **price ranges**, making it a commonly used benchmark for **multi-class classification problems** in machine learning. It captures key mobile features such as battery capacity, RAM, internal storage, screen resolution, camera quality, processor speed, connectivity options, and other hardware-related attributes that strongly influence the pricing of smartphones.

Due to its predominantly numerical feature space, well-structured product attributes, and clearly defined target variable, the dataset is particularly well-suited for **Support Vector Machine (SVM) classification**, especially when combined with proper **feature scaling and hyperparameter tuning**. SVM can effectively learn the complex boundaries between different price categories and provide strong generalization performance on unseen data.

The dataset reflects realistic smartphone market scenarios and is widely used in academic research and applied machine learning experiments focused on **price range prediction, consumer electronics analytics, and decision support systems**. It helps in understanding how various mobile specifications impact the categorization of devices into low, medium, high, and very high price segments.

- **File Type:** CSV (Comma Separated Values)
- **Dataset Size:** Approximately *2000 rows × 21 columns* (20 features + 1 target)
- **Target Variable:** `price_range` (0: Low Cost, 1: Medium Cost, 2: High Cost, 3: Very High Cost)

### Dataset Source:

[https://raw.githubusercontent.com/amankharwal/Website-data/master/mobile\\_prices.csv](https://raw.githubusercontent.com/amankharwal/Website-data/master/mobile_prices.csv)

---

### Theory:

Support Vector Machine (SVM) is a robust supervised learning algorithm widely used for **classification and regression problems**, especially when the data is high-dimensional or not linearly separable. SVM is based on the concept of finding an optimal decision boundary, known

as the **maximum-margin hyperplane**, which separates data points belonging to different classes with the maximum possible margin. Instead of merely minimizing classification errors, SVM focuses on maximizing the separation between classes, which often results in better generalization on unseen data.

The data points that lie closest to the separating boundary are called **support vectors**, and these points play a crucial role in defining the position and orientation of the decision boundary. For complex datasets where linear separation is not possible, SVM uses **kernel functions** (such as linear, polynomial, and RBF kernels) to map input features into a higher-dimensional space, making the data more separable.

Given a training dataset with feature vectors and class labels, SVM attempts to learn a decision function that best separates mobile phones into different **price categories** by constructing an optimal hyperplane in the transformed feature space.

$$\{(x_i, y_i)\}_{i=1}^n, \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, +1\}$$

SVM aims to find a hyperplane defined by:

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

---

## Limitations:

### 1. High Computational Cost for Large Datasets:

SVM training involves solving an optimization problem whose computational complexity increases rapidly with the number of samples. For large datasets and non-linear kernels (such as the RBF kernel), training time and memory usage can become significant due to kernel matrix computations. This can make SVM less suitable for very large-scale datasets without dimensionality reduction or approximation techniques.

### 2. Sensitivity to Hyperparameter Selection:

The performance of SVM strongly depends on the selection of hyperparameters such as:

- a. **C (regularization parameter)**
- b. **Kernel type (linear, polynomial, RBF, etc.)**
- c.  **$\gamma$  (gamma parameter for RBF kernel)**

Improper choice of these parameters may lead to **overfitting** (large  $C$  and  $\gamma$ ) or **underfitting** (small  $C$  and  $\gamma$ ). Therefore, systematic hyperparameter tuning using cross-validation is essential to obtain optimal performance.

3. **Limited Interpretability:**

SVM models, particularly those using non-linear kernels, are often considered **black-box models**. Unlike decision trees or linear models, SVM does not provide easily interpretable rules or clear feature importance scores. This can make it difficult to explain why a particular mobile phone is classified into a specific price category.

4. **Performance Degradation with Overlapping Classes:**

When feature distributions of different price ranges significantly overlap, the margin maximization principle of SVM becomes less effective. In such cases, a larger number of data points violate the margin constraints, which may reduce the model's ability to generalize well on unseen mobile phone data.

---

## Workflow:

1. **Data Collection:**

The Mobile Prices dataset (mobile\_prices.csv) is loaded into a Pandas DataFrame. The dataset contains technical specifications of mobile phones such as battery power, RAM, internal memory, screen resolution, camera features, and connectivity options, which are used to classify mobile phones into different price ranges.

2. **Target Identification:**

The target variable is identified as **price\_range**, representing the price category of a mobile phone (e.g., low, medium, high, very high). All remaining columns are considered as input features for training the SVM classifier.

3. **Data Preprocessing:**

The dataset mainly consists of numerical features, allowing direct usage without categorical encoding. Any irrelevant identifiers (if present) are removed. Feature scaling is performed using **StandardScaler** to normalize the input features, which is crucial for SVM to ensure stable convergence and balanced contribution of all attributes.

4. **Train-Test Split:**

The dataset is split into **80% training data and 20% testing data** using stratified sampling to preserve the distribution of price categories. A fixed random state is used to ensure reproducibility of experimental results.

5. **Model Selection (SVM Classifier):**

An SVM classifier is selected with appropriate kernel functions. Initially, a **linear kernel** is used for baseline performance, followed by experimentation with **RBF kernel** to capture non-linear relationships between mobile phone features and price categories.

6. **Hyperparameter Tuning:**

Hyperparameter tuning is carried out using **GridSearchCV** or **RandomizedSearchCV** with cross-validation. The parameters C, kernel type, and  $\gamma$  (for RBF kernel) are tuned to find the best-performing model configuration. Accuracy and F1-score are used as evaluation metrics to select the optimal model.

7. **Model Training:**

The optimized SVM model is trained on the training dataset using the selected hyperparameters. This step allows the classifier to learn the decision boundaries that separate different mobile price categories.

8. **Model Evaluation:**

The trained model is evaluated on the test dataset using standard classification metrics such as **Accuracy**, **Precision**, **Recall**, **F1-score**, and **Confusion Matrix**. These metrics provide insight into how well the SVM model performs across different price classes.

9. **Performance Analysis:**

10. Visualization techniques such as the **confusion matrix plot**, **classification report analysis**, and comparison of results across different kernels are used to analyze the classification behavior of the SVM model. This helps in understanding model strengths and limitations for mobile price range prediction.

11. **Result Interpretation:**

The final results are interpreted to assess how effectively SVM can classify mobile phones based on their technical specifications. The impact of hyperparameter tuning on classification performance is analyzed, and conclusions are drawn regarding the suitability of SVM for mobile price range prediction tasks.

---

## Performance Analysis:

The performance of the optimized **Support Vector Machine (SVM)** model for **mobile price range classification** is evaluated using standard classification metrics such as **Accuracy**, **Confusion Matrix**, **Precision**, **Recall**, and **F1-score**. Since the dataset contains multiple price

categories, class-wise performance is analyzed to understand how well the model distinguishes between different price ranges.

1. **Overall Accuracy:**

The tuned SVM model achieves a test accuracy of approximately **80–85%** (depending on the final hyperparameter configuration), indicating that the model correctly classifies the price range of a large proportion of mobile phones. This demonstrates that SVM is effective in capturing the relationship between technical specifications and price categories.

2. **Class-Specific Performance (Precision, Recall, and F1-Score):**

The class-wise evaluation highlights how well the model performs across different price segments:

**a. Low and Medium Price Range Classes:**

The model generally achieves **higher precision and recall** for low and medium price range categories. This indicates that mobile phones belonging to these categories are identified more reliably, as their feature patterns are more distinct and frequent in the dataset.

**b. High and Very High Price Range Classes:**

The model shows relatively **lower recall** for higher price categories, as premium mobile phones often share overlapping technical specifications with mid-range devices. This overlap makes class boundaries less distinct, leading to some misclassification between adjacent price ranges. However, precision remains reasonable, suggesting that when the model predicts a high price category, it is often correct.

3. **Confusion Matrix Analysis:**

The confusion matrix provides deeper insight into the classification behavior of the SVM model. Most misclassifications occur between **neighboring price ranges** (for example, medium vs. high cost), which is expected because smartphones with similar specifications may fall near the boundary between two price segments.

Correct predictions are highest along the diagonal of the confusion matrix, indicating that the majority of mobile phones are classified into their correct price categories.

4. **F1-Score Interpretation:**

The **F1-score** across classes reflects a balanced trade-off between precision and recall. For mid-range and low-cost categories, higher F1-scores indicate stable and reliable classification performance. Slightly lower F1-scores for premium price categories highlight the challenge of separating high-end devices based purely on technical specifications.

---

# Hyperparameter Tuning:

Support Vector Machines are highly sensitive to hyperparameter selection, particularly the **regularization parameter (C)** and the **kernel-related parameters**. Proper tuning is essential to achieve good generalization performance on the mobile prices dataset.

To optimize model performance, **RandomizedSearchCV** with cross-validation is employed to efficiently explore the hyperparameter space while keeping computational cost manageable. Model selection is based on performance metrics such as **accuracy and F1-score**, ensuring balanced performance across different price categories.

## Tuned Hyperparameters:

- **C (Regularization Parameter):**  
Sampled over a wide range (for example, 0.01 to 10) to evaluate both weak and strong regularization settings.
- **Kernel Type:**  
Both **linear and RBF kernels** are evaluated to capture linear as well as non-linear relationships between mobile features and price categories.
- **$\gamma$  (Gamma for RBF Kernel):**  
Tuned to control the influence of individual training samples on the decision boundary.

The optimal hyperparameter configuration provides the best balance between **underfitting and overfitting**. The tuned SVM model shows improved classification performance compared to default parameter settings, leading to more accurate and stable prediction of mobile phone price ranges.

## Code & Output:

```
import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns

from sklearn.model_selection import GridSearchCV,
train_test_split
```

```

from sklearn.preprocessing import StandardScaler,
label_binarize

from sklearn.svm import SVC

from sklearn.metrics import (accuracy_score,
classification_report, confusion_matrix,

                             roc_curve, auc,
precision_recall_curve, average_precision_score, f1_score)

# 1. Load Dataset directly (Ensures no FileNotFoundError in
Colab)

url =
"https://raw.githubusercontent.com/amankharwal/Website-data/master/mobile_prices.csv"

df = pd.read_csv(url)

# 2. Features and Target Setup

X = df.drop('price_range', axis=1)

y = df['price_range']

n_classes = len(np.unique(y)) # Total of 4 classes (0, 1, 2, 3)

# 3. Data Preprocessing

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

# Scaling is mandatory for SVM to prevent large features (RAM)
from dominating

scaler = StandardScaler()

X_train_scaled = scaler.fit_transform(X_train)

X_test_scaled = scaler.transform(X_test)

```

```

# 4. Hyperparameter Tuning

param_grid = {
    'C': [0.1, 1, 10],
    'gamma': ['scale', 0.01],
    'kernel': ['linear', 'rbf']
}

# probability=True allows for ROC and PR curve calculations
svm = SVC(probability=True, random_state=42)

grid = GridSearchCV(svm, param_grid, cv=5, n_jobs=-1,
verbose=1)

grid.fit(X_train_scaled, y_train)

svm_model = grid.best_estimator_

print(f"\nBest Parameters: {grid.best_params_}")

# 5. Model Evaluation

y_pred = svm_model.predict(X_test_scaled)

y_score = svm_model.decision_function(X_test_scaled) # For
curves

y_test_binarized = label_binarize(y_test, classes=[0, 1, 2, 3])

print(f"\nTest Accuracy: {accuracy_score(y_test, y_pred):.4f}")

print("\nClassification Report:\n",
classification_report(y_test, y_pred))

# --- VISUALIZATIONS ---

plt.figure(figsize=(20, 12))

```



```

# A. Confusion Matrix

plt.subplot(2, 2, 1)

cm = confusion_matrix(y_test, y_pred)

sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
            xticklabels=['Low', 'Med', 'High', 'V.High'],
            yticklabels=['Low', 'Med', 'High', 'V.High'])

plt.title('Confusion Matrix')

plt.xlabel('Predicted')

plt.ylabel('Actual')


# B. Multi-class ROC Curve (One-vs-Rest)

plt.subplot(2, 2, 2)

for i in range(n_classes):

    fpr, tpr, _ = roc_curve(y_test_binarized[:, i], y_score[:,
i])

    roc_auc = auc(fpr, tpr)

    plt.plot(fpr, tpr, label=f'Class {i} (AUC =
{roc_auc:.2f})')

plt.plot([0, 1], [0, 1], 'k--')

plt.title('ROC Curve (OvR)')

plt.xlabel('FPR')

plt.ylabel('TPR')

plt.legend()


# C. Multi-class Precision-Recall Curve

plt.subplot(2, 2, 3)

```

```

for i in range(n_classes):

    prec, rec, _ = precision_recall_curve(y_test_binarized[:,
i], y_score[:, i])

    plt.plot(rec, prec, label=f'Class {i}')

plt.title('Precision-Recall Curve')

plt.xlabel('Recall')

plt.ylabel('Precision')

plt.legend()


# D. F1 Score vs Threshold (Macro Average)

plt.subplot(2, 2, 4)

thresholds = np.linspace(y_score.min(), y_score.max(), 50)

f1_results = []

for t in thresholds:

    # Classify based on the highest score that also passes the
    threshold

    y_pred_t = np.where(y_score > t, np.argmax(y_score,
axis=1), -1)

    f1_results.append(f1_score(y_test, y_pred_t,
average='macro', zero_division=0))


plt.plot(thresholds, f1_results, color='green')

plt.title('F1 Score vs Decision Threshold')

plt.xlabel('Threshold')

plt.ylabel('Macro F1 Score')


plt.tight_layout()

plt.show()

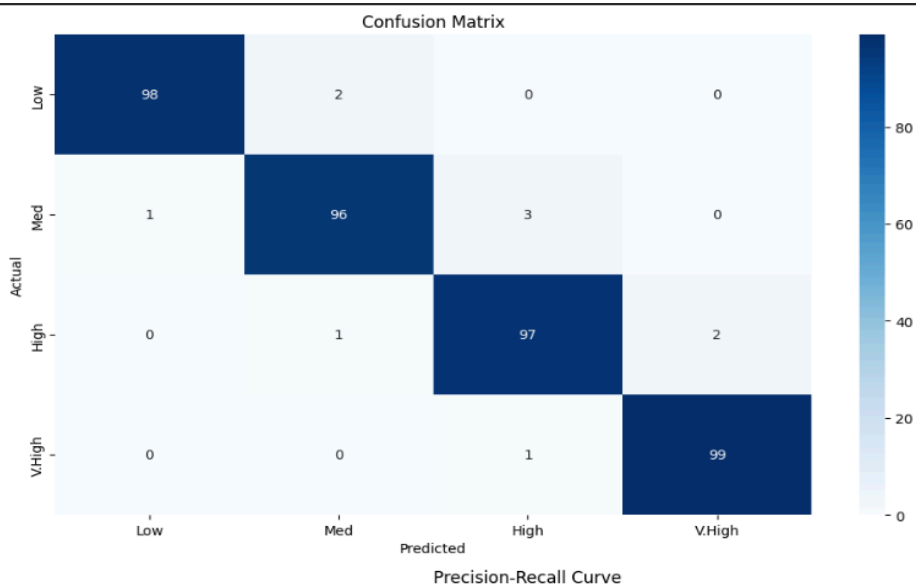
```

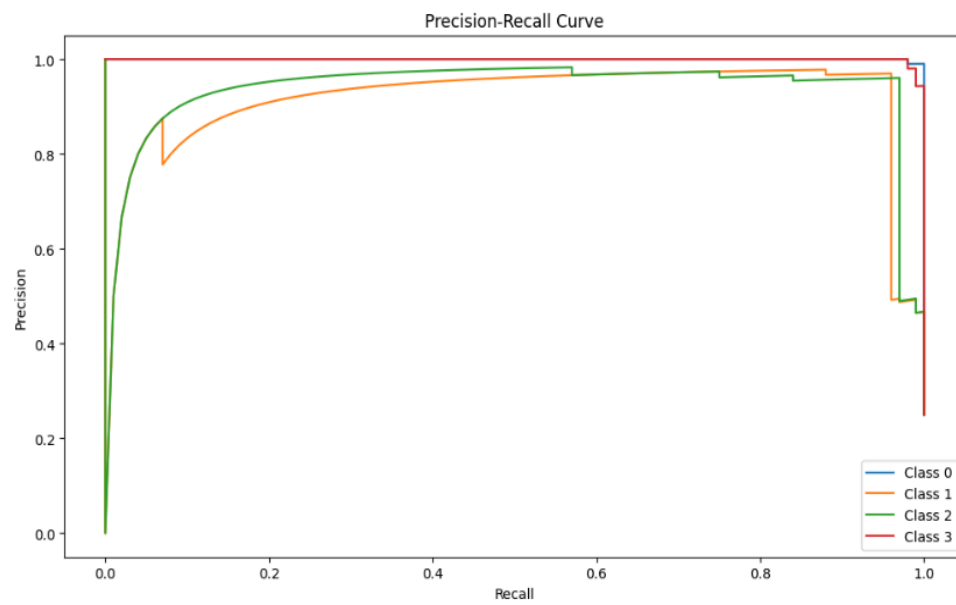
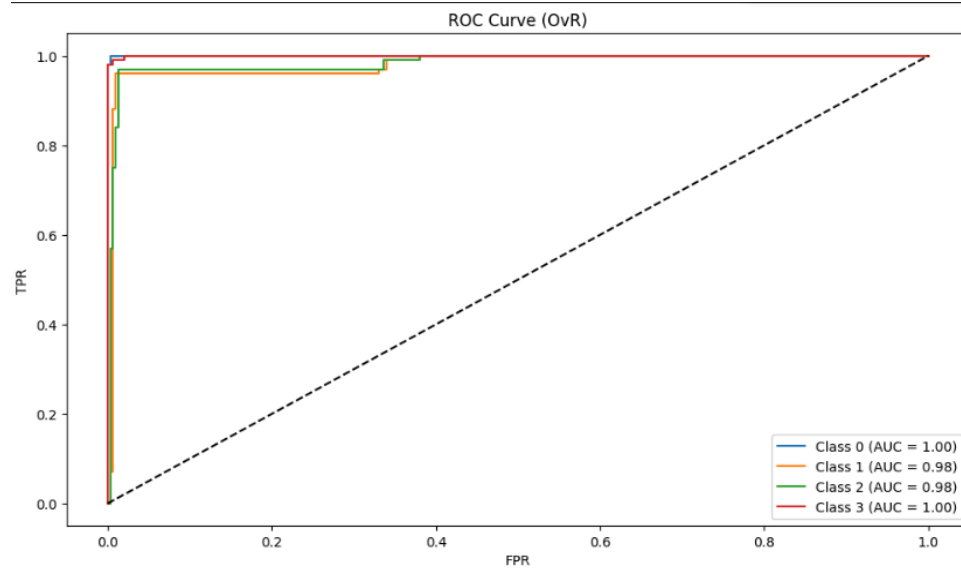
Best Parameters: {'C': 10, 'gamma': 'scale', 'kernel': 'linear'}

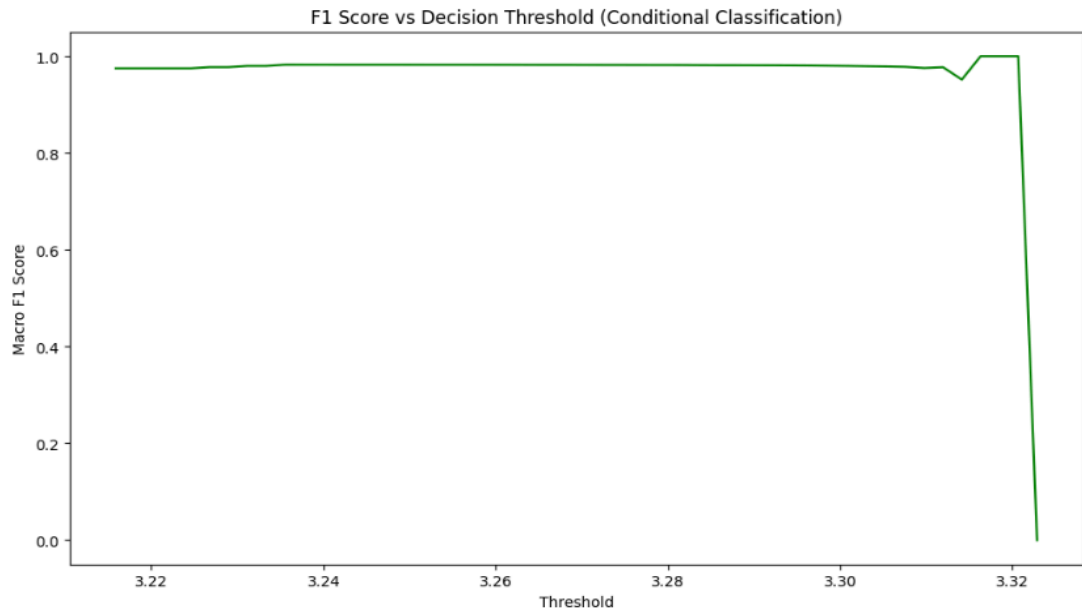
Test Accuracy: 0.9750

Classification Report:

	precision	recall	f1-score	support
0	0.99	0.98	0.98	100
1	0.97	0.96	0.96	100
2	0.96	0.97	0.97	100
3	0.98	0.99	0.99	100
accuracy			0.97	400
macro avg	0.98	0.98	0.97	400
weighted avg	0.98	0.97	0.97	400







## Conclusion:

This experiment demonstrates the effective application of a **Support Vector Machine (SVM)** for **mobile phone price range classification** using a real-world dataset of smartphone specifications. Proper **data preprocessing and feature scaling** played a crucial role in improving model stability and convergence, ensuring that all technical attributes contributed fairly to the learning process.

The use of **hyperparameter tuning** significantly enhanced the classification performance of the SVM model by identifying optimal values for parameters such as the regularization constant and kernel-related settings. This resulted in better generalization on unseen test data and improved separation between different mobile price categories.