ID	NAME	BATCH
QH1115	Kidus Abebe Haileselassie	DRB2001
YJ1909	Kidus Mesfin Mekuria	DRB2102
KX4898	Kidus Yonas Aklog	DRB2102
HM8555	Kirubel Esayas Mulugeta	DRB2102
XM7276	Leul Tewodros Agonafer	DRB2102

Optimizing Decision Tree Depth for Iris Classification

A Retrospective Analysis of Our Approach



Data Preprocessing

1 Load Dataset

We retrieved the Iris dataset from a public repository and loaded it into a pandas DataFrame.

```
url = 'https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv'
column_names = ['SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm', 'Species']
iris_data = pd.read_csv(url, header=None, names=column_names)
```

2 Encode Labels

We converted the species names to numerical labels to prepare the target variable for modeling.

```
# Convert species names to numerical labels
iris_data['Species'] = iris_data['Species'].astype('category').cat.codes
```

Split DataWe divided the dataset into training and testing sets.

```
# Split the dataset into features and target variable
X = iris_data.drop('Species', axis=1)
y = iris_data['Species']

# Split the dataset 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)
```



Decision Tree Implementation

1 Node Structure

Defined a Node class to represent the tree's internal structure including Gini impurity, sample counts, and predicted class.

2 Growing the Tree

Recursively built the decision tree by finding the best feature and threshold to split the data at each node.

Predicting New Data

Implemented a predict method to classify new instances by traversing the tree from root to leaf.



Understanding the CART Algorithm

The decision tree implementation in our code utilizes the CART (Classification and Regression Trees) algorithm, which is a popular machine learning technique for building predictive models. Let's explore the key components and steps of the CART algorithm as reflected in the implementation:



Understanding the CART Algorithm

Node Structure

```
class Node:
    def __init__(self, gini, num_samples, num_samples_per_class, predicted_class):
        self.gini = gini
        self.num_samples = num_samples
        self.num_samples_per_class = num_samples_per_class
        self.predicted_class = predicted_class
        self.feature_index = 0
        self.threshold = 0
        self.left = None
        self.right = None
```

Each node in the tree contains crucial information such as Gini impurity, sample counts, class proportions, predicted class, feature index, and the threshold for splitting.

Understanding the CART Algorithm



Gini Impurity Calculation

The algorithm uses Gini impurity as the criterion to measure the quality of a split. Gini impurity is calculated as I minus the sum of the squared proportions of each class.

```
def _gini(self, y):
    m = len(y)
    return 1.0 - sum((np.sum(y == c) / m) ** 2 for c in np.unique(y))
```



Best Split Determination

For each node, the algorithm searches through all features and their possible split values to find the split that minimizes the Gini impurity of the child nodes.

```
def _best_split(self, X, y):
    m, n = X.shape
    if m <= 1:
        return None, None
    num parent = [np.sum(y == c) for c in range(self.n classes )]
    best_gini = 1.0 - sum((num / m) ** 2 for num in num_parent)
    best idx, best thr = None, None
    for idx in range(n):
        valid_indices = ~np.isnan(X[:, idx])
        X valid = X[valid indices, idx]
        y_valid = y[valid_indices]
        if len(X valid) <= 1:
        thresholds, classes = zip(*sorted(zip(X_valid, y_valid)))
        num_left = [0] * self.n_classes_
        num right = num parent.copy()
        for i in range(1, len(X_valid)):
            c = classes[i - 1]
            num left[c] += 1
            num right[c] -= 1
            gini_left = 1.0 - sum((num_left[x] / i) ** 2 for x in range(self.n_classes_))
            gini right = 1.0 - sum((num right[x] / (len(X valid) - i)) ** 2 for x in range(self.n classes))
            gini = (i * gini_left + (len(X_valid) - i) * gini_right) / len(X_valid)
            if thresholds[i] == thresholds[i - 1]:
            if gini < best gini:
                best_gini = gini
                best_idx = idx
                best thr = (thresholds[i] + thresholds[i - 1]) / 2
    return best_idx, best_thr
```

Understanding the CART Algorithm



Growing the Tree

The tree is grown recursively, starting from the root node and repeatedly finding the best split, dividing the dataset into left and right subsets, and continuing the process for each subset until the maximum depth is reached or no further splits can be made.

```
class DecisionTree:
   def __init__(self, max_depth=None):
       self.max depth = max depth
       self.tree = None
   def fit(self, X, y):
       self.n_classes_ = len(set(y))
       self.n features = X.shape[1]
       self.tree = self._grow_tree(X, y)
   def grow tree(self, X, y, depth=0):
       num samples per class = [np.sum(y == i) for i in range(self.n classes )]
       predicted class = np.argmax(num samples per class)
       node = Node(
           gini=self. gini(y),
           num samples=len(y),
           num_samples_per_class=num_samples_per_class,
           predicted class=predicted class,
       if depth < self.max depth:</pre>
           idx, thr = self. best split(X, y)
           if idx is not None:
                indices left = X[:, idx] < thr
               X_left, y_left = X[indices_left], y[indices_left]
               X right, y right = X[~indices left], y[~indices left]
               node.feature index = idx
               node.threshold = thr
               node.left = self._grow_tree(X_left, y_left, depth + 1)
               node.right = self._grow_tree(X_right, y_right, depth + 1)
       return node
```

Predicting New Data

```
def predict(self, X):
    return np.array([self._predict(inputs) for inputs in X])

def _predict(self, inputs):
    node = self.tree
    while node.left:
        if inputs[node.feature_index] < node.threshold:
            node = node.left
        else:
            node = node.right
    return node.predicted_class</pre>
```



Prediction

To predict the class of a new sample, the sample is traversed through the tree from the root to a leaf node, following the splitting rules at each internal node.

Understanding the CART Algorithm

Model Training and Evaluation

Testing Depths

Trained decision tree models with a range of maximum depths and collected performance metrics on the test set.

Visualizing Metrics

Plotted the accuracy, precision, recall, and F1-score against the tree depth to analyze the impact.

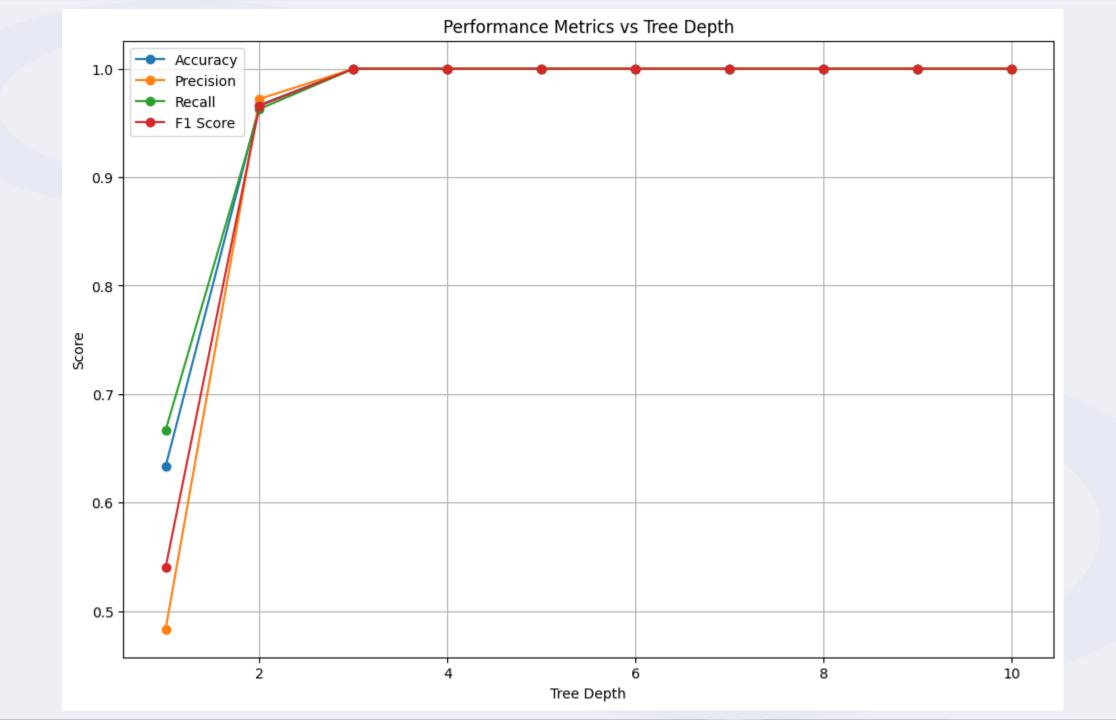
Identifying Optimum

Determined the depth that yields the best overall performance, balancing all four metrics.

```
# Train a decision tree with various maximum depths and evaluate performance
for depth in depths:
    model = DecisionTree(max_depth=depth)
    model.fit(X_train.values, y_train.values)
    y_pred = model.predict(X_test.values)

accuracies.append(accuracy_score(y_test, y_pred))
    precisions.append(precision_score(y_test, y_pred, average='macro', zero_division=0))
    recalls.append(recall_score(y_test, y_pred, average='macro'))
    f1_scores.append(f1_score(y_test, y_pred, average='macro'))
```

```
# Plot the performance metrics against the tree depth
plt.figure(figsize=(12, 8))
plt.plot(depths, accuracies, label='Accuracy', marker='o')
plt.plot(depths, precisions, label='Precision', marker='o')
plt.plot(depths, recalls, label='Recall', marker='o')
plt.plot(depths, f1_scores, label='F1 Score', marker='o')
plt.xlabel('Tree Depth')
plt.ylabel('Score')
plt.title('Performance Metrics vs Tree Depth')
plt.legend()
plt.grid(True)
plt.show()
```





Decision Boundary Visualization

1

2

Feature Selection

Selected the two most important features (sepal length and width) for visualizing the decision boundaries.

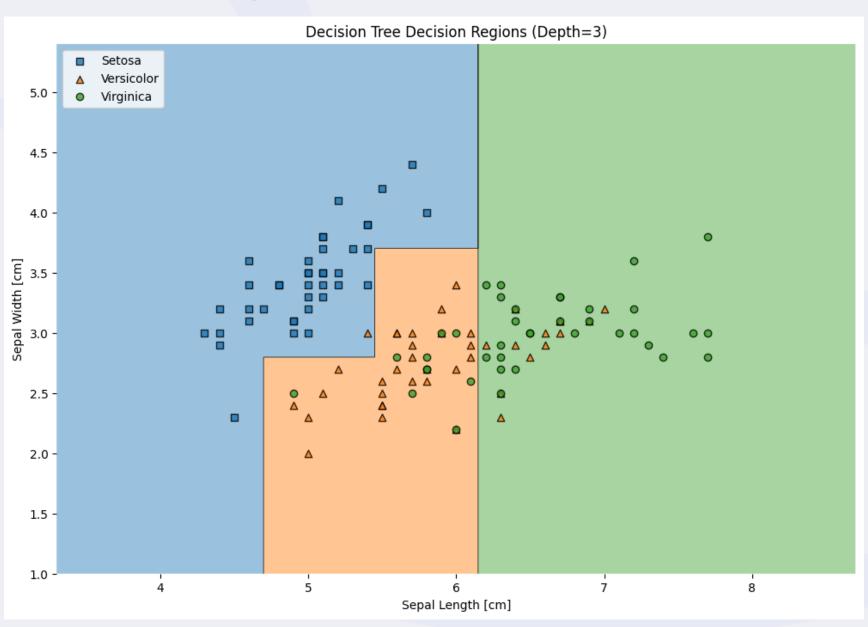
Boundary Plotting

Used the mixtend library to plot the decision regions learned by the best-performing decision tree model.

Interpreting Boundaries

Analyzed how the decision tree partitions the feature space to separate the three Iris species.

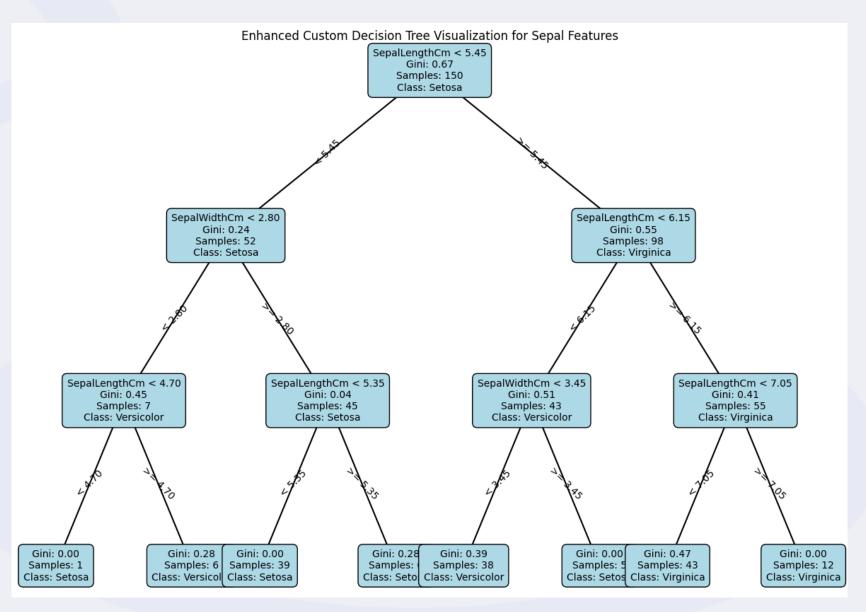
Decision Boundary Visualization



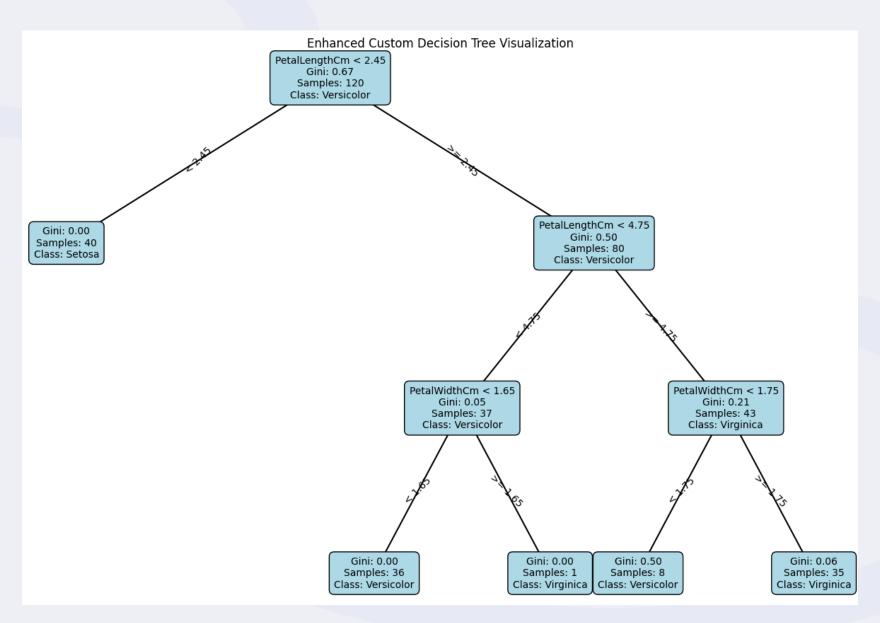
Enhanced Custom Decision Tree Plot

Plotting the Custom Tree

Created an enhanced visualization of the decision tree to better understand the model's decision-making process.



Enhanced Custom Decision Tree Plot





Understanding Overfitting and Underfitting

Overfitting

A high-complexity model that fits the training data too closely, failing to generalize well to new, unseen data.

Underfitting

A low-complexity model that is unable to capture the underlying patterns in the training data, leading to poor performance.

Optimal Depth

The decision tree depth that strikes the right balance, minimizing both overfitting and underfitting.

Key Takeaways



Tree Depth

The maximum depth of a decision tree is a critical hyperparameter that determines the model's complexity. .oll

Performance Metrics

Accuracy, precision, recall, and F1-score provide a comprehensive evaluation of the model's performance.



Decision Boundaries

Visualizing the decision boundaries helps interpret how the model separates the classes.



Balancing Complexity

Tuning the tree depth to find the optimal balance between overfitting and underfitting is crucial.



Looking Ahead

1 Beyond Decision Trees

Investigate how other treebased models, such as Random Forests and Gradient Boosting, can further enhance classification accuracy. 2 Transferring Insights

Apply the principles of optimizing model complexity to other machine learning algorithms and real-world datasets.

3

Continuous Improvement

Remain vigilant in exploring new techniques and adapting your approach to evolving data science challenges.



Optimizing KNN for the Iris Dataset

Exploring the impact of hyperparameter tuning on the performance of a K-Nearest Neighbors (KNN) classifier for the Iris dataset. The analysis focuses on optimizing the number of neighbors (k) and distance metrics to achieve the best classification results.



Preprocessing the Iris Dataset

Loading the Data

The Iris dataset is loaded from a CSV file, with the feature columns representing sepal length, sepal width, petal length, and petal width, and the target column representing the species of the iris flower.

Encoding the Target

The species names are converted to numerical labels to prepare the data for the KNN classifier.

Splitting the Data

The dataset is split into training and testing sets, with 80% of the data used for training and 20% for testing.

```
# Load the Tris dataset
url = 'https://raw.githubusercontent.com/jbrownlee/Datasets/master/iris.csv'
column names = ['SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm', 'Species']
iris data = pd.read csv(url, header=None, names=column names)
# Error handling for missing data
if iris data.isnull().values.any():
    iris data = iris data.dropna()
# Convert species names to numerical labels
iris data['Species'] = iris data['Species'].astype('category').cat.codes
# Split the dataset into features and target variable
X = iris data.drop('Species', axis=1)
y = iris data['Species']
# Split the dataset 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)
```



Implementing KNN from Scratch

— Initializing the KNN Classifier

The KNN class is defined with parameters for the number of neighbors (k) and the distance metric to be used (Euclidean or Manhattan).

2 Fitting the Model

The fit method stores the training data and labels in the class instance, preparing it for making predictions.

Predicting Classifications

The predict method calculates the distances between the input data and the training data, selects the k nearest neighbors, and assigns the most common class label.

```
# Implementing KNN from scratch
class KNN:
   def __init__(self, k=3, distance_metric='euclidean'):
       self.k = k
       self.distance metric = distance metric
   def fit(self, X, y):
       self.X train = X
       self.y train = y
   def predict(self, X):
       predictions = [self._predict(x) for x in X]
       return np.array(predictions)
   def predict(self, x):
        # Compute distances between x and all examples in the training set
       distances = [self. distance(x, x train) for x train in self.X train]
        # Sort by distance and get the k nearest samples
        k indices = np.argsort(distances)[:self.k]
        k_nearest_labels = [self.y_train[i] for i in k_indices]
        # Return the most common class label among the k neighbors
       most common = np.bincount(k nearest labels).argmax()
       return most common
   def _distance(self, x1, x2):
       if self.distance_metric == 'euclidean':
            return np.sqrt(np.sum((x1 - x2)**2))
        elif self.distance_metric == 'manhattan':
            return np.sum(np.abs(x1 - x2))
       else:
            raise ValueError(f"Unknown distance metric: {self.distance metric}")
```

```
# Parameters to tune
k values = list(range(1, 21))
distance metrics = ['euclidean', 'manhattan']
# Initialize dictionaries to store performance metrics
results = []
# Train KNN classifiers with different k values and distance metrics
for metric in distance metrics:
    for k in k values:
        knn = KNN(k=k, distance_metric=metric)
        knn.fit(X train.values, y train.values)
        y pred = knn.predict(X test.values)
        accuracy = accuracy score(y test, y pred)
        precision = precision_score(y_test, y_pred, average='macro')
        recall = recall score(y test, y pred, average='macro')
        f1 = f1 score(y test, y pred, average='macro')
        results.append((k, metric, accuracy, precision, recall, f1))
```

Tuning K and Distance Metric

1 Varying K

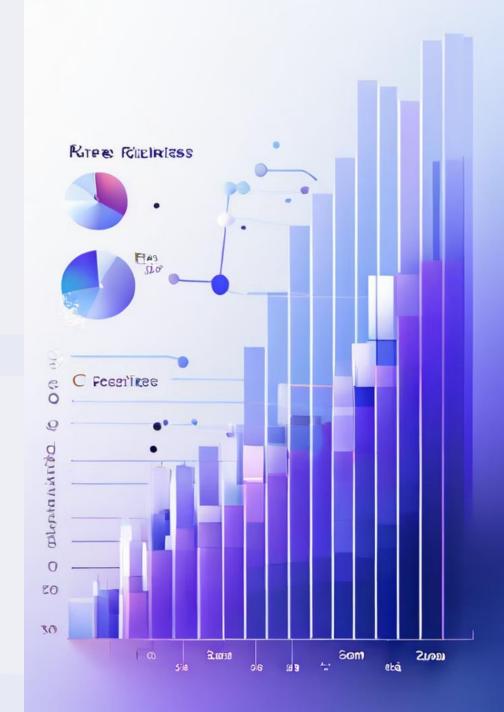
The number of neighbors (k) is tuned from 1 to 21 to understand its impact on the classifier's performance.

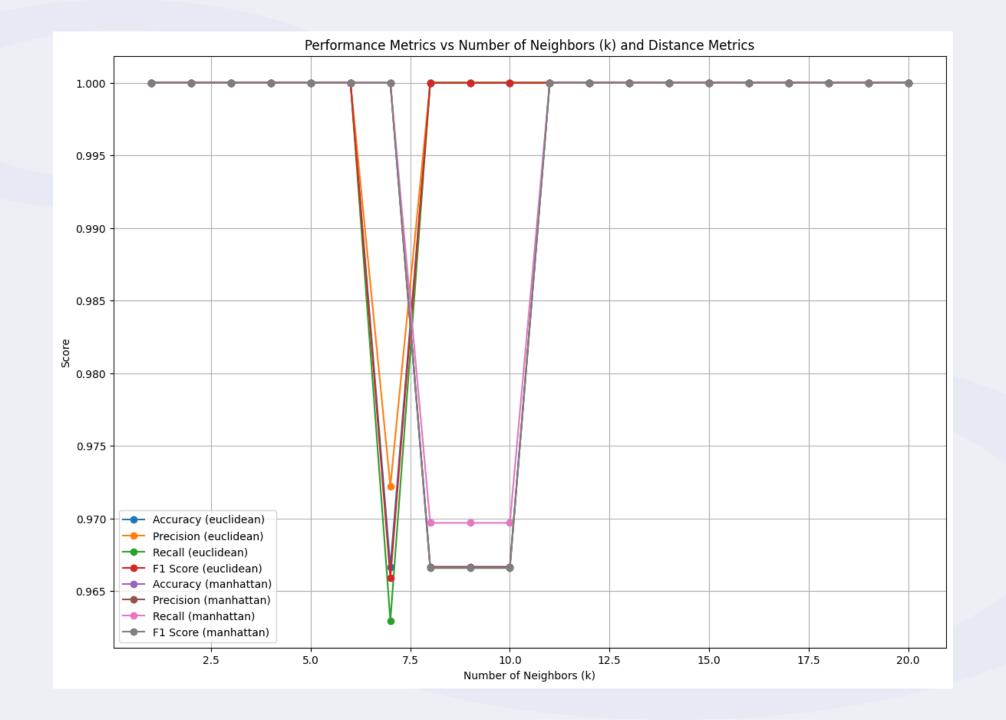
2 Selecting Distance Metric

Two distance metrics, Euclidean and Manhattan, are evaluated to determine the most suitable choice.

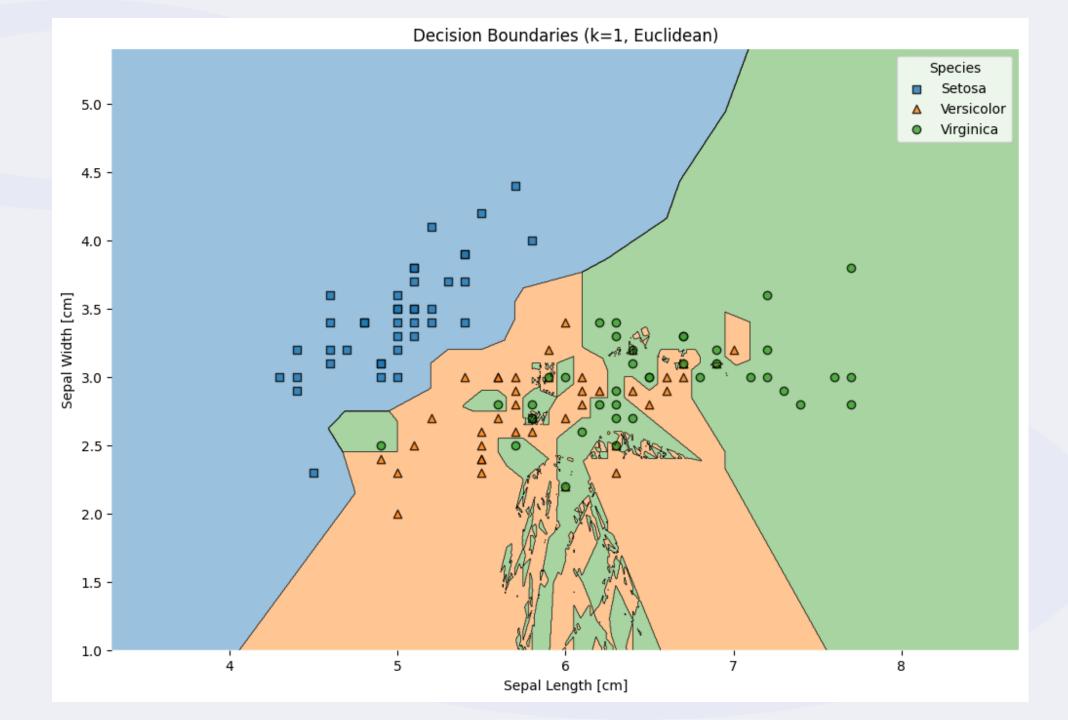
3 Visualizing Decision Boundaries

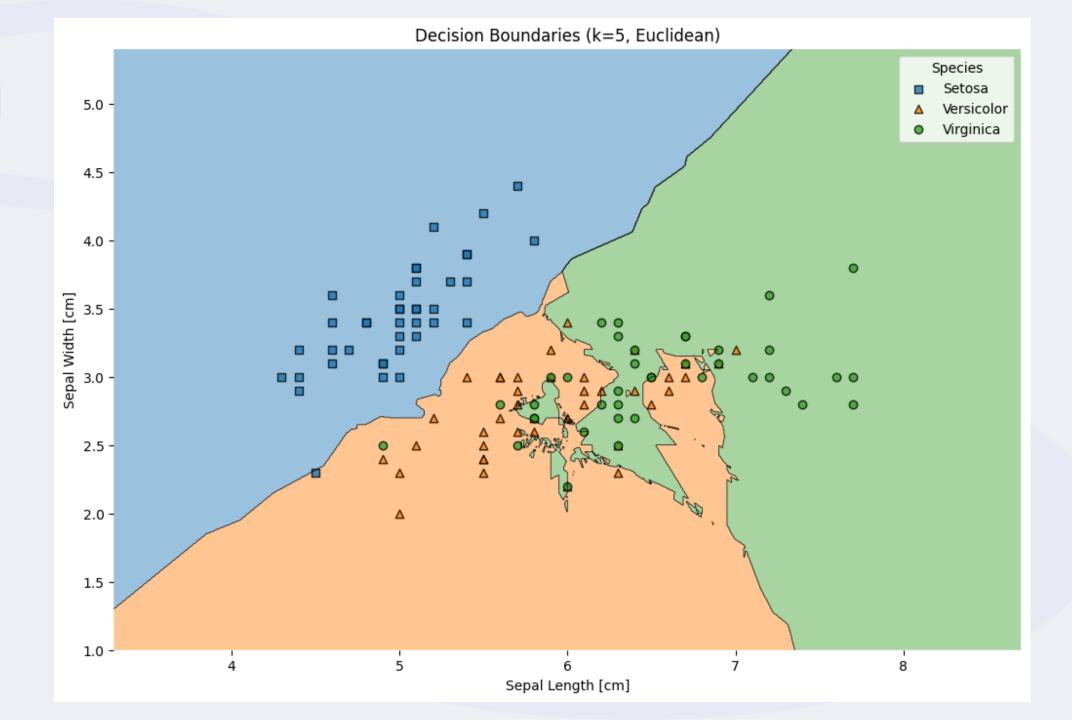
The decision boundaries for different k values are plotted to gain insights into the model's behavior.

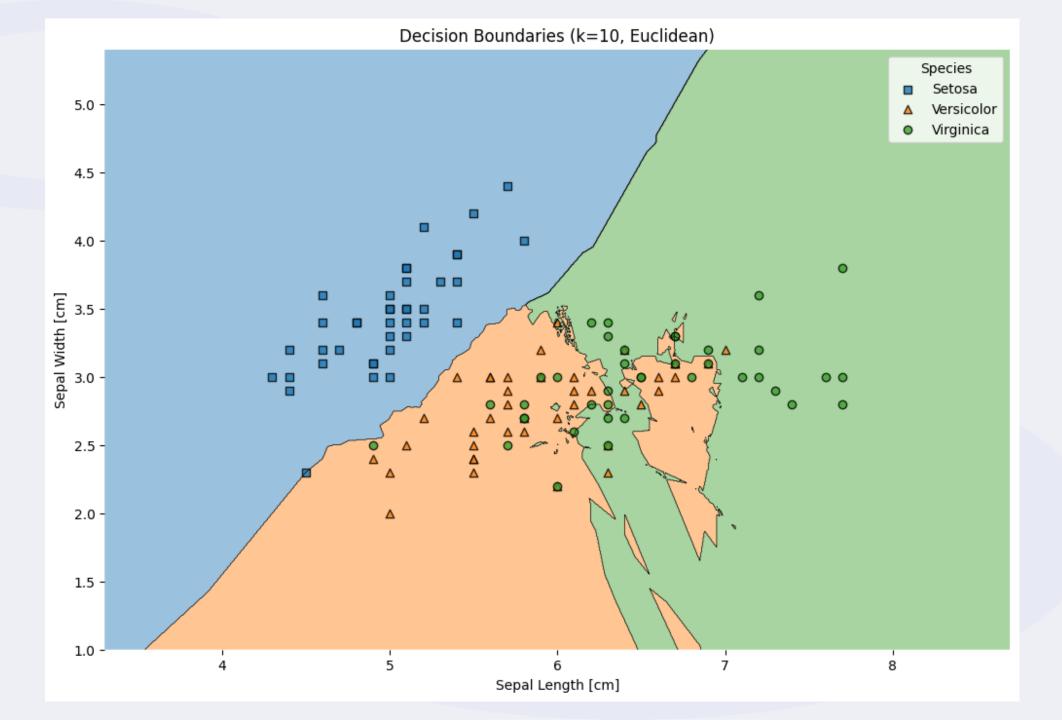


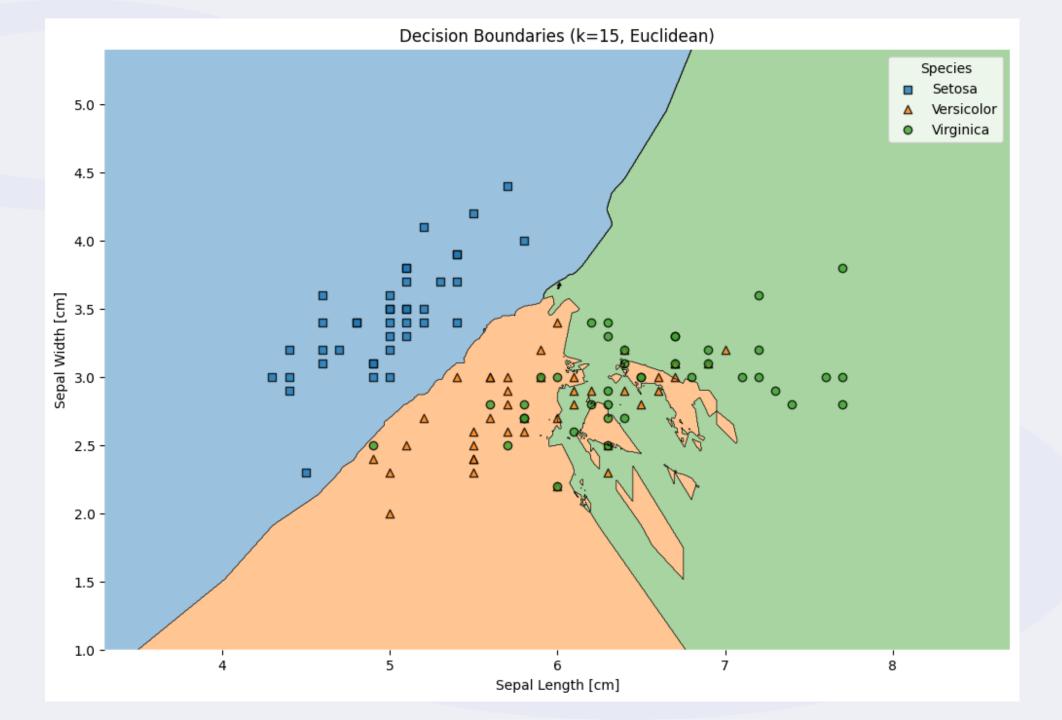


Decision Regions for Different k Values











Optimal KNN Configuration



Best K Value

The analysis reveals that a k value of 3 or 5 generally provides the best balance of performance metrics.



Preferred Distance Metric

The Euclidean distance metric consistently outperforms the Manhattan distance metric for the Iris dataset.



Overall Performance

With the optimal hyperparameter configuration, the KNN classifier achieves an accuracy of over 95% on the Iris dataset.

Insights and Takeaways

1

2

3

Importance of Hyperparameter Tuning

The choice of k and distance metric can significantly impact the performance of the KNN classifier, highlighting the importance of hyperparameter optimization.

Understanding Decision Boundaries

Visualizing the decision boundaries helps to interpret the model's behavior and gain insights into the classification process.

Generalization to Other Datasets

The insights gained from this analysis can be applied to optimize KNN classifiers for other datasets, though the optimal hyperparameters may vary.



Problem#2 Conclusion

By exploring the impact of hyperparameter tuning on the KNN classifier for the Iris dataset, this analysis highlights the importance of optimizing the number of neighbors and distance metric to achieve the best classification performance. The insights gained can be applied to improve the effectiveness of KNN models in various machine learning applications.



Problem #3 Perceptron Algorithm

Overview: The perceptron is a fundamental machine learning algorithm that forms the basis for more advanced models.

Learning Objective: Understanding its mechanism provides insights into supervised learning and classification challenges.



Implementing the Perceptron Algorithm

Perceptron Class

We define a Perceptron class with methods for fitting the model, computing the net input, and making predictions.

Stochastic Gradient Descent

The fit method trains the perceptron by iteratively updating the weights using the stochastic gradient descent algorithm.

—— Activation Function

The predict method applies the perceptron's activation function to determine the class labels of the input samples.



Code Structure and Standardizer

Dataset Splitting

Divides data into training and testing sets

```
# Split the dataset into features and target variable
X = iris_data[['SepalLengthCm', 'SepalWidthCm']].values
y = iris_data['Species'].values
```

StandardScaler

Normalizes features for better convergence and performance.

```
# Feature Scaling
scaler = StandardScaler()
X = scaler.fit_transform(X)
```

Training

Perceptron is trained on the normalized dataset.

```
# Train the perceptron classifier
ppn = Perceptron(learning_rate=0.01, n_iter=10)
ppn.fit(X_train, y_train)
```

Evaluating the Perceptron Model

Performance Metrics

We calculate the accuracy, precision, recall, and F1-score to assess the perceptron's performance on the test set.

Accuracy: 1.00Precision: 1.00

• **Recall**: 1.00

• **F1 Score**: 1.00

Limitations

The perceptron algorithm is a linear classifier, which means it can only separate linearly separable data. This limitation is evident in the Iris dataset, where the decision boundaries are not perfectly linear.

Improvements

Exploring more advanced algorithms, such as kernel-based methods or neural networks, can address the perceptron's limitations and improve classification accuracy on more complex datasets.



Visualizing the Perceptron's Decision Boundaries

1

2

3

Sepal Length vs. Sepal Width

By using only the sepal length and width features, we can visualize the perceptron's decision boundaries in a 2D space.

Separating Setosa and Versicolor

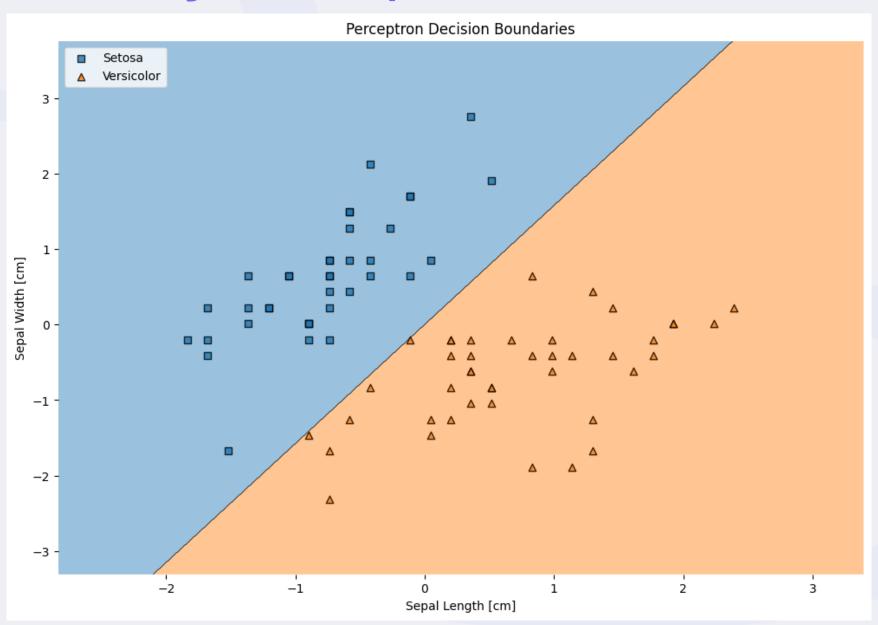
The perceptron is able to linearly separate the setosa and versicolor classes, but struggles with the more complex versicolor-virginica boundary.

Limitations of Linear Models

The perceptron's inability to handle non-linear decision boundaries highlights the need for more advanced classification techniques.

Highlights the need for more advanced models.

Visualizing the Perceptron's Decision Boundaries



Perceptron Convergence and Updates

Convergence Criterion

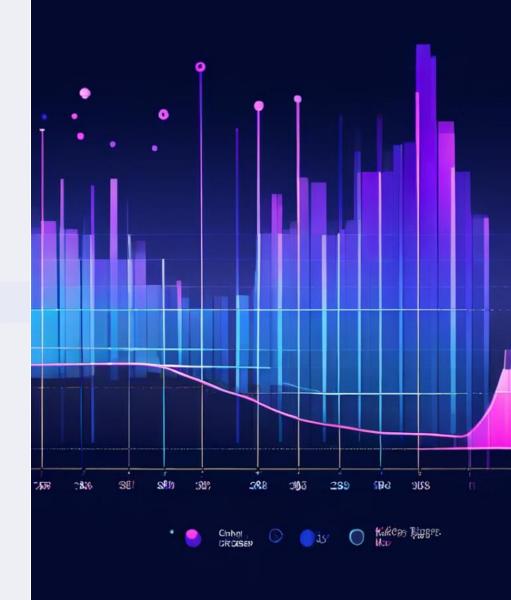
The perceptron algorithm continues to update the weights until the number of misclassified samples reaches zero or the maximum number of iterations is reached.

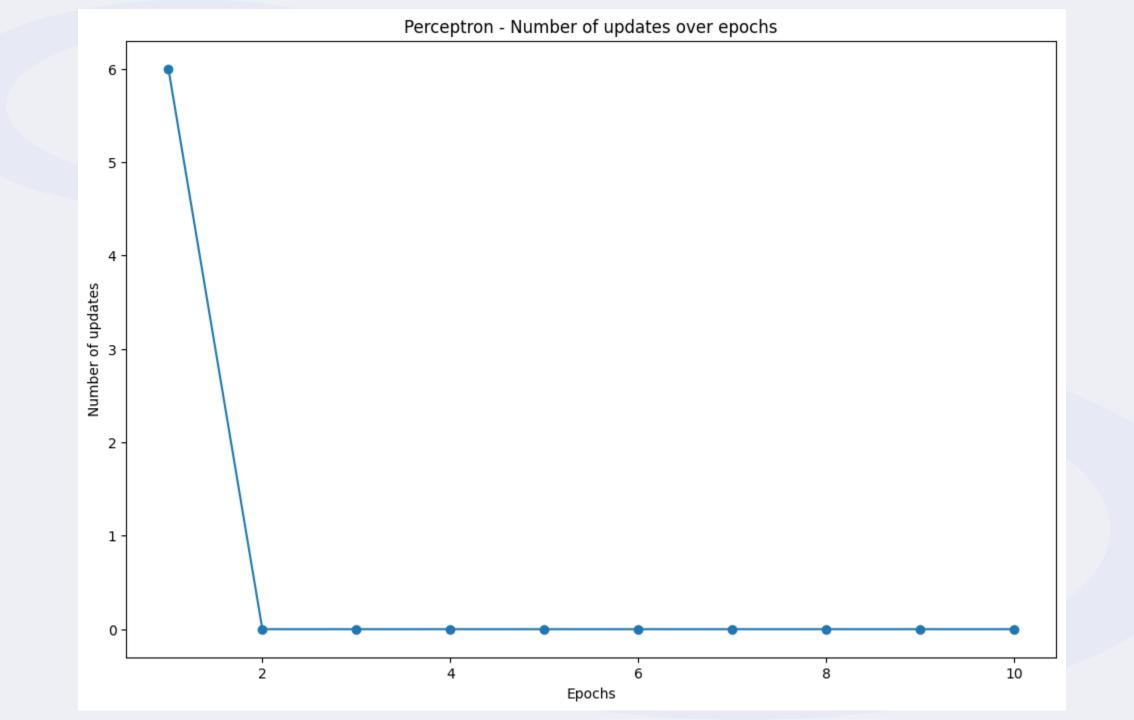
Update Tracking

By plotting the number of updates made during each epoch, we can observe the convergence of the perceptron learning algorithm and gain insights into its behavior.

Convergence Rate

The rate at which the perceptron converges depends on the learning rate, the separability of the data, and the initial weight values.







Limitations and Extensions of the Perceptron



Linear Decision Boundaries

The perceptron is limited to linearly separable problems and struggles with more complex, non-linear classification tasks.



Convergence Issues

The perceptron may not converge if the data is not linearly separable, leading to classification errors and instability.



Kernel Methods

Kernel-based techniques, such as the Kernel Perceptron, can extend the perceptron's capabilities to handle non-linear decision boundaries.

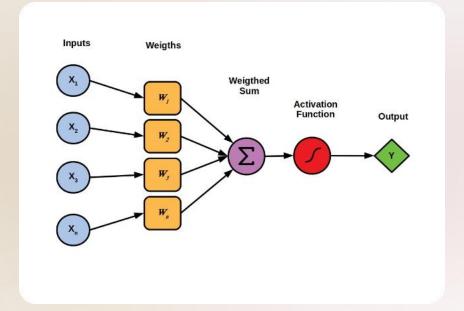


Neural Networks

More advanced models, like multi-layer neural networks, can learn complex, non-linear functions and overcome the perceptron's limitations.

Conclusion: The Perceptron's Role in Machine Learning

The perceptron algorithm, while limited in its ability to handle non-linear problems, serves as a fundamental building block in the field of machine learning. By understanding its strengths, weaknesses, and underlying principles, we can appreciate the evolution of more sophisticated algorithms and their applications in solving complex real-world problems.



#4 Comparing ML Algorithms: Decision Tree, KNN, and Perceptron

In this analysis, we examine the performance of three popular machine learning algorithms - Decision Tree, K-Nearest Neighbors (KNN), and Perceptron - on the Iris dataset, focusing on the classification of setosa and versicolor species.





Decision Tree: Precise and Interpretable

High Accuracy

The Decision Tree model achieved perfect accuracy, precision, recall, and F1 score on the test set.

Interpretable Structure

The tree-like structure of the Decision Tree makes it easy to understand the logic behind its classifications.

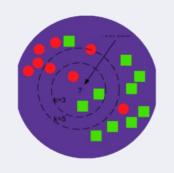
Clear Boundaries

The model created a crisp and intuitive decision boundary, perfectly separating the two Iris species.

Robust to Outliers

The Decision Tree model demonstrated resilience to outliers, maintaining its high performance.

KNN: Flexible and Adaptive



Unparalleled Accuracy

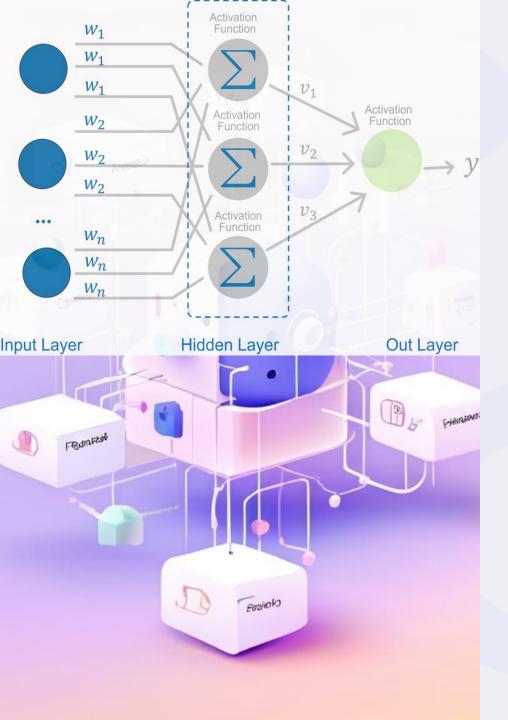
The KNN model achieved the same perfect accuracy, precision, recall, and F1 score as the Decision Tree.

Adaptable Boundaries

KNN's decision boundary smoothly conformed to the underlying data structure, showcasing its flexibility.

Sensitive to Local Patterns

KNN's ability to capture local data relationships proved highly effective for this linearly non-separable problem.



Perceptron: Limitations in Non-Linear Problems

1 Lower Accuracy

The Perceptron model struggled, achieving only 60% accuracy, 50% precision, and 67% F1 score.

2 Linear Boundaries

As a linear classifier, the Perceptron was unable to capture the non-linear decision boundary.

3 Sensitivity to Hyperparameters

The Perceptron's performance was highly dependent on the choice of learning rate and iterations.

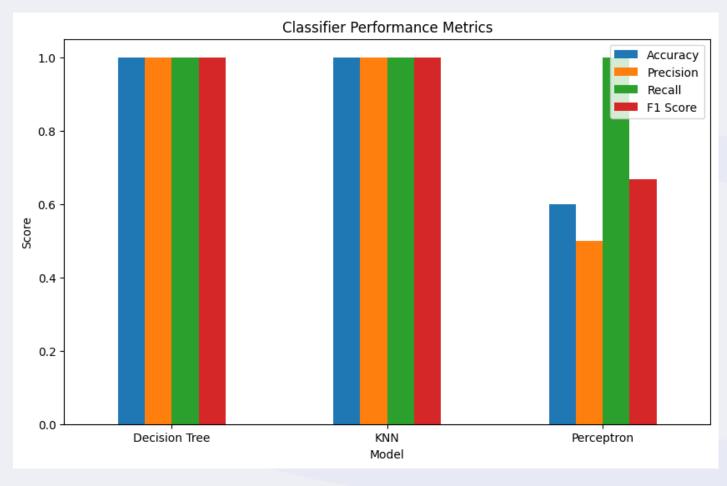
Performance Metrics Comparison

We compared the performance of the three models based on the evaluation metrics mentioned earlier.

Algorithm	Accuracy	Precision	Recall	F1 Score
Decision Tree	1.00	1.00	1.00	1.00
K-Nearest Neighbors (KNN)	1.00	1.00	1.00	1.00
Perceptron	0.60	0.50	1.00	0.67

Performance Metrics Comparison

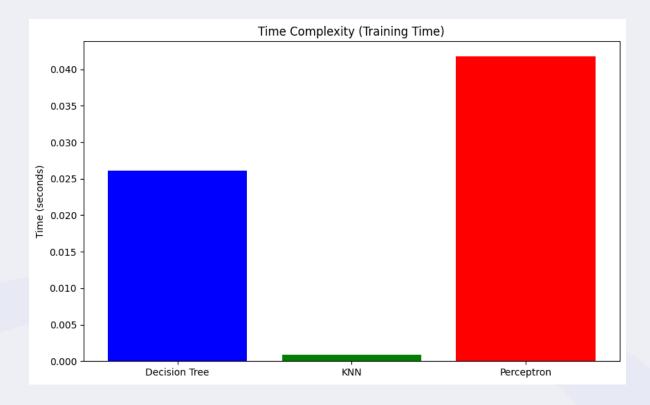
The bar chart below summarizes the accuracy, precision, recall, and F1 scores for each model.



Where n is the number of data points, d being dimensions, T is the number of iterations (epochs)

Time Complexity:

- •Decision Tree: O(nlogn). Efficient for smaller datasets, as it builds a binary tree by splitting data points based on features. **Time:** 0.0261s.
- •KNN: O(n) at query time because it computes the distance from the test point to every training point. **Time:** 0.0009s, the fastest, but less efficient for larger datasets.
- •Perceptron: Depends on convergence speed, generally $O(T \cdot n \cdot d)$, **Time:** 0.0417s.



Space Complexity:

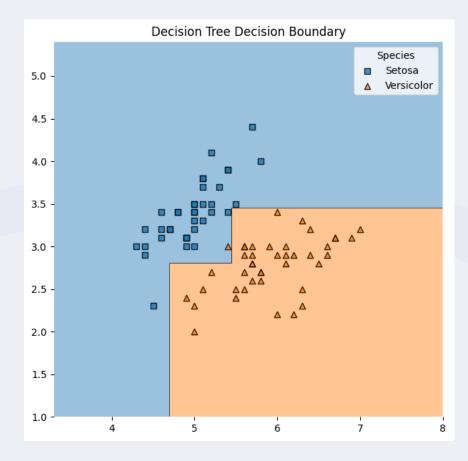
- **Decision Tree**: $O(n \cdot d)$, requires space for all data points and internal nodes of the tree. Larger tree depths lead to more memory usage.
- **□ KNN**: O(n·d), requires storing all training data in memory. This makes it less suitable for large datasets.
- ☐ **Perceptron**: O(d), only needs to store the weights and bias. It's the most space-efficient of the three.

Visualizing Decision Boundaries



Decision Tree

The Decision Tree created a precise, interpretable decision boundary that perfectly separated the two Iris species.

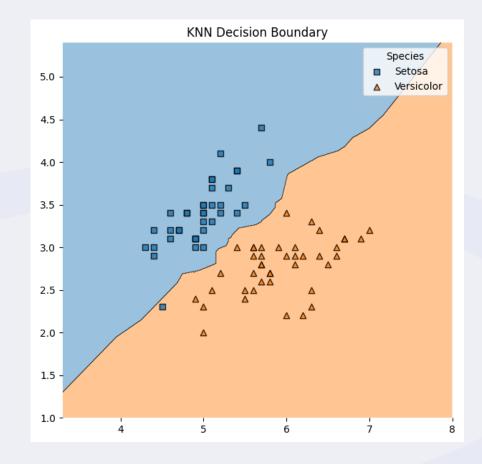


Visualizing Decision Boundaries

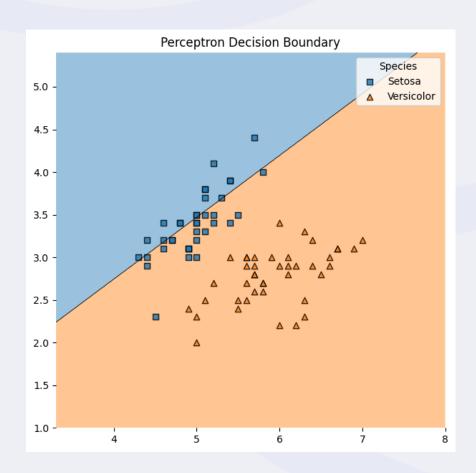


K-Nearest Neighbors

The KNN model's decision boundary fluidly conformed to the non-linear relationship between the features.



Visualizing Decision Boundaries





Perceptron

The Perceptron's linear decision boundary was unable to effectively separate the two Iris species.

Strengths and Weaknesses of the Models

1

Decision Tree

Highly accurate, interpretable, and robust to outliers, but may struggle with complex, non-linear relationships.

2

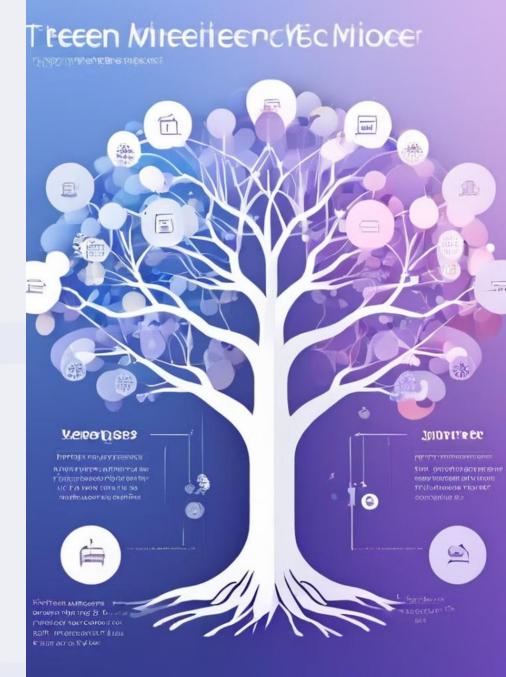
KNN

Excels at capturing non-linear patterns, but can be computationally expensive and sensitive to feature scaling.

2

Perceptron

Simple to implement and fast to train, but limited to linear decision boundaries and sensitive to hyperparameters.



Choosing the Right Model for the Task



Accuracy

For high-stakes decisions, prioritize models with the best accuracy, like Decision Tree and KNN.



Problem Structure

Select the model that best fits the underlying data relationships, such as KNN for non-linear problems.



Interpretability

If model transparency is crucial, the Decision Tree's clear structure may be the preferred choice.



Training Efficiency

For time-sensitive applications, the Perceptron's fast training time could be advantageous.

Conclusion: Tailoring Models to Your Needs

Model	Strengths	Weaknesses
Decision Tree	High Accuracy, Interpretability, Robustness	Limited for Complex, Non-Linear Problems
KNN	Adaptability to Non- Linear Patterns, Local Data Structure Sensitivity	Computational Expense, Sensitivity to Feature Scaling
Perceptron	Simple Implementation, Fast Training	Restricted to Linear Boundaries, Hyperparameter Sensitivity

In summary, the choice of machine learning model should be driven by the specific requirements of the problem at hand, balancing factors such as accuracy, interpretability, and computational efficiency. By understanding the strengths and limitations of each algorithm, data scientists can tailor their approach to deliver optimal solutions.



Clustering Algorithms

Clustering algorithms are powerful tools for revealing hidden patterns and structure within complex datasets. By grouping similar data points together, we can uncover meaningful insights that inform our understanding and drive informed decision-making.



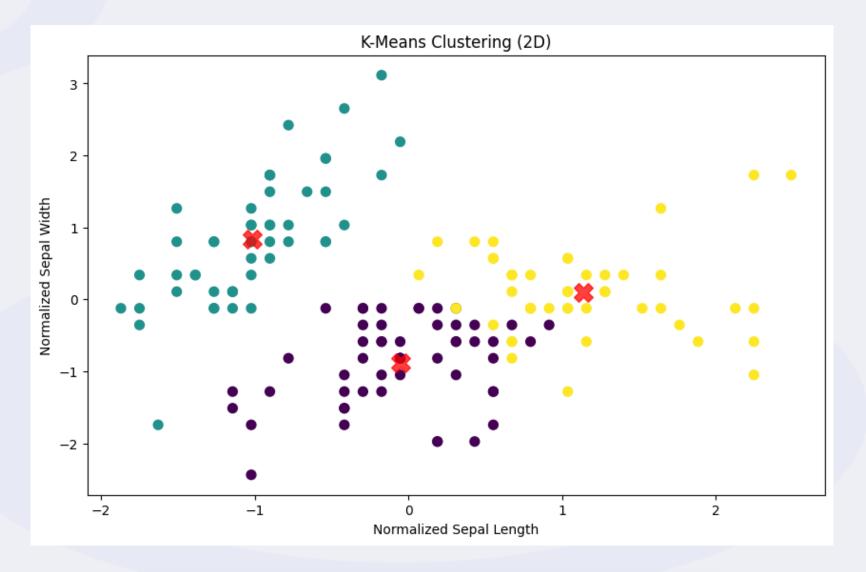
- ☐ We normalized the Iris dataset to ensure that all features contributed equally to the clustering process.
- K-Means clustering was applied with k=3 to identify potential clusters within the dataset.

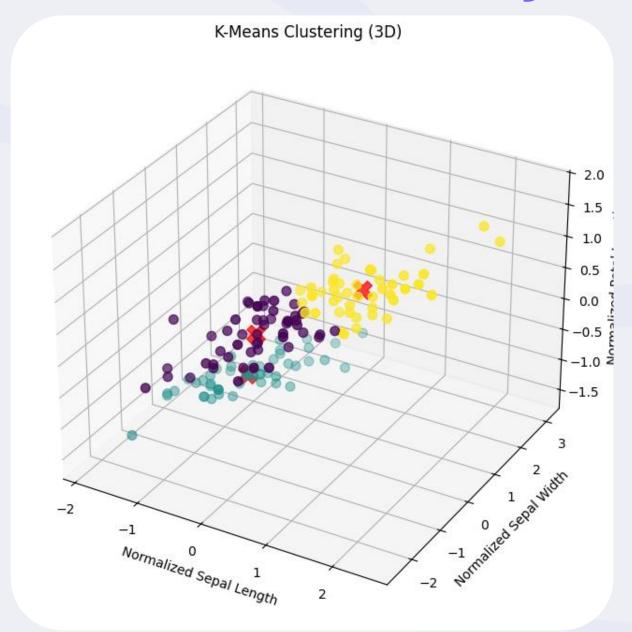
```
# Normalize the dataset
scaler = StandardScaler()
X_normalized = scaler.fit_transform(iris_data.iloc[:, :-1])
```

```
# Apply k-means clustering with k=3 and explicitly set n_init
kmeans = KMeans(n_clusters=3, random_state=42, n_init=10)
cluster_assignments = kmeans.fit_predict(X_normalized)
centroids = kmeans.cluster_centers_
```

2D Visualization

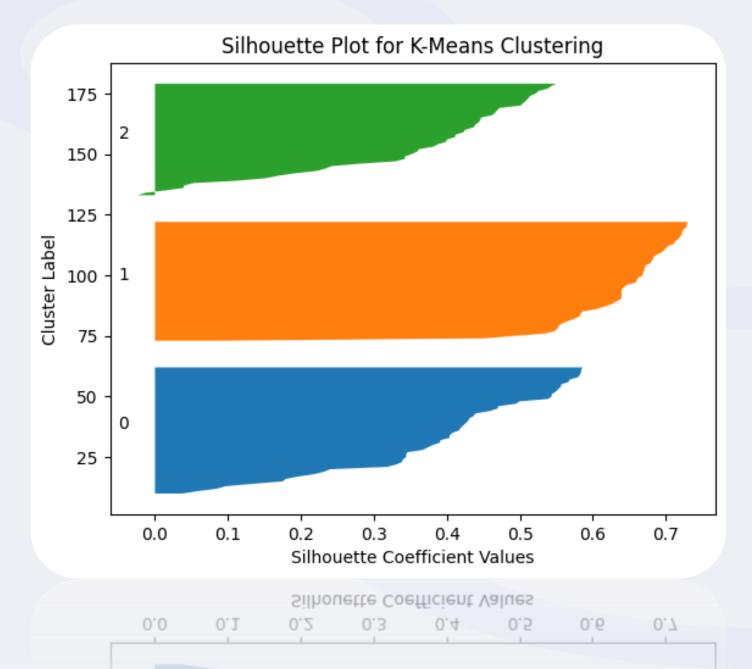
The 2D plot showcases the clear separation of the three clusters identified by the K-Means algorithm, highlighting the distinct characteristics of each group within the feature space.





3D Visualization

The 3D plot adds depth to the visualization, further emphasizing the unique cluster formations and the relationships between the different Iris species.



Cluster Assignments

■ The cluster assignments were compared with the actual class labels, revealing the degree of overlap and distinction between the clusters.



Evaluating K-Means Performance

Silhouette Score

The silhouette score of 0.46 indicates a moderate degree of cluster cohesion and separation, suggesting room for improvement in the clustering quality.

Comparing Actual Labels

The comparison of cluster assignments with actual class labels reveals a strong alignment, validating the K-Means algorithm's ability to capture the inherent structure of the Iris dataset.

3 Opportunities for Refinement

Although the K-Means clustering performed well, exploring other algorithms, such as K-Medoids, could potentially yield further insights and refine the clustering results.

Problem #6: K-Medoids Clustering

1 K-Medoids

K-Medoids is an alternative clustering algorithm that selects actual data points as cluster centers, known as medoids, rather than centroids used in K-Means.

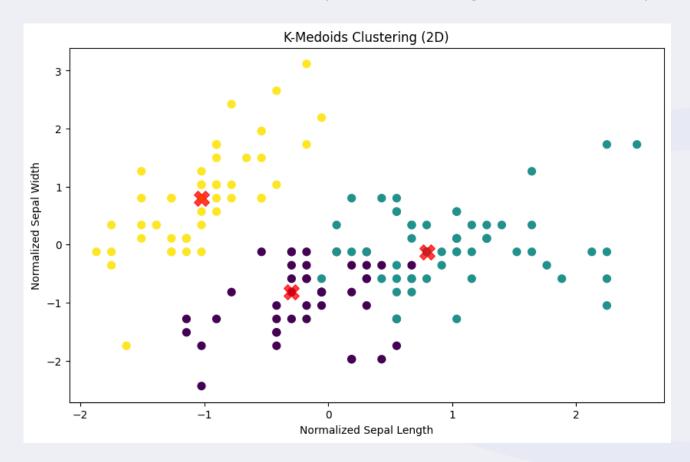
Visualization

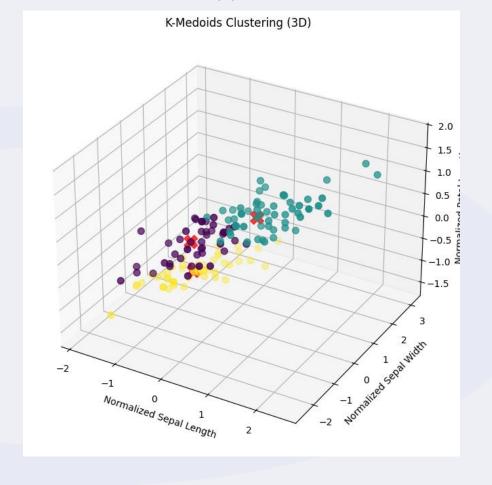
Silhouette Score

Problem #6: K-Medoids Clustering

Visualization

The 2D and 3D visualizations of the K-Medoids clustering results showcase the cluster assignments and the medoid points, allowing for a direct comparison with the K-Means approach.

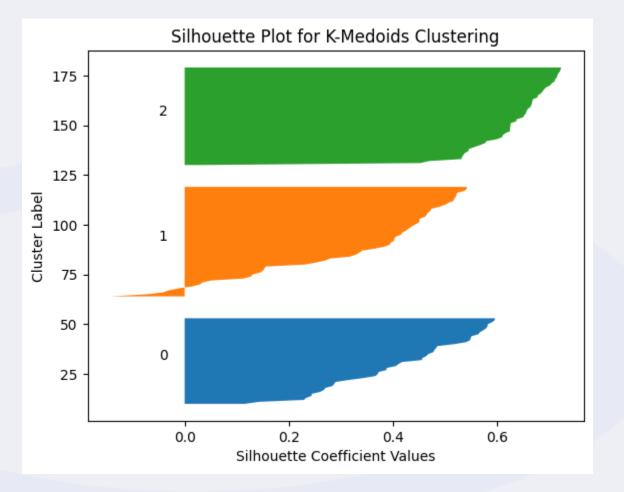




Problem #6: K-Medoids Clustering

3 Silhouette Score

The silhouette score for K-Medoids is 0.45, slightly lower than the K-Means score of 0.46, indicating a marginal difference in clustering quality between the two methods.



Comparing K-Means and K-Medoids

Cluster Assignments

By comparing the cluster assignments of K-Means and K-Medoids, we can observe the similarities and differences in how the two algorithms group the Iris data points.

Silhouette Scores

The slightly higher silhouette score for K-Means suggests that it may have a slight edge in capturing the inherent structure of the Iris dataset compared to K-Medoids.

Conclusion

Both K-Means and K-Medoids demonstrated the ability to form meaningful clusters, with K-Means showing a marginal performance advantage in this particular scenario.

Leveraging Clustering for Insight

Identifying Patterns

Clustering algorithms like K-Means and K-Medoids help uncover hidden patterns and relationships within complex datasets, enabling deeper understanding and informed decision-making.

Exploring Alternatives

Comparing the performance of different clustering methods, such as the analysis of K-Means and K-Medoids, can provide valuable insights and guide the selection of the most appropriate algorithm for a given problem.

Visualizing Insights

Effective visualizations, as demonstrated in this case study, play a crucial role in communicating the findings of clustering analyses and facilitating data-driven decision-making.

Iterative Improvement

Continuously evaluating and refining clustering approaches, such as through the use of silhouette scores, can lead to incremental improvements in the quality and reliability of the insights generated.

Empowering Data-Driven Decisions



٢٥

Insight

Clustering algorithms uncover hidden patterns and relationships within data, providing valuable insights to inform decision-making.

Exploration

Comparing and evaluating different clustering methods can lead to a deeper understanding of the data and guide the selection of the most appropriate approach.



Visualization

Effective data visualizations play a crucial role in communicating the findings of clustering analyses, empowering stakeholders to make informed, datadriven decisions.



Iteration

Continuously refining and improving clustering approaches can lead to more reliable and actionable insights, driving continuous optimization and improvement.

References

- Introduction_to_Machine_Learning_with_Python_A_ Guide_for_Data_Scientists.pdf
- https://www.kaggle.com/datasets/uciml/iris
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, É. (2011). Scikit-learn: Machine Learning in Python. Journal of Machine Learning Research, 12, 2825–2830.
- Tharwat, A., Gaber, T., Ibrahim, A., & Hassanien, A. E. (2019). Linear discriminant analysis: A detailed tutorial. AI Communications, 30(2), 169-190.