```
1  pip install ucimlrepo
```

**CS613 Project -Francis Kurian**

# Support Vector Machines (SVM)

Support Vector Machines (SVM) are a class of supervised machine learning algorithms used for classification and regression tasks. SVM aims to find a decision boundary that maximizes the margin between different classes while minimizing classification errors. The "support vectors" are the data points closest to the decision boundary, and they play a crucial role in determining the boundary's position.

## ▾ Project Summary:

The project demonstrates the use of SVM in a classification task using the heart disease dataset. Here is a summary of the project's main components:

**Data Preparation:** The heart disease dataset is fetched and prepared. Missing values are handled, and the data is split into features (x) and the target variable (y). Additionally, the data is divided into training and testing sets.

**Base Kernel Comparison:** The code compares SVM models with different default kernels ('linear,' 'poly,' 'rbf,' 'sigmoid') by calculating accuracy for each kernel. The accuracy is presented as a percentage on a bar chart, allowing for a visual comparison of model performance.

**Kernel and Degree Tuning:** The code then proceeds to explore SVM models with various kernel types and degrees. It calculates the mean error and accuracy for different combinations of kernels and degrees and presents the results using line plots. This analysis helps identify the impact of kernel choice and degree on model performance.

```
 1 import warnings
 2 import pandas as pd
 3 import numpy as np
 4 from sklearn.svm import SVC
 5 from sklearn.metrics import accuracy_score
 6 import matplotlib.pyplot as plt
 7 from ucimlrepo import fetch_ucirepo
 8 from sklearn.model_selection import train_test_split
 9 from sklearn.preprocessing import StandardScaler
10 warnings.simplefilter(action='ignore', category=FutureWarning)
11 # Fetch the heart disease dataset
12 heart_disease = fetch_ucirepo(id=45)
13
14 # Prepare the data
15 X = heart_disease.data.features
16 y = heart_disease.data.targets
17
18 combined_df = pd.concat([X, y], axis=1)
19 combined_df = combined_df.dropna()
20 combined_df.reset_index(drop=True, inplace=True)
21
22 # Split the data into features (x) and the target variable (y)
23 x = combined_df.iloc[:, 0:13].values
24 y = combined_df['num'].values
25
26 # Split the data into training and testing sets
27 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=0)
28
29 # Standardize the features
30 scaler = StandardScaler()
31 x_train = scaler.fit_transform(x_train)
32 x_test = scaler.transform(x_test)
33 kernels = ['linear', 'poly', 'rbf', 'sigmoid']
```
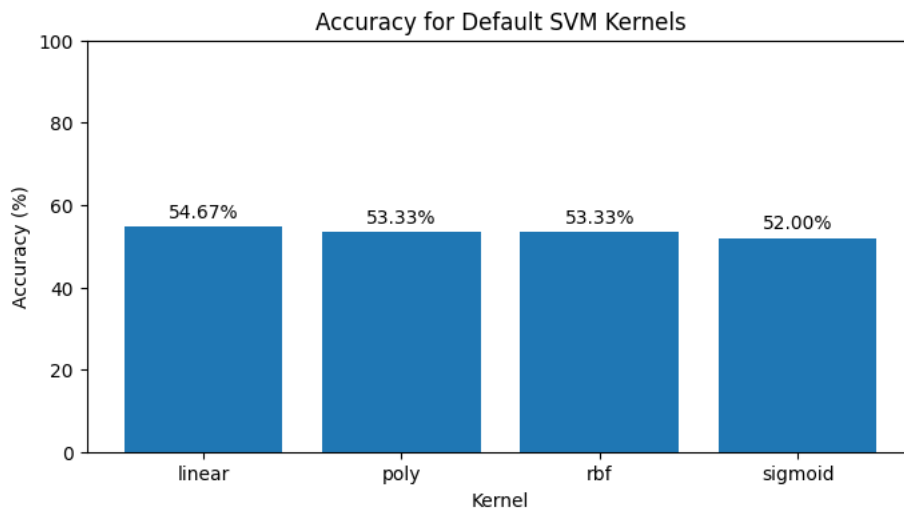
## ▾ Baseline SVM Model by Kernel Types

```
1 # Create a list to store accuracy for different kernels
2 default_accuracy = []
3
4 # Create a DataFrame to store results for default settings
5 default_results_df = pd.DataFrame(columns=['Kernel', 'Accuracy'])
```

```
 6
 7 # Evaluate each kernel with default settings
 8 for kernel in kernels:
 9     default_svm = SVC(kernel=kernel)
10     default_svm.fit(x_train, y_train)
11     default_pred = default_svm.predict(x_test)
12
13     # Calculate accuracy for default settings and convert to percentage
14     default_acc = accuracy_score(y_test, default_pred) * 100
15
16     # Append accuracy to the list and the DataFrame
17     default_accuracy.append(default_acc)
18     default_results_df = default_results_df.append({'Kernel': kernel, 'Accuracy': default_acc}, ignore_index=True)
19
20 # Visualize accuracy for default settings in percentage with labeled bars
21 plt.figure(figsize=(8, 4))
22 bars = plt.bar(default_results_df['Kernel'], default_results_df['Accuracy'])
23 plt.title('Accuracy for Default SVM Kernels')
24 plt.xlabel('Kernel')
25 plt.ylabel('Accuracy (%)')
26 plt.ylim(0, 100)
27
28 # Label the bars with their corresponding percentage values
29 for bar in bars:
30     height = bar.get_height()
31     plt.annotate(f'{height:.2f}%',  # Format the label as a percentage with two decimal places
32                  xy=(bar.get_x() + bar.get_width() / 2, height),
33                  xytext=(0, 3),  # Offset for the label position
34                  textcoords='offset points',
35                  ha='center', va='bottom')
36
37 plt.show()
38
```



## Mean Error Comparison by Kernel & Degrees

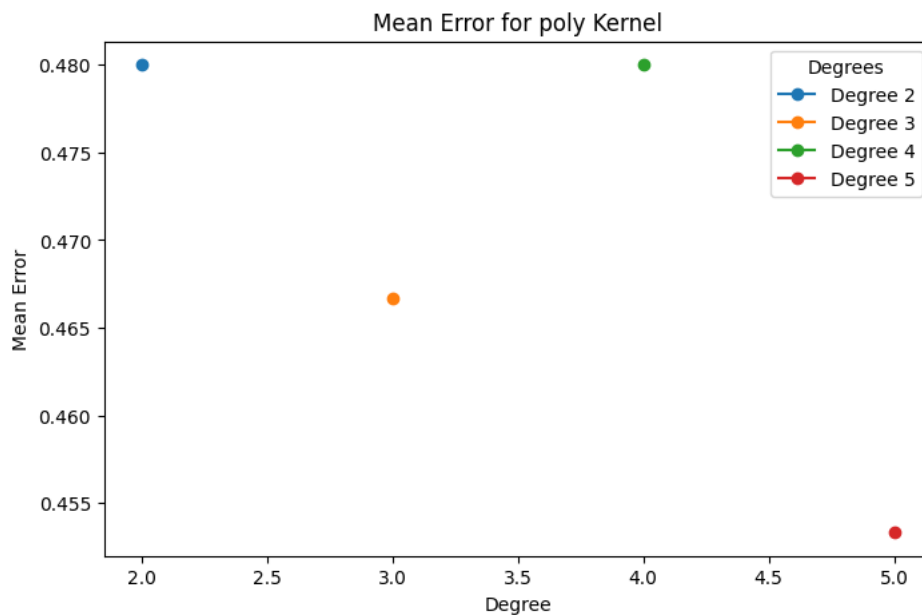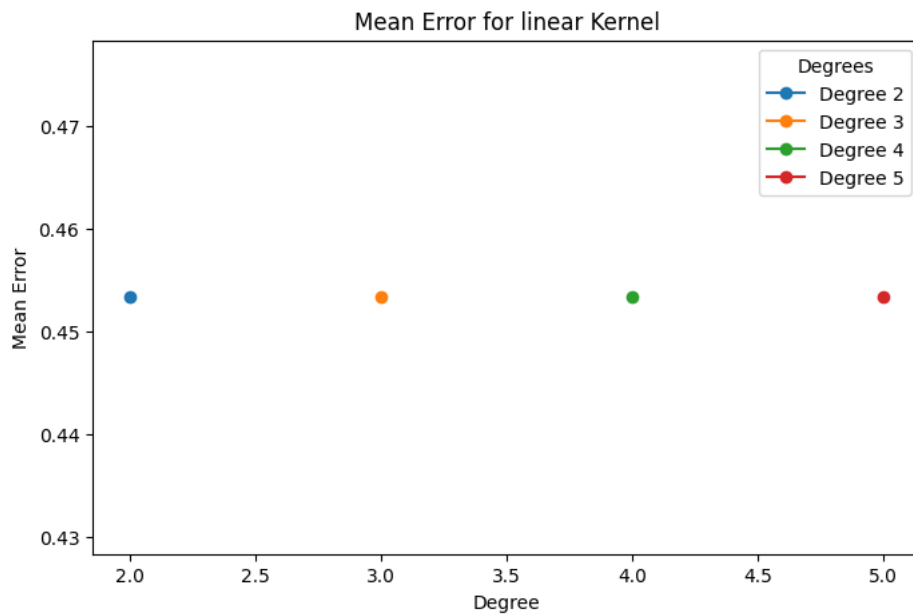- Degrees Hyperparameter is only applicable to Kernel =polynomial.

```
 1 # Create empty lists to store error and accuracy
 2 error = []
 3 accuracy = []
 4
 5 # Create a DataFrame to store results
 6 results_df = pd.DataFrame(columns=['Kernel', 'Degree', 'Mean Error', 'Accuracy'])
 7
 8 # Define the list of kernels and degrees to iterate over
 9
10 degrees = [2, 3, 4, 5]
11
12 for kernel in kernels:
13     for degree in degrees:
14         svm = SVC(kernel=kernel, degree=degree)
15         svm.fit(x_train, y_train)
```

```
16          pred_i = svm.predict(x_test)
17
18          # Calculate mean error and accuracy
19          mean_error = np.mean(pred_i != y_test)
20          acc = accuracy_score(y_test, pred_i)
21
22          # Append values to the lists
23          error.append(mean_error)
24          accuracy.append(acc)
25
26          # Append values to the DataFrame
27          results_df = results_df.append({'Kernel': kernel, 'Degree': degree, 'Mean Error': mean_error, 'Accuracy': acc},
28                                          ignore_index=True)
29
30 results_df.to_csv("svm_results.csv", index=False)
31
```

```
 1 for kernel in kernels:
 2     plt.figure(figsize=(8, 5))
 3     for degree in degrees:
 4         subset = results_df[(results_df['Kernel'] == kernel) & (results_df['Degree'] == degree)]
 5         plt.plot(subset['Degree'], subset['Mean Error'], marker='o', label=f'Degree {degree}')
 6
 7     plt.title(f'Mean Error for {kernel} Kernel')
 8     plt.xlabel('Degree')
 9     plt.ylabel('Mean Error')
10     plt.legend(title='Degrees')
11     plt.show()
```

## Mean Error for linear Kernel



## Mean Error for poly Kernel



## Mean Error for rbf Kernel
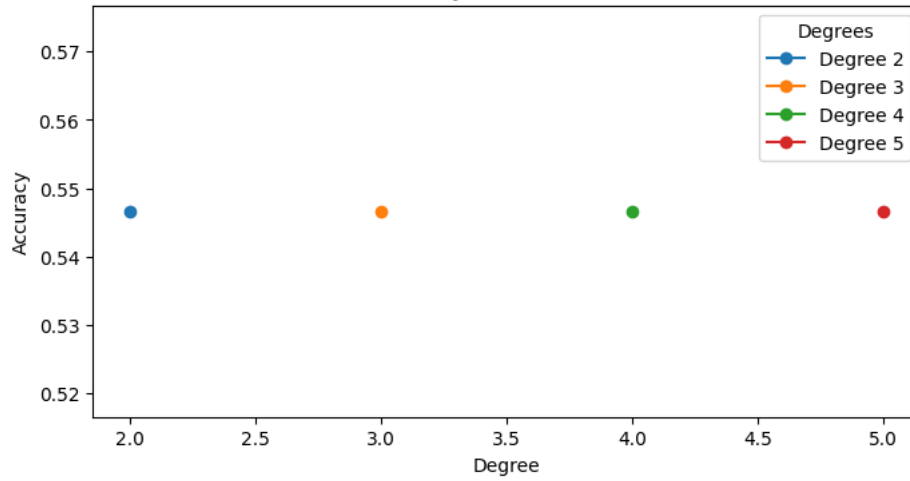
## ▾ Accuracy Comparison by Kernel & Degrees

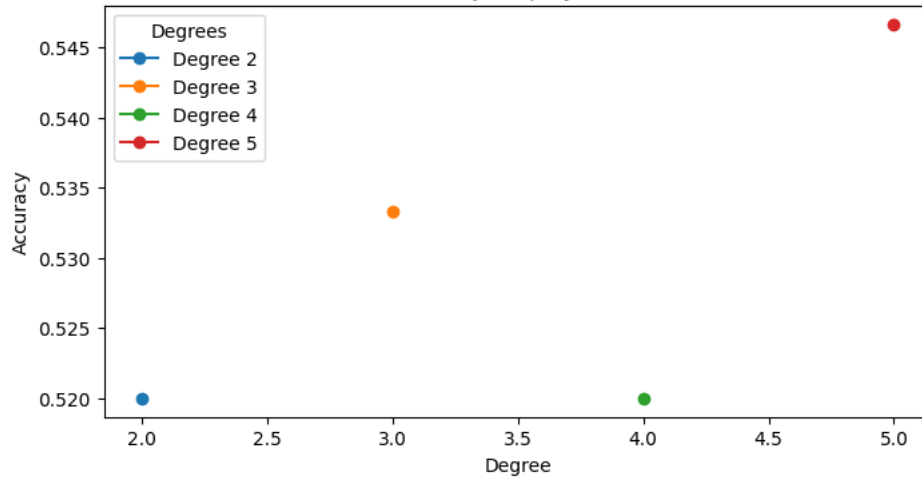- Degrees Hyperparameter is only applicable to Kernel =polynomial.

```
1 # Visualize accuracy for different kernels and degrees
2 for kernel in kernels:
3     plt.figure(figsize=(8, 4))
4     for degree in degrees:
5         subset = results_df[(results_df['Kernel'] == kernel) & (results_df['Degree'] == degree)]
6         plt.plot(subset['Degree'], subset['Accuracy'], marker='o', label=f'Degree {degree}')
7
8     plt.title(f'Accuracy for {kernel} Kernel')
9     plt.xlabel('Degree')
10    plt.ylabel('Accuracy')
11    plt.legend(title='Degrees')
12    plt.show()
```
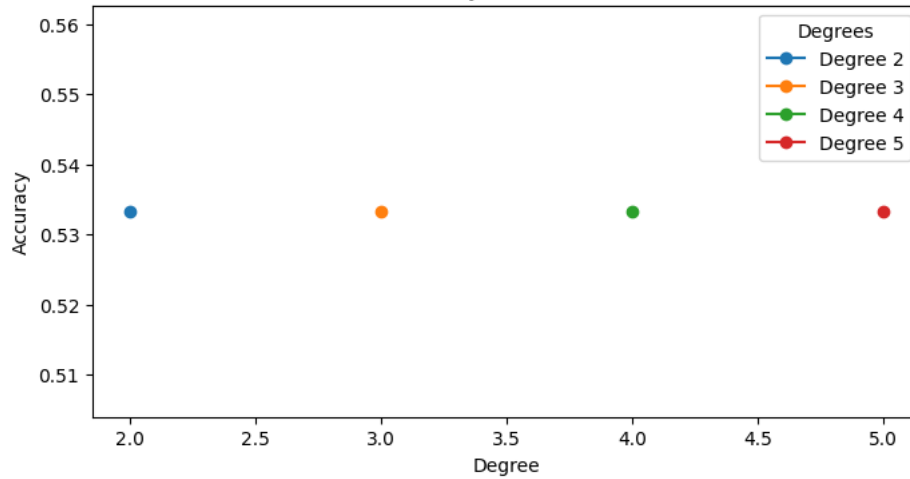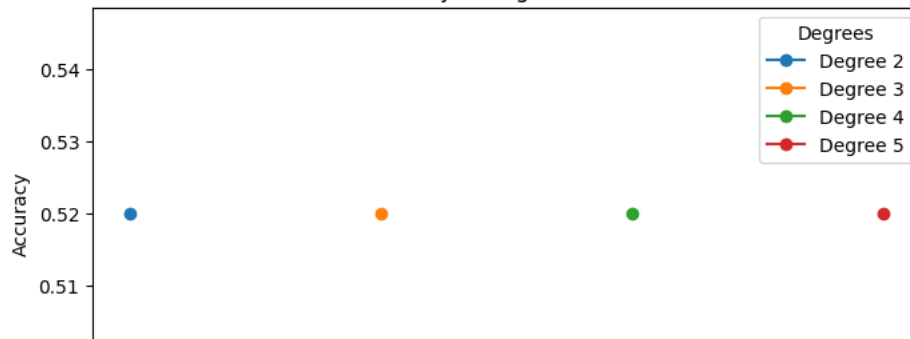
Accuracy for linear Kernel

Accuracy for poly Kernel

Accuracy for rbf Kernel

Accuracy for sigmoid Kernel

# Summary

The project serves as a comprehensive example of SVM usage in a classification task, considering different kernels and degrees. It showcases how to evaluate SVM models and provides visualizations to aid in understanding the performance of different kernel configurations.

- Linear an Poly Kernel with degree 5 achieve the highest level of accuracy -54%
- Degree Hyper Parameter is only applicable to Kernel =Polynomial