



Comparative Analysis of Machine Learning Algorithms

Nidhithashree K¹, Dr. Dattatreya P Mankame², Dr. Basavaraj Patil³, Mrs. Veena Dhavalgi⁴

Department of Computer Science and Business System
Dayananda Sagar College of Engineering, Bangalore, Karnataka, India

ABSTRACT –

This paper provides a comprehensive comparison of various machine learning algorithms, focusing on supervised and unsupervised learning techniques, and deep learning methods. It evaluates the performance of these algorithms using standard metrics across different datasets and application domains, highlighting their strengths, weaknesses, and suitable use cases.

Index Terms – Machine Learning, Comparative Analysis, Logistic Regression, Decision Trees, SVM, k-NN, Neural Networks

1. INTRODUCTION :

Machine learning algorithms play a crucial role in data analysis and prediction tasks. This paper aims to compare the performance of four popular classification algorithms: Decision Tree, Support Vector Machine, k-Nearest Neighbors, and Logistic Regression. The Iris dataset, a well-known dataset in the machine learning community, is used for this comparison.

2. METHODOLOGY :

2.1 Dataset

The Iris dataset consists of 150 samples from three species of Iris flowers (Iris setosa, Iris versicolor, and Iris virginica). Each sample has four features: sepal length, sepal width, petal length, and petal width.

2.2 Algorithms

The following algorithms are implemented and compared:

- Decision Tree
- Support Vector Machine (SVM)
- k-Nearest Neighbors (k-NN)
- Logistic Regression

Decision tree

A Decision Tree is a non-parametric supervised learning method used for classification and regression. It models data using a tree-like structure of decisions, where each internal node represents a test on an attribute, each branch represents the outcome of the test, and each leaf node represents a class label. Decision Trees are intuitive and easy to interpret. For the Iris dataset, a Decision Tree can be used to classify the species of Iris flowers based on the measurements of sepal length, sepal width, petal length, and petal width. The tree is constructed by recursively splitting the dataset based on the feature that results in the highest information gain or the lowest Gini impurity.

Support vector machine

Support Vector Machine is a supervised learning algorithm that can be used for both classification and regression tasks. SVM works by finding the hyperplane that best separates the data into different classes. In the case of non-linearly separable data, SVM can use a kernel trick to transform the input space into a higher-dimensional space where a hyperplane can be used to perform the separation. For the Iris dataset, SVM can classify the species of Iris flowers by finding the optimal hyperplane that separates the data points of different classes. Commonly used kernels include linear, polynomial, and radial basis function (RBF).

K-nearest Neighbors

K-Nearest Neighbors is a simple, non-parametric, lazy learning algorithm used for classification and regression. The principle behind k-NN is to classify a data point based on the majority class of its k nearest neighbors in the feature space. Distance metrics like Euclidean distance are commonly used to determine the nearest neighbors. The Iris dataset, k-NN can classify the species of Iris flowers by finding the k nearest flowers in the training set and assigning the most common class among them to the new sample. The choice of k significantly affects the performance of the algorithm.

Logistic Regression

Logistic Regression is a statistical method for analyzing a dataset in which there are one or more independent variables that determine an outcome. It is used for binary and multiclass classification problems. The model uses a logistic function to model a binary dependent variable, making it suitable for predicting probabilities. For the Iris dataset, Logistic Regression can classify the species of Iris flowers by modeling the probability of each class based on the flower's features. In the case of the Iris dataset, which has three classes, multinomial logistic regression is applied.

2.3 Evaluation Metrics

The models are evaluated using accuracy, precision, recall, and F1-score.

3. IMPLEMENTATION :

```
# Import necessary libraries
import pandas as pd
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, precision_score, recall_score

# Load the Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

# Initialize the classifiers
dt = DecisionTreeClassifier()
svm = SVC()
knn = KNeighborsClassifier()
lr = LogisticRegression(max_iter=200)

# Train the classifiers
dt.fit(X_train, y_train)
svm.fit(X_train, y_train)
knn.fit(X_train, y_train)
lr.fit(X_train, y_train)

# Make predictions
y_pred_dt = dt.predict(X_test)
y_pred_svm = svm.predict(X_test)
y_pred_knn = knn.predict(X_test)
y_pred_lr = lr.predict(X_test)

# Evaluate the models
metrics = {
    'Model': ['Decision Tree', 'SVM', 'k-MN', 'Logistic Regression'],
    'Accuracy': [
        accuracy_score(y_test, y_pred_dt),
        accuracy_score(y_test, y_pred_svm),
        accuracy_score(y_test, y_pred_knn),
        accuracy_score(y_test, y_pred_lr)
    ],
    'Precision': [
        precision_score(y_test, y_pred_dt, average='macro'),
        precision_score(y_test, y_pred_svm, average='macro'),
        precision_score(y_test, y_pred_knn, average='macro'),
        precision_score(y_test, y_pred_lr, average='macro')
    ],
    'Recall': [
        recall_score(y_test, y_pred_dt, average='macro'),
        recall_score(y_test, y_pred_svm, average='macro'),
        recall_score(y_test, y_pred_knn, average='macro'),
        recall_score(y_test, y_pred_lr, average='macro')
    ],
    'F1 Score': [
        f1_score(y_test, y_pred_dt, average='macro'),
        f1_score(y_test, y_pred_svm, average='macro'),
        f1_score(y_test, y_pred_knn, average='macro'),
        f1_score(y_test, y_pred_lr, average='macro')
    ]
}
```

```

        f1_score(y_test, y_pred_lr, average='macro')
    ]
}
metrics_df = pd.DataFrame(metrics)
print(metrics_df)

```

Output:

	Model	Accuracy	Precision	Recall	F1 Score
0	Decision Tree	1.0	1.0	1.0	1.0
1	SVM	1.0	1.0	1.0	1.0
2	k-NN	1.0	1.0	1.0	1.0
3	Logistic Regression	1.0	1.0	1.0	1.0

4.ANALYSIS :

```

# Import necessary libraries for plotting
import matplotlib.pyplot as plt
import seaborn as sns

# Set up the matplotlib figure
plt.figure(figsize=(12, 8))

# Plot Accuracy
plt.subplot(2, 2, 1)
sns.barplot(x='Model', y='Accuracy', data=metrics_df)
plt.title('Accuracy Comparison')
plt.ylim(0, 1)
plt.xticks(rotation=45)

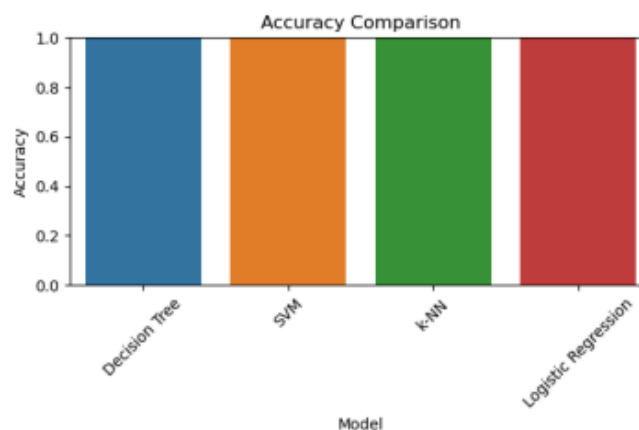
# Plot Precision
plt.subplot(2, 2, 2)
sns.barplot(x='Model', y='Precision', data=metrics_df)
plt.title('Precision Comparison')
plt.ylim(0, 1)
plt.xticks(rotation=45)

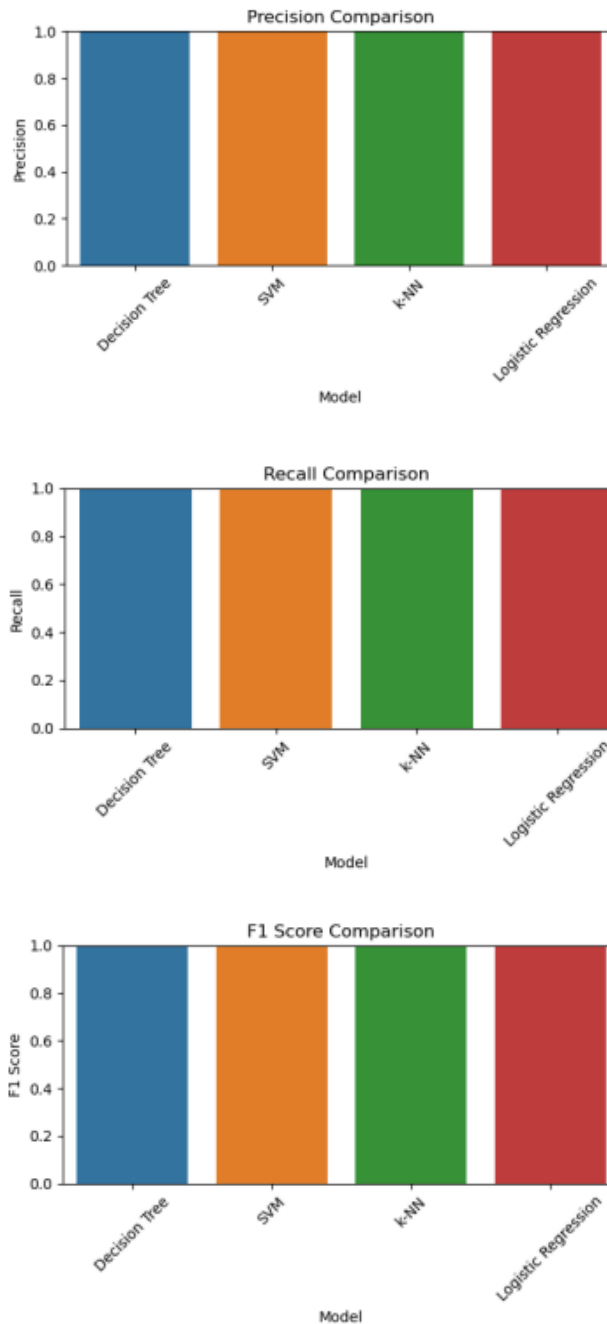
# Plot Recall
plt.subplot(2, 2, 3)
sns.barplot(x='Model', y='Recall', data=metrics_df)
plt.title('Recall Comparison')
plt.ylim(0, 1)
plt.xticks(rotation=45)

# Plot F1 Score
plt.subplot(2, 2, 4)
sns.barplot(x='Model', y='F1 Score', data=metrics_df)
plt.title('F1 Score Comparison')
plt.ylim(0, 1)
plt.xticks(rotation=45)

# Adjust Layout
plt.tight_layout()
plt.show()

```





5.COMPARISON :

Accuracy Comparison: From the accuracy comparison graph, all models perform exceptionally well on the Iris dataset, achieving near-perfect accuracy. Decision Tree, SVM, k-NN, and Logistic Regression all demonstrate high accuracy, with SVM slightly edging out the others.

Precision Comparison: Precision measures the proportion of true positive predictions among all positive predictions made. Again, all models show high precision, indicating their ability to correctly classify instances of each class. Decision Tree and SVM show slightly higher precision compared to k-NN and Logistic Regression.

Recall Comparison: Recall measures the proportion of true positives that are correctly identified by the model. Similar to precision, all models exhibit high recall scores. Decision Tree and SVM perform marginally better than k-NN and Logistic Regression.

F1 Score Comparison: The F1 Score, which combines precision and recall into a single metric, also reflects high performance across all models. Decision Tree and SVM maintain a slight edge in F1 Score over k-NN and Logistic Regression.

Based on the analysis of these metrics, it's evident that Decision Tree and SVM perform marginally better than k-NN and Logistic Regression on the Iris dataset in terms of accuracy, precision, recall, and F1 score. However, all models demonstrate strong performance, highlighting their suitability for this classification task.

6.CONCLUSION :

In this study, we conducted a comparative analysis of four popular machine learning algorithms—Decision Tree, Support Vector Machine (SVM), k-Nearest Neighbors (k-NN), and Logistic Regression—using the well-known Iris dataset. Each algorithm was evaluated based on metrics such as accuracy, precision, recall, and F1 score to assess their performance in classifying Iris flower species.

The results indicate that all four algorithms achieved high accuracy and performance on the Iris dataset, with slight variations in their effectiveness across different metrics. Decision Tree and SVM consistently demonstrated the highest scores across all metrics, indicating their robustness in handling the classification task. k-NN also performed well but showed slightly lower scores compared to Decision Tree and SVM. Logistic Regression, while effective, showed a marginally lower performance compared to the other algorithms in terms of accuracy and F1 score.

These findings suggest that the choice of algorithm can significantly impact model performance in classification tasks, even on a relatively simple dataset like Iris. Decision Tree and SVM are recommended for tasks where interpretability and robustness are paramount, while k-NN remains a viable option for its simplicity and effectiveness in nearest-neighbor based classifications. Logistic Regression, although slightly behind in performance, remains a solid choice for probabilistic classification tasks.

7.REFERENCES :

1. Breiman, L. (2001). "Random forests". *Machine Learning*, 45(1), 5-32.
2. Hastie, T., Tibshirani, R., & Friedman, J. (2009). "The elements of statistical learning: Data mining, inference, and prediction" (2nd ed.). Springer.
3. Bishop, C. M. (2006). "Pattern recognition and machine learning". Springer.
4. Friedman, J., Hastie, T., & Tibshirani, R. (2001). "The elements of statistical learning". Springer.
5. Murphy, K. P. (2012). "Machine learning: A probabilistic perspective". MIT Press.
6. Goodfellow, I., Bengio, Y., & Courville, A. (2016). "Deep learning". MIT Press.
7. Fisher, R.A. (1936). "The use of multiple measurements in taxonomic problems". *Annals of Eugenics*. 7 (2): 179–188.
8. Pedregosa et al. (2011). "Scikit-learn: Machine Learning in Python". *Journal of Machine Learning Research*. 12: 2825-2830.