

Dataset Description

The dataset used in this analysis appears to be related to text classification. Although the specific details of the dataset are not provided, we can infer that it likely involves a binary or multi-class classification problem based on the performance metrics and the classifiers employed. The text data has been preprocessed using two main techniques: Bag of Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF), which are commonly used to convert text into numerical features for machine learning models.

Objective

The primary objective of this analysis is to evaluate and compare the performance of various machine learning classifiers on a text classification task. The models are trained and tested using two different feature extraction methods (BoW and TF-IDF). The classifiers compared include:

1. Support Vector Classifier (SVC)
2. K-Neighbors Classifier (KNN)
3. Logistic Regression
4. Decision Tree Classifier
5. Random Forest Classifier
6. Multinomial Naive Bayes (MultinomialNB)

Methodology

1. **Data Preprocessing:**
 - Text data is converted into numerical features using two methods: Bag of Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF).
2. **Model Training and Evaluation:**
 - Each classifier is trained on the training dataset and evaluated on both the training and test datasets.
 - Performance metrics used for evaluation include Accuracy, Precision, Recall, and F1 Score.

Results

Bag of Words (BoW) Feature Extraction:

1. **Support Vector Classifier (SVC):**
 - **Train:** Accuracy = 97.94%, Precision = 98.77%, Recall = 96.37%, F1 Score = 97.56%
 - **Test:** Accuracy = 94.56%, Precision = 97.07%, Recall = 89.74%, F1 Score = 93.26%
2. **K-Neighbors Classifier (KNN):**
 - **Train:** Accuracy = 92.46%, Precision = 97.05%, Recall = 84.93%, F1 Score = 90.59%
 - **Test:** Accuracy = 90.58%, Precision = 95.85%, Recall = 81.03%, F1 Score = 87.82%
3. **Logistic Regression:**

- **Train:** Accuracy = 97.28%, Precision = 98.29%, Recall = 95.30%, F1 Score = 96.77%
- **Test:** Accuracy = 94.67%, Precision = 96.57%, Recall = 90.50%, F1 Score = 93.43%
- 4. **Decision Tree Classifier:**
 - **Train:** Accuracy = 99.96%, Precision = 99.96%, Recall = 99.93%, F1 Score = 99.95%
 - **Test:** Accuracy = 93.73%, Precision = 93.26%, Recall = 91.65%, F1 Score = 92.45%
- 5. **Random Forest Classifier:**
 - **Train:** Accuracy = 99.96%, Precision = 99.95%, Recall = 99.95%, F1 Score = 99.95%
 - **Test:** Accuracy = 94.13%, Precision = 94.40%, Recall = 91.42%, F1 Score = 92.89%
- 6. **Multinomial Naive Bayes (MultinomialNB):**
 - **Train:** Accuracy = 93.70%, Precision = 90.59%, Recall = 95.13%, F1 Score = 92.80%
 - **Test:** Accuracy = 91.99%, Precision = 88.48%, Recall = 92.99%, F1 Score = 90.68%

TF-IDF Feature Extraction:

1. **Support Vector Classifier (SVC):**
 - **Train:** Accuracy = 98.71%, Precision = 99.33%, Recall = 97.64%, F1 Score = 98.48%
 - **Test:** Accuracy = 94.51%, Precision = 96.44%, Recall = 90.43%, F1 Score = 93.34%
2. **K-Neighbors Classifier (KNN):**
 - **Train:** Accuracy = 68.82%, Precision = 95.34%, Recall = 28.14%, F1 Score = 43.45%
 - **Test:** Accuracy = 65.57%, Precision = 92.35%, Recall = 20.76%, F1 Score = 33.90%
3. **Logistic Regression:**
 - **Train:** Accuracy = 95.21%, Precision = 97.56%, Recall = 91.02%, F1 Score = 94.17%
 - **Test:** Accuracy = 93.52%, Precision = 95.55%, Recall = 88.92%, F1 Score = 92.11%
4. **Decision Tree Classifier:**
 - **Train:** Accuracy = 99.96%, Precision = 99.96%, Recall = 99.94%, F1 Score = 99.95%
 - **Test:** Accuracy = 92.90%, Precision = 91.58%, Recall = 91.73%, F1 Score = 91.66%
5. **Random Forest Classifier:**
 - **Train:** Accuracy = 99.96%, Precision = 99.95%, Recall = 99.95%, F1 Score = 99.95%
 - **Test:** Accuracy = 93.80%, Precision = 93.48%, Recall = 91.82%, F1 Score = 92.64%
6. **Multinomial Naive Bayes (MultinomialNB):**
 - **Train:** Accuracy = 94.17%, Precision = 92.59%, Recall = 93.82%, F1 Score = 93.20%

- **Test:** Accuracy = 91.65%, Precision = 89.87%, Recall = 90.57%, F1 Score = 90.22%

Conclusion

- **Support Vector Classifier (SVC)** generally performs well on both BoW and TF-IDF, showing high accuracy and balanced precision and recall.
- **K-Neighbors Classifier (KNN)** exhibits poor performance with TF-IDF, indicating this method is not suitable for this classifier in this context.
- **Logistic Regression** shows robust performance across both feature extraction methods, making it a reliable choice.
- **Decision Tree and Random Forest Classifiers** show signs of overfitting, as indicated by near-perfect scores on the training set but lower performance on the test set.
- **Multinomial Naive Bayes** performs reasonably well, especially with BoW, but its performance is slightly lower with TF-IDF.

Based on this analysis, **SVC and Logistic Regression** are the most effective classifiers for this text classification task, with consistent high performance across different feature extraction methods. The results highlight the importance of selecting appropriate classifiers and feature extraction techniques for text classification problems.