# EMPLOYEE CHURN PREDICTION USING ML MODELS

By AFIYA AFSHEEN
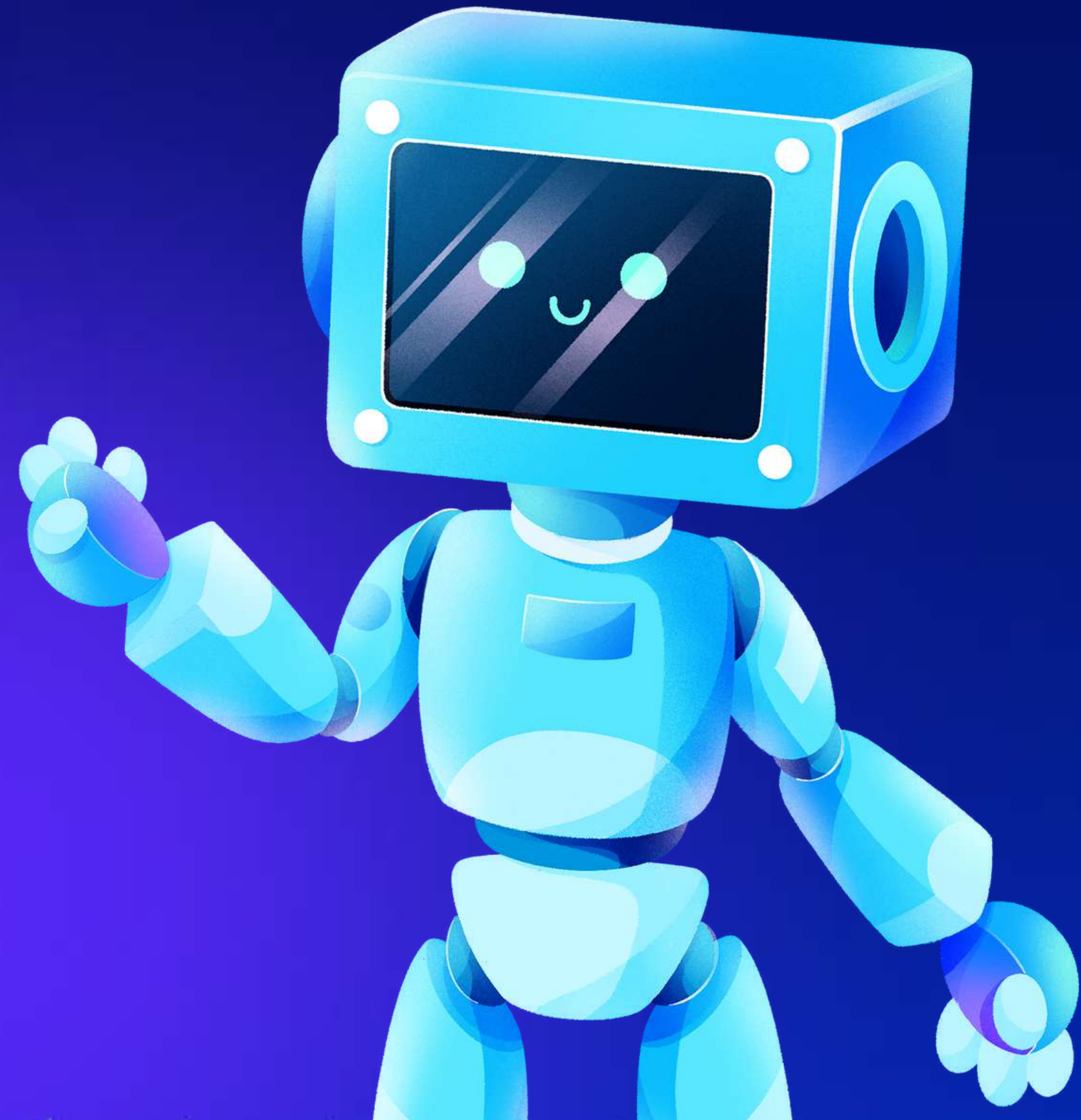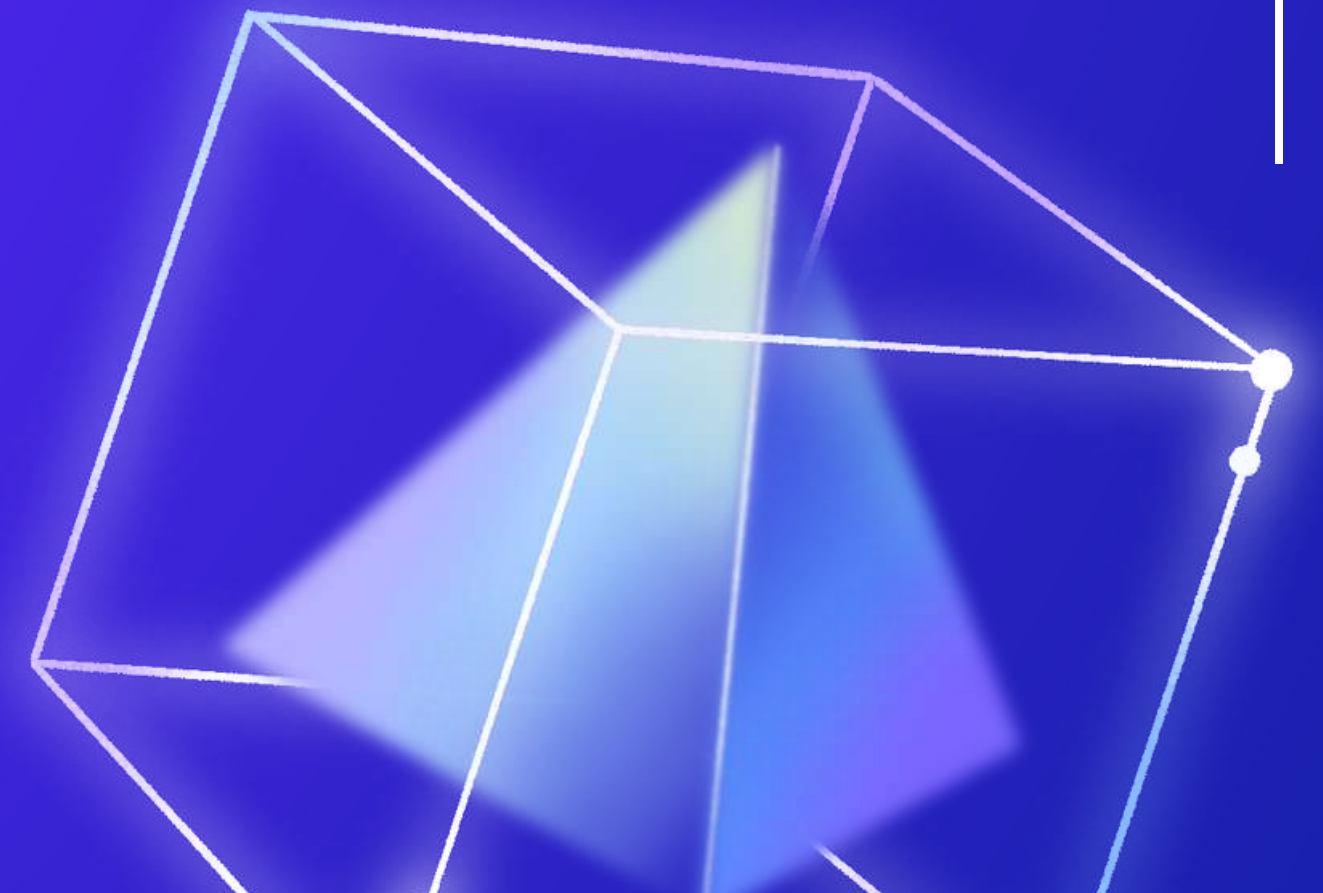
# TABLE OF CONTENTS
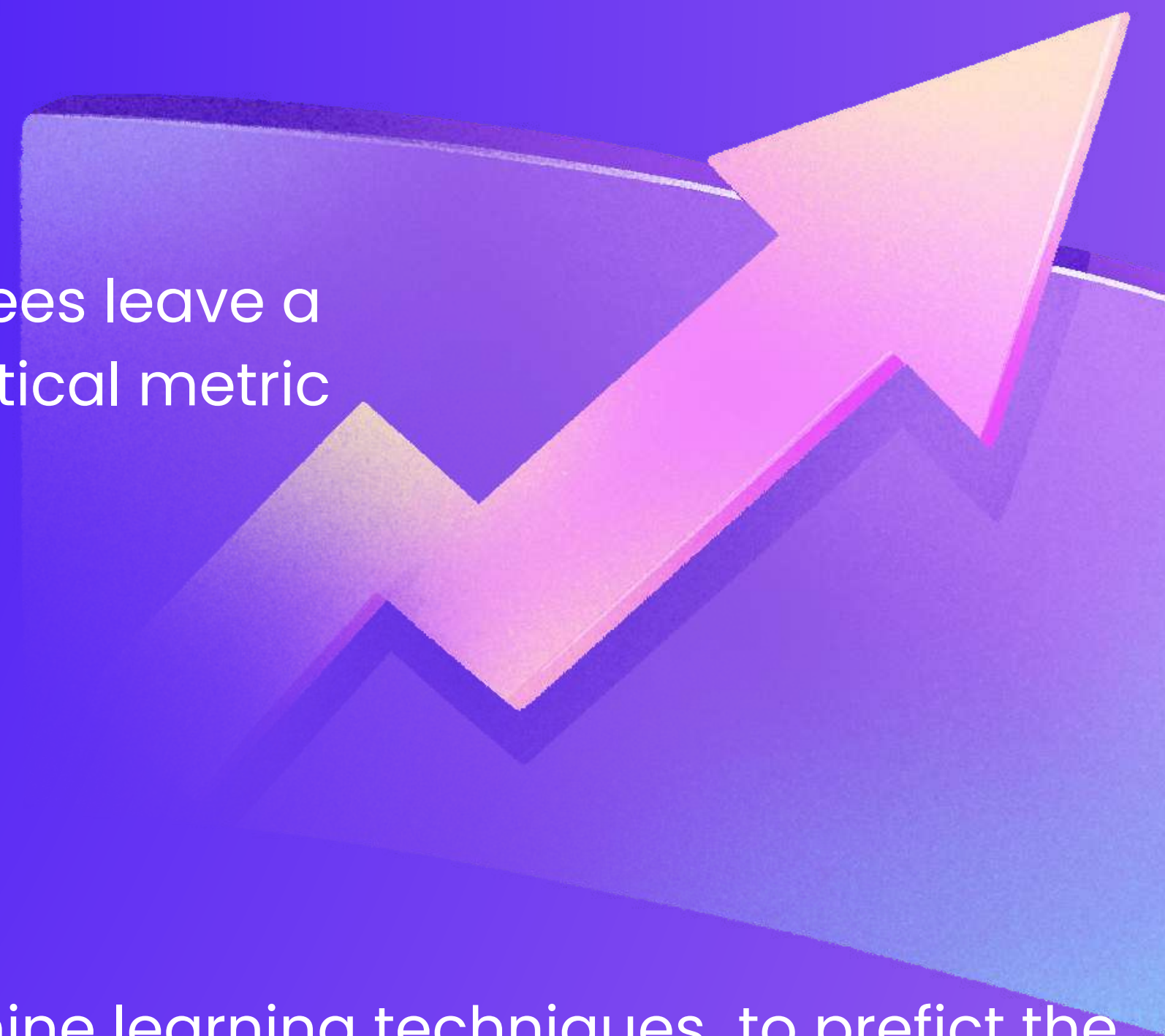
- Introduction
- Project Flow
- About Dataset
- Basic Cleaning
- Methodology-wise Basic DataFrame
- **Visualization**
- **Correlation Matrix**
- Class Imbalnce ,Outliers,Skewness.
- DataFrame_2
- Dataframe_1 (v/s) DataFrame_2
- Feature Selection
- Conclusion

# INTRODUCTION

Employee turnover, or the rate at which employees leave a company and are replaced by new hires, is a critical metric for organizations across industries.

In this model we will be deplloying supervised Machine learning techniques  to prefict the turnover

# FLOW OF THE PROJECT
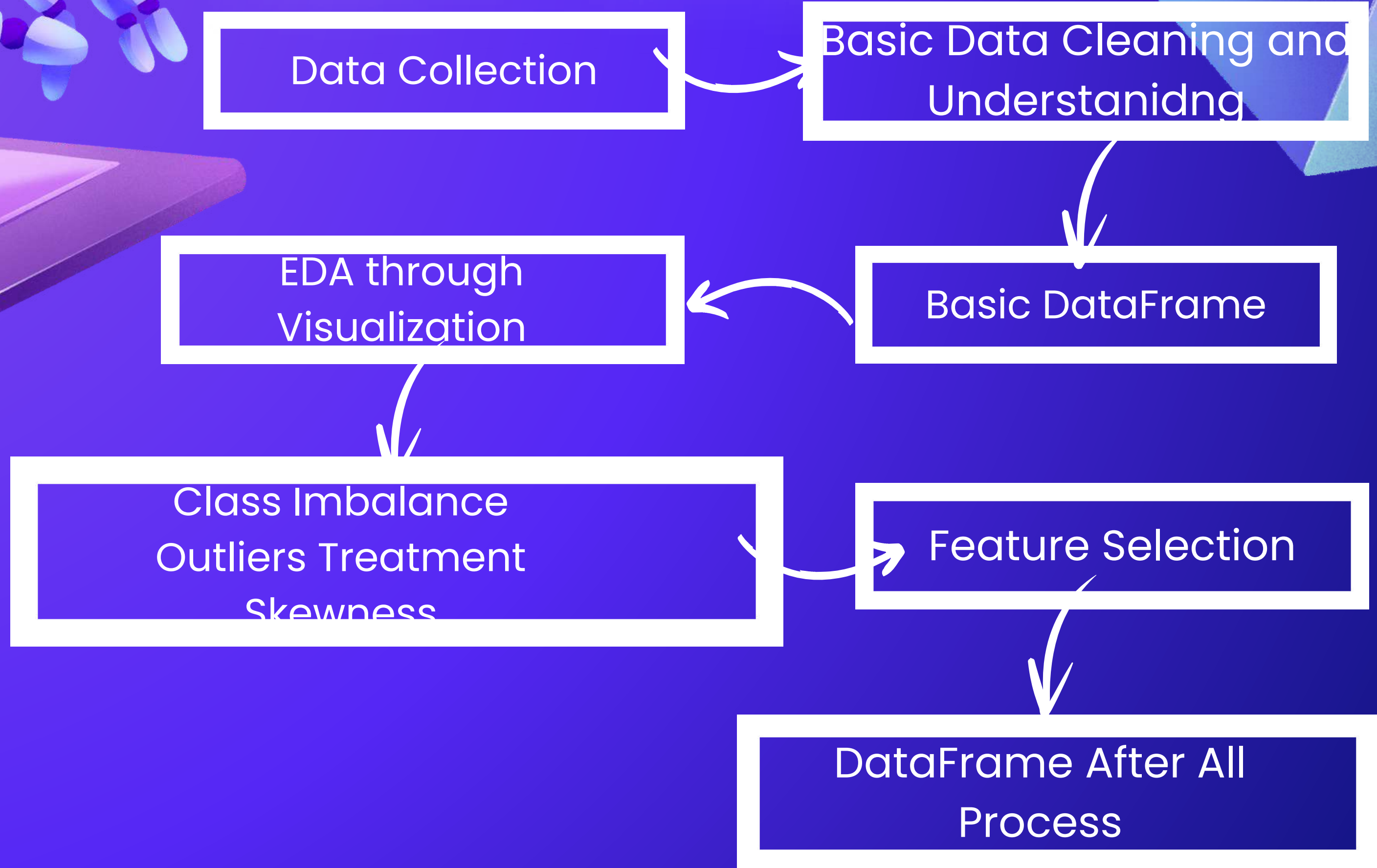
Data Collection → Basic Data Cleaning and Understanidng → Basic DataFrame → EDA through Visualization → Class Imbalance Outliers Treatment Skewness → Feature Selection → DataFrame After All Process

# ABOUT DATASET

DataSet: Real dataset shared from Edward Abushkin's blog used to predict an Employee's risk of quitting.

Columns: 10

| Features | "department" , "promoted" ,"review" ,"projects" ,"salary","tenure","satisfaction" ,"bonus" ,"avg_hrs_month" , "left" |
|---|---|
| Target | "Left" |
| No.of Rows | 9540 |

## Problem Statement

Predict if employee will leave company.

# BASIC CLEANING:

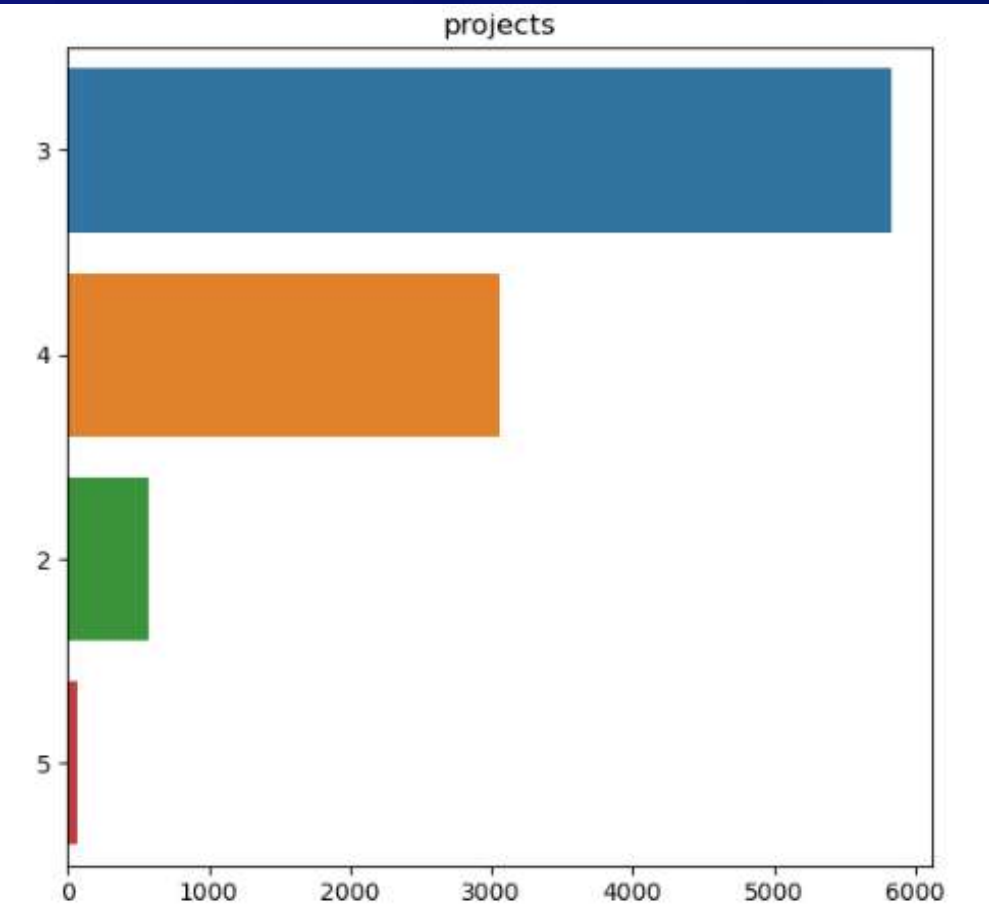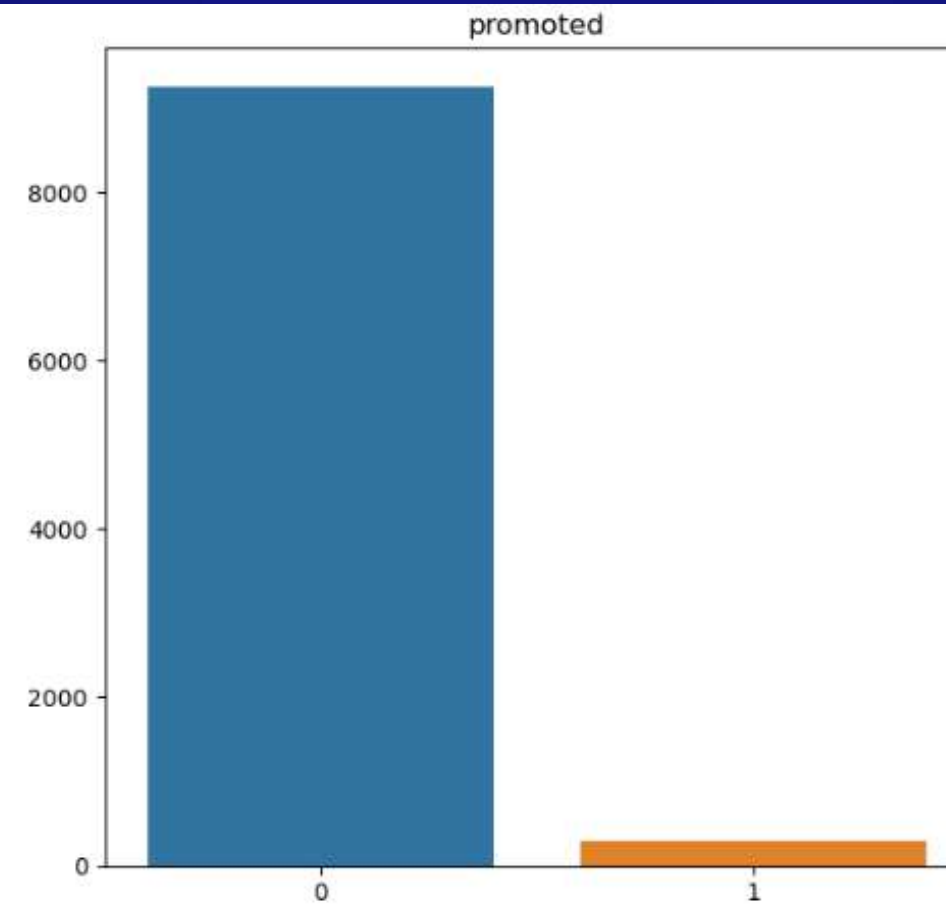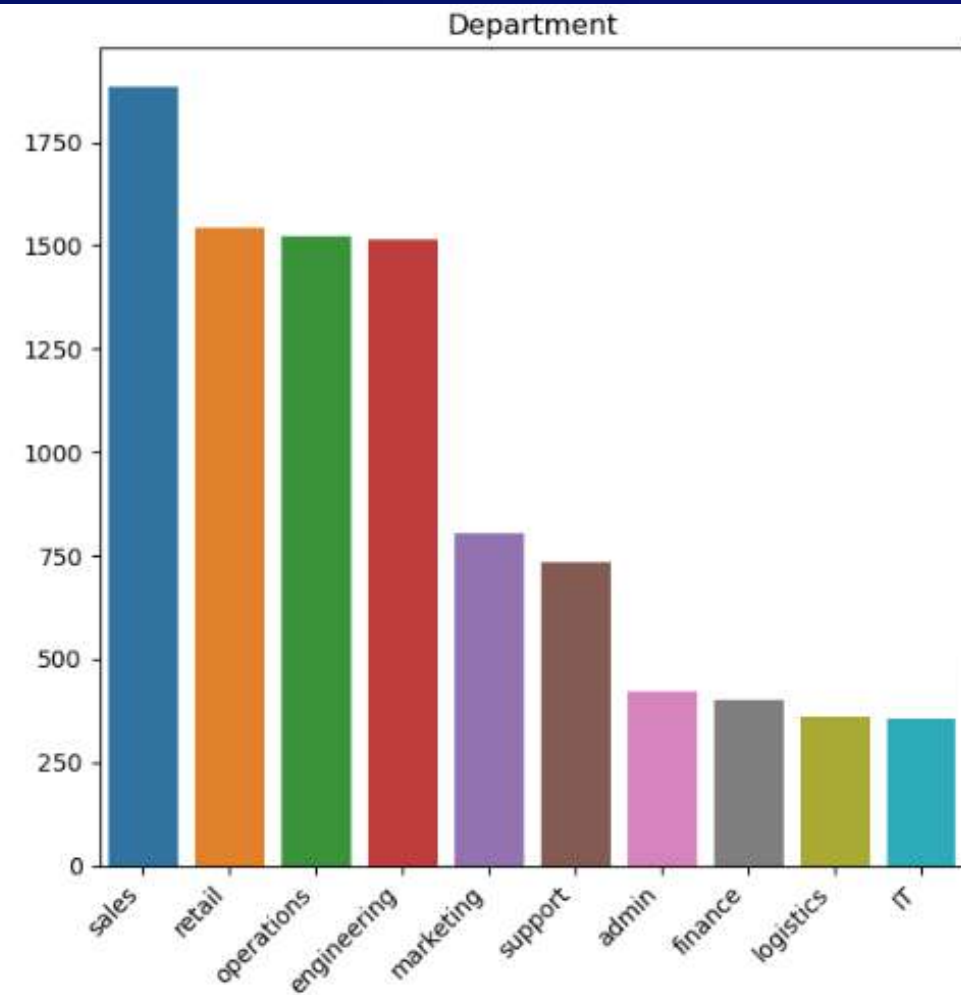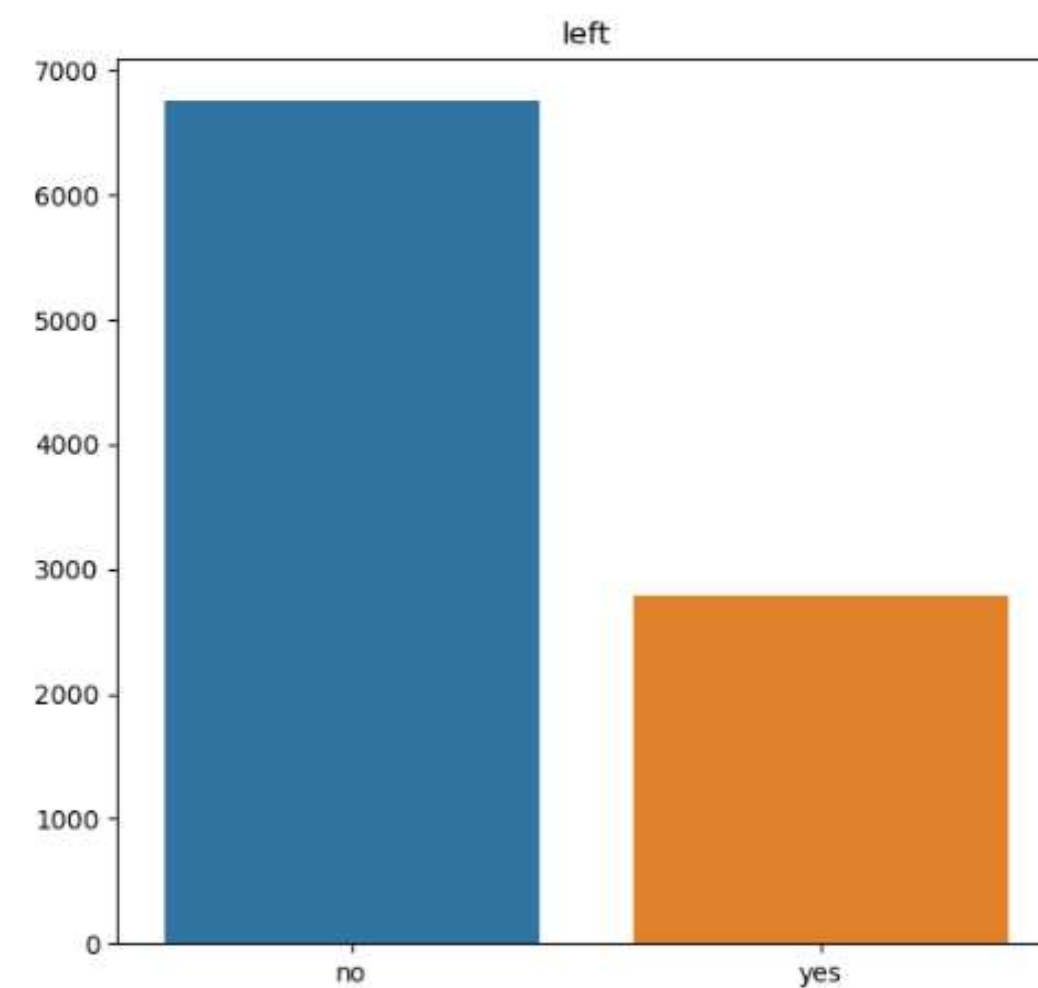| Null Values | O |
|---|---|
| DATA TYPE OBJECT | "department","salary","left","Promoted","bonus |
| DATA TYPE NUMERIC | "review","projects","tenure",satisfaction","avg_hrs_month". |

DATA TYPE CONVERSION USING LABEL ENCODER

# BASIC DATAFRAME

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Logistic Regression | 72.327044 | 0.726810 | 0.966165 | 0.264624 |
| Decision Tree | 81.813417 | 0.860844 | 0.874436 | 0.694811 |
| Random Forest | 85.691824 | 0.872702 | 0.927820 | 0.747456 |
| KNN | 73.322851 | 0.772757 | 0.874436 | 0.481142 |
| SVM | 69.706499 | 0.697065 | 1.000000 | 0.000000 |
| XGBoost | 85.534591 | 0.879683 | 0.918045 | 0.748634 |

01

# CLASS IMBALANCE

**03**

TARGET COLUMN "LEFT "

HAD CLASS IMBALANCE

ACTION TAKEN : APPLIED SMOTE TECHNIQUE

"LEFT" - "YES" IF THE EMPLOYEE ENDED UP LEAVING, "NO" OTHERWISE.

| Instance | No.of "NO" | **No.of"YES"** |
|----------|-----------|-----------|
| Before | 6756 | 2784 |
| After | 6756 | 6756 |

# OUTLIERS

## Boxplots of Numerical Columns



NUMERICAL COLUMNS WHICH REQUIRE OUTLIERS

TREATMENT :

'REVIEW','PROJECTS','TENURE",'AVG_HRS_MONTH'

| Type of Distribution | Method Used | Columns |
|---|---|---|
| Normal | Standard Deviation | "review","project","tenure","avg_hrs_month" |
| Skewed | IQR | NA |

## 05 SKEWNESS

| Column | Skewness | Skewness Treatment Required |
|---|---|---|
| department | −0.526152 | NO |
| promoted | 6.533033 | NO |
| review | 0.153407 | NO |
| projects | 0.023403 | NO |
| salary | −1.237002 | NO |
| tenure | −0.099847 | NO |
| satisfaction | 0.099041 | NO |
| bonus | 1.554054 | NO |
| avg_hrs_month | −0.119797 | NO |

# DATAFRAME_2

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Logistic Regression | 53.015168 | 0.691834 | 0.664693 | 0.558414 |
| Decision Tree | 80.096189 | 0.842502 | 0.827535 | 0.802496 |
| Random Forest | 85.423603 | 0.884384 | 0.871947 | 0.85482 |
| KNN | 74.953755 | 0.797524 | 0.667654 | 0.764522 |
| SVM | 49.981502 | 0.498144 | 0.993338 | 0.000000 |
| XGBoost | 83.3518 | 0.862324 | 0.874625 | 0.834680 |
| AdaBoost | 69.367370 | 0.811979 | 0.792746 | 0.675549 |

# DATAFRAME_1 V/S DATAFRAME_2

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Logistic Regression | 72.327044 | 0.726810 | 0.966165 | 0.264624 |
| Decision Tree | 81.813417 | 0.860844 | 0.874436 | 0.694811 |
| Random Forest | 85.691824 | 0.872702 | 0.927820 | 0.747456 |
| KNN | 73.322851 | 0.772757 | 0.874436 | 0.481142 |
| SVM | 69.706499 | 0.697065 | 1.000000 | 0.000000 |
| XGBoost | 85.534591 | 0.879683 | 0.918045 | 0.748634 |

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Logistic Regression | 53.015168 | 0.691834 | 0.664693 | 0.558414 |
| Decision Tree | 80.096189 | 0.842502 | 0.827535 | 0.802496 |
| Random Forest | 85.423603 | 0.884384 | 0.871947 | 0.85482 |
| KNN | 74.953755 | 0.797524 | 0.667654 | 0.764522 |
| SVM | 49.981502 | 0.498144 | 0.993338 | 0.000000 |
| XGBoost | 83.3518 | 0.862324 | 0.874625 | 0.834680 |
| AdaBoost | 69.367370 | 0.811979 | 0.792746 | 0.675549 |

# FEATURE SELECTION

We see promoted Salary projects and Bonus are least important

| Importance | Column |
|------------|--------|
| 0.049497 | department |
| 0.003577 | promoted |
| 0.260457 | review |
| 0.021433 | projects |
| 0.018584 | salary |
| 0.065972 | tenure |
| 0.301302 | satisfaction |
| 0.012011 | bonus |
| 0.267167 | avg_hrs_month |

# DATA FRAME AFTER FEATURE SELECTION

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Logistic Regression | 53.570107 | 0.530918 | 0.629164 | 0.487127 |
| Decision Tree | 78.727340 | 0.794852 | 0.777202 | 0.789454 |
| Random Forest | 83.980762 | 0.848187 | 0.831236 | 0.841450 |
| KNN | 80.133185 | 0.821853 | 0.768320 | 0.808556 |
| SVM | 46.281909 | 0.495896 | 0.983716 | 0.004115 |
| XGBoost | 83.166852 | 0.829756 | 0.834077 | 0.833150 |
| AdaBoost | 68.183500 | 0.693796 | 0.769800 | 0.661684 |

| Model | BEFORE Accuracy | Precision | Recall | F1-Score | AFTER Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|---|---|---|---|
| Logistic Regression | 53.015168 | 0.691834 | 0.664693 | 0.558414 | 53.570107 | 0.530918 | 0.629164 | 0.487127 |
| Decision Tree | 80.096189 | 0.842502 | 0.827535 | 0.802496 | 78.727340 | 0.794852 | 0.777202 | 0.789454 |
| Random Forest | 85.423603 | 0.884384 | 0.871947 | 0.85482 | 83.980762 | 0.848187 | 0.831236 | 0.841450 |
| KNN | 74.953755 | 0.797524 | 0.667654 | 0.764522 | 80.133185 | 0.821853 | 0.768320 | 0.808556 |
| SVM | 49.981502 | 0.498144 | 0.993338 | 0.000000 | 46.281909 | 0.495896 | 0.983716 | 0.004115 |
| XGBoost | 83.3518 | 0.862324 | 0.874625 | 0.834680 | 83.166852 | 0.829756 | 0.834077 | 0.833150 |
| AdaBoost | 69.367370 | 0.811979 | 0.792746 | 0.675549 | 68.183500 | 0.693796 | 0.769800 | 0.661684 |

# PROJECT

## DEMERITS

- Error Chances
- Dependence on Data Quality

## FACTS

The Accurecy which ever atmost we get we cannot tell it 100% that employee will leave or stay as this thing have much more dimensions in reality

## BASED ON DATASET

**Conclusion** : Random Forest is model which is giving the best accurecy outof Linear, Decison tree , Random forest , SVM,KNN,,XGBoost and Ada boost