# Capstone Project

# Analysis on Engineering Graduate Salary

**Group-6**
**Team Members**

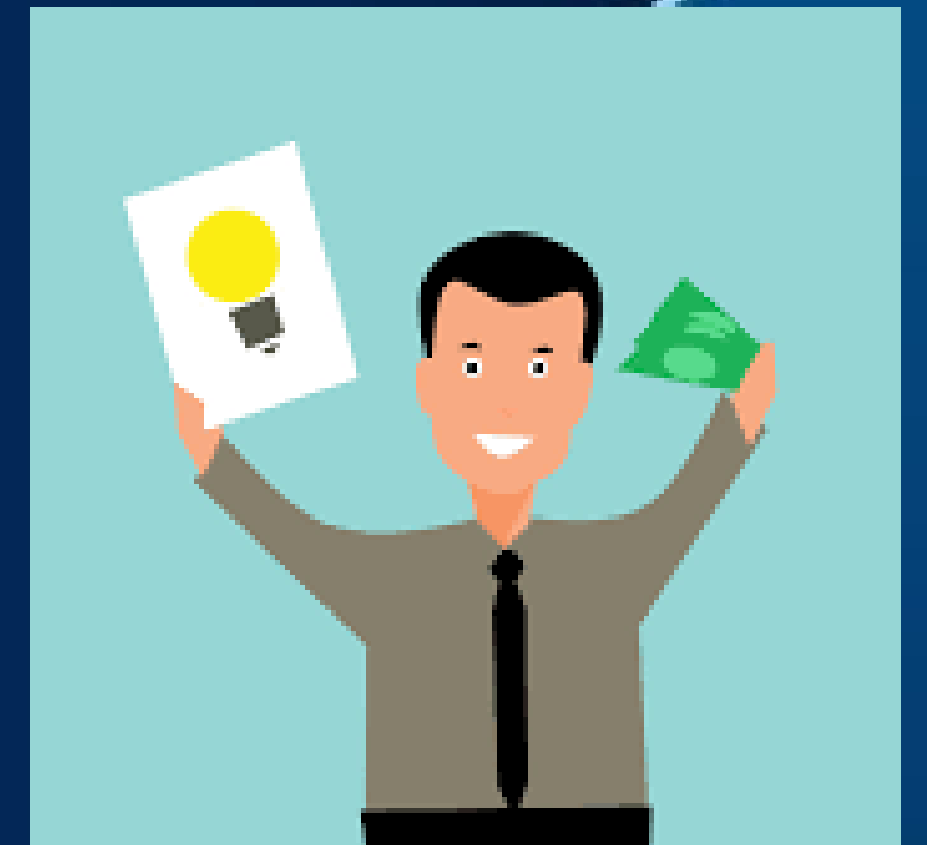1.UudhhayKiirran

2.Bhavana Reddy

3.Jashwanth

4.SaiPrasanna

# Introduction

**Dataset : Engineering_Graduate Salary**

**Objective : Our objective is to determine the salary of an engineering graduate .**

**Technical contents:**

- **Data importing**
- **Data exploration**
- **Data Preprocessing**
- **Data Modelling**

# Data Imported:

| ID | Gender | DOB | 10percentage | 10board | 12graduation | 12percentage | 12board | CollegeID | CollegeTier | Degree | ... | MechanicalEngg | ElectricalEngg | TelecomEngg | CivilEngg | conscientiousness | agreeableness | extraversion | nueroticism | openess_to_experience | Salary |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 604399 | f | 22-10-1990 | 87.80 | cbse | 2009 | 84.00 | cbse | 6920 | 1 | B.Tech/B.E. | ... | -1 | -1 | -1 | -1 | -0.1590 | 0.3789 | 1.2396 | 0.14590 | 0.2889 | 445000 |
| 988334 | m | 15-05-1990 | 57.00 | cbse | 2010 | 64.50 | cbse | 6624 | 2 | B.Tech/B.E. | ... | -1 | -1 | -1 | -1 | 1.1336 | 0.0459 | 1.2396 | 0.52620 | -0.2859 | 110000 |
| 301647 | m | 21-08-1989 | 77.33 | maharashtra state board,pune | 2007 | 85.17 | amravati divisional board | 9084 | 2 | B.Tech/B.E. | ... | -1 | -1 | 260 | -1 | 0.5100 | -0.1232 | 1.5428 | -0.29020 | -0.2875 | 255000 |
| 582313 | m | 04-05-1991 | 84.30 | cbse | 2009 | 86.00 | cbse | 8195 | 1 | B.Tech/B.E. | ... | -1 | -1 | -1 | -1 | -0.4463 | 0.2124 | 0.3174 | 0.27270 | 0.4805 | 420000 |
| 339001 | f | 30-10-1990 | 82.00 | cbse | 2008 | 75.00 | cbse | 4889 | 2 | B.Tech/B.E. | ... | -1 | -1 | -1 | -1 | -1.4992 | -0.7473 | -1.0697 | 0.06223 | 0.1864 | 200000 |

- This dataset contains 2,998 rows and 33 columns.
- The target variable is salary.
- This dataset contians numerical values.
- The columns of the dataset describe about Engineering graduates.

**Checking for Null values**

```
[5]  data.isnull().sum()

     Gender                 0
     DOB                    0
     10percentage           0
     10board                0
     12graduation           0
     12percentage           0
     12board                0
     CollegeID              0
     CollegeTier            0
     Degree                 0
     Specialization         0
     collegeGPA             0
     CollegeCityID          0
     CollegeCityTier        0
     CollegeState           0
     GraduationYear         0
     English                0
     Logical                0
     Quant                  0
     Domain                 0
     ComputerProgramming    0
     ElectronicsAndSemicon  0
     ComputerScience        0
     MechanicalEngg         0
     ElectricalEngg         0
     TelecomEngg            0
     CivilEngg              0
     conscientiousness      0
     agreeableness          0
     extraversion           0
     nueroticism            0
     openess_to_experience  0
     Salary                 0
     dtype: int64
```
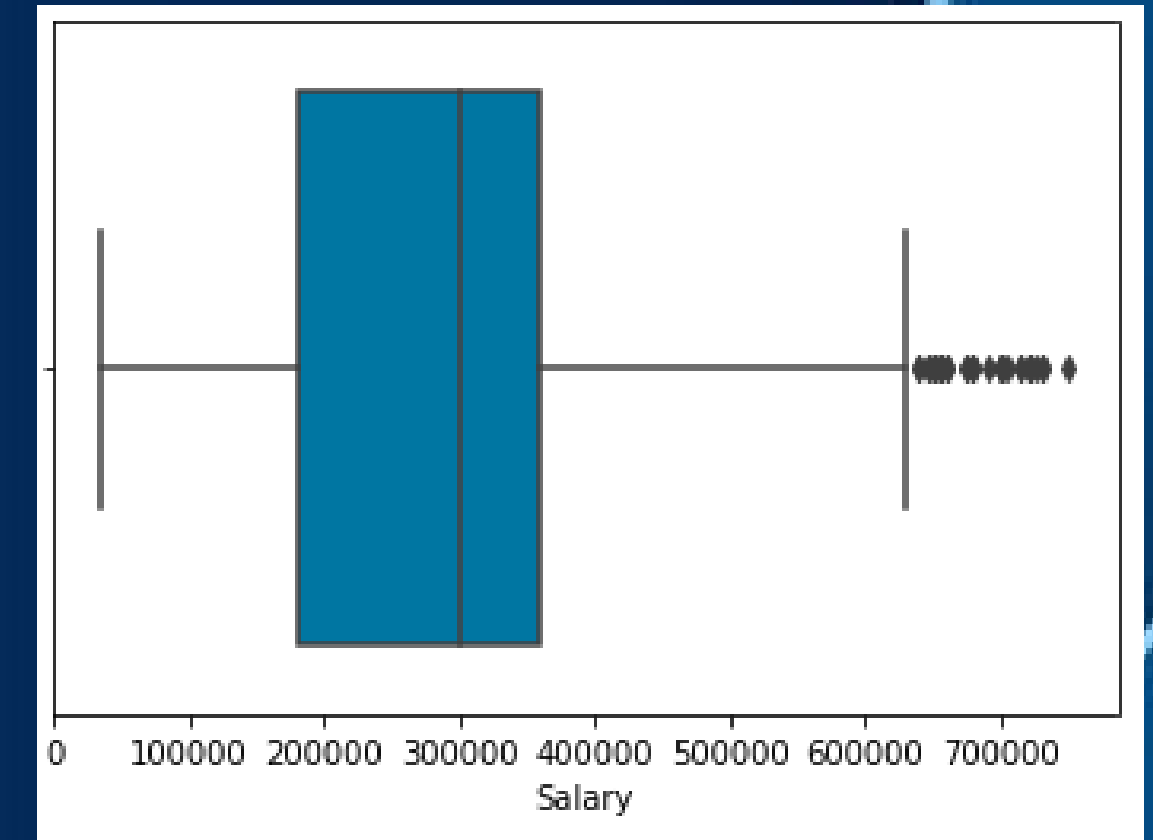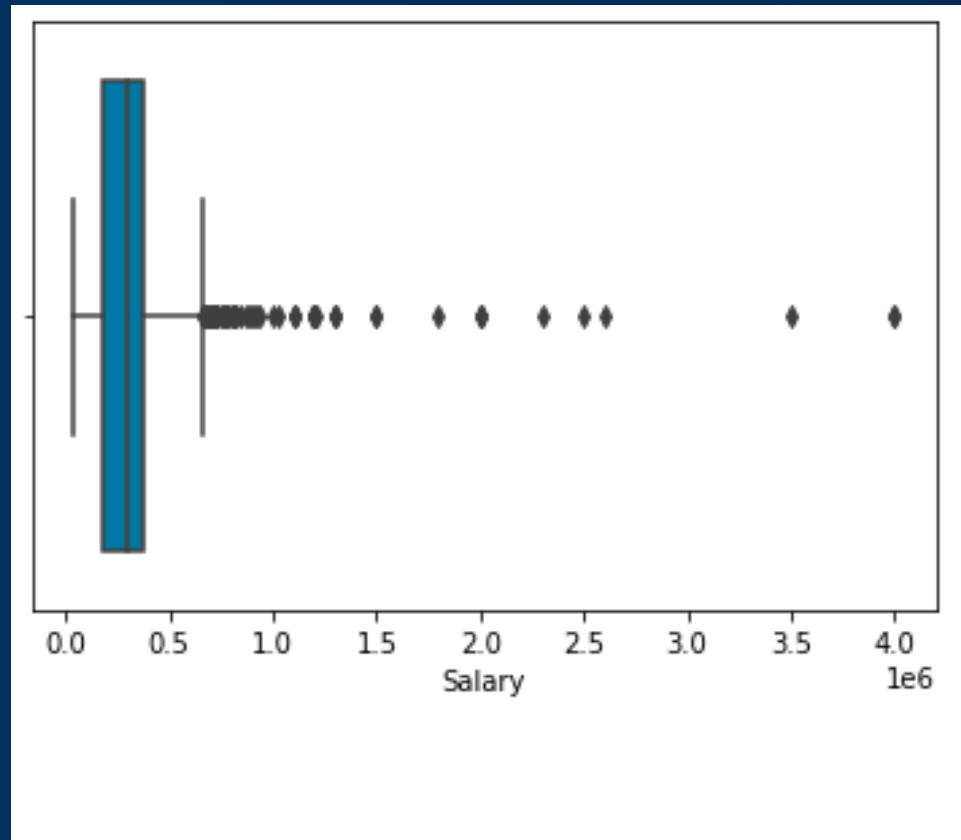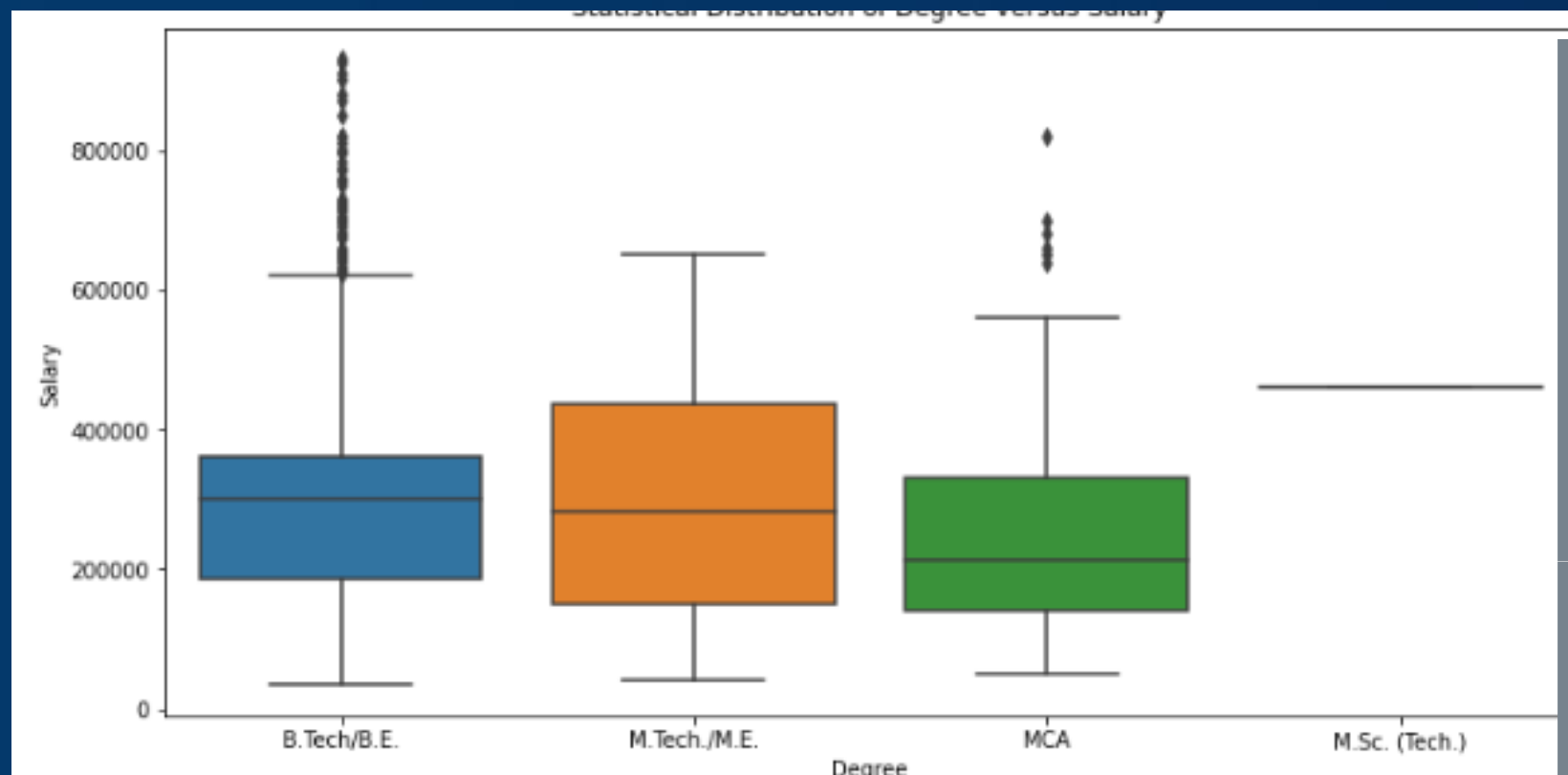
- The first box plot represents the outlers in dataset.
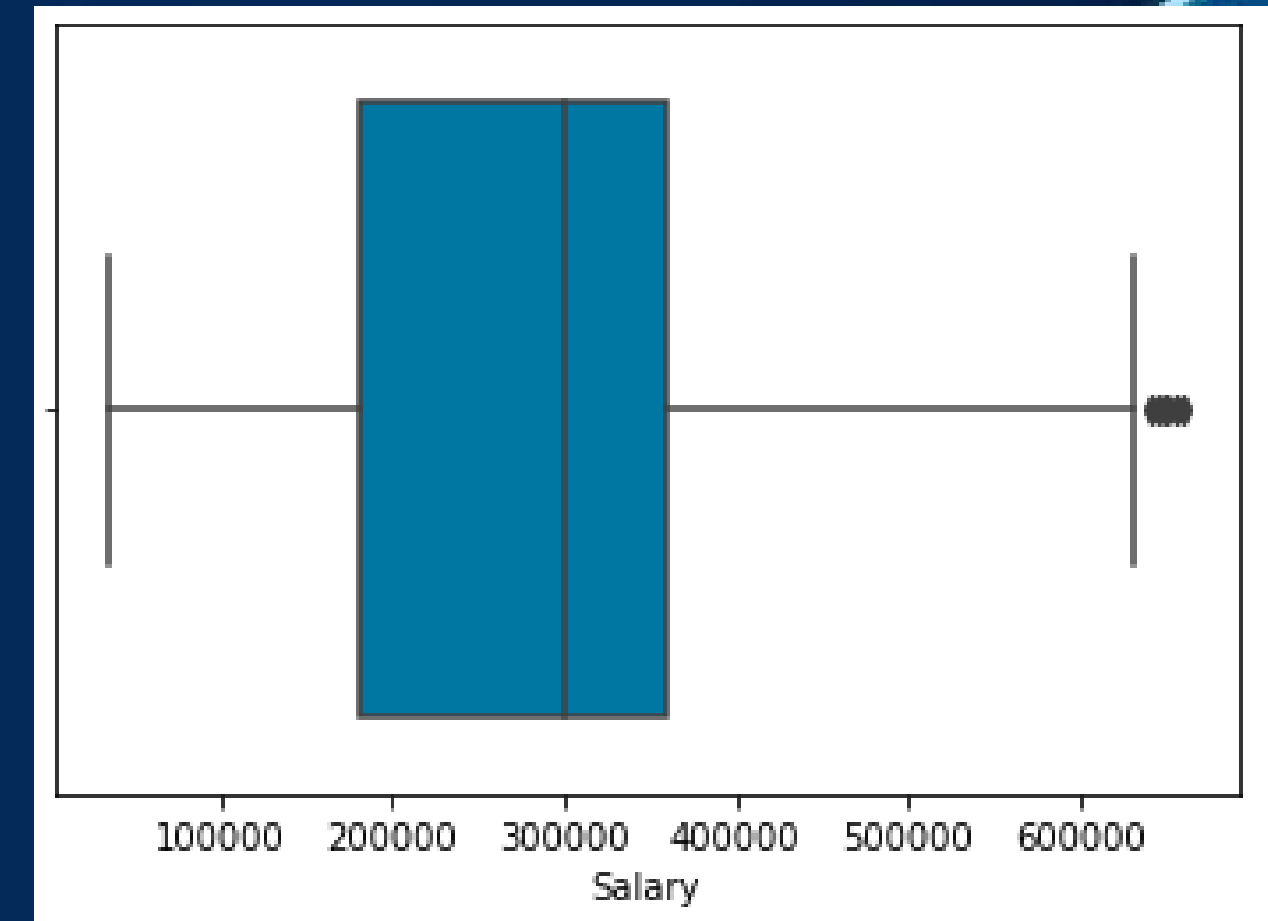- We have used hample method to remove the outliers.

We used again zscore method to decrease the number of outliers.



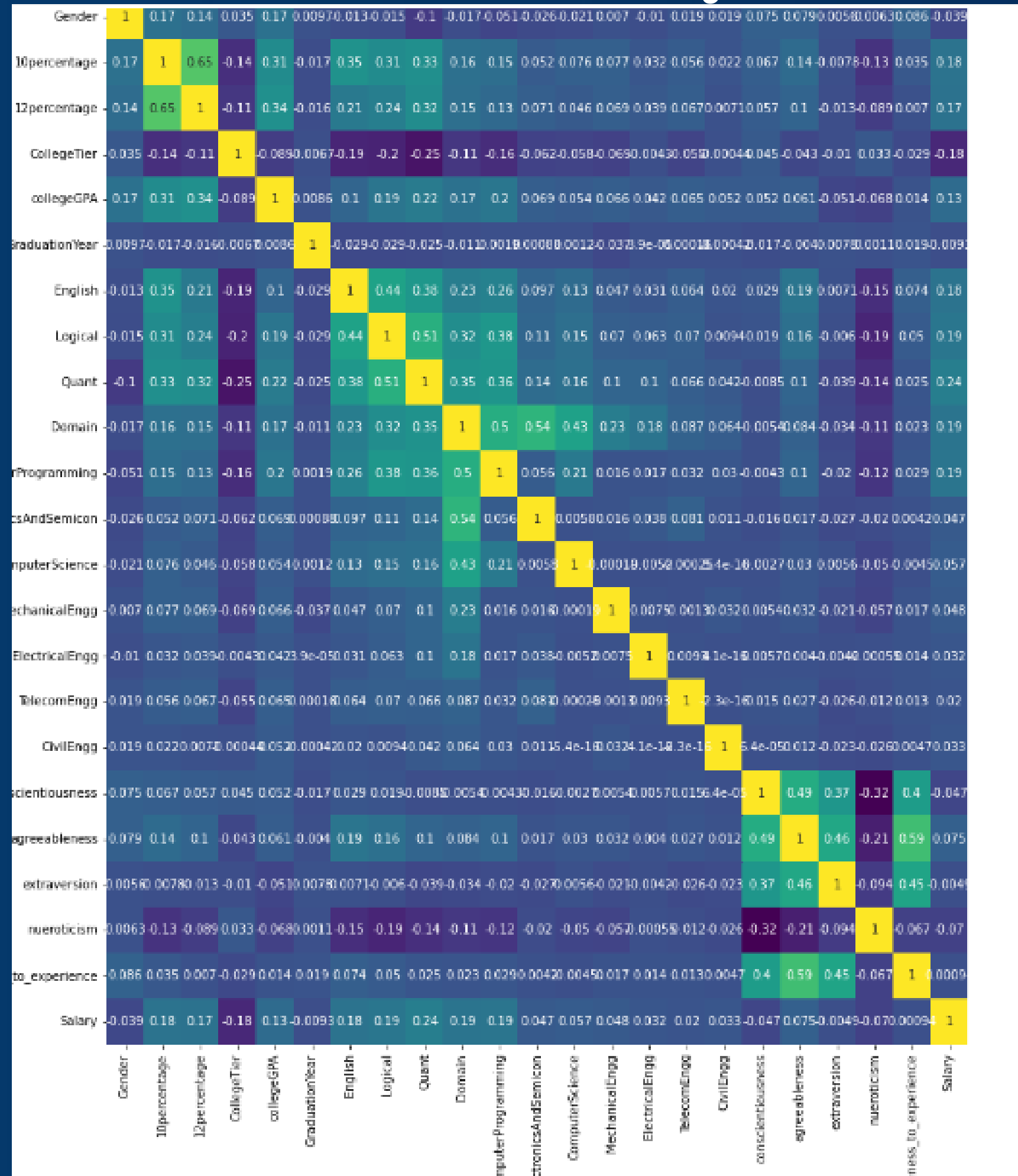# Exploratory Data Analysis:

## Analysis of Salary with Degree

## Analysis of Salary with Gender



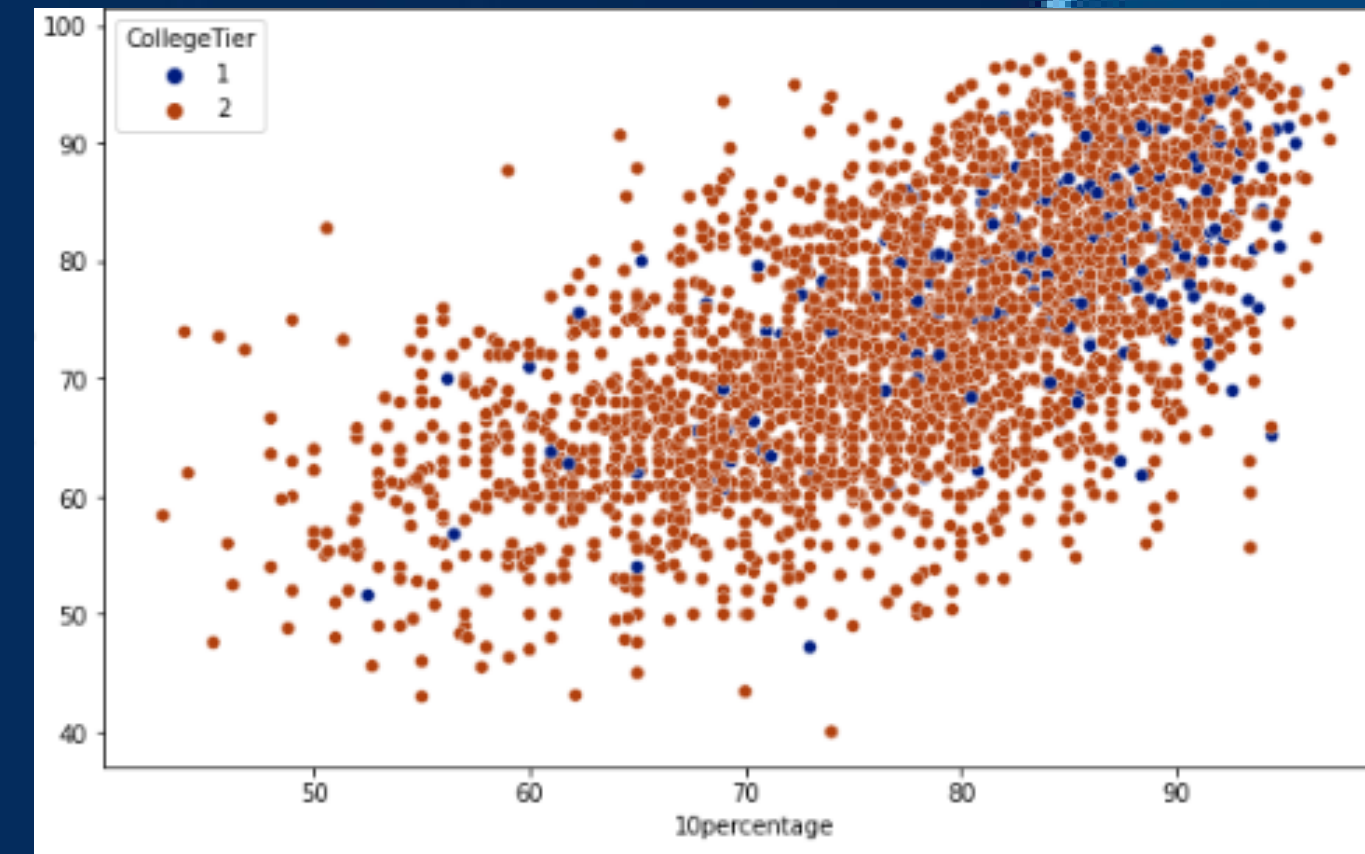Average salary is highest for BE/B.tech graduates as compared to any other degree graduates.
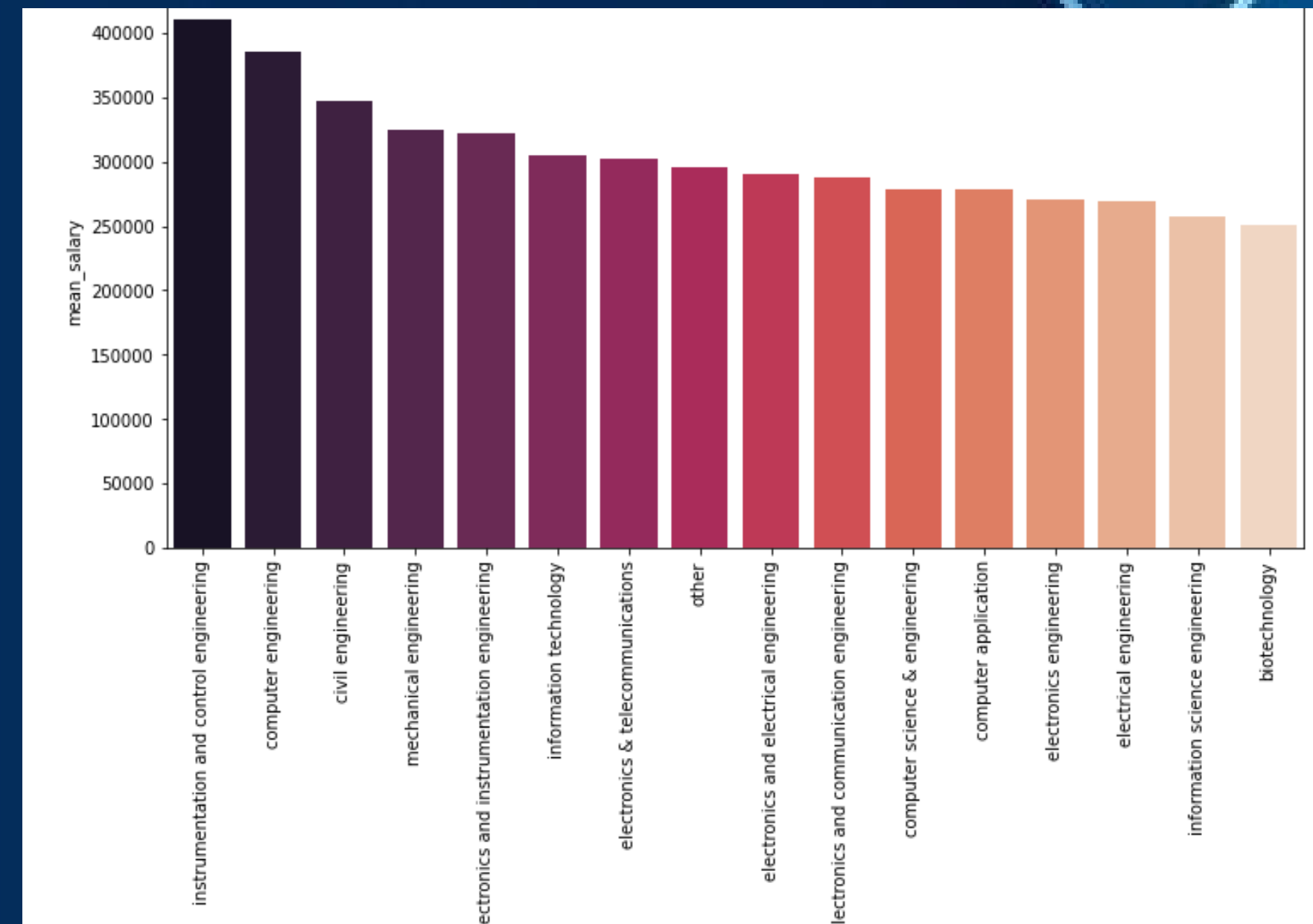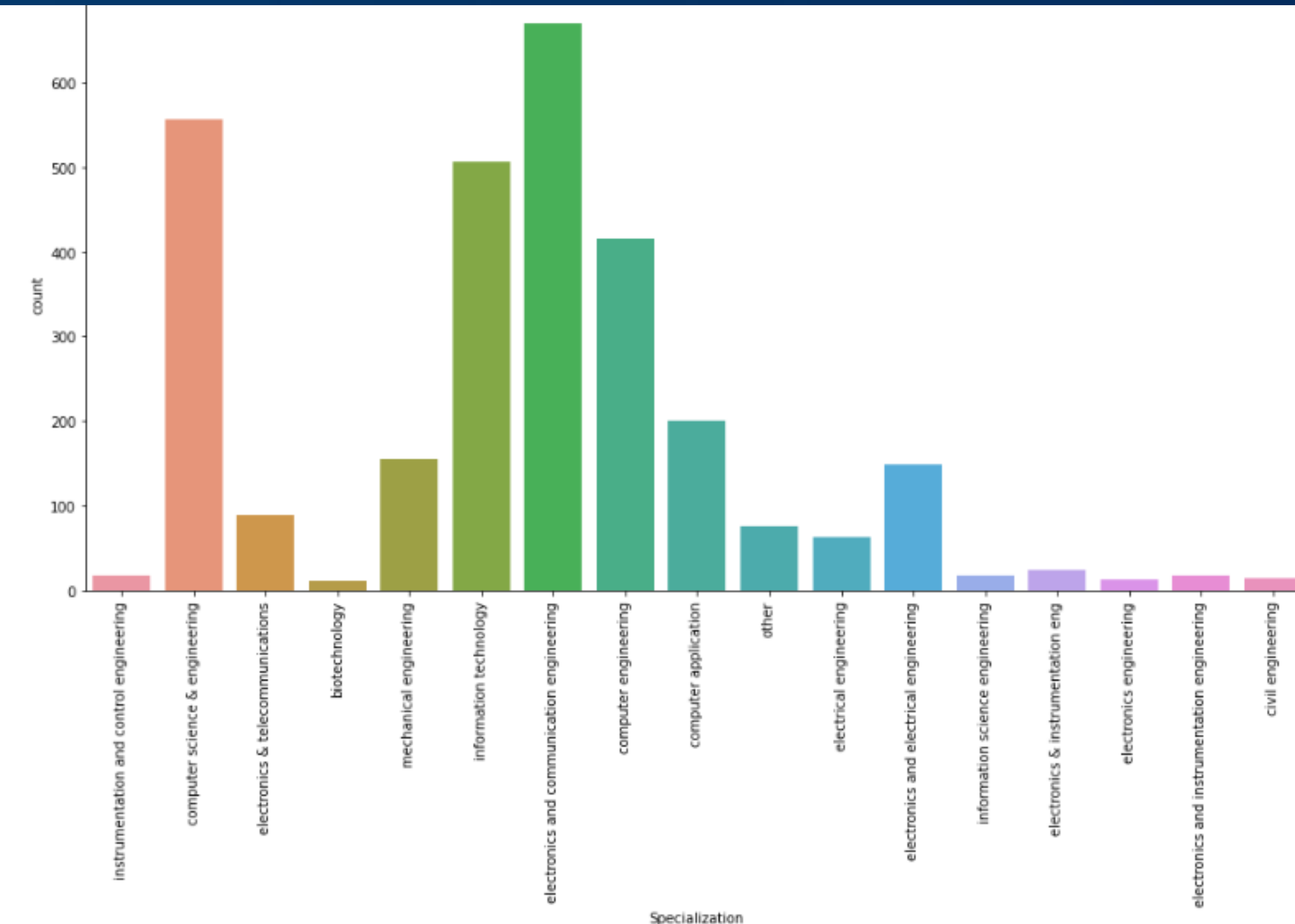
# Correlation Analysis



## Analysis of 10th and 12th percentage by collage tier



According to correlation and scatterplot we can see 10th and 12th are positively correlated and this is the case of multicollinearity.
So we decide to keep only one

# Analysation of Salary with Specialisation



ICE Engineer, Computer Engineer and Electronics Engineer having highest mean salary

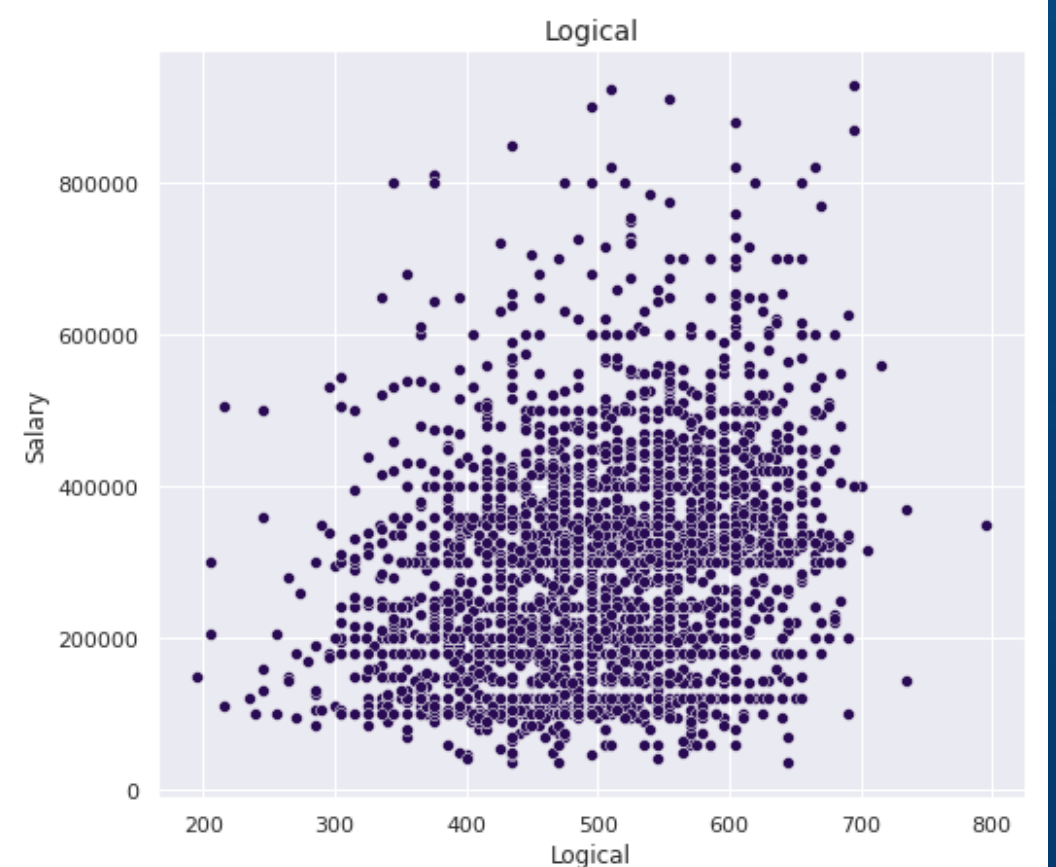# Analysis of Collage GPA and Salary

## Analysis of AMCAT subjects and Salary



We have to remove the outliers present as we can view them in the plot .

**Here the on the basis of the above data visualised we observe that gender,degree, collage GPA,Specialisation are useful for analysing meaningful insight for the variable salary.**

We have visulaised the Salary relation with AMCAT subjects through scatterplot above .

As AMCAT exams are conducted as an entrance for the jobs  and it plays a  key role  but we found from the above graphs there is no impact on the  the salary of the graduate.

# Linear regression

Linear regression analysis is used to predict the value of a variable based on the value of another variable.

- For 60-40 ratio we got the least MAE value i.e: 91973.1911

| Linear Regression | |
| --- | --- |
| **Train-Test Proportion** | **MAE** |
| 70-30 | 137905.4249 |
| 80-20 | 93409.354 |
| 60-40 | 91973.1911 |
| 75-25 | 93001.6601 |

# Neural Networks

In the neural networks we are working with respect to two optimisers : Adam ,SGD So we applied it with various experiments and architectures

We find for 70-30 ratioat 20--10--5--1 architecture and Adam as optimiser at 2000epoch we have least MAE value: 91597.9844

| Train-Test Proportion | Architecture | Optimizer | epochs | MAE |
|---|---|---|---|---|
| 80-20 | 24--12--6--3--1 | Adam | 500 | 94020.5156 |
| 80-20 | 24--12--8--4--1 | Adam | 600 | 94340.7344 |
| 80-20 | 24--12--8--4--2 | SGD | 600 | nan |
| 80-20 | 24--10--5--1 | Adam | 400 | 93993.6484 |
| `80-20 | 24--10--5--2 | SGD | 400 | nan |
| 80-20 | 24--8--3--1 | Adam | 1000 | 93821.4453 |
| 80-20 | 24--8--3--2 | SGD | 1000 | nan |
| 80-20 | 18--10--5--1 | Adam | 800 | 93947.5234 |
| 80-20 | 18--10--5--2 | SGD | 800 | nan |
| 80-20 | 15--7--1 | Adam | 400 | 93815.2109 |
| 80-20 | 15--7--1 | SGD | 400 | nan |
| 70-30 | 8--5--3--1 | Adam | 500 | 91953.3516 |
| 70-30 | 8--5--2--1 | SGD | 400 | nan |
| 70-30 | 20--10--5--1 | Adam | 2000 | 91597.9844 |
| 70-30 | 20--10--5--1 | SGD | 1000 | nan |
| 70-30 | 6--4--2--1 | Adam | 3000 | 91706.4141 |
| 70-30 | 6--4--2--1 | Adam | 1500 | 93134.6562 |
| 70-30 | 6--4--2--1 | Adam | 600 | 91989.6328 |
| 70-30 | 12--6--3--1 | Adam | 600 | 91960.9766 |
| 70-30 | 6--3--1 | Adam | 400 | 92000.8594 |
| 60-40 | 8--4--2--1 | Adam | 60 | 93275.0938 |
| 80-20 | 10--6--3--1 | Adam | 60 | 94089.2656 |
| 80-20 | 8--6--3--1 | SGD | 60 | nan |
| 70-30 | 6--4--2--1 | Adam | 600 | 91911.7812 |
| 70-30 | 24--12--8--1 | Adam | 500 | 91876.1797 |
| 70-30 | 8--5--1 | Adam | 800 | 91968.3516 |
| 70-30 | 5--2--1 | Adam | 400 | 291649.625 |
| 60-40 | 5--2--1 | Adam | 1000 | 92912.625 |
| 60-40 | 10--3--1 | Adam | 800 | 92916.8906 |
| 60-40 | 20--10--1 | Adam | 800 | 92865.4688 |
| 60-40 | 24--18--11--5--2--1 | Adam | 800 | 92603.4453 |
| 60-40 | 24--20--15--10--5--3--1 | Adam | 1200 | 92270.1953 |

# Bagging

## Bagging Regressor

For 70-30 ratio we got the least MAE value i.e: 92739.88662

| Bagging Regressor | |
|---|---|
| Train-Test Proportion | MAE |
| 80-20 | 94824.04762 |
| 70-30 | 92739.88662 |
| 60-40 | 93150.40816 |
| 75-25 | 93758.92517 |

## Decision Tree

For 75-25 ratio we got the least MAE value i.e: 95019.70234

| Decision Tree | |
|---|---|
| Train-Test Proportion | MAE |
| 80-20 | 95177.45818 |
| 70-30 | 95133.15039 |
| 60-40 | 96296.33257 |
| 75-25 | 95019.70234 |

## Random Forest

For 70-30 ratio we got the least MAE value i.e: 92889.95465

## Boosting

### Adaboost

For 70-30 ratio we got the least MAE value i.e: 92201.85737

**Random Forest**

| Train-Test Proportion | MAE |
|---|---|
| 60-40 | 93183.97109 |
| 80-20 | 94383.46939 |
| 70-30 | 92889.95465 |
| 75-25 | 94236.12245 |

**Adaboosting**

| Train-Test Proportion | MAE |
|---|---|
| 60-40 | 92492.67825 |
| 75-25 | 93917.87464 |
| 80-20 | 93830.78736 |
| 70-30 | 92201.85737 |

## Gradient Boost

For 80-20 ratio we got the least MAE value i.e: 91502.9275

**Gradient Boosting**

| Train-Test Proportion | MAE |
|---|---|
| 80-20 | 91502.9275 |
| 70-30 | 91654.51334 |
| 60-40 | 93531.56278 |
| 75-25 | 92380.45955 |

## XGboost

For 70-30 ratio we got the least MAE value i.e: 90676.26513

**XGboosting**

| Train-Test Proportion | MAE |
|---|---|
| 80-20 | 91894.45148 |
| 70-30 | 90676.26513 |
| 60-40 | 92553.95049 |
| 75-25 | 92353.62725 |

# Comparision Table

| Algorithm | MAE | Ratio |
|---|---|---|
| XGboosting | 90676.26513 | 80-20 |
| Gradient Boosting | 91502.9275 | 80-20 |
| Neural Networks | 91597.9844 | 70-30 |
| Linear Regression | 91973.1911 | 60-40 |
| Adaboosting | 92201.85737 | 70-30 |
| Bagging Regressor | 92739.88662 | 70-30 |
| Random Forest | 92889.95465 | 70-30 |
| Decision Tree | 95019.70234 | 75-25 |

# Conclusion

The Engineering Graduate salary dataset is a regression based data. So we have applied all the algorithms of machine learning with respect to regression such as linear regression, Neural Network, Decision tree, Random forest and Boosting.

Among all the algorithms in boosting XGboost techinque with 80-20 performed well an gave least MAE value when comapred to others.

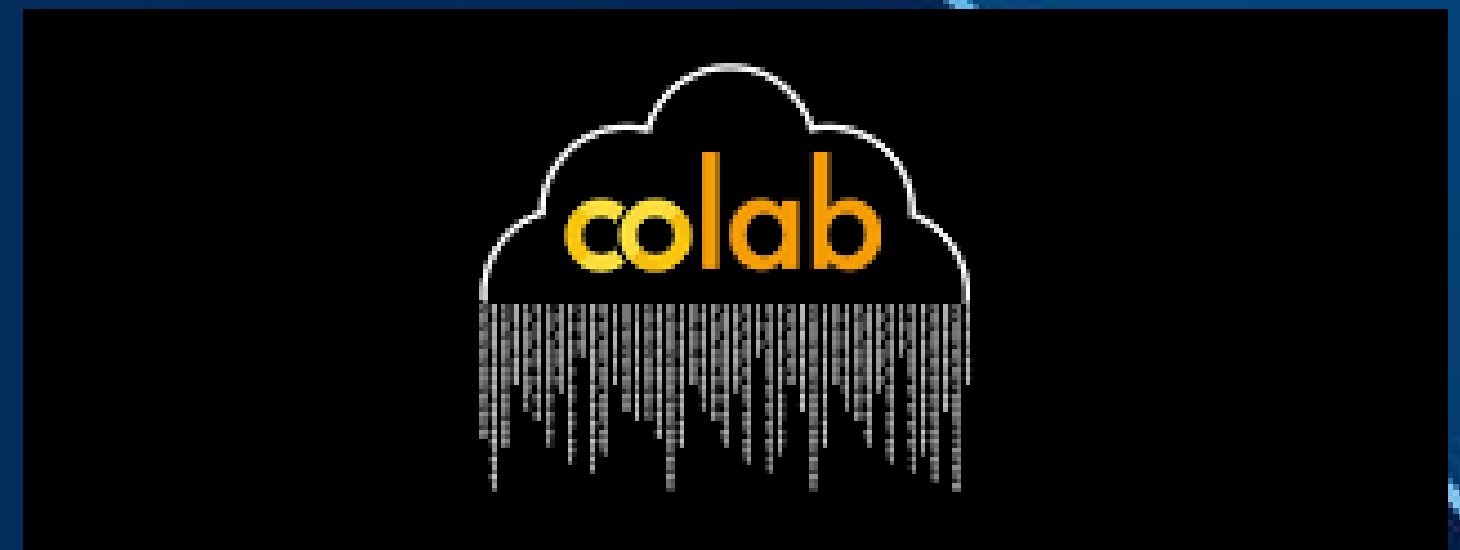So we can conclude boosting algorithm fits good and can be used for further usage of model .

# Team members and their Roles

- N.Bhavana Reddy-Coding,PPT ,EDA
- UudhhayKiirran-Coding,PPT
- Sai Prasanna-PPT
- Jashwanth-PPT

Click on the icon of github  and colab  for
more details of the project .