# Sentiment Analysis
# (Random Forest Classifier)

## 1. Classification and Workflow

The primary goal of a **classification model** is to accurately predict a dependent variable (DV) using one or more independent variables (IVs).

- **Dependent Variable (DV):** The variable being predicted or classified. Commonly referred to as Y, target, output, or response variable. In classification, the DV typically represents distinct categories or classes.
- **Independent Variables (IVs):** The predictors or features used to forecast the DV. Also known as X variables, features, predictors, explanatory variables, or covariates. These variables help explain patterns or correlations with the DV and serve as the model's input.

A standard classification workflow consists of several key steps:

| | |
|---|---|
| **Project Scoping** | <ul><li>Define the business objective, success criteria, and stakeholders.</li><li>Understand the broader context and the desired outcome of the model.</li></ul> |
| **Preliminary Quality Assurance** | <ul><li>Align with stakeholders on scope and goals.</li><li>Review data access, governance policies, and modeling framework.</li></ul> |
| **Data Profiling** | <ul><li>Understand the structure and quality of your data.</li><li>Assess variable types, null values, distributions, and key statistics.</li></ul> |
| **Feature Engineering** | <ul><li>**One-Hot Encoding:** Convert categorical variables into binary indicator variables.</li><li>**Scaling:** Normalize numeric features to a common range (e.g., 0–1 or z-scores).</li><li>**Log Transformation:** Reduce skewness in data for more stable modeling.</li><li>**Discretization:** Bucket continuous variables into discrete bins.</li><li>**Date Extraction:** Decompose date/time fields into parts (e.g., day, month, weekday).</li><li>**Boolean Logic:** Create interaction features using logical conditions.</li></ul> |
| **Data Splitting** | <ul><li>Divide the dataset into training and testing subsets.</li><li>**Training Set:** Used to fit and optimize the model.</li><li>**Test Set:** Used to evaluate model performance on unseen data.</li><li>This step is essential to prevent overfitting, where the model performs well on training data but poorly on new, unseen data.</li></ul> |
| **Model Training** | <ul><li>Train one or more classification models using the training dataset.</li><li>Apply the trained models to the test data to assess their generalization performance.</li></ul> |
| **Model Selection & Hyperparameter Tuning** | <ul><li>Choose the best-performing model based on evaluation metrics (e.g., accuracy, ROC-AUC, F1-score). Fine-tune model hyperparameters to optimize performance and reduce bias or variance.</li><li>Plan for monitoring and re-tuning the model over time to handle concept drift (changes in data patterns).</li></ul> |

A **Confusion Matrix** is a tabular representation of actual vs. predicted classifications on a test dataset. It is one of the most widely used tools for evaluating the performance of classification models.

- **Accuracy:** Measures the overall correctness of the model. When the correct predictions of both classes are equally important
- **Precision:** Indicates the proportion of predicted positives that were correct. When false positives are more harmful than false negatives
- **Recall:** Measures how well the model identifies all actual positives. When false negatives are more harmful than false positives

## 2. Decision Trees and Random Forests

A **Decision Tree** is a supervised machine learning algorithm that uses a sequence of binary (true/false) rules to classify data. It is suitable for both binary and multiclass classification problems.

- Each **internal node** in the tree represents a decision rule based on one independent variable (IV).
- Each **split** divides the dataset into two groups, ideally separating observations to favor one class in each subset.
- The **tree** continues to split the data recursively, aiming to create purer subsets concerning the target variable.

Key considerations in building decision trees:

| | |
|---|---|
| **Entropy** | • Measures the impurity or uncertainty in a dataset.<br>• A perfectly homogeneous set has **entropy = 0**, while maximum uncertainty has higher entropy. |
| **Information Gain** | • The reduction in entropy after a dataset is split on an attribute. It helps determine which variable provides the most effective split.<br>• The process of selecting splits is automated. The algorithm evaluates multiple variables and thresholds to maximize information gain. |
| **Splitting** | • Variables that better reduce entropy (uncertainty) are preferred for splitting<br>• For continuous variables, the algorithm tests various thresholds and chooses the split point that maximizes information gain. |
| **Stopping** | • Excessive splitting can lead to overfitting, where the model fits the training data too well but performs poorly on unseen data.<br>• This is controlled using hyperparameters, such as: **maximum tree depth, minimum samples per leaf, and minimum samples required to split a node** |
| **Root Node** | • First split may not be the best overall. A locally optimal split may not lead to the most accurate model globally.<br>• To address this, ensemble methods like **Random Forests** are often used. |

A **Random Forest** is an ensemble learning method that combines multiple decision trees to produce a more robust and accurate prediction.

- Each tree is trained on a random subset of the training data using bootstrapping. At each node, a random subset of features is considered for splitting. The model combines the predictions from all trees.
- **Advantages:** Reduces overfitting (by averaging multiple models, variance is minimized without significantly increasing bias), improves accuracy (ensemble methods generally outperform individual models, especially on complex or noisy datasets), handles missing values, and maintains high performance even when some data is uninformative or irrelevant.
- For **classification**, the majority vote is used. For **regression**, the average of all tree predictions is taken.
- Hyperparameter tuning involves systematically adjusting these settings to enhance model performance. Effective tuning can improve prediction accuracy and help prevent overfitting or underfitting. Common optimization techniques include **grid search** (which exhaustively tests all combinations of specified hyperparameter values), **random search** (which tests a random subset of combinations for quicker results), and **Bayesian optimization** (which utilizes probabilistic models to intelligently select hyperparameters).

Breakdown of hyperparameters:

| | |
|---|---|
| **max_depth** | • Tree complexity control<br>• Limits tree growth to reduce overfitting |
| **max_features** | • Feature selection per split<br>• Introduces diversity among trees |
| **min_samples_leaf** | • Minimum samples at leaf nodes<br>• Larger sample size increases generalization |
| **n_estimators** | • Number of trees<br>• Fewer trees are faster but less stable |

## 3. Sentiment Analysis

**Sentiment Analysis is a Natural Language Processing (NLP) technique** used to identify and extract the emotional tone conveyed in text-based data. It is most commonly implemented as a classification task that categorizes text into sentiment classes such as positive, negative, or neutral.

- While rooted in NLP, sentiment analysis typically involves **transforming unstructured text into structured numerical inputs**, enabling machine learning models to interpret and analyze the data effectively.
- Sentiment models often employ a Bag of Words (BoW) approach, where text is converted into numerical representations based on the frequency or presence (1/0) of specific keywords.

Different aspects of sentiment analysis:

| | |
|---|---|
| **Text Preprocessing and Quality Assurance** | <ul><li>Removing punctuation, capitalization, and special characters</li><li>Correcting spelling and grammatical errors</li><li>Ensuring proper text encoding (e.g., UTF-8)</li><li>Applying lemmatization or stemming to reduce words to their root form</li><li>Removing stop words (common words such as "a", "the", "is", "or", "of", etc.) that do not carry sentiment</li></ul> |
| **Feature Engineering: Bag of Words** | <ul><li>**Tokenization:** Split text into individual words or tokens</li><li>**Vectorization:** Create features (columns) representing each word, using binary flags or term frequencies</li></ul> |
| **Feature Engineering: TF-IDF Vectorization** | <ul><li>**Tokenization:** Split cleaned text into individual words (tokens)</li><li>**Term Frequency (TF):** Measures how frequently a word appears in a specific document</li><li>**Inverse Document Frequency (IDF):** Scale down the importance of commonly used words that appear across many documents</li><li>**TF-IDF Score:** Multiply TF and IDF to highlight terms that are both frequent in a document and rare across the corpus</li><li>**Vector Representation:** Convert each document into a numerical vector based on the TF-IDF scores of its terms</li></ul> |
| **Language Nuances** | <ul><li>Use larger and more diverse training datasets with accurately labeled sentiment</li><li>Enhance preprocessing and feature engineering with techniques such as n-grams, part-of-speech tagging, or dependency parsing</li><li>Consider advanced models like **transformers (e.g., BERT),** which understand context and semantic relationships more deeply than traditional vectorization techniques</li></ul> |

A standard sentiment analysis workflow consists of several key steps:

| | |
|---|---|
| **Data** | <ul><li>**Location:** https://github.com/f2005636/sentiment-analysis</li><li>**File:** input.csv (0-negative review, 1-positive review)</li></ul> |
| **Clean the Text** | <ul><li>Convert to lowercase</li><li>Remove numbers and punctuation</li><li>Remove extra spaces</li><li>Remove common stop words (like "the", "and", "is")</li><li>Apply stemming (e.g., "running" becomes "run")</li></ul> |
| **Convert Text to Numbers** | <ul><li>The cleaned text is converted into numerical features using **TF-IDF**</li><li>It is a technique that highlights important words</li></ul> |
| **Model Selection** | <ul><li>Grid Search is used to find the best hyperparameters for a **Random Forest classifier** based on the AUC-ROC score.</li><li>The model is trained using the selected hyperparameters.</li><li>The model performance is evaluated using a **confusion matrix** and **AUC-ROC**</li></ul> |