

IBM HR Analytics

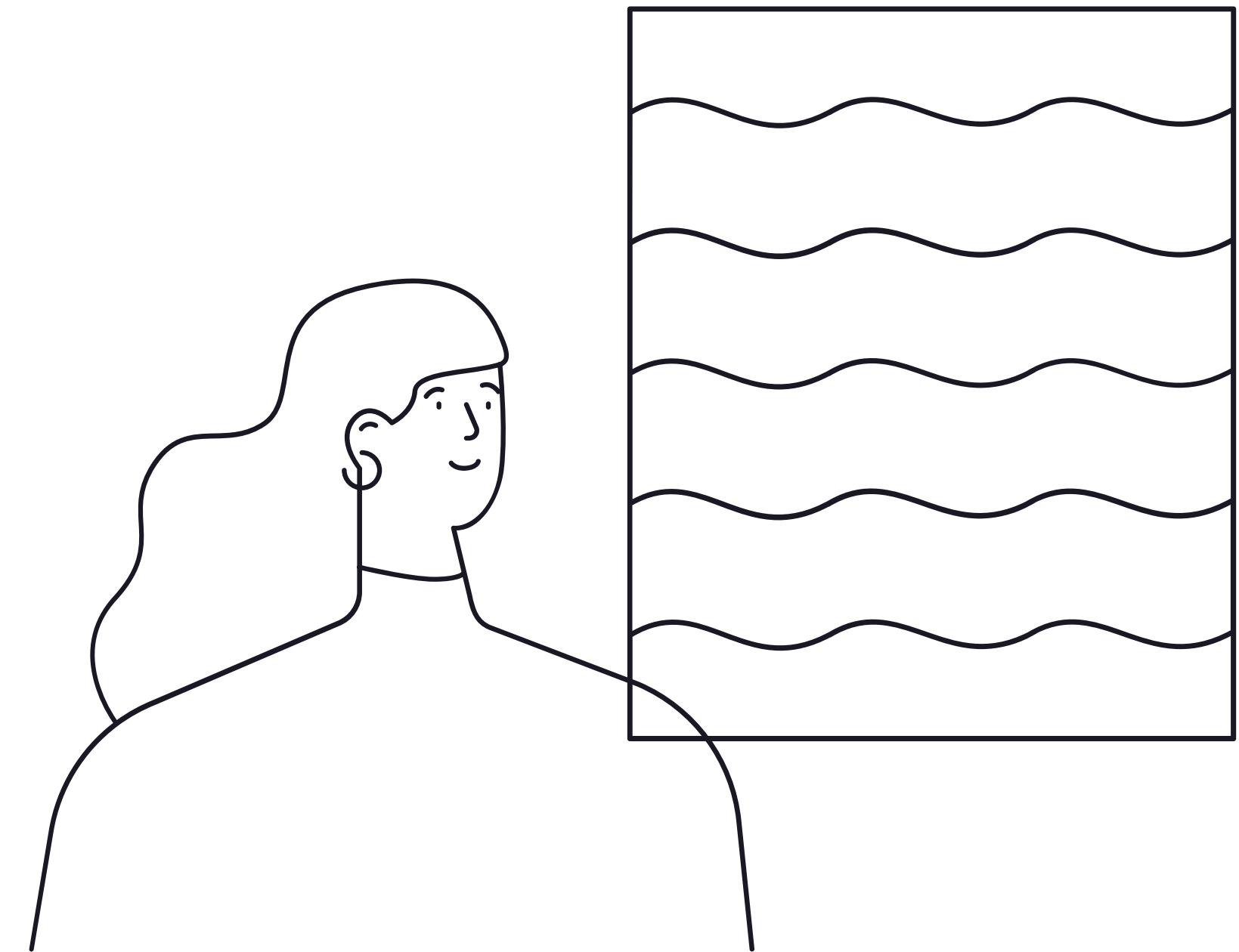
Team 4: Emory's Next Top Models
Neha Bansal
Emmy Fortunato
Boping Zhang

ATTRITION

- Predict the likelihood that an existing employee will quit

BUSINESS VALUE

- Expensive in terms of money and time to train new employees
- Loss of experienced employees
- Impact on profit because of business disruption



PREDICTIVE MODEL

THE DATA



The business question that we want to answer is:

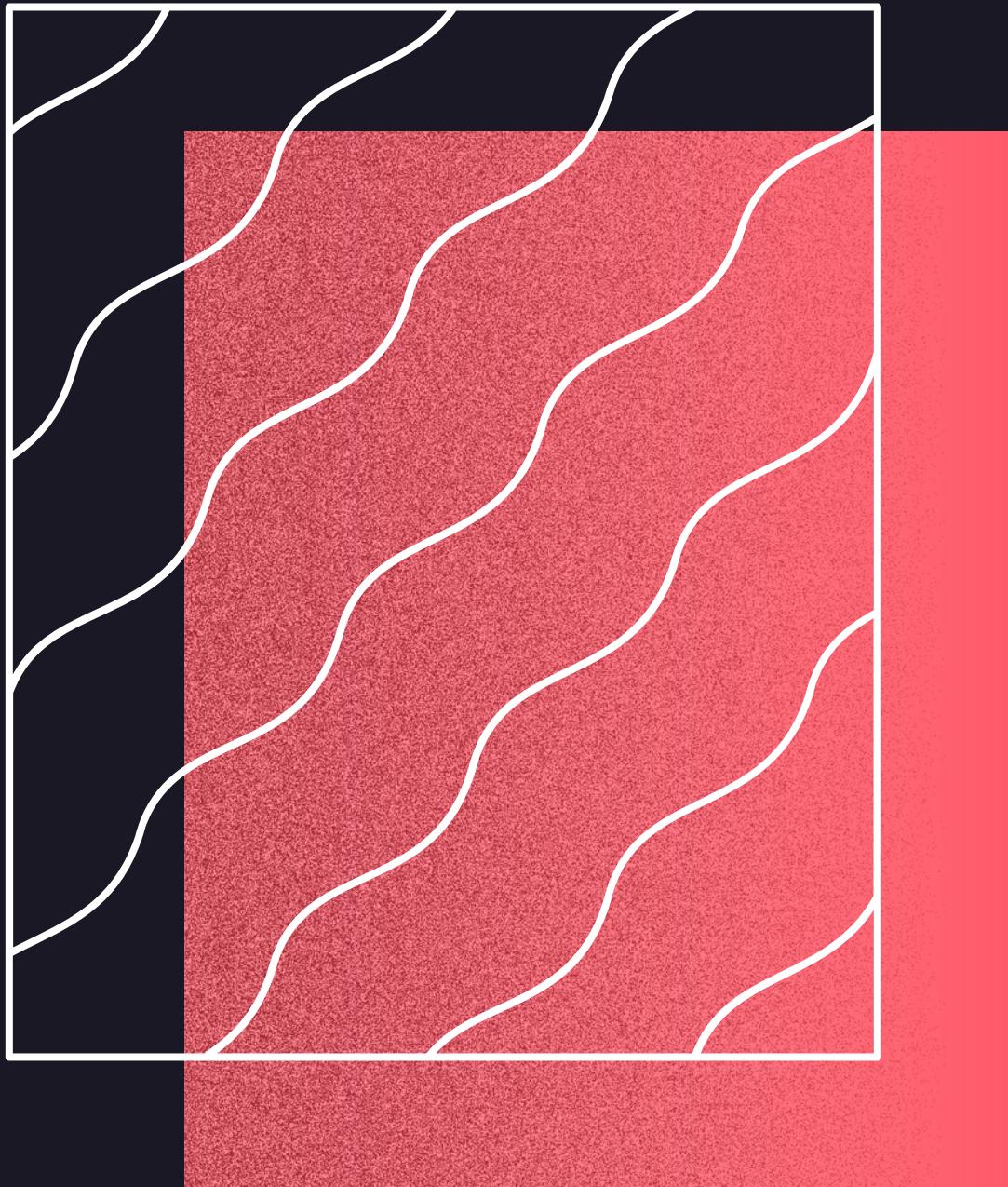
- Will employee "X" quit?

So the data science questions that we want to answer are:

- Which variables are the best predictors of employee attrition at IBM?
- Based on these variables, can we expect that employee "X" will voluntarily leave IBM?

To answer these questions, we will use the data available to us to build a predictive model.

TARGET VARIABLE



In the case of our dataset, Attrition is a binomial variable, with a value of 0 indicating that the employee has not left the company and a value of 1 indicating that the employee has voluntarily left the company.

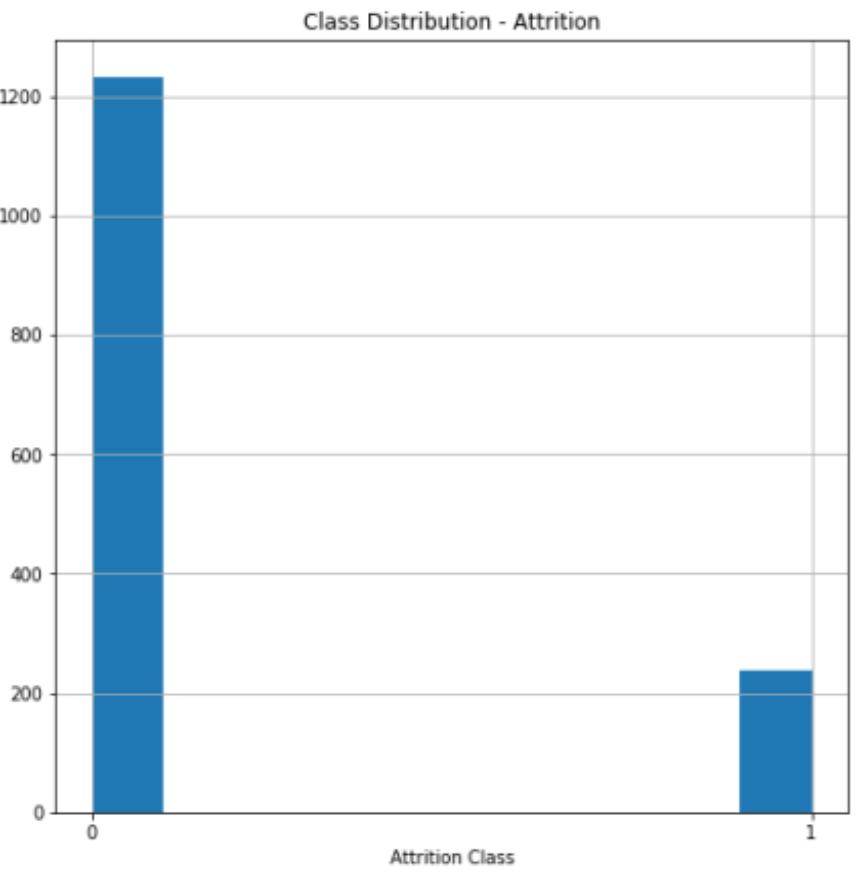
FEATURES

The X-variables in the dataset are made up of:

- numerical variables: measures of time in the workforce, wages earned, and some demographics
- nominal categorical variables: objective classifications by job, education, and demographic
- ordinal categorical variables: subjective rankings of job-related engagement, satisfaction, and performance

DATA PREPARATION

```
#3  # Ordinal Variable: BusinessTravel
#4  unique_travel = X['BusinessTravel'].unique()
#5  print(unique_travel)
#6  size_mapping = { 'Travel_Frequently': 3,
#7      'Travel_Rarely': 2,
#8      'Non-Travel': 1}
#9  X['BusinessTravel'] = X['BusinessTravel'].map(size_mapping)
#10
#11 # Binary Variable: Gender
#12 X["Gender"] = X["Gender"].replace(["Male"], 1)
#13 X["Gender"] = X["Gender"].replace(["Female"], 0)
#14 unique_gender = X['Gender'].unique()
#15 print(unique_gender)
#16
#17 # Binary Variable: Over18
#18 X["Over18"] = X["Over18"].replace(["Y"], 1)
#19 X["Over18"] = X["Over18"].replace(["N"], 0)
#20 unique_over18 = ibm_hr["Over18"].unique()
#21 print(unique_over18) #can remove!
#22
#23 # Binary Variable: Overtime
#24 X["Overtime"] = X["Overtime"].replace(["Yes"], 1)
#25 X["Overtime"] = X["Overtime"].replace(["No"], 0)
#26 unique_overtime = X["Overtime"].unique()
#27 print(unique_overtime)
```

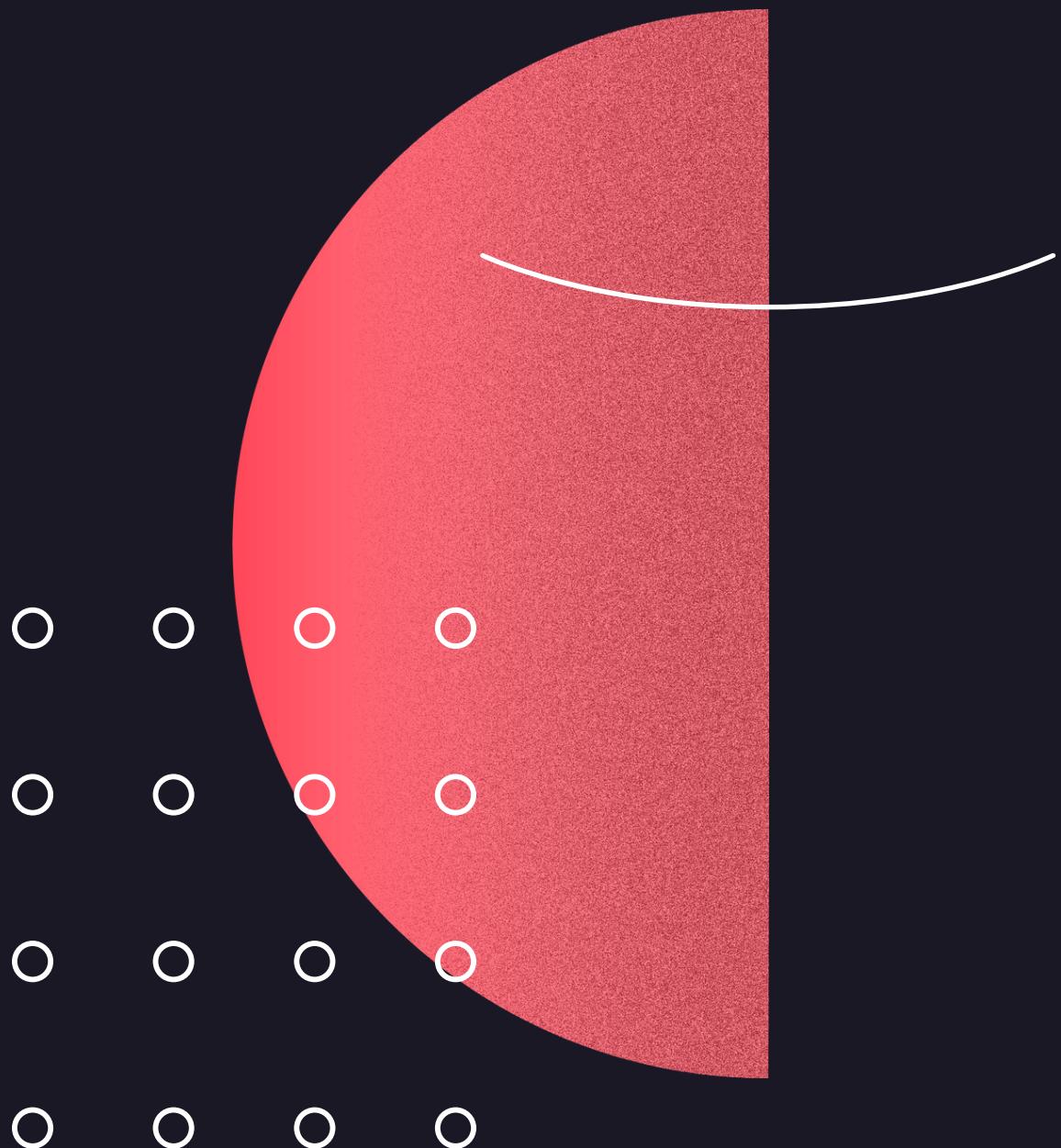


Our approach to preparing the data involved a few key steps:

- 1) Transform text variables to numeric variables.
- 2) Remove uninformative attributes.
- 3) Reduce multicollinearity.
- 4) Balance the data.
- 5) Conduct feature engineering.



DATA ENGINEERING



We add 5 new estimated features based on existing features:

- Annual Salary = $12 * \text{Monthly Income}$
- Annual Compensation = $12 * \text{Monthly Rate}$
- Bonus = Annual Compensation - Annual Salary
- Fidelity = Number of Company Worked / Total Working Years
- Performance of Job Satisfaction = Average of (Job Satisfaction + Job Involvement)

Model: Decision Tree Classifier

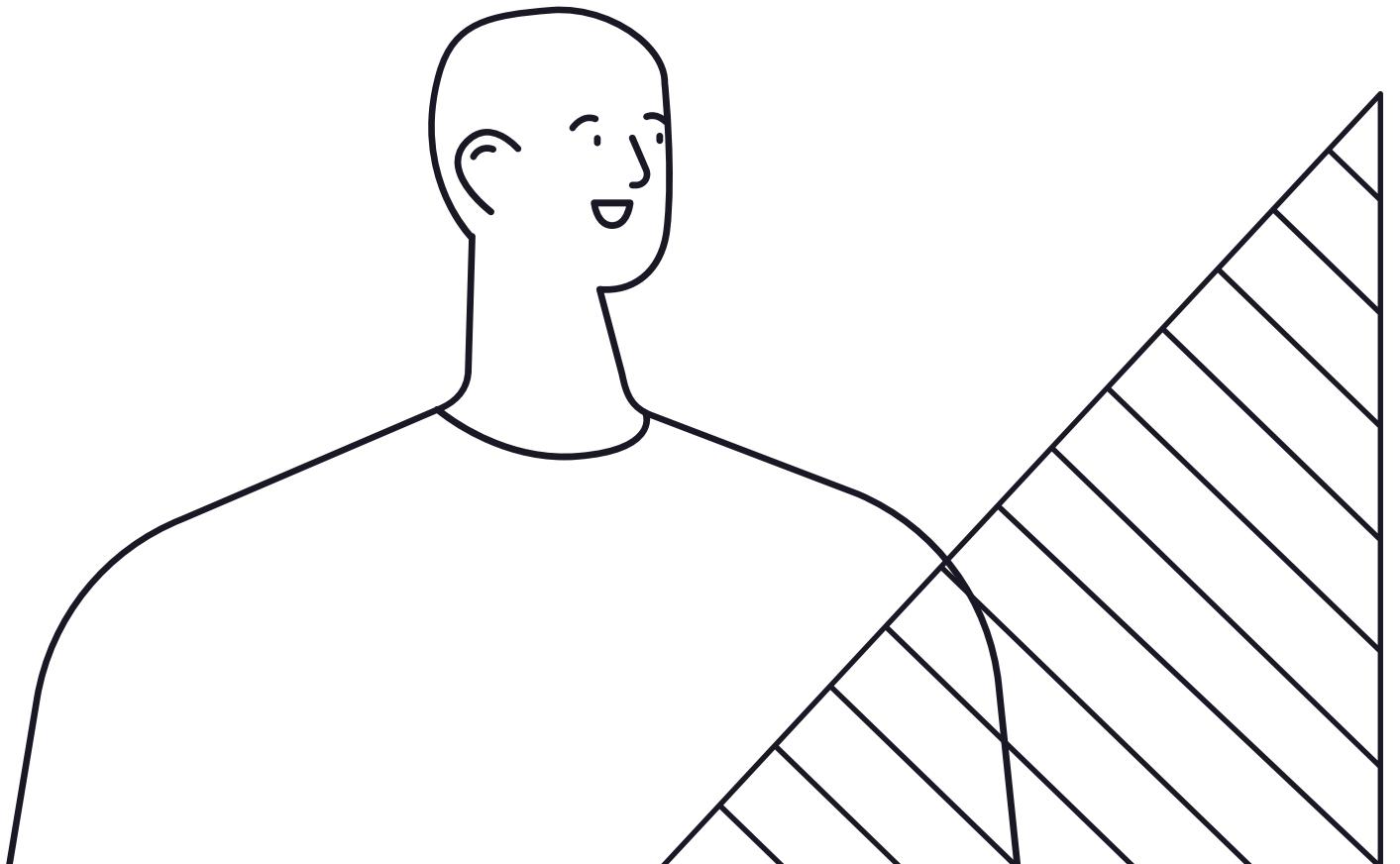
MODEL

Advantages:

- Consider all possible outcomes of a decision
- Assign specific values to each problem
- Easy to use
- Easy to understand for non-professionals
- The deeper branch can identify if an employee will leave the company

Disadvantage:

- Overfitting

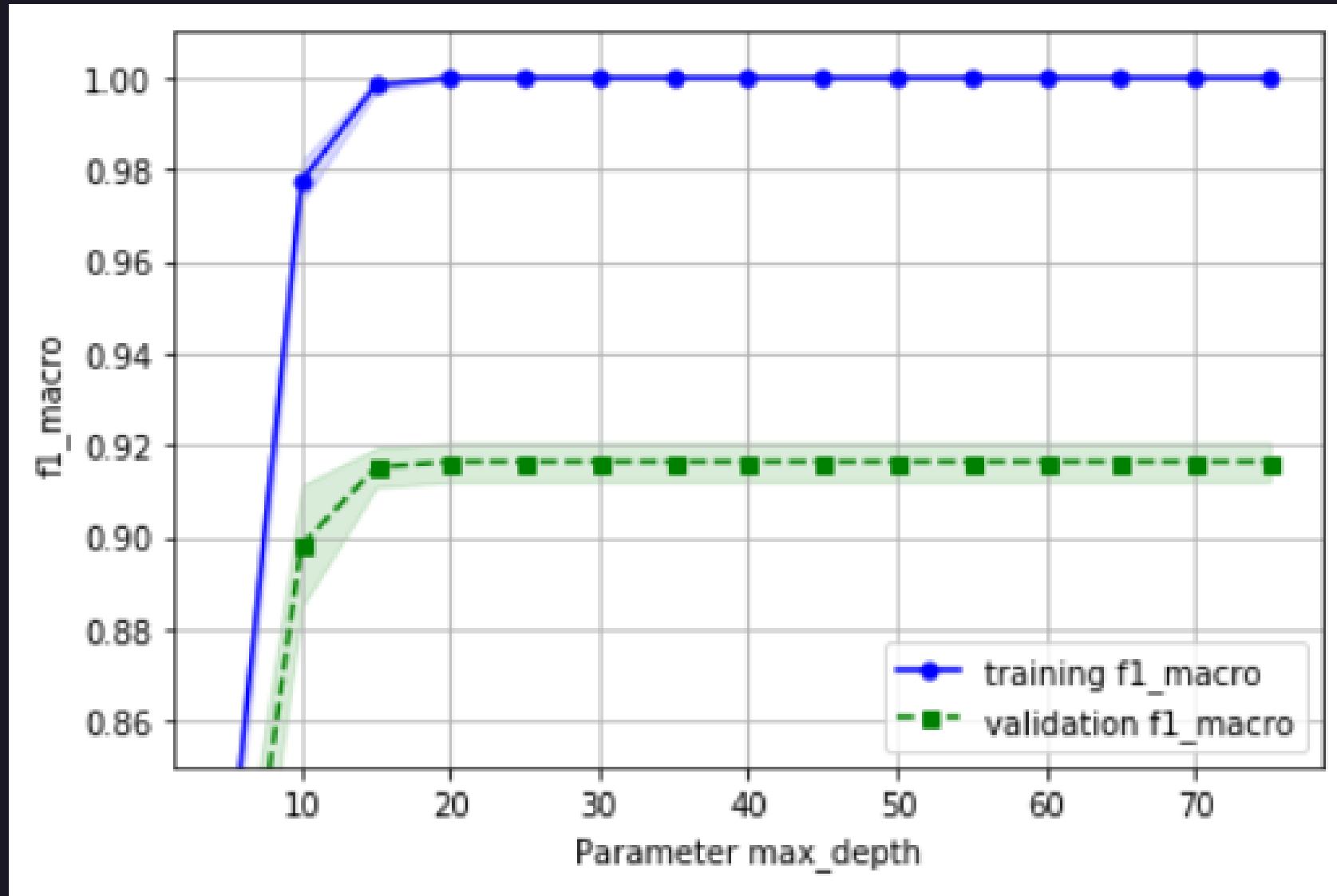


EVALUATION

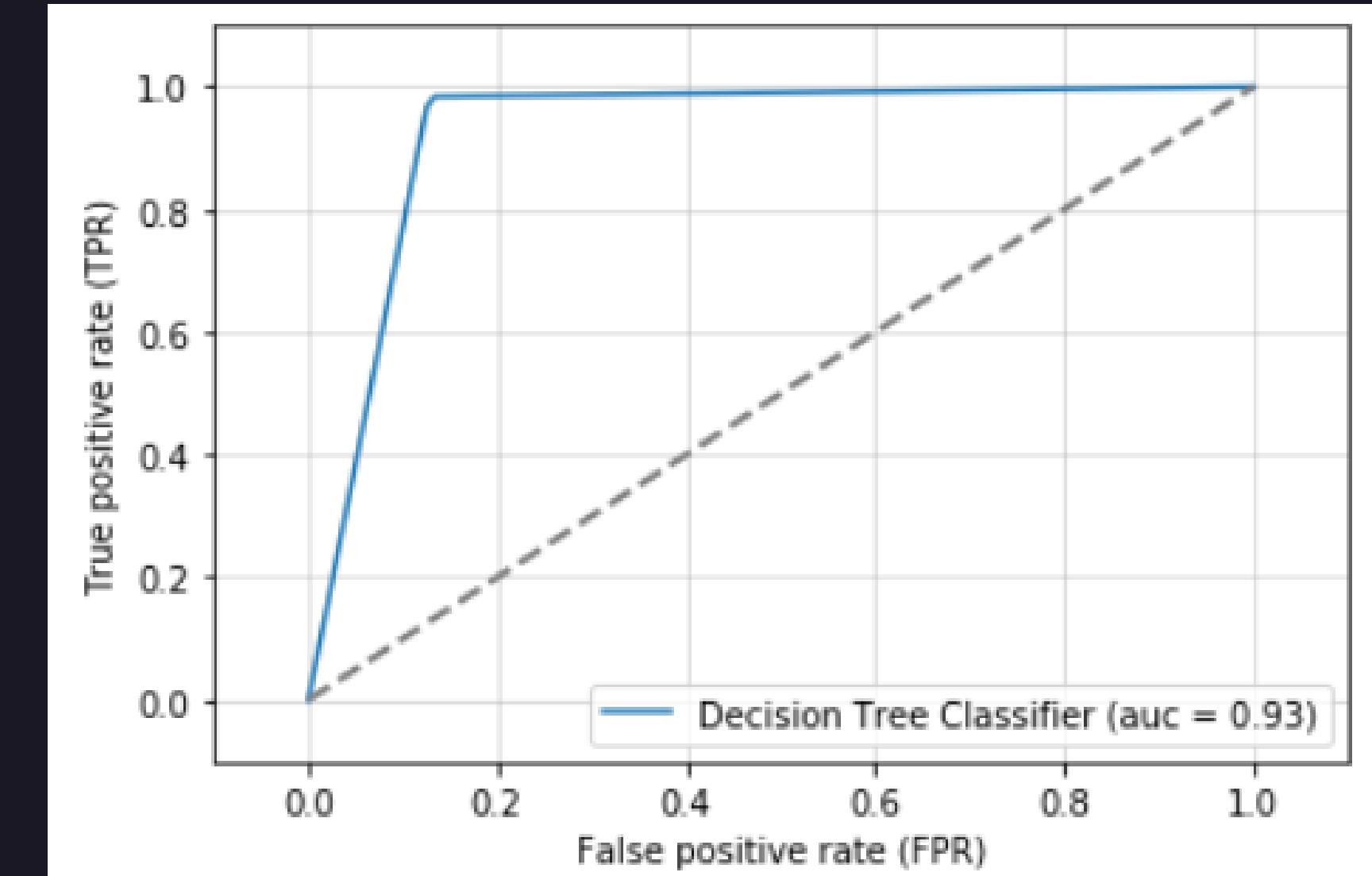
Non-nested f1_macro: 93.4%

Nested f1_macro: 92.07%

FITTING GRAPH



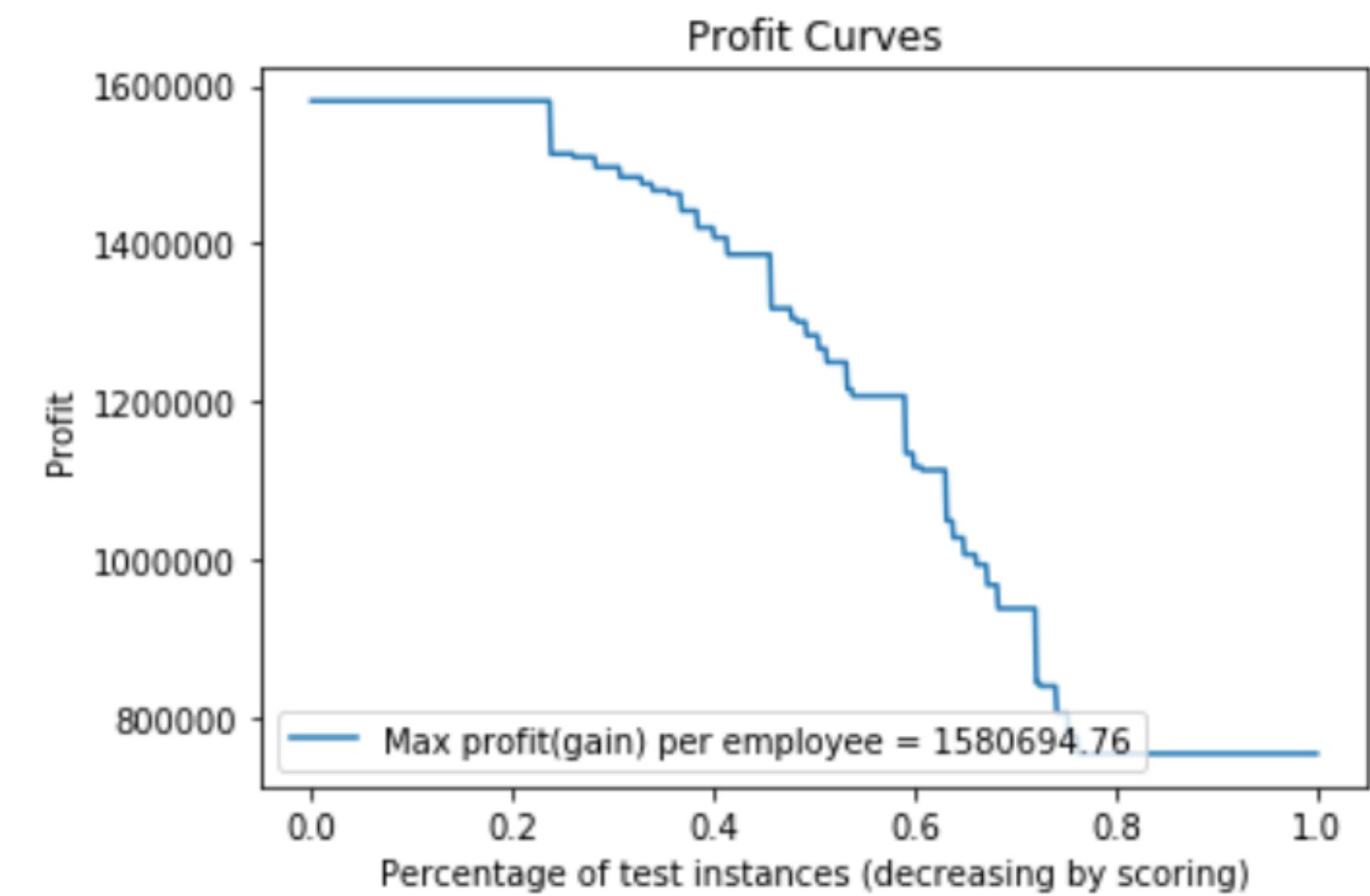
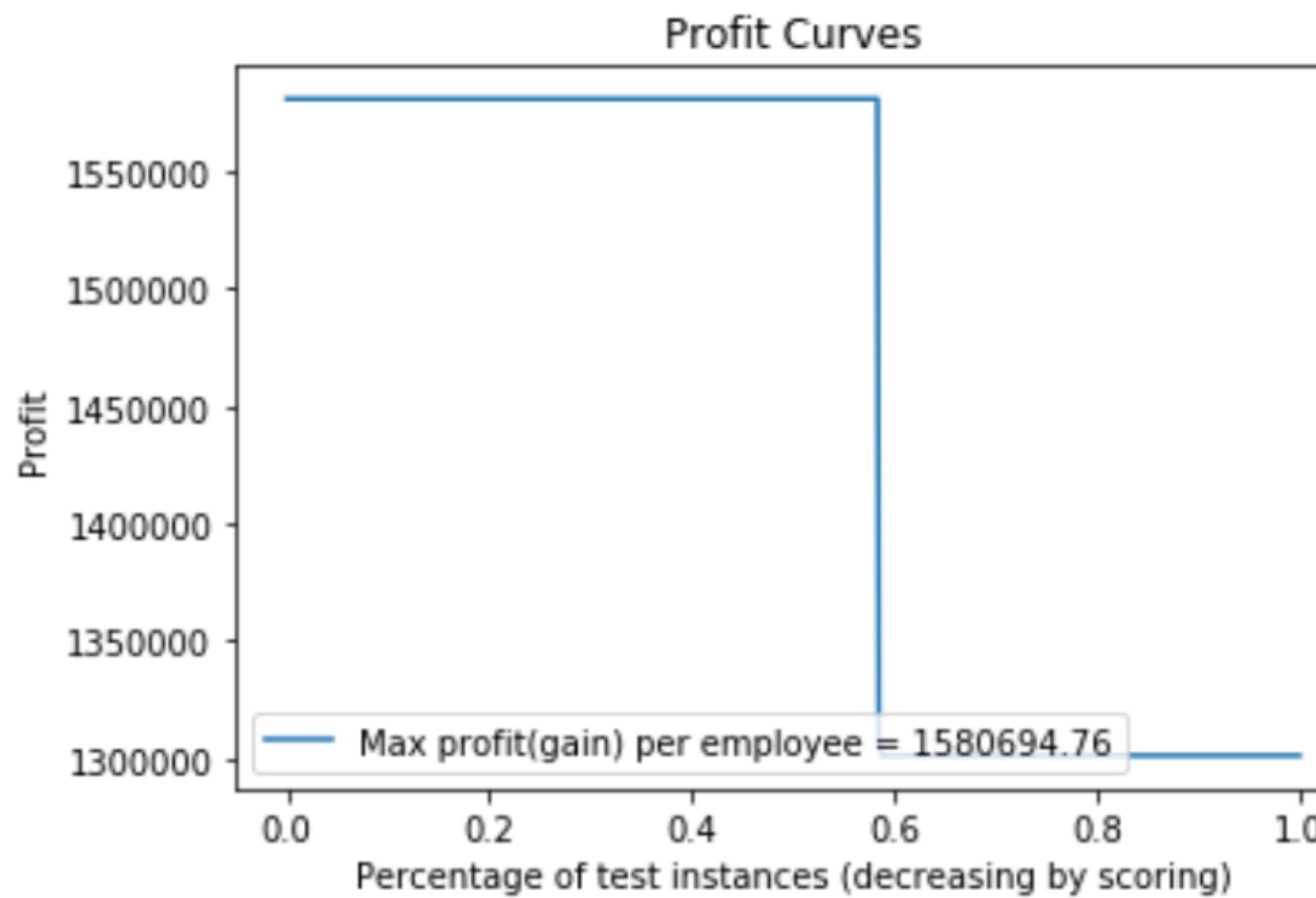
ROC CURVE



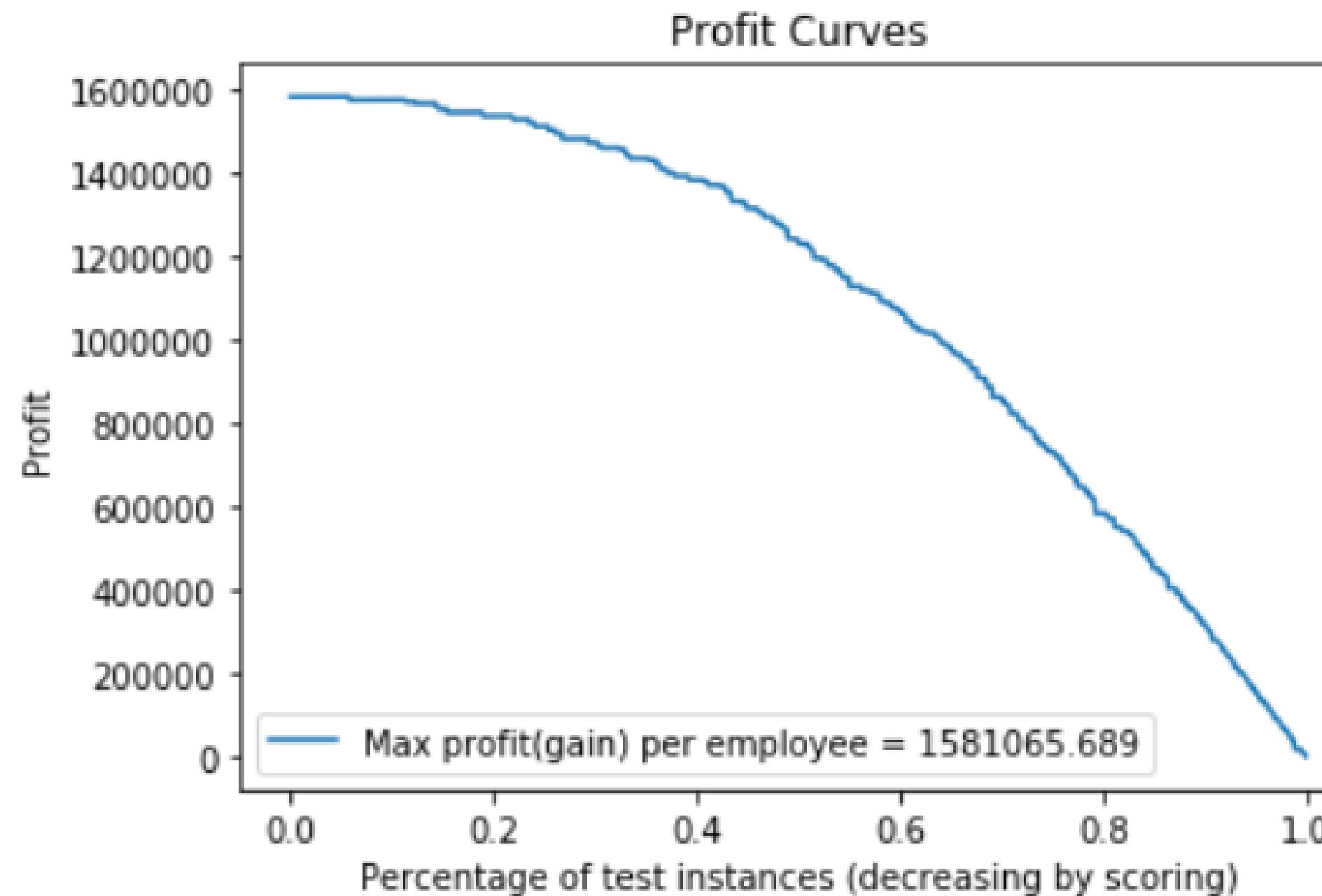
COST BENEFIT MATRIX

		p	n	
		Y	3,167,772.95	-6383.43
		N	0	0
Salary promotion offered to employee, and employee accepts				Salary promotion was offered (one month salary promotion has been awarded to employee) but employee does not accept
Salary promotion not offered but employee would have accepted if it was offered				Salary promotion not offered and employee wouldn't have accepted if it was offered

PROFIT CURVE - DECISION TREE

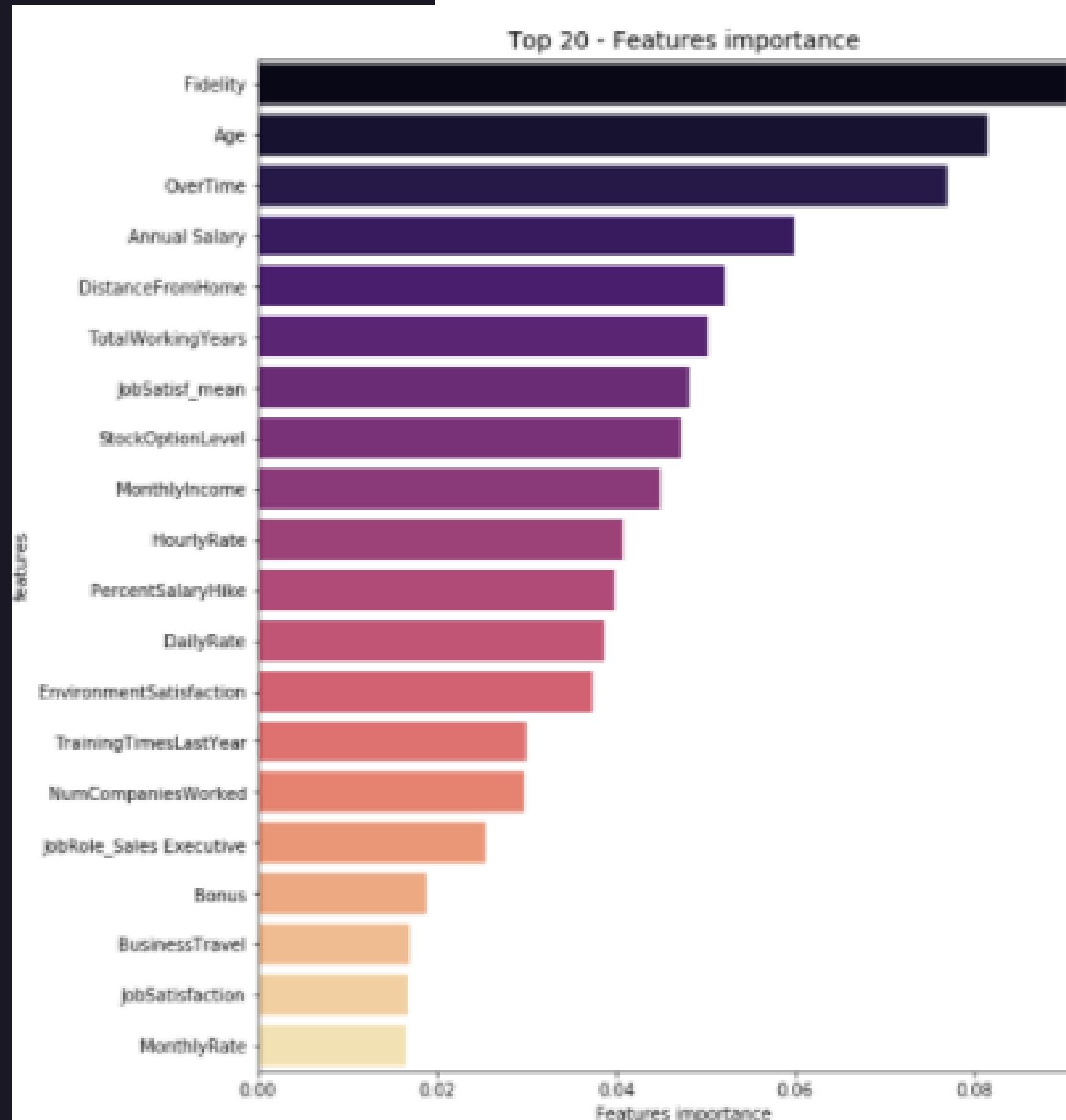


PROFIT CURVE - LOGISTIC REGRESSION



FEATURE IMPORTANCE

In order to know which features influence attrition more compared to others features. We used the (feature_importances) to identify the Top 20 essential features.



KEY TAKEAWAYS FOR HR

	Attrition	No Attrition
Fidelity	Lower	Higher
Annual Salary	Very low and mid-high	Mid and Very high
Age	27- 32	33 - 40
Performance of Job Satisfaction	Low and mid	High
Distance From Home	Far and near	Mid and near
Stock Option Level (4 is the highest)	1 and 3	2 and 4

MODEL DEPLOYMENT

Ways that HR can utilize the model:

- Most predictive attributes/key takeaways
- Actionable, intentional follow-up after Employee Pulse Surveys and Performance Evaluations
- Proactively plan for turnover

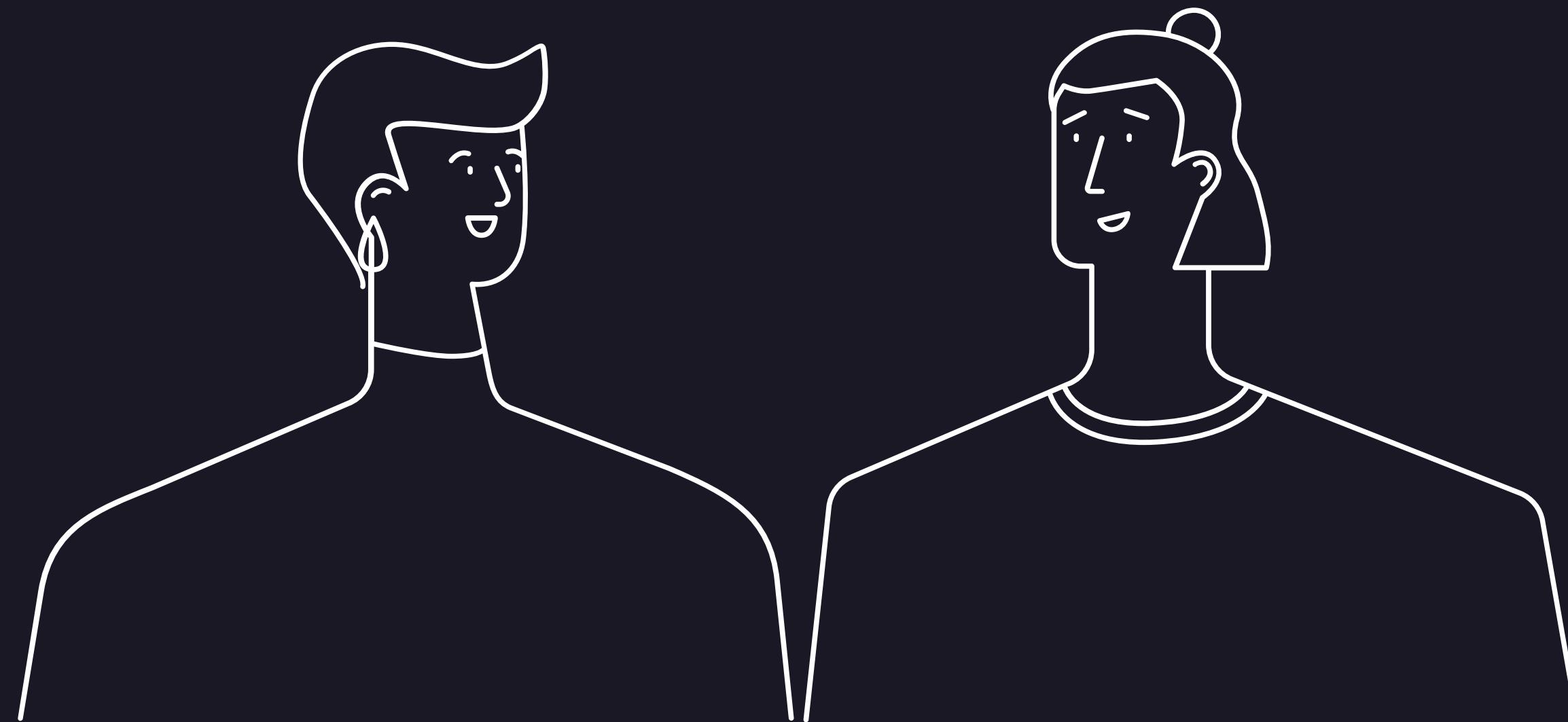
A word of caution:

HR folks need to be aware that although the model is informative, it is not 100% accurate; it is important that those taking action based on model results are mindful to not treat employees differently based on what the model predicts about their future choices.

ETHICS

When it comes to building an ethical model, we want to be particularly mindful of variables that could be used in a discriminatory fashion. In the case of the data provided to us, variables that we may consider removing in favor of a more ethical model are:

- Age
- Gender
- TotalWorkingYears (as a proxy for Age)

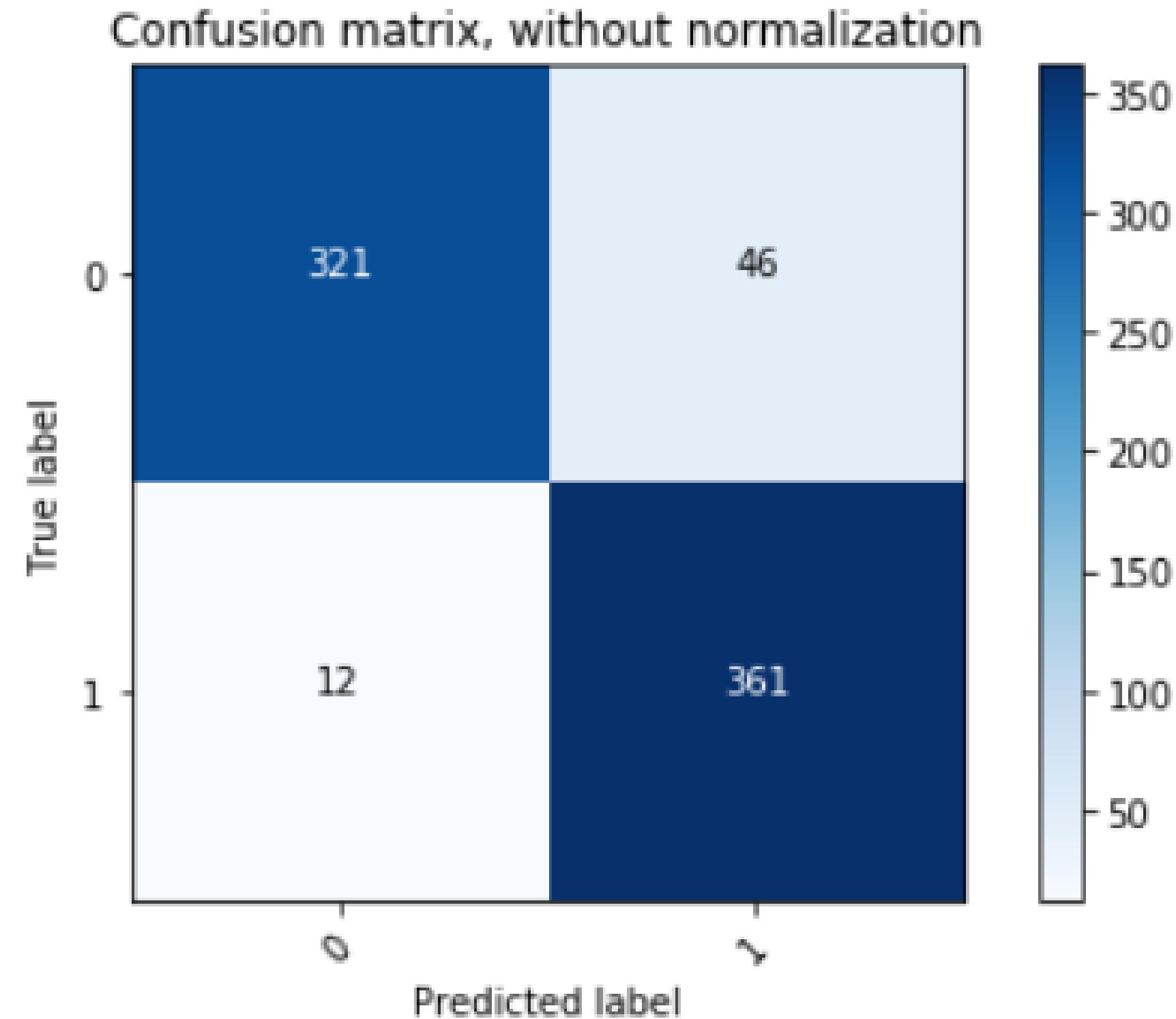




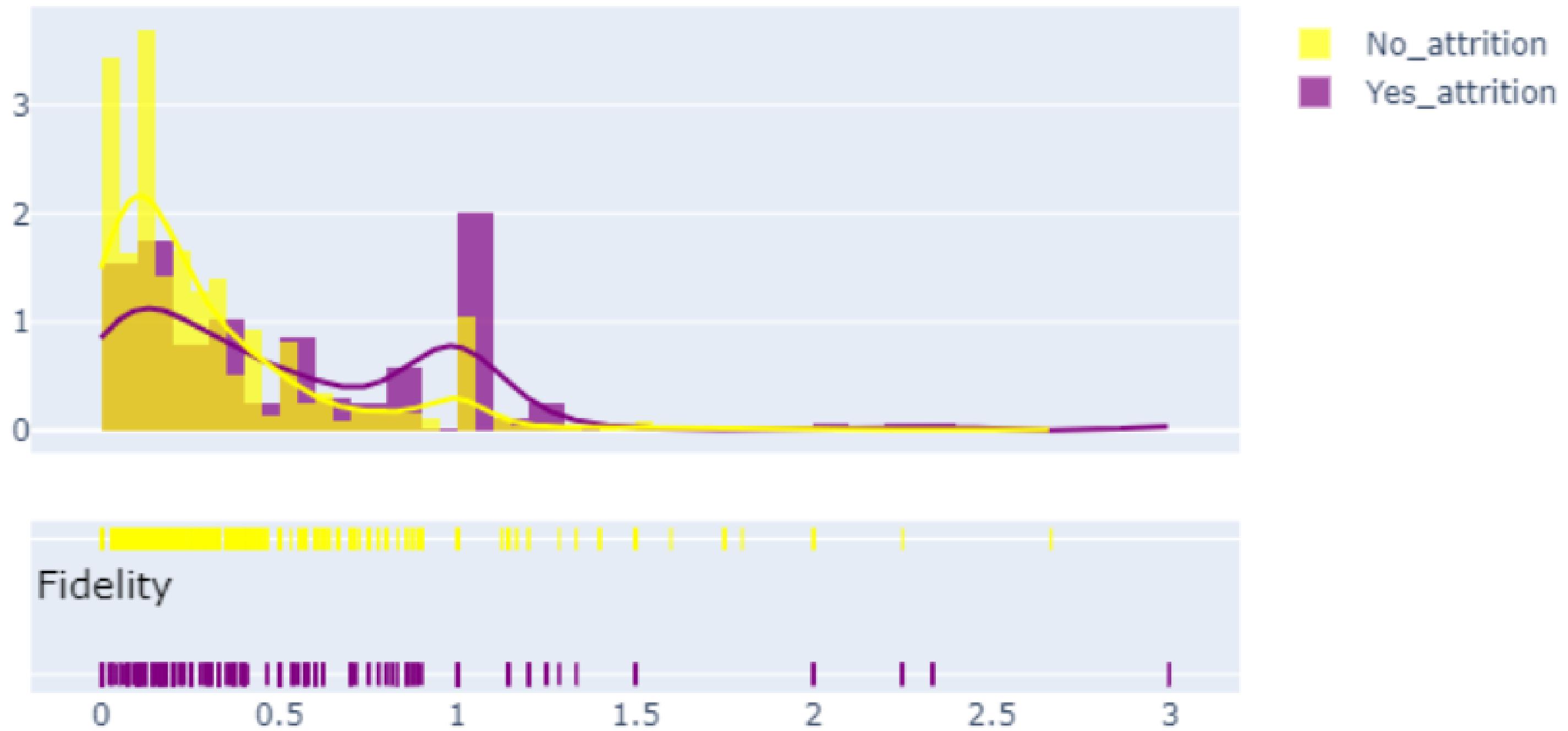
THANK YOU

APPENDIX

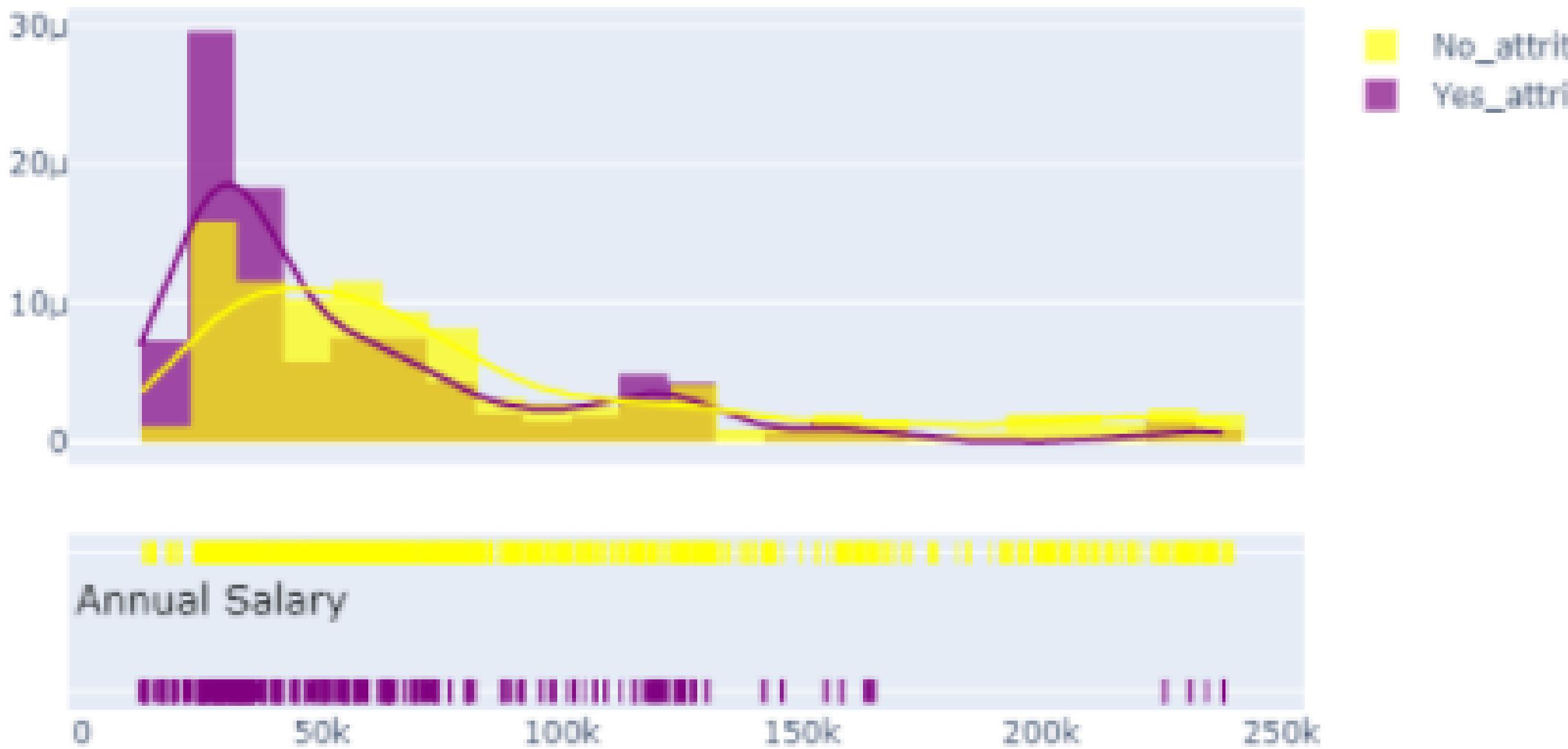
CONFUSION MATRIX



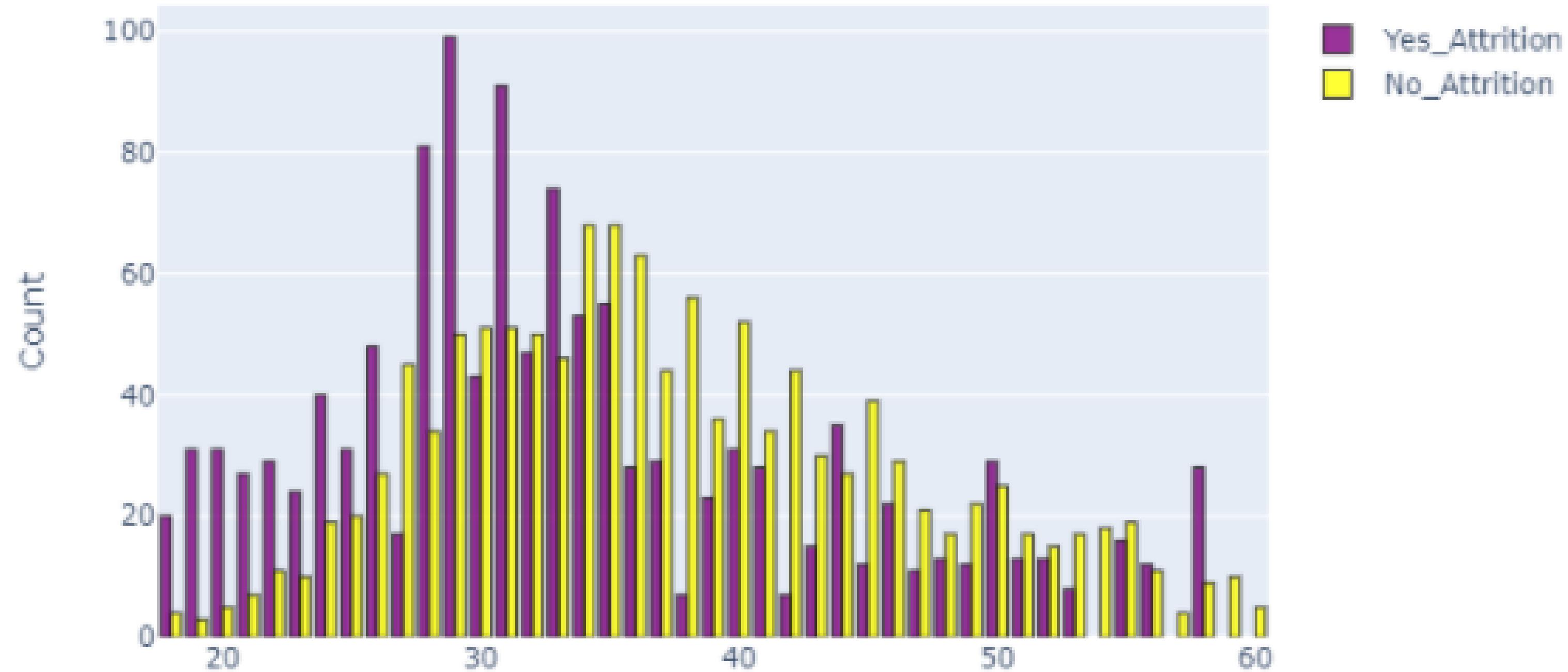
DISTRIBUTION OF FIDELITY (BASED ON WEIGHT)



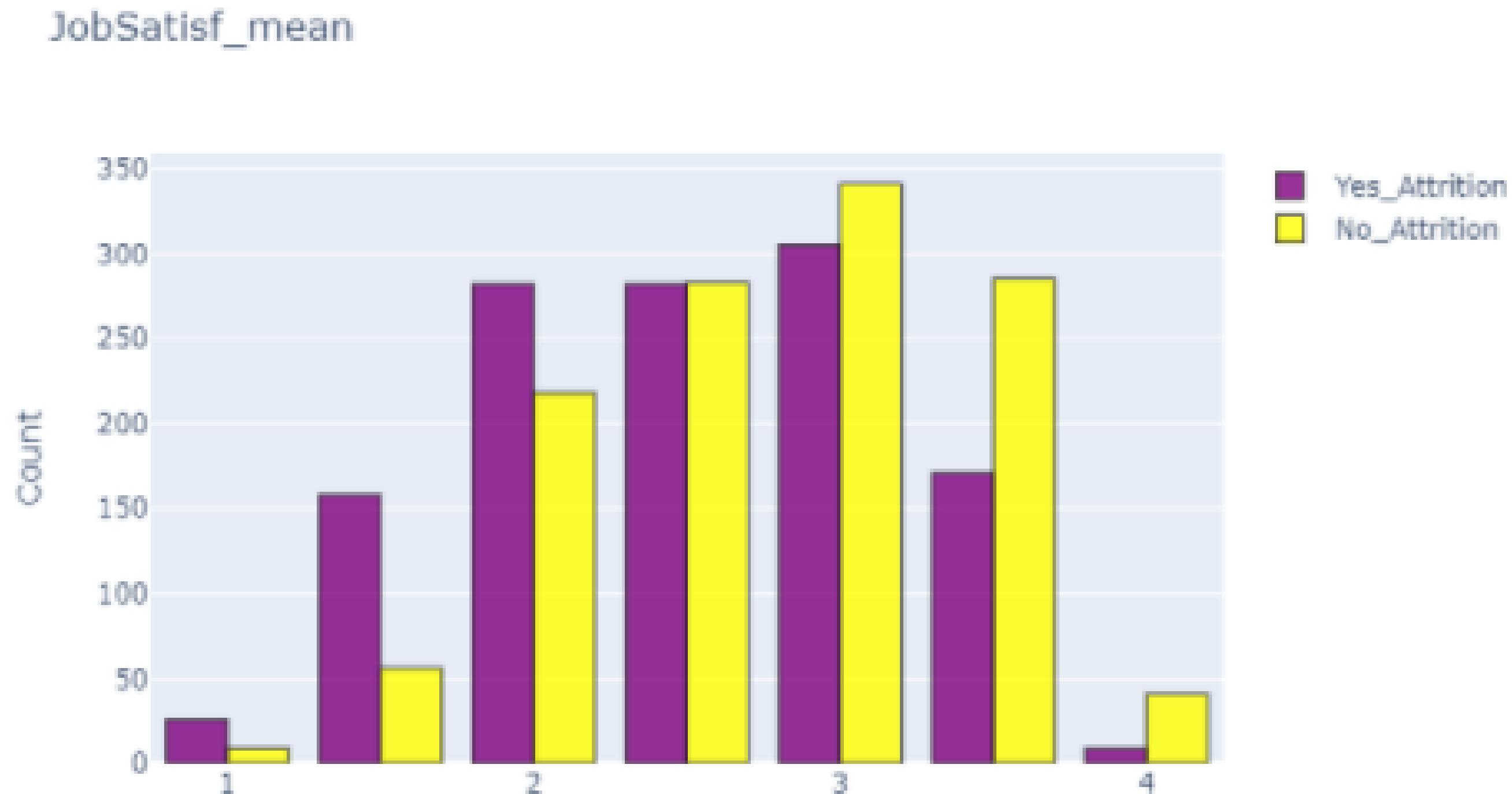
DISTRIBUTION OF ANNUAL SALARY (BASED ON WEIGHT)



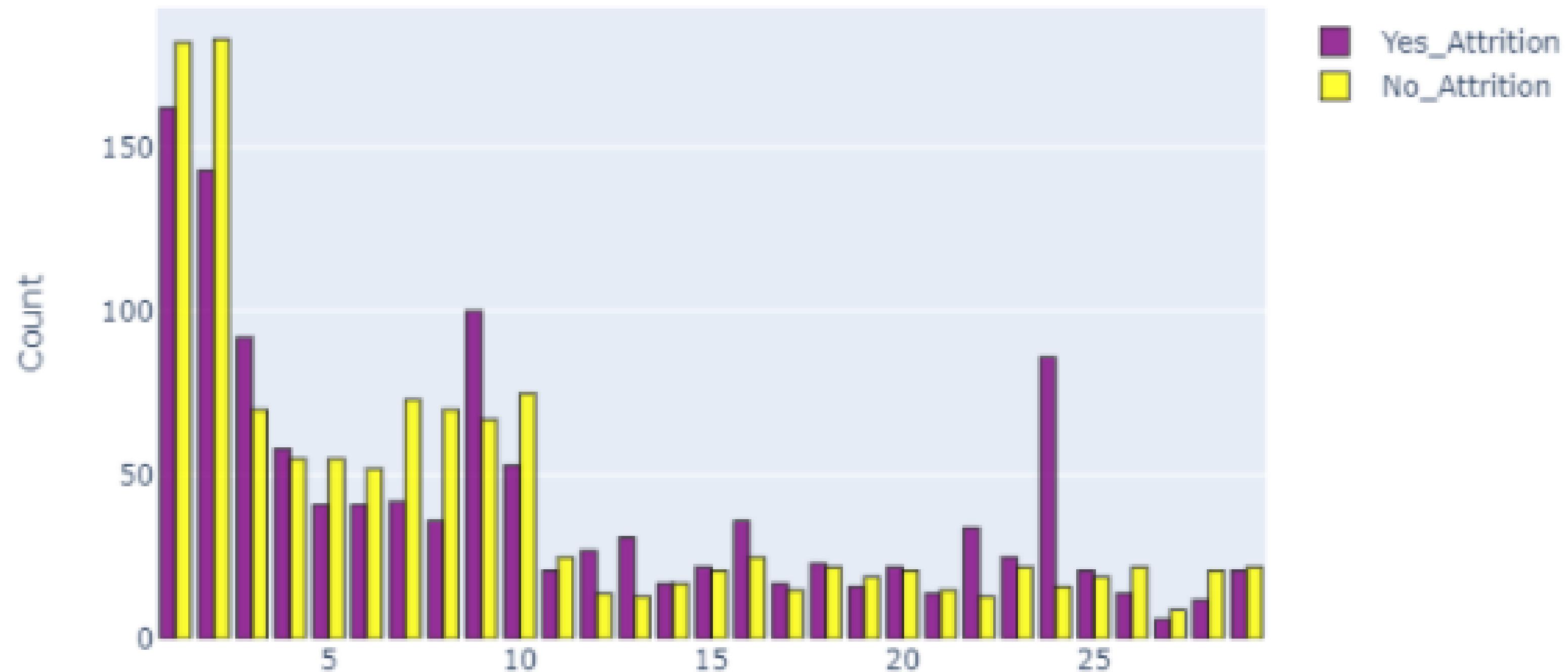
COUNT OF ATTRITION SEGMENTED BY AGE



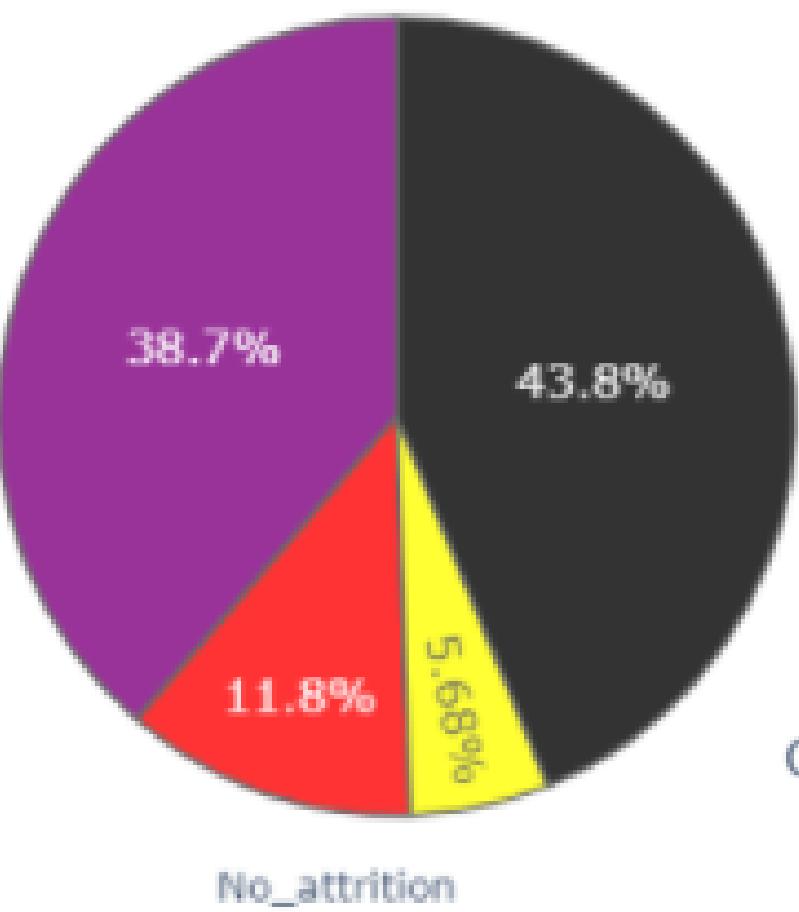
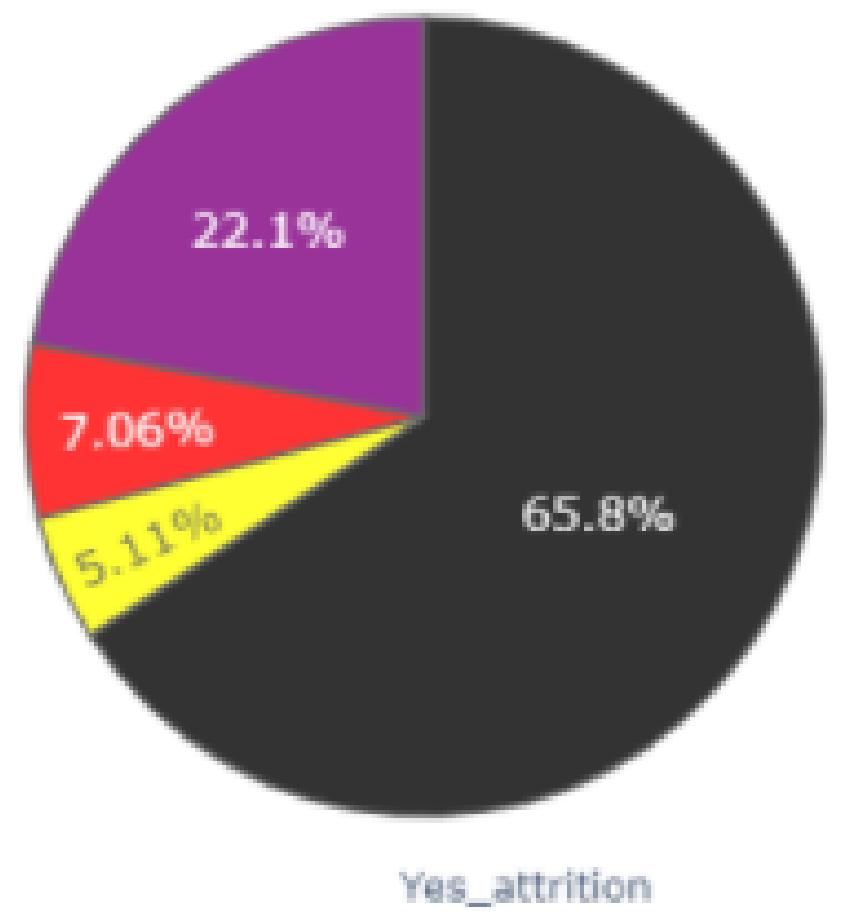
COUNT OF ATTRITION SEGMENTED BY JOB SATISFACTION AND JOB INVOLVEMENT



COUNT OF ATTRITION SEGMENTED BY DISTANCE FROM HOME

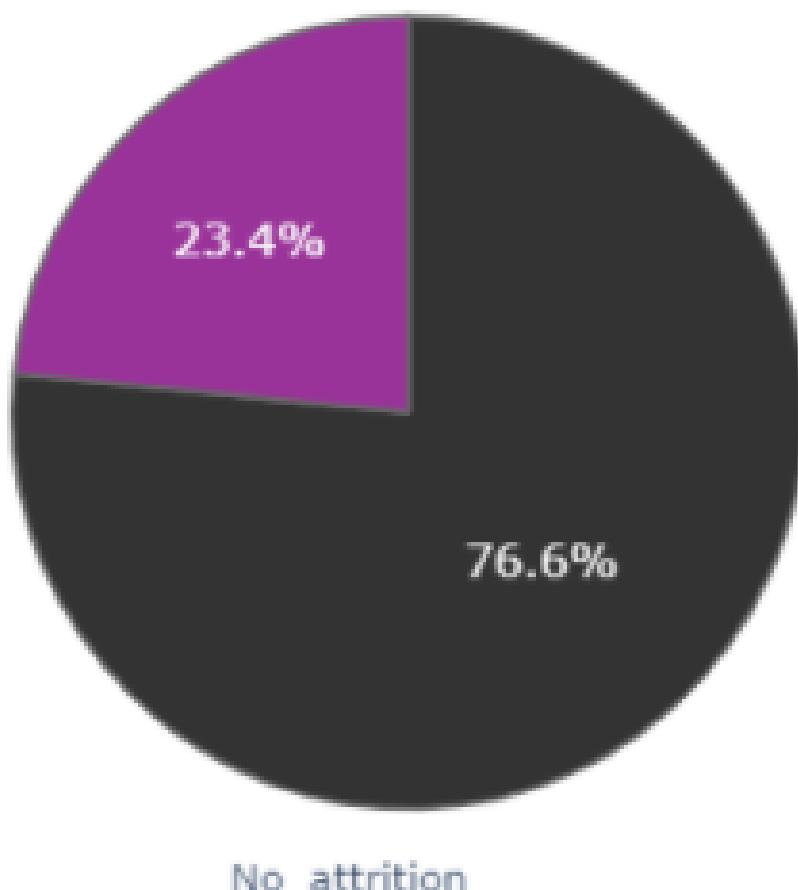
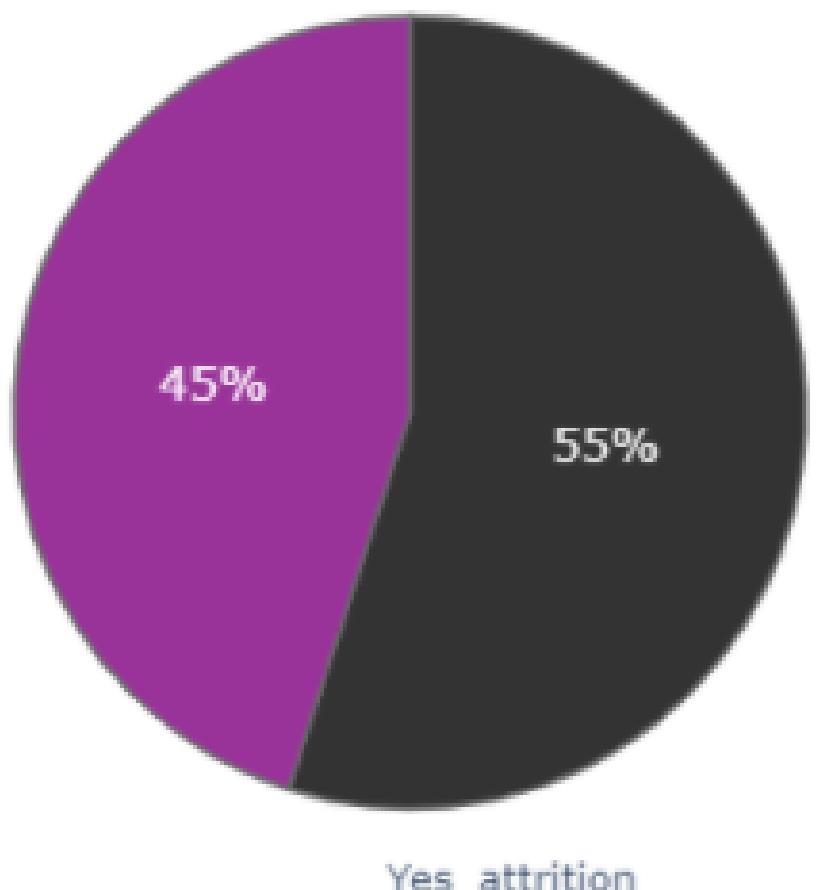


StockOptionLevel distribution in employees attrition



0
1
3
2

OverTime distribution in employees attrition



1
0