

DIFFERENCE BETWEEN CATBOOST AND XGBOOST.

Feature	CatBoost	XGBoost
Algorithm	Uses ordered boosting to reduce overfitting	Uses gradient boosting with decision trees
Handling Categorical Features	Handles categorical variables natively without encoding	Requires one-hot or label encoding for categorical features
Training Speed	Faster when handling categorical data	Slower with categorical features due to encoding
Overfitting Prevention	Uses ordered boosting to prevent target leakage	Uses regularization (L1, L2) to reduce overfitting
Performance on Small Datasets	Performs well due to efficient feature handling	Requires careful hyperparameter tuning for small datasets

WHAT IS THE REASON FOR THE MODEL TO PREDICT THE PARTICULAR OUTPUT(YES OR NO).

A model predicts "Yes" or "No" based on the learned patterns in the dataset. It considers:

- **Feature Importance:** Some features have more influence on predictions.
- **Decision Boundaries:** It calculates probabilities using learned weights.
- **Thresholding:** If the probability is above a threshold (e.g., 0.5), it predicts "Yes"; otherwise, "No".
- **Feature Contributions:** Using SHAP or LIME, we can analyze how much each feature contributes to the final decision.

ANALAYSE THE COLUMN(MONTHLY RATE, MONTHLY INCOME, HOURLY RATE, DAILY INCOME), TO UNDERSTAND THE COLUMN.

The issue occurs because:

- **Irregular Work Hours:** Employees may work variable hours, so hourly and daily earnings don't always sum up logically.
- **Missing Data or Errors:** Possible data entry mistakes or missing records.

How to Overcome This Problem?

- **Check Data Consistency:** Ensure correct data entry and calculations.
- **Feature Engineering:** Create derived features like "Effective Hourly Rate" (Monthly Income / Hours Worked).
- **Normalization:** Scale the data to reduce inconsistencies.

WHAT IS CORRELATION? WHAT DID YOU FIND IN CORRELATION MAPPING?

Correlation measures the relationship between two numerical variables.

Values range from **-1 (strong negative correlation)** to **1 (strong positive correlation)**.

In correlation mapping:

- A strong correlation suggests dependent relationships (e.g., high "Monthly Income" correlates with "Monthly Rate").
- No correlation suggests independent variables.

WHAT DOES VALUE COUNT MEANS AND HOW DOES IT WORK.

`value_counts()` is a Pandas method used to count occurrences of unique values in a column.

WHAT METHOD DID YOU USE TO REMOVE THE COLUMNS FROM THE DATASET.

`value_counts()` can help identify columns with a single dominant value, which may not be useful for modeling.

WHAT ARE OUTLIERS. HOW DO YOU HANDLE OUTLIERS.

Outliers are extreme values that deviate significantly from most data points.

Handling Outliers:

- **Z-score method:** Remove values beyond 3 standard deviations.
- **IQR (Interquartile Range) method:** Remove values outside $1.5 * IQR$.
- **Transformation:** Use log or Box-Cox transformations to reduce skewness.
- **Capping:** Replace extreme values with the upper/lower percentile.

WHAT IS STANDARDISATION? WHY WE NEED IT.

Standardization scales data to have a mean of 0 and a standard deviation of 1.

Why Needed?

- Ensures all features contribute equally to model performance.
- Improves convergence speed in gradient-based optimization (e.g., logistic regression, neural networks).

WHAT IS ONE HOT ENCODING. WHAT IS LABEL ENCODING. WHEN WE NEED TO USE EACH ONE.

One-Hot Encoding (OHE): Converts categorical variables into multiple binary columns.

Label Encoding: Assigns unique numerical values to each category.

WHAT IS OVERSAMPLING.

Oversampling is increasing the number of samples in the minority class to balance the dataset.

Example: If Class "No" has 500 samples and "Yes" has 100, oversampling increases "Yes" to 500.

Methods: Random Oversampling, SMOTE.

WHAT IS UNDERSAMPLING.

Undersampling reduces the number of majority class samples to balance the dataset.

Example: If "No" has 500 samples and "Yes" has 100, undersampling reduces "No" to 100.

Risk: Loss of valuable information.

DIFFERENCE BETWEEN OVERSAMPLING AND UNDERSAMPLING.

Aspect	Oversampling	Undersampling
Method	Adds synthetic/minority samples	Removes majority samples
Goal	Balance the dataset by increasing minority samples	Balance by reducing majority samples
Risk	Overfitting	Loss of information

WHAT IS SMOTE?

Synthetic Minority Over-sampling Technique (SMOTE) is an advanced oversampling method.

Instead of duplicating samples, it creates **synthetic samples** by interpolating between existing minority class samples.

WHAT IS PRECISION, RECALL. DIFFERENCE BETWEEN THEM.

Metric	Formula	Meaning	When to Focus on it?
Precision	$TP / (TP + FP)$	Measures how many positive predictions were correct	When false positives are costly (e.g., fraud detection)
Recall	$TP / (TP + FN)$	Measures how many actual positives were correctly identified	When false negatives are costly (e.g., medical diagnosis)