

# **Progress Report on Developing a Machine Learning Model for Turkish Constitutional Court Decisions**

## **1. Summary of Modeling Work Completed**

The goal of this project was to develop a robust classification model for predicting the outcomes of Turkish Constitutional Court decisions. To achieve this, we followed a comprehensive machine learning pipeline, which included extensive preprocessing, feature engineering, and experimentation with five machine learning algorithms: Decision Trees, Logistic Regression, Random Forest, Naive Bayes, and a Deep Learning-based Neural Network approach. Below is a detailed summary of the steps we completed:

### **Key Steps Completed:**

#### **1. Data Collection and Preprocessing:**

- Acquired a dataset consisting of textual descriptions and corresponding outcomes of Turkish Constitutional Court decisions.
- Performed rigorous data cleaning to remove inconsistencies, including duplicate entries and missing values, ensuring a high-quality dataset.
- Preprocessed textual data by implementing tokenization, stop-word removal, and stemming to standardize text content.
- Converted textual data into numerical representations using both TF-IDF vectorization for sparse feature representation and Word Embedding techniques (e.g., Word2Vec) to capture semantic meanings.

#### **2. Feature Engineering:**

- Extracted key features from text, such as word frequency counts and semantic embeddings.
- Incorporated additional metadata features, such as the type of case, decision length, and temporal information, to enrich the dataset.

#### **3. Algorithm Implementation:**

- Systematically implemented and tested the following five machine learning algorithms to determine their effectiveness for classification:
  - **Decision Trees:** Provided interpretable models with hierarchical splitting rules.
  - **Logistic Regression:** Offered a probabilistic linear approach for binary classification.

- **Random Forest:** Leveraged ensemble learning to enhance prediction robustness.
- **Naive Bayes:** Utilized probabilistic assumptions to handle text-based features efficiently.
- **Neural Networks:** Designed and trained deep learning architectures capable of capturing complex patterns in textual data.

#### 4. **Model Evaluation:**

- Employed a diverse set of metrics, including accuracy, precision, recall, F1-score, and AUC-ROC, to assess model performance comprehensively.
- Conducted cross-validation to ensure the reliability and generalizability of the trained models across different subsets of the data.
- Visualized evaluation results through confusion matrices and ROC curves to identify strengths and areas for improvement.

## 2. Challenges Encountered and Solutions

### Challenge 1: Imbalanced Dataset

- **Problem:** The dataset exhibited a significant imbalance in the distribution of decision outcomes, which led to biased model predictions favoring the majority class.
- **Solution:** To address this, we employed both oversampling and undersampling techniques. The Synthetic Minority Oversampling Technique (SMOTE) was used to generate synthetic samples for the minority class, while undersampling reduced instances of the majority class. This combination ensured a balanced dataset and improved model fairness.

### Challenge 2: High Dimensionality of Text Data

- **Problem:** The textual nature of the data resulted in high-dimensional feature spaces, especially when using TF-IDF or embedding techniques. This increased computational complexity and risked overfitting.
- **Solution:** Dimensionality reduction techniques were applied to mitigate these challenges. Principal Component Analysis (PCA) was utilized to reduce TF-IDF feature space dimensions, retaining the most informative components. For embeddings, we limited vector size and applied truncation techniques, streamlining model input without significant loss of information.

### Challenge 3: Hyperparameter Optimization

- **Problem:** Determining the optimal hyperparameters for each algorithm was computationally intensive and time-consuming, given the wide parameter search space.
- **Solution:** We implemented systematic optimization approaches such as grid search for exhaustive exploration of parameter combinations and random search for quicker, stochastic sampling. Additionally, Bayesian optimization was tested for certain algorithms, further accelerating convergence to optimal hyperparameters.

### Challenge 4: Overfitting in Neural Networks

- **Problem:** During training, the Neural Network model showed signs of overfitting, where the model performed well on the training data but poorly on the validation set.
- **Solution:** Several techniques were employed to prevent overfitting. Dropout regularization was applied to randomly deactivate neurons during training, reducing reliance on specific neurons. Early stopping monitored validation loss during training and halted training when performance plateaued, preventing overtraining. Additionally, we introduced L2 regularization to penalize large weights, ensuring a more generalized model.

## 3. Algorithm Details

### 3.1 Decision Trees

- **Concept:** Decision Trees are hierarchical models that recursively split data into subsets based on feature values. Each node represents a decision rule, while leaves denote outcome predictions. The tree's depth and splitting criteria define its complexity.
- **Advantages:**
  - Highly interpretable and visually intuitive.
  - Handles both numerical and categorical data seamlessly.
  - Requires minimal data preprocessing.
- **Implementation:**
  - Used Gini impurity and entropy as splitting criteria to determine the best splits.
  - Pruned the tree to avoid overfitting by setting a maximum depth and minimum samples per leaf.
- **Performance:**
  - Demonstrated moderate accuracy, particularly in binary classifications.
  - Prone to overfitting when the tree's depth was excessive, which was mitigated through regularization techniques.

### 3.2 Logistic Regression

- **Concept:** Logistic Regression is a linear model used for binary classification. It predicts probabilities using a sigmoid function, outputting values between 0 and 1, which are interpreted as class probabilities.
- **Advantages:**
  - Computationally efficient and robust to multicollinearity.
  - Performs well when the relationship between features and the target is linear.
  - Easy to implement and interpret.
- **Implementation:**
  - Applied L2 regularization to prevent overfitting and improve generalization.
  - Optimized hyperparameters such as the regularization strength (C) using grid search.
- **Performance:**

- Performed well on simpler, linearly separable cases.
- Achieved competitive precision and recall for binary classifications but struggled with non-linear relationships.

### 3.3 Random Forest

- **Concept:** Random Forest is an ensemble learning technique that constructs multiple Decision Trees during training and combines their predictions through majority voting for classification tasks.
- **Advantages:**
  - Reduces the risk of overfitting inherent in individual Decision Trees.
  - Capable of handling high-dimensional and complex datasets effectively.
  - Robust against noise and missing data.
- **Implementation:**
  - Experimented with different numbers of trees (`n_estimators`) and maximum tree depths.
  - Used grid search to optimize parameters such as `max_depth` and `min_samples_split`.
- **Performance:**
  - Delivered high accuracy and robust predictions.
  - Balanced precision and recall, excelling in datasets with complex feature interactions.

### 3.4 Naive Bayes

- **Concept:** Naive Bayes is a probabilistic classifier based on Bayes' theorem. It assumes that features are independent, a simplification that works well in certain domains, especially text classification.
- **Advantages:**
  - Fast to train and predict, making it suitable for large datasets.
  - Performs well with sparse data, such as text represented by TF-IDF or bag-of-words.
  - Naturally handles imbalanced datasets.
- **Implementation:**

- Used Multinomial Naive Bayes, ideal for text classification tasks involving word counts or frequencies.
- Applied Laplace smoothing to handle zero probabilities.
- **Performance:**
  - Provided strong baseline results with minimal computational overhead.
  - Worked particularly well for text data but struggled with datasets where feature independence assumptions were violated.

### 3.5 Neural Networks

- **Concept:** Neural Networks are deep learning models that consist of layers of interconnected neurons. They excel in capturing non-linear and complex patterns in data through hierarchical feature extraction.
- **Advantages:**
  - Highly flexible and capable of learning intricate relationships.
  - Performs well on large and unstructured datasets, such as text and images.
  - Can adapt to various problem types with appropriate architecture tuning.
- **Implementation:**
  - **Architecture:** Designed a model with an input layer (text embeddings), multiple hidden layers with ReLU activations, and a softmax output layer for multiclass classification.
  - **Optimization:** Used the Adam optimizer for efficient gradient-based updates.
  - **Regularization:** Included dropout layers to reduce overfitting and applied early stopping based on validation loss.
  - **Loss Function:** Employed Categorical Crossentropy to compute the classification loss.
- **Performance:**
  - Achieved the best accuracy among all algorithms tested.
  - Required significant computational resources and careful tuning of hyperparameters such as learning rate, batch size, and number of epochs.

## 4. Code Documentation and Workflow Overview

The codebase for this project was structured into modular components, each responsible for a distinct aspect of the machine learning pipeline. Each module was documented extensively to ensure maintainability and ease of understanding for future contributors.

#### 4.1 Preprocessing Workflow

- **Purpose:** The preprocessing phase prepared the raw data for machine learning by performing necessary cleaning, transformation, and feature extraction steps.
- **Functionality:**
  - Handled data cleaning processes, including removal of duplicates, null entries, and irrelevant features.
  - Tokenized text data into individual terms to facilitate analysis.
  - Applied stemming and lemmatization techniques to reduce words to their root forms, ensuring consistency across textual representations.
  - Extracted features using both TF-IDF vectorization, capturing the importance of terms in documents, and word embeddings like Word2Vec to represent semantic similarities between words.

#### 4.2 Model Training Workflow

- **Purpose:** This module encapsulated the logic for implementing and training various machine learning models.
- **Functionality:**
  - Trained five distinct algorithms: Decision Trees, Logistic Regression, Random Forest, Naive Bayes, and Neural Networks, each configured with optimal parameters.
  - Evaluated each model against a predefined validation dataset using metrics like accuracy, precision, recall, and F1-score.
  - Saved trained models for reuse, enabling seamless integration with evaluation and deployment workflows.

#### 4.3 Neural Network Workflow

- **Purpose:** Focused specifically on deep learning models, this module handled complex data representations and model architectures.
- **Functionality:**
  - Designed a custom architecture with multiple hidden layers, dropout for regularization, and ReLU activations for efficient learning.

- Trained the model on text embeddings to capture nuanced semantic features of the data.
- Monitored validation metrics and implemented early stopping to prevent overfitting, ensuring model robustness.
- Saved the final trained model and its weights for deployment or further experimentation.

#### 4.4 Evaluation and Analysis Workflow

- **Purpose:** This phase was dedicated to assessing model performance and deriving insights from the results.
- **Functionality:**
  - Computed comprehensive performance metrics, including confusion matrices, AUC-ROC curves, and classification reports, providing a detailed view of each model's strengths and weaknesses.
  - Visualized results through plots and charts, making it easier to communicate findings.
  - Compared models against benchmarks to identify the best-performing approach for the given task.

#### Example Workflow Execution

##### Logistic Regression Implementation

```
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, roc_auc_score

# Splitting data
features_train, features_test, labels_train, labels_test = train_test_split(features, labels, test_size=0.2,
random_state=42)

# Initializing and training the model
logistic_model = LogisticRegression(penalty='l2', solver='liblinear', random_state=42)
logistic_model.fit(features_train, labels_train)
```



```
# Making predictions
predicted_labels = logistic_model.predict(features_test)

# Evaluating the model
print(classification_report(labels_test, predicted_labels))
print("AUC-ROC:", roc_auc_score(labels_test, logistic_model.predict_proba(features_test)[:, 1]))
```

## Code Documentation Highlights

- **Preprocessing Functions:** Clearly outlined with detailed docstrings explaining input parameters, processing steps, and outputs.
- **Model Training Logic:** Annotated with comments clarifying hyperparameter choices, algorithmic trade-offs, and performance expectations.
- **Evaluation Metrics:** Included inline explanations for interpreting output metrics and assessing their relevance to the project's goals.

## 5. Future Work

- **Hyperparameter Fine-Tuning:** Further optimize model parameters for enhanced performance.
- **Explainability:** Use SHAP or LIME to interpret model decisions.
- **Deployment:** Develop an API to make the model accessible for real-time predictions.
- **Scalability:** Experiment with distributed frameworks like TensorFlow Serving or PyTorch's TorchServe for large-scale deployment.