

Data Preprocessing Steps:

- 1. Loading Data:** The dataset is loaded using pandas' `read_csv` function.
- 2. Data Cleaning:** Rows with label 'O' are removed as they seem to represent non-hate speech comments.
- 3. Encoding Labels:** Labels are encoded using `LabelEncoder` to convert them into numerical format.
- 4. Text Preprocessing:**
 - Tokenization: The comments are tokenized using NLTK's `word_tokenize` function.
 - Cleaning: '@user' mentions are removed from the comments using regular expressions.
 - Stemming: Porter stemming is applied to reduce words to their root form.
- 5. TF-IDF Vectorization:** Text data is transformed into numerical vectors using TF-IDF vectorization.

Model Architecture and Parameters:

- 1. Logistic Regression:**
 - Regularization: L1 and L2 penalties.
 - Optimization Algorithm: 'liblinear' and 'saga'.
- 2. Naive Bayes:**
 - Multinomial Naive Bayes classifier.
- 3. Random Forest:**
 - Hyperparameters:
 - Number of estimators: 50, 100, 150.
 - Maximum depth of trees: 5, 10, 15.
 - Minimum samples split: 2, 5, 10.
 - Minimum samples leaf: 1, 2, 4.
 - Maximum features: 'auto', 'sqrt', 0.5.
- 4. XGBoost:**
 - Objective: Binary logistic.
 - Evaluation Metric: Log loss.

Training Process and Hyperparameters:

Grid Search: Used to find the best hyperparameters for Logistic Regression, Random Forest, and Support Vector Classifier models.

Cross-Validation: 5-fold cross-validation is used for model evaluation during grid search.

Scoring Metric: F1-score weighted is used as the scoring metric.

Tokenization and Cleaning: NLTK's word tokenizer is used for tokenization and comments are cleaned using regular expressions to remove '@user' mentions.

Stemming: Porter stemming is used to reduce words to their root form.

Feature Extraction: TF-IDF vectorization is used to convert text data into numerical vectors.

Evaluation Results and Analysis:

1. Logistic Regression:

- Best Parameters: {'C': 1, 'penalty': 'l2', 'solver': 'liblinear'}.
- Validation Scores: Accuracy=0.844, F1-Score=0.839.

2. Naive Bayes:

- Training Scores: Accuracy=0.802, F1-Score=0.794.
- Validation Scores: Accuracy=0.776, F1-Score=0.767.

3. Random Forest:

- Best Parameters: {'max_depth': 15, 'max_features': 0.5, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 150}.
- Training Scores: Accuracy=0.955, F1-Score=0.954.
- Validation Scores: Accuracy=0.854, F1-Score=0.848.

4. XGBoost:

- Training Scores: Accuracy=0.916, F1-Score=0.914.
- Validation Scores: Accuracy=0.852, F1-Score=0.848.

5. Support Vector Classifier:

- Best Parameters: {'C': 1, 'gamma': 'scale', 'kernel': 'rbf'}.
- Training Scores: Accuracy=0.965, F1-Score=0.965.
- Validation Scores: Accuracy=0.876, F1-Score=0.871.

6. Other Ensemble Models: Various ensemble methods like Gradient Boosting, AdaBoost, Extra Trees, and Bagging Classifiers are also trained and evaluated.

7. Voting Classifier:

- Ensemble of Random Forest, Extra Trees, and Bagging Classifiers.
- Validation Scores: Accuracy=0.863, F1-Score=0.859.

Overall, the models perform well, with Random Forest, Support Vector Classifier, and Voting Classifier achieving the highest validation scores. Random Forest and Support Vector Classifier exhibit strong performance with F1-Scores around 0.85 and 0.87, respectively.