

## 1. Support Vector Machine (SVM)

The provided notebook uses a **Support Vector Machine (SVM)**, which is a powerful supervised learning algorithm used for both classification and regression tasks. The main goal of an SVM classifier is to find the optimal hyperplane that best separates data points into different classes. In our notebook, the SVM is used to classify different species of Iris flowers based on their physical measurements.

---

## 2. Data Preprocessing

The data preprocessing steps were performed as follows:

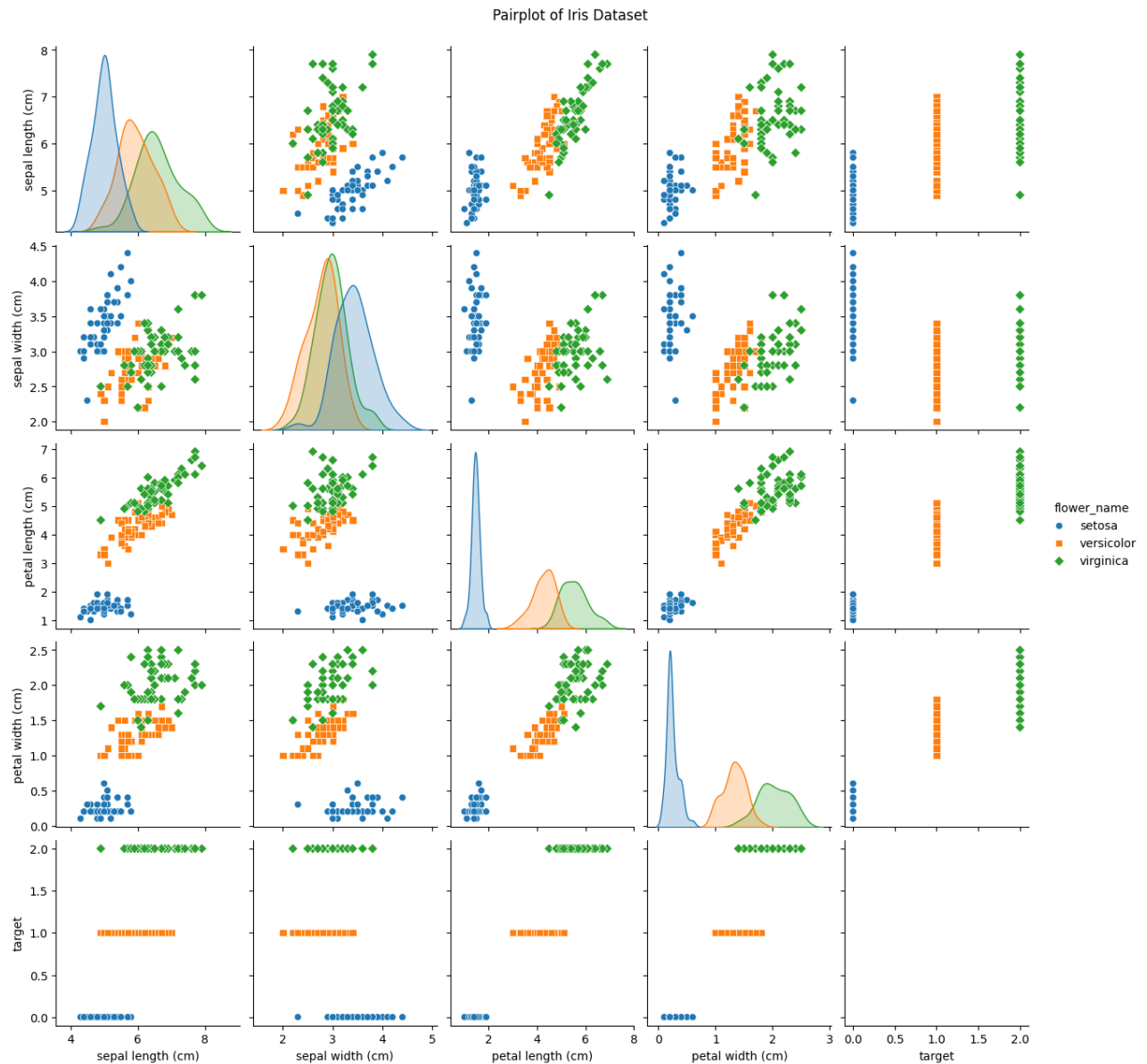
- The **Iris dataset** was loaded from `sklearn.datasets`.
  - A pandas DataFrame was created using the feature data (`iris.data`) and feature names (`iris.feature_names`).
  - The target values (`iris.target`) were added as a new column named `target`.
  - A new column, `flower_name`, was created to map the numerical targets (0, 1, 2) to their corresponding flower names (`setosa`, `versicolor`, `virginica`) for better readability.
  - The features (X) were defined by dropping the `target` and `flower_name` columns from the DataFrame.
  - The target variable (y) was defined as the `target` column.
  - The data was split into training and testing sets using `train_test_split` with a **20% test size** and a `random_state` of 10 for reproducibility.
- 

## 3. Data Visualization

Data visualization was crucial for understanding the relationships between the features and the classes. The notebook generated a pairplot and two separate scatter plots.

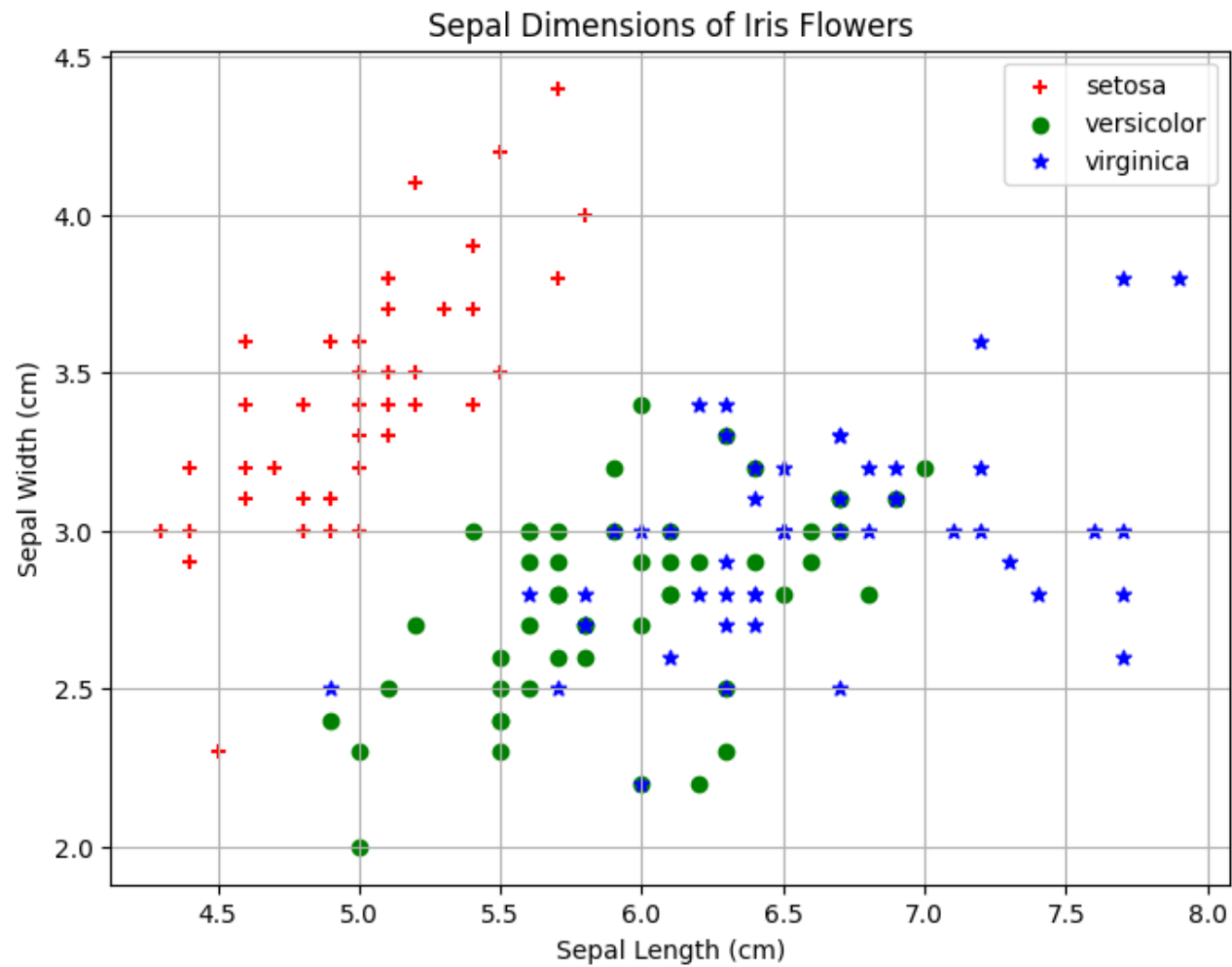
- **Pairplot Analysis:** The pairplot visualizes the relationship between all pairs of features. The diagonal plots show the distribution of each feature. The most notable observation is the clear separation of the '**setosa**' species from the other two species in almost all feature combinations. The '**versicolor**' and '**virginica**' species show some overlap, particularly in the sepal dimensions, suggesting they are harder

to distinguish than the *setosa* species.



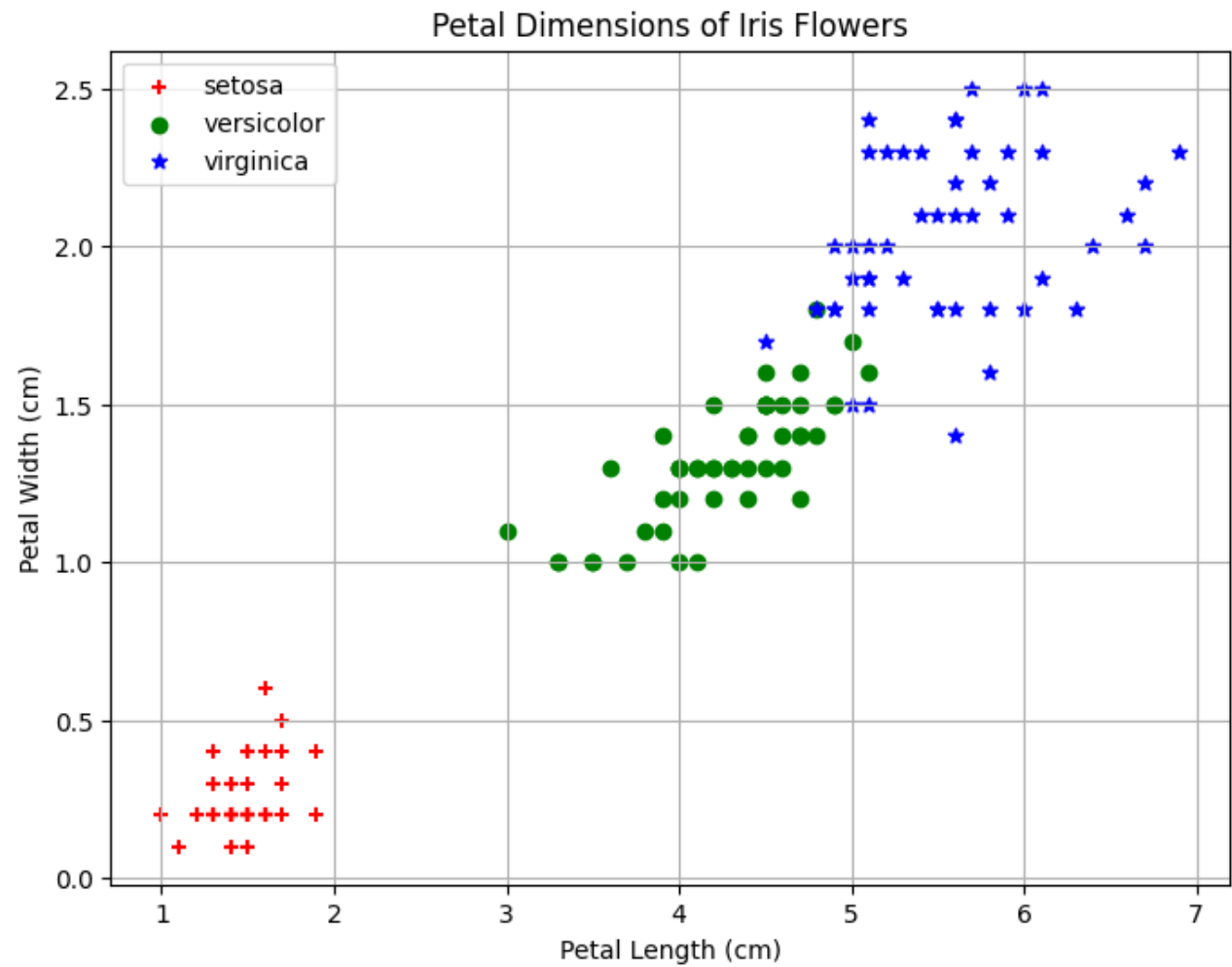
- **Sepal Dimensions Plot:** This scatter plot shows the relationship between sepal length and sepal width. While the *setosa* species is well-separated, there's significant overlap between *versicolor* and *virginica*. This visualizes why a linear model might struggle to perfectly classify these two species

using only sepal measurements.



- **Petal Dimensions Plot:** This plot displays the relationship between petal length and petal width. This is the most informative plot, as it shows a clear separation of all three species. The *setosa* species is distinctly clustered, and while *versicolor* and *virginica* are closer, they are largely linearly separable. This suggests that the petal dimensions are the most effective features for distinguishing

between the Iris species.



4. Model Training

- The SVM model was trained using the `SVC` (Support Vector Classifier) class from scikit-learn.
- The notebook experimented with different `kernel` types (`linear` and `poly`) and hyperparameters (`C` and `gamma`) to find the best performing model.
- **Radial Basis Function (RBF)** kernel is used, which is a common choice for non-linear data classification.
- The model was fit to the training data (`X_train` and `y_train`).

5. Evaluation Metrics

The model's performance was evaluated using several key metrics:

- **Accuracy:** The model achieved an **accuracy of 1.00** on the test set. This indicates that the model correctly predicted all the species in the test set.
- Model Accuracy: 1.00
- **Classification Report:** This report provides a more detailed breakdown of performance for each class.

- **Precision:** For all three species, the precision is **1.00**. This means that of all the instances predicted as a certain species, 100% of them were correct.
- **Recall:** The recall for all three species is also **1.00**. This means that the model correctly identified 100% of all actual instances of each species.
- **F1-Score:** The F1-score is the harmonic mean of precision and recall. A score of **1.00** for all classes indicates a perfect balance between precision and recall.
- **Support:** This shows the number of instances of each species in the test set (10 for *setosa*, 13 for *versicolor*, and 7 for *virginica*).
- **Confusion Matrix:** The confusion matrix visually confirms the perfect classification. The matrix shows that all predictions fall along the diagonal, with zero off-diagonal values.

The evaluation metrics collectively show that the SVM model was able to perfectly classify the Iris test data, which is a testament to both the separability of the dataset's classes and the effectiveness of the SVM algorithm for this particular problem.

Classification Report:				
	precision	recall	f1-score	support
setosa	1.00	1.00	1.00	10
versicolor	1.00	1.00	1.00	13
virginica	1.00	1.00	1.00	7
accuracy			1.00	30
macro avg	1.00	1.00	1.00	30
weighted avg	1.00	1.00	1.00	30

## 6. Why the Model Achieved 100% Accuracy

The perfect score on the test data is a direct result of the specific characteristics of the Iris dataset. While a 100% score on a small test set could indicate overfitting, in this case, it is more likely due to the inherent quality and separability of the data itself.

- **Linear Separability:** The SVM model's perfect performance is primarily due to the clear linear separability of the different Iris species based on their features. The *setosa* species is particularly well-separated from the other two species by all features, especially petal dimensions.
- **Effective Features:** The petal length and petal width features are highly effective at distinguishing between all three species. As seen in the plot, the classes form distinct, non-overlapping clusters that can be easily separated by a hyperplane. The SVM algorithm is designed to find this optimal boundary, and in this case, it was able to do so flawlessly.
- **Small Sample Size:** The test set is small, consisting of only 30 samples. While this doesn't invalidate the results, it does mean that the model's performance on a larger, more varied dataset might differ. The

perfect score is a promising indicator but should be interpreted with this small sample size in mind.