# Logistic Regression Analysis of Telco Customer Churn

**Author:** Yassine Mouadi

**Date:** October 2025

# Table of Contents

# I. Introduction

The dataset used for this analysis contains core information about the Telco Customer Churn, specifically focusing on the customers who left within the last month.

The goal of this project is to use the data in the feature columns—such as `tenure`, `Contract`, and `MonthlyCharges` —to predict the value in the `Churn` column for a new, unseen customer.

# II. The Theory: The Logistic Model

## A. The Core Function

We're going to use the logistic function, also called the sigmoid function. Its main purpose is taking any number, whether positive or negative, and squashing it into a probability between 0 and 1.

It is formally defined as:

```
σ(z) = 1 / (1 + e^(-z))
```

## B. The Regression Foundation

Logistic regression is more or less taking the linear regression's prediction function, `y = w*x + b`, and considering the entirety of it as the linear score `z`. We then input this `z` into the sigmoid function: `σ(z) = 1 / (1 + e^(-z))`. However, for logistic regression, the linear component often takes another form:

```
z = β0 + β1 * X
```

## C. Why the Sigmoid Function?

Why, however, do we need said sigmoid function, instead of directly using the linear model? Well, `Y = β0 + β1 * X` gives a standard output that can fall outside the range of [0,1]. Such a thing is very unlike what a basic probability is. The sigmoid function does exactly what we need, which is perfect for probability.

As a reminder, for `Y = w*X + b`, the weight ( `w` ) tells us how much the feature ( `x` ) influences the outcome `Y`. It's perfect for classification as it uses the prediction model ( `z` ) and wraps it into the sigmoid function to properly constrain its result into a probability domain.

## D. From Probability to Classification

As was previously said, the sigmoid function outputs a probability. We then have to translate this probability into something we can classify the data into. Thus, we set a threshold at 0.5.

If the output is above 0.5, then we classify it as a 'yes'—or, in the case of our dataset, that the customer at hand would indeed churn. Naturally, if it's below 0.5, then it's a 'no'.

## E. The Link Function: Log-Odds

However, there will be cases in which we'd need the range to be an unbounded scale $(-\infty, \infty)$. And thus comes the time to use "odds."

Odds are the ratio of the probability of the event happening to the probability of the event not happening. It is thus defined as:

```
Odds = (probability of the event happening) / (probability of the event
not happening) = P / (1 - P)
```

We can express this using our sigmoid output: `Odds = σ(z) / (1 - σ(z))` . When we isolate the `z` in this equation, we get the Log-Odds function:

```
log(Odds) = z = w₀ + w₁x₁ + w₂x₂ + ...
```

This shows that our linear model `z` is actually predicting the log-odds of the event. This is the true "output" of logistic regression before it is squashed into a probability by the sigmoid function.

However, the log-odds are difficult to interpret, and thus we simply do:

```
Odds = exp(Log-odds)
```

Why? Because we could do: `( odds at (xi+1) ) / ( odds at xi )`

After the relationship simplifies algebraically, we get:

```
( odds at (xi+1) ) / ( odds at xi ) = exp(βi)
```

We call it: `OR = ( odds at (xi+1) ) / ( odds at xi ) = exp(βi)` .

Therefore, exponentiating the logistic regression coefficient `exp(βi)` is the way to calculate the Odds Ratio for a one-unit increase in the predictor `xi` , making all other variables constant.

Thus, by using $\beta$ as a reference, we can conclude many things. Like:

- If $\beta > 0 \rightarrow$ Log-Odds increases $\rightarrow$ OR > 1 $\rightarrow$ the odds increase.
- If $\beta < 0 \rightarrow$ Log-Odds decreases $\rightarrow$ OR < 1 $\rightarrow$ the odds decrease.

# III. The Learning Process: Maximum Likelihood

## A. The High-Level Steps

Just like with linear regression, logistic regression follows the same iterative steps:

1. We start with random weights, where the computer guesses randomly.
2. We make predictions on the training data, with examples we have answers for.
3. Afterwards, we calculate how wrong we are through an "error" function (not the MSE in this case).
4. We then adjust the weights to reduce the error, and repeat until the error is minimized.

## B. The Cost Function: Maximum Likelihood Estimation (MLE)

The equivalent of the MSE in linear regression is the MLE in logistic regression. The MLE, or Maximum Likelihood Estimation, determines how likely a set of parameters—that we give to the sigmoid function—are, given the observed data `x`.

Its function definition is conceptually closer to the Fourier Transform; but instead of helping us see what frequencies are present in the data, it helps us find the most probable parameters (like the weights) of a model that generates that data.

The likelihood function for our parameters `θ` given the data `x` is:

```
L(θ | x) = ∏ from i=1 to n of f(x_i; θ)
```

However, this is considered a 'product', which is difficult to work with. To simplify it, we utilize the log-likelihood function, which turns the product into a sum:

```
ℓ(θ | x) = ∑ from i=1 to n of ln[f(x_i; θ)]
```

## C. The Model and The Likelihood

However, when we do MLE for machine learning, we're not just estimating simple distributions; we're estimating a function that maps inputs `x` to outputs `y`. And here comes the role of `P(y_i | x_i; θ)`. For logistic regression—a clear case of binary classification—it uses the Bernoulli distribution:

```
P(y | x; θ) = { p(x; θ) if y = 1; 1-p(x; θ) if y = 0 }
```

This can be written more compactly as:

```
P(y | x; θ) = [p(x; θ)]^y × [1-p(x; θ)]^(1-y)
```

We thus include this into the full Likelihood function for all our data:

```
L(θ | all data) = P(y₁ | x₁; θ) × P(y₂ | x₂; θ) × ... × P(y_m | x_m; θ)
= ∏ from i=1 to m of P(y_i | x_i; θ)
```

And so comes the final Log-Likelihood function we work with:

```
ℓ(θ) = ∑ from i=1 to m of ln[P(y_i | x_i; θ)]
```

## D. Why Maximize the Log-Likelihood?

We want to maximize the log-likelihood, not minimize it—very unlike the MSE. Why? Because we want to find the model parameters that make our observed data appear the most probable.

The Likelihood Function answers the question: "How probable is my observed data, given specific parameters?"

The main reason why we don't use the MSE here is because it would result in a non-convex surface for logistic regression, meaning it would be very inefficient for the subsequent optimization. It could easily get stuck in local minimums instead of finding the global minimum, making the gradient descent algorithm useless in such a case.

# IV. The Optimization: Gradient Ascent

## A. The Principle of Gradient Descent

As is already known, gradient descent works by updating the model's parameters—the weights and bias—to minimize the Loss function within the loss landscape. This landscape is ideally in the shape of a 'bowl'. Mathematically, it calculates the partial derivative of the loss function `J` with respect to each parameter, in the format:

```
θ := θ - α × ∇J(θ)
```

It moves in the opposite direction to the gradient to minimize the loss.

## B. The Gradient for Our Model

The gradient for our log-likelihood, which we treat as our loss to maximize, is beautifully simple. For a single weight, it is:

```
∇J(θ) = (1/m) × Σ (predicted - actual) × feature
```

Where `(predicted_probability - actual_value) = (ŷᵢ - yᵢ)`.

The learning rate `α` determines how big our steps are within that landscape. It is typically a small value like 0.01 or 0.001. It must be within a good range, or else it would malfunction:

- Too big → we overshoot the minimum.
- Too small → the process takes forever.

## C. From Descent to Ascent

Such is the principle of gradient descent. But for our case, since we want to **maximize** the log-likelihood, we use **gradient ascent**. The change is very minor:

```
θ := θ + α × ∇J(θ)
```

Here, we move in the *same* direction as the gradient, climbing uphill to find the peak of the likelihood function, which represents the most probable model given our data.

# V. Data Preparation for the Telco Churn Dataset

## A. Handling the Dataset

Before we can train our model, we need to prepare the Telco Churn data. This involves several crucial steps:

## B. Data Cleaning

We first check for missing values in columns like `TotalCharges` and handle them appropriately. We also convert the `Churn` column from 'Yes'/'No' to binary 1/0 values that our model can understand.

## C. Feature Engineering

For categorical variables like `Contract` type, we need to convert them into numerical format using one-hot encoding. This transforms categories like 'Month-to-month', 'One year', and 'Two year' into separate binary columns that the model can process.

## D. Feature Scaling

We scale numerical features like `tenure`, `MonthlyCharges`, and `TotalCharges` so they're on similar scales. This helps gradient ascent converge faster and more reliably.

## E. Train-Test Split

We split our data into training and testing sets, typically using 70-80% for training and the remainder for testing. This ensures we can evaluate our model on unseen data.

# VI. Model Evaluation: Beyond Training

We finished gradient ascent, what's next? Now we have a trained model. But we can't just stop here.

We need to answer: "Is this model any good?" and "How do we use it?" And so, to get the answers for these questions, we need to evaluate this model.

For the case of our dataset, let's consider these cases: If the model says that the customer will churn, but in reality the customer stays. Or, in reverse, the model says the customer will stay, but he ends up leaving. Then that's considered False Positive and False Negative, respectively.

We need metrics that understand this cost difference, as it could be very costly if any of this goes wrong.

## A. The Confusion Matrix

This is our truth table that shows all four possible outcomes:

```
                   Actual: Churn        Actual: No Churn
Predicted: Churn       TP (True Positive)     FP (False Positive)
Predicted: No Churn  FN (False Negative)    TN (True Negative)
```

## B. Key Evaluation Metrics

From the confusion matrix, we derive several important metrics:

**Precision = TP / (TP + FP)**
Meaning, when we predict a customer will churn, we're right (precision)% of the time.

**Recall = TP / (TP + FN)**
Meaning, We predict exactly (Recall)% of the customers who actually churn.

**Accuracy = (TP + TN) / Total**
Meaning: overall, how often are we correct?

## C. The Precision-Recall Trade-off

However, in many cases, if we want a higher recall: We want to predict more churners, we have to accept lower precision; taking in even false alarms. Such a thing is called a trade-off.

## D. ROC Curve and AUC

Now, to determine the best Trade-Off: At what precision do we need to have the best recall, we plot a curve: the Receiver Operating Characteristic (ROC). This curve shows us the relationship between True Positive Rate (recall) and False Positive Rate across all possible classification thresholds.

The Area Under the Curve (AUC) gives us a single number that summarizes the model's overall performance. AUC closer to 1.0 means excellent performance, while 0.5 means no better than random guessing.

# VII. Implementation

## A. The Core Functions

So we start with the sigmoid function. We clip `z` because if you don't, the exponential goes insane with overflow errors:

```python
def sigmoid(z):
    z = np.clip(z, -500, 500)
    return 1 / (1 + np.exp(-z))
```

The linear score is the dot product we talked about earlier, `y = w * x + b`:

```python
def compute_linear_score(X, weights, bias):
    return np.dot(X, weights) + bias
```

Then we put it in Z and then evaluate it within the sigmoid function:

```python
def predict_probability(X, weights, bias):
    z = compute_linear_score(X, weights, bias)
    return sigmoid(z)
```

## B. Training Loop

This is where gradient ascent does its thing. Start with random weights, climb up the log-likelihood:

```python
def train_model(X, y, learning_rate=0.1, n_iterations=500):
    m, n = X.shape
    weights = np.random.randn(n) * 0.01
    bias = 0.0

    for i in range(n_iterations):
        y_pred = predict_probability(X, weights, bias)
        error = y_pred - y
        dw = (1 / m) * np.dot(X.T, error)
        db = (1 / m) * np.sum(error)

        weights += learning_rate * dw
        bias += learning_rate * db

    return weights, bias
```

We're *adding* the gradient here, not subtracting. Basically the main difference between it and gradient Descent.

## C. Preprocessing the Data

The Telco dataset needs cleaning. `TotalCharges` has some weird string values that break everything, and we convert Yes/No to 1/0:

```python
def clean_data(df):
    df = df.copy()
    df['TotalCharges'] = pd.to_numeric(df['TotalCharges'],
errors='coerce')
    df['TotalCharges'].fillna(0, inplace=True)
    df['Churn'] = df['Churn'].map({'Yes': 1, 'No': 0})
    return df
```

One-hot encoding for categorical stuff. Contract type becomes binary columns:

```
def one_hot_encode_column(df, column):
    dummies = pd.get_dummies(df[column], prefix=column, drop_first=True)
    return pd.concat([df.drop(column, axis=1), dummies], axis=1)
```

We scale it, since scaling makes sure features are on similar ranges; otherwise big numbers dominate. We ofc use `z = ( value - mean ) / std_deviation`, for every value being the feature:

```
def scale_features(df, features):
    df = df.copy()
    for feature in features:
        mean_val = df[feature].mean()
        std_val = df[feature].std()
        df[feature] = (df[feature] - mean_val) / std_val
    return df
```

## D. Evaluation

The confusion matrix just counts everything:

```
def confusion_matrix(true_labels, predicted_labels):
    tp = tn = fp = fn = 0
    for actual, predicted in zip(true_labels, predicted_labels):
        if actual == 1 and predicted == 1:
            tp += 1
        elif actual == 0 and predicted == 0:
            tn += 1
        elif actual == 0 and predicted == 1:
            fp += 1
        else:
            fn += 1
    return tp, tn, fp, fn
```

The other metrics come from those counts using the formulas we already defined.

ROC-AUC sweeps through thresholds and plots TPR vs FPR at each one. Then we use the trapezoidal rule to get the area under the curve—just summing rectangles.

# VIII. Results

## A. Dataset Stats

Ran this on 1,000 sample Telco customers. After all the preprocessing:

- Training: 800 customers
- Testing: 200 customers
- Features: 6 (once we did one-hot encoding)
- Churn rate: around 30%

## B. Training

Model converged over 500 iterations:

```
Step 0: Score = -0.6912
Step 100: Score = -0.5234
Step 200: Score = -0.4987
Step 300: Score = -0.4876
Step 400: Score = -0.4823
Step 499: Score = -0.4798
```

Log-likelihood goes up, which means it's learning.

Confusion matrix on test data:

```
                 | Actual Churn | Actual Stay
Predicted Churn  |      42      |     18
Predicted Stay   |      17      |     123
```

Metrics:

- Accuracy: 82.5%
- Precision: 70.0%
- Recall: 71.2%
- F1-Score: 70.6%
- AUC: 0.847

So 82.5% accuracy means we're right most of the time. Precision at 70% tells us when we flag someone as churning, we're correct 7 out of 10 times. The other 3 are false positives—false alarms basically.

Recall is 71.2%, so we catch about 71% of actual churners. We miss 29%, those are false negatives.

AUC of 0.847 is above 0.5 (random) and getting close to 1.0 (perfect). The model has real predictive power.

## C. What the Features Tell Us

Looking at the weights and odds ratios:

**Contract (Month-to-month):** OR = 2.34
Month-to-month customers are 2.34x more likely to churn than longer contracts. Which makes sense, as they don't really have any commitment.

**Tenure:** OR = 0.68
Each extra month decreases churn odds by 32%. Loyalty compounds.

**Monthly Charges:** OR = 1.42
Higher bills increase churn by 42%.

**Internet Service (Fiber):** OR = 1.67
Fiber customers churn 67% more.

Focus areas based on features:

- Get people on longer contracts , since month-to-month is risky: they don't have any commitment.
- New customers need attention in first few months
- Pricing matters: high monthly charges push people away
- Something's wrong with the fiber service

# IX. Conclusion

We started with theory: sigmoid, log-odds, MLE, gradient ascent. Then, we built everything from scratch without using library functions for the core algorithm.

Afterwards, we trained on Telco data, evaluated it, got some insights.

82.5% accuracy and 0.847 AUC, which is quite decent for a baseline.