

# Text Classification - XGBoost

Promothesh Chatterjee\*

---

\*Copyright© by Promothesh Chatterjee. All rights reserved. No part of this note may be reproduced, distributed, or transmitted in any form or by any means, including photocopying, recording, or other electronic or mechanical methods, without the prior written permission of the author.

**Introduction to XGBoost** XGBoost (**Extreme Gradient Boosting**) is an advanced implementation of gradient boosting, designed for efficiency, flexibility, and high performance. It is widely used in machine learning competitions and real-world applications due to its speed and predictive accuracy. XGBoost enhances traditional gradient boosting algorithms with optimizations and additional features, making it suitable for handling large datasets and complex tasks. XGBoost primarily uses decision trees as its base models. A **decision tree** is a machine learning model that makes decisions based on a series of binary (yes/no) questions, where each internal node represents a feature, each branch represents a decision rule, and each leaf node represents an outcome or prediction. It recursively splits the data into subsets based on the most significant feature at each step, aiming to create a model that accurately predicts the target variable.

## Key Concepts and Features

### 1. Gradient Boosting:

- **Concept:** Gradient Boosting builds an ensemble (a group) of weak models, typically decision trees, in a sequential manner. Each new model is trained to correct the errors made by the previous models. Specifically, it employs an ensemble of decision trees that are sequentially added to the model. Each tree in the ensemble is trained to correct the errors made by the previous trees, following the principles of gradient boosting. This iterative process helps improve the overall predictive performance of the model.

### 2. Mathematical Foundations:

- **Gradient Descent:** Gradient Descent is an optimization algorithm used to minimize a loss function by iteratively moving towards the minimum value of the function. In the context of XGBoost, it is used to minimize the error of the model.
- **Gradient Boosting:** Gradient Boosting combines boosting with gradient descent. Each new model is trained to minimize the residual errors (differences between actual values and predictions) of the ensemble's current predictions.

**Gradient Boosting and Gradient Descent** Gradient Boosting enhances model performance by combining gradient descent and boosting. Here's a detailed explanation:

## Gradient Descent: A Detailed and Intuitive Explanation

**Basic Concept** Gradient descent is an optimization algorithm used to minimize a function by iteratively moving towards the lowest point, or minimum, of that function. It's widely used in machine learning to find the optimal parameters of a model that minimize the loss function (a measure of error).

## Key Components

1. **Function to Minimize:** Imagine you have a surface, like a hilly landscape, representing the loss function. The goal is to find the lowest point in this landscape, which corresponds to the minimum loss.

2. **Parameters:** These are the variables of the function you're trying to optimize. In the landscape analogy, these would be your coordinates (latitude and longitude) on the surface.
3. **Learning Rate:** This is a step size, representing how far you move in each iteration. It's crucial to choose an appropriate learning rate to ensure efficient and accurate convergence.

## How Gradient Descent Works

1. **Start with Initial Parameters:** Begin at a random point on the surface (initial values for your parameters).
2. **Compute the Gradient:** Calculate the slope of the surface at that point. The gradient is a vector that points in the direction of the steepest ascent (uphill).
3. **Move in the Opposite Direction:** Since you want to minimize the function, move in the opposite direction of the gradient (downhill).
4. **Update Parameters:** Adjust the parameters slightly in the direction that reduces the loss. The amount by which you adjust the parameters is determined by the learning rate.
5. **Iterate:** Repeat the process until you reach a point where the slope is zero or near zero, indicating you've found a minimum.

**Analogy: Hiking Down a Hill** Imagine you're hiking down a hill to reach the lowest point in a valley:

1. **Starting Point:** You start at a random position on the hill.
2. **Look at the Slope:** You observe the steepness and direction of the slope at your current position.
3. **Decide on Step Size:** You decide how big a step to take. If you take too large a step, you might overshoot the lowest point or end up on a different slope. If your steps are too small, it will take a long time to get there.
4. **Take a Step Downhill:** You move a little bit downhill in the direction that decreases your elevation the most.
5. **Repeat:** You continue this process, adjusting your direction and step size, until you find yourself at the bottom of the valley where the slope is flat.

## Types of Gradient Descent

1. **Batch Gradient Descent:** Uses the entire dataset to compute the gradient. It's like having a complete map of the terrain and planning your entire route in one go. This is accurate but can be slow for large datasets.
2. **Stochastic Gradient Descent (SGD):** Uses one data point at a time to compute the gradient. This is like taking steps based on only the terrain immediately around you. It's faster but can be noisy.
3. **Mini-batch Gradient Descent:** A compromise between the two, using small random subsets (mini-batches) of the dataset to compute the gradient. It balances speed and accuracy.

## 2. Gradient Boosting:

- **Residuals:** Each new model in gradient boosting is trained to predict the residual errors of the current model ensemble. Residual errors are the differences between the actual values and the predicted values.
- **Additive Model:** Models are added sequentially, with each new model improving the overall prediction by focusing on the residuals.
- **Objective Function:** The loss function (e.g., mean squared error for regression, log loss for classification) is minimized using gradient descent.
- **Regularization:** Techniques like shrinkage (learning rate) and tree pruning are used to prevent overfitting, ensuring that the model generalizes well to new data.

## Technical Details of XGBoost

### 1. Objective Functions:

- XGBoost supports various objective functions, which define the optimization goal. Common objectives include:
  - **Regression:** `reg:squarederror` for mean squared error.
  - **Classification:** `binary:logistic` for binary classification using logistic regression.
  - **Ranking:** `rank:pairwise` for ranking tasks.

### 2. Regularization:

- **L1 Regularization (Lasso):** Adds a penalty equal to the absolute value of the coefficients' magnitude. This helps in feature selection by shrinking some coefficients to zero.
- **L2 Regularization (Ridge):** Adds a penalty equal to the square of the coefficients' magnitude. This distributes the error term across all features, preventing any one feature from dominating.

### 3. Tree Pruning:

- XGBoost uses a max depth parameter to limit the depth of trees and prevent overfitting.
- **Minimum Child Weight:** A parameter that sets the minimum sum of instance weight (hessian) needed in a child, further preventing overfitting.

### 4. Handling Missing Values:

- XGBoost automatically learns the best imputation values for missing data, allowing it to handle datasets with missing values effectively.

#### 5. Column Block Structure:

- Utilizes a compressed column block structure for storing datasets, which helps optimize operations by reducing the memory footprint and speeding up computations.

#### 6. Parallel and Distributed Computing:

- XGBoost supports parallel processing on a single machine and distributed computing on clusters, making it scalable for large datasets.

### Practical Application: XGBoost in R

Below is a practical implementation of XGBoost for a binary classification problem using a synthetic dataset. This example uses text data and several numerical and categorical features to predict customer churn.

### Customer Churn Analysis

#### Step-by-Step Explanation of the Code

##### 1. Load Libraries:

- Install and load the necessary libraries for data manipulation, text processing, and modeling.

```
library(xgboost)
library(Matrix)
library(quanteda)
library(tidyverse)
```

##### 2. Set Seed for Reproducibility:

- Set a random seed to ensure the reproducibility of the results.

```
set.seed(123)
```

##### 3. Generate Synthetic Data:

- Create a list of possible review sentences.
- Generate 500 random reviews by sampling and combining sentences.
- Generate random values for tenure, churn, age, gender, location, usage, payment history, and customer service interactions.
- Combine these into a tibble.

```

# Define a set of possible sentences for the reviews
sentences <- c("I love the service.",
               "The customer service was terrible.",
               "I am unhappy with the pricing.",
               "The quality of the service was great.",
               "I had an awful experience.",
               "The customer support was excellent.",
               "The service was satisfactory.",
               "I am not happy with the service.",
               "The pricing was very reasonable.",
               "I had a bad experience with the customer support.")

# Generate 500 random reviews
reviews <- map_chr(1:500, ~ paste(sample(sentences, 5, replace = TRUE), collapse = " "))

# Generate a random tenure between 1 and 60 months
tenure <- sample(1:60, 500, replace = TRUE)

# Generate a binary churn variable, with a 20% churn rate
churn <- case_when(
  runif(500) < 0.2 ~ 1,
  TRUE ~ 0
)

# Generate demographic information
age <- sample(18:75, 500, replace = TRUE) # Age
gender <- sample(c("Male", "Female"), 500, replace = TRUE) # Gender
location <- sample(c("Urban", "Suburban", "Rural"), 500, replace = TRUE) # Location

# Generate usage data
usage <- runif(500, min = 1, max = 100) # Usage could represent hours of use, data consumed, etc.

# Generate financial data
payment_history <- sample(c("On Time", "Late"), 500, replace = TRUE, prob = c(0.7, 0.3)) # Payment history

# Generate customer service interactions
num_interactions <- rpois(500, lambda = 2) # Number of customer service interactions

# Combine into a tibble (tidyverse version of a data frame)
data <- tibble(reviews, tenure, churn, age, gender, location, usage, payment_history, num_interactions)

```

#### 4. Text Preprocessing:

- Tokenize the reviews and create a document-feature matrix (DFM).

- Convert the DFM to a dense matrix and then to a sparse matrix.

```
# Text preprocessing
corpus <- corpus(data$reviews)
tokens <- tokens(corpus, remove_punct = TRUE)
dfm <- dfm(tokens, tolower = TRUE)
dense_dfm <- as.matrix(dfm) # Convert dfm to dense matrix
sparse_dfm <- Matrix(dense_dfm, sparse = TRUE) # Convert dense matrix to sparse matrix
```

## 5. One-Hot Encoding for Categorical Variables:

- Convert categorical variables (gender, location, payment history) to dummy variables and then to a sparse matrix.

```
# One-hot encoding for categorical variables using model.matrix
dummies <- model.matrix(~ gender + location + payment_history - 1, data = data)
sparse_dummies <- as(dummies, "dgCMatrix")
```

## 6. Combine Data:

- Combine numerical data, dummy variables, and the sparse matrix of the text data into a single sparse matrix.

```
# Combine numerical data, dummies, and sparse matrix
numeric_data <- data %>% select(tenure, age, usage, num_interactions)
sparse_numeric_data <- as(as.matrix(numeric_data), "dgCMatrix")
combined_data <- cbind(sparse_numeric_data, sparse_dummies, sparse_dfm)
```

## 7. Convert to DMatrix for XGBoost:

- Convert the combined data to the DMatrix format required by XGBoost.

```
# Convert to DMatrix for XGBoost
dall <- xgb.DMatrix(data = combined_data, label = data$churn)~
```

## 8. Train the XGBoost Model:

- Define the parameters for the XGBoost model.
- Train the model using the combined data.

```
# Train the XGBoost model
params <- list( objective = "binary:logistic", eval_metric = "logloss" )
xgb_model <- xgb.train(params, dall)
```

## 9. Make Predictions and Evaluate the Model:

- Make predictions on the same dataset.
- Convert the probabilities to binary predictions using a threshold of 0.5.
- Create and print a confusion matrix to evaluate the model's performance.

```

# Make predictions and evaluate the model
preds <- predict(xgb_model, dall)
preds_binary <- ifelse(preds > 0.5, 1, 0)

# Confusion matrix
conf_matrix <- table(preds_binary, data$churn,
                     dnn = c("Predicted", "Actual"))

# Print the confusion matrix
print(conf_matrix)
# Convert to data frame for better presentation
conf_matrix_df <- as.data.frame.matrix(conf_matrix)

# Replace row and column names with appropriate labels
rownames(conf_matrix_df) <- c("Predicted No Churn", "Predicted Churn")
colnames(conf_matrix_df) <- c("Actual No Churn", "Actual Churn")

# Load knitr and display the table with a caption
library(knitr)
kable(conf_matrix_df, caption = "Confusion Matrix")

```

Output

Table 1: Confusion Matrix

	Actual No Churn	Actual Churn
Predicted No Churn	404	0
Predicted Churn	0	96

## Summary

This example demonstrates how to preprocess text and numerical data, combine them into a format suitable for machine learning, and train an XGBoost model for binary classification. The model is evaluated using a confusion matrix, which provides insights into its performance. The process highlights the flexibility and efficiency of XGBoost in handling various types of data and performing complex predictive tasks.