



# Salary Predictions of Data Professions

By:Pratiksha  
Saheb



# Problem Statement

Salaries in the field of data professions vary widely based on factors such as experience, job role, and performance. Accurately predicting salaries for data professionals is essential for both job seekers and employers.

# Methodology



Exploratory Data Analysis



Data Preprocessing



Feature Engineering



Machine Learning Model



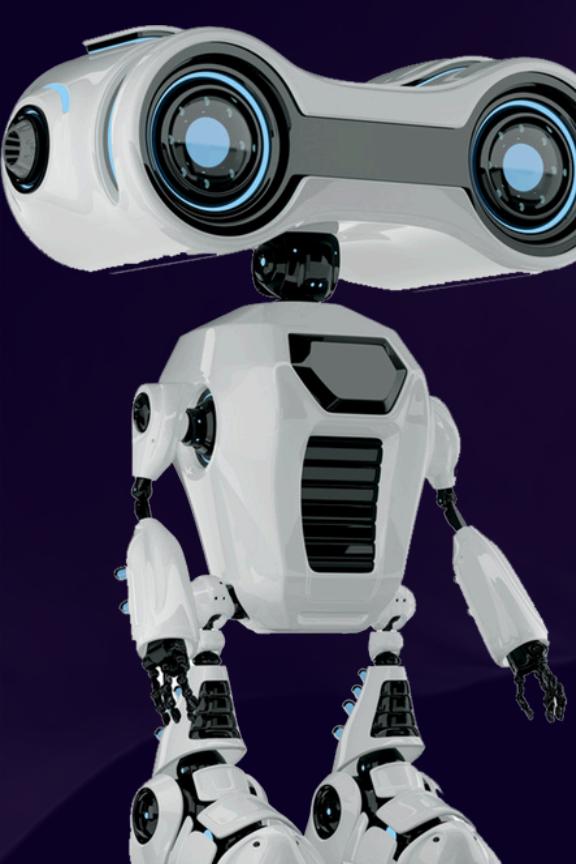
Model Evaluation

# Exploratory Data Analysis

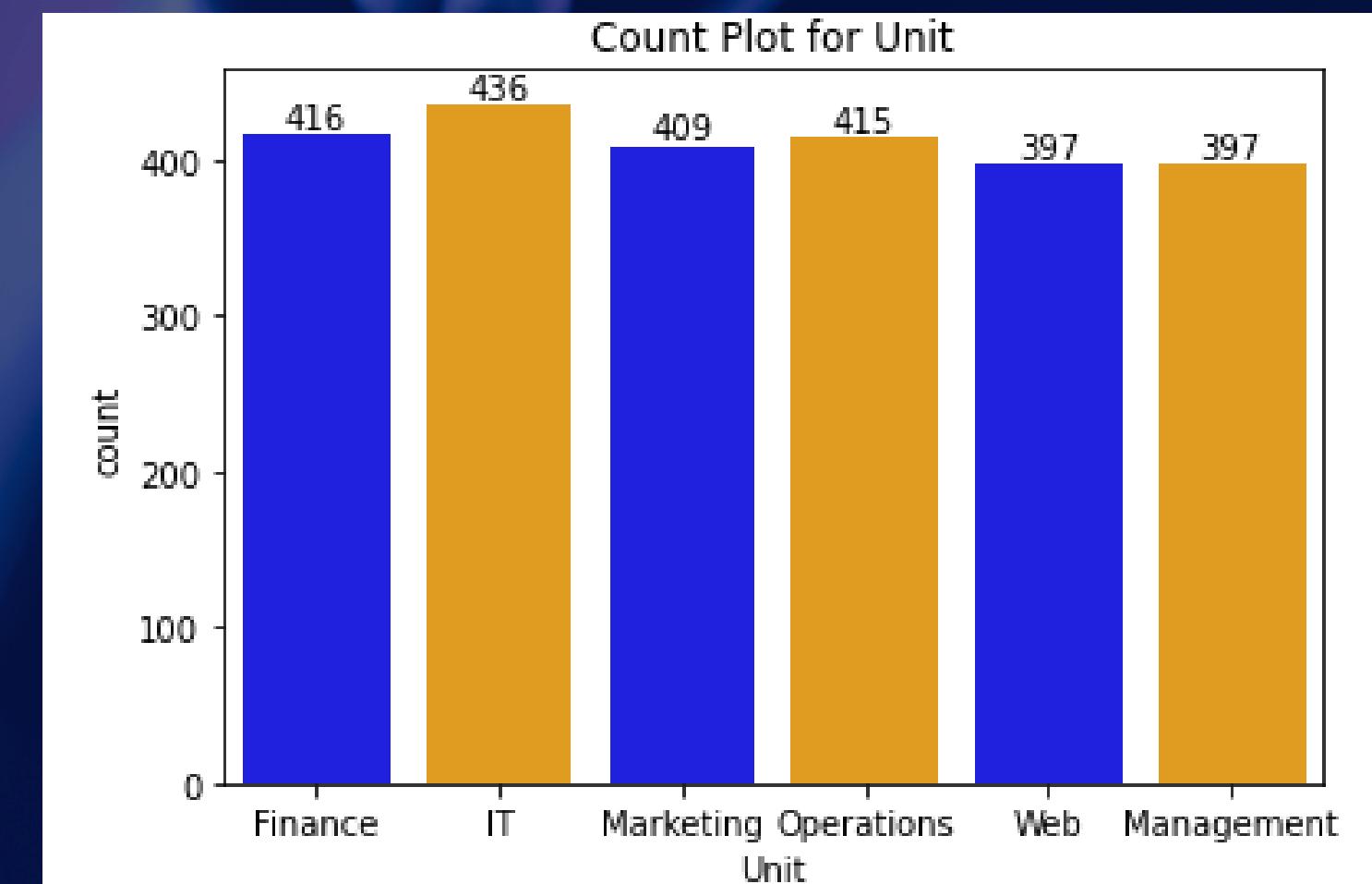
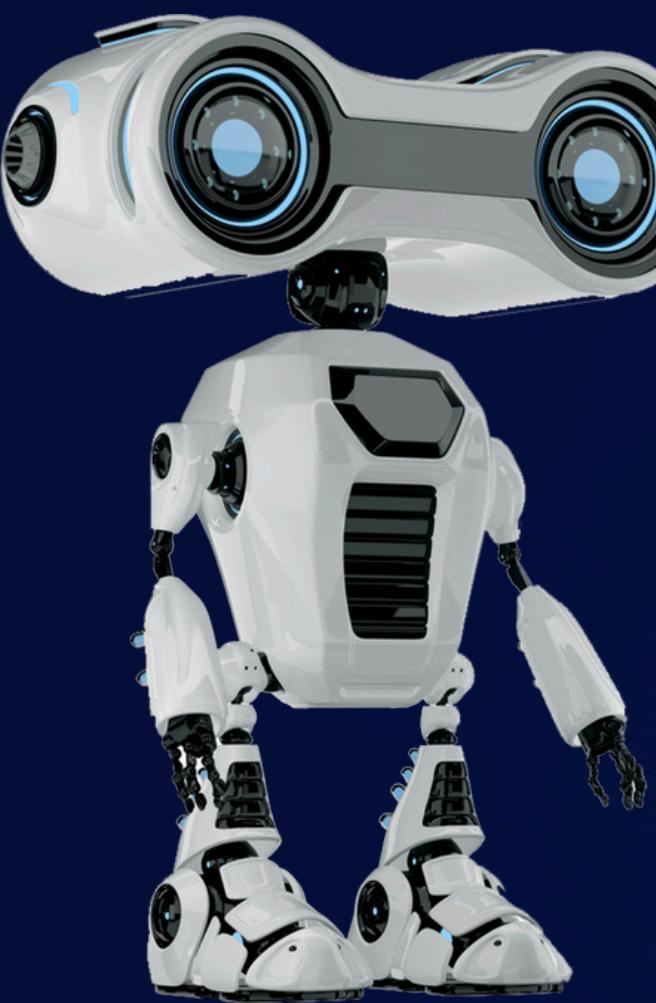
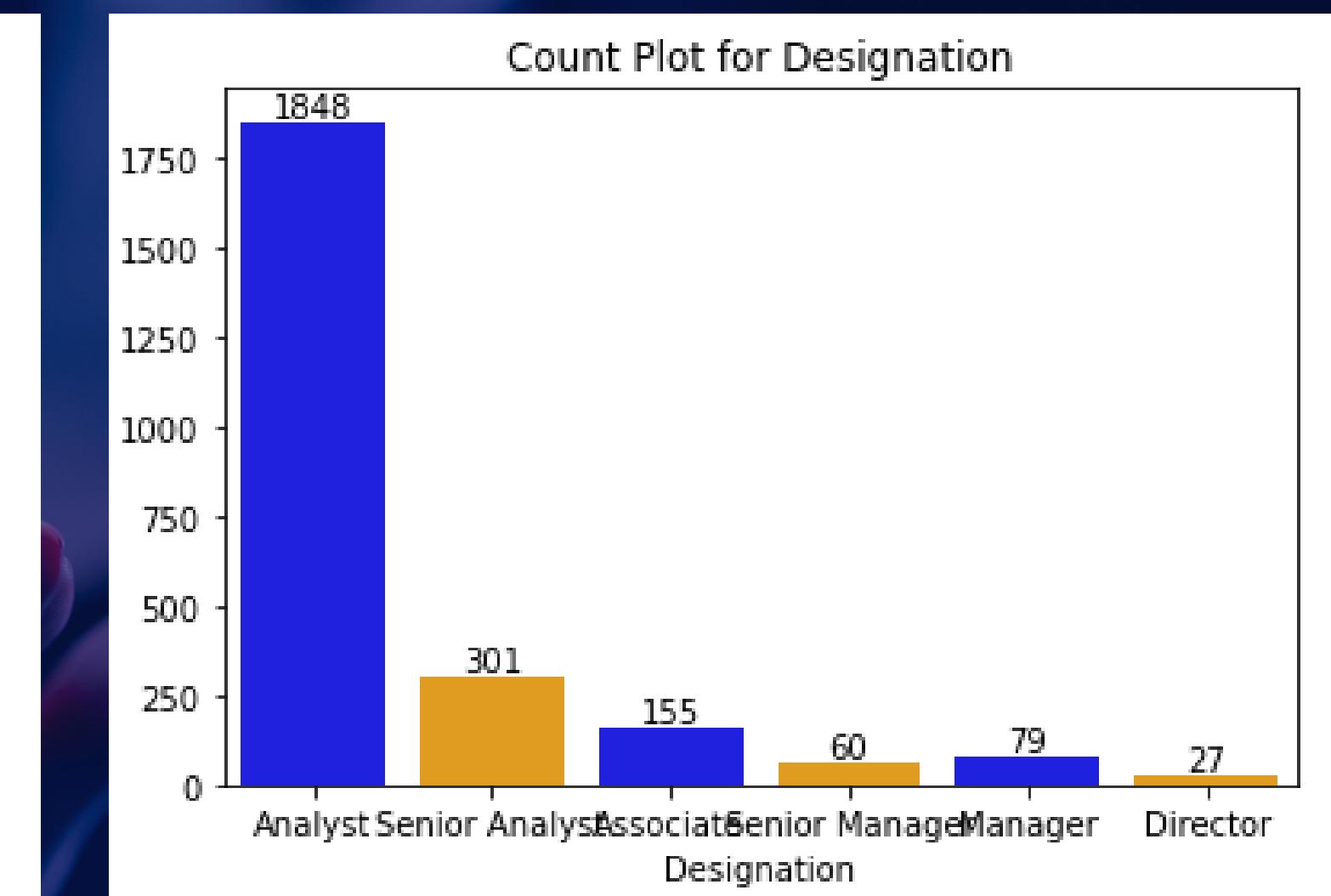
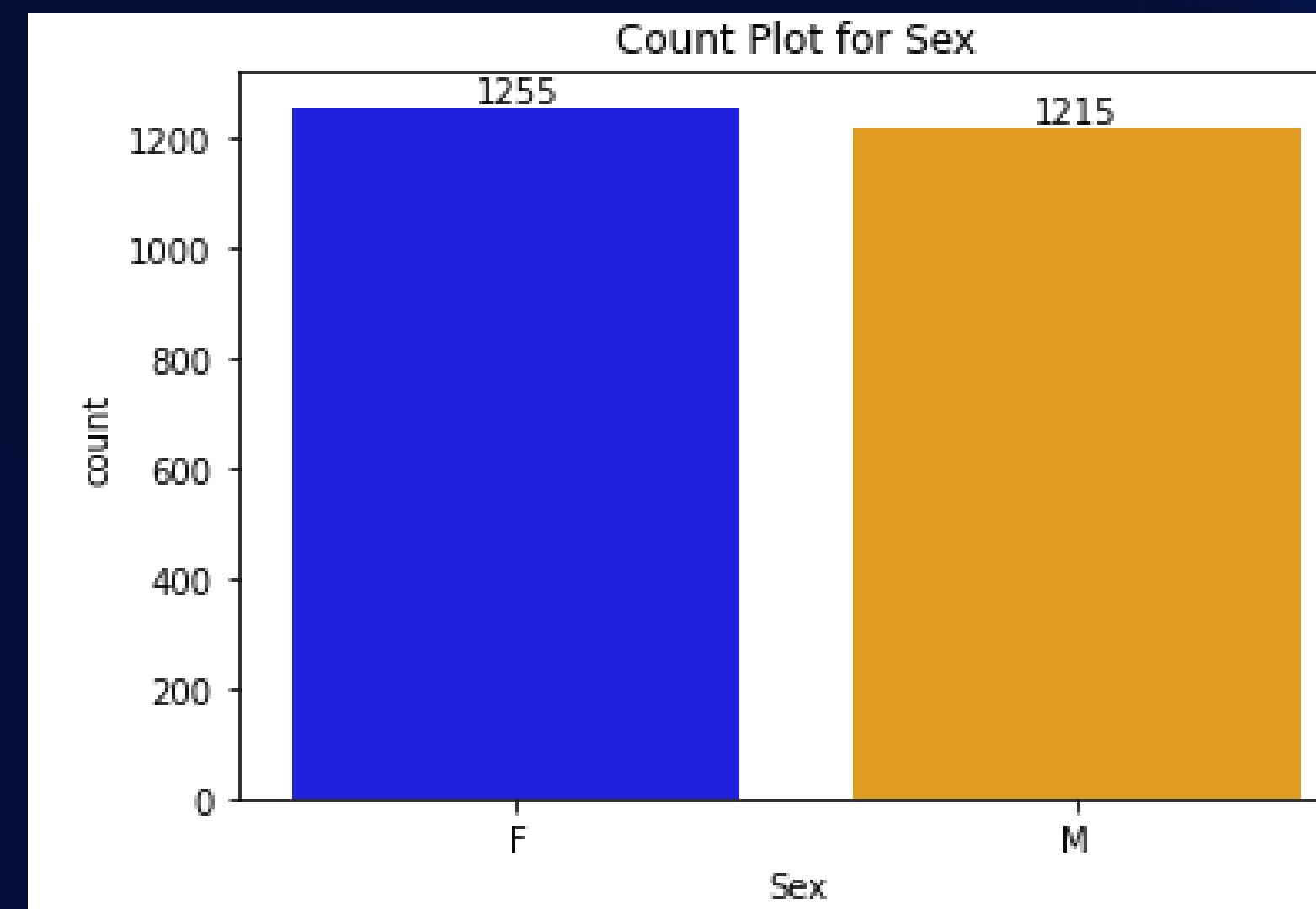
```
1 df = pd.read_csv('Salary Prediction of Data Professions.csv')
```

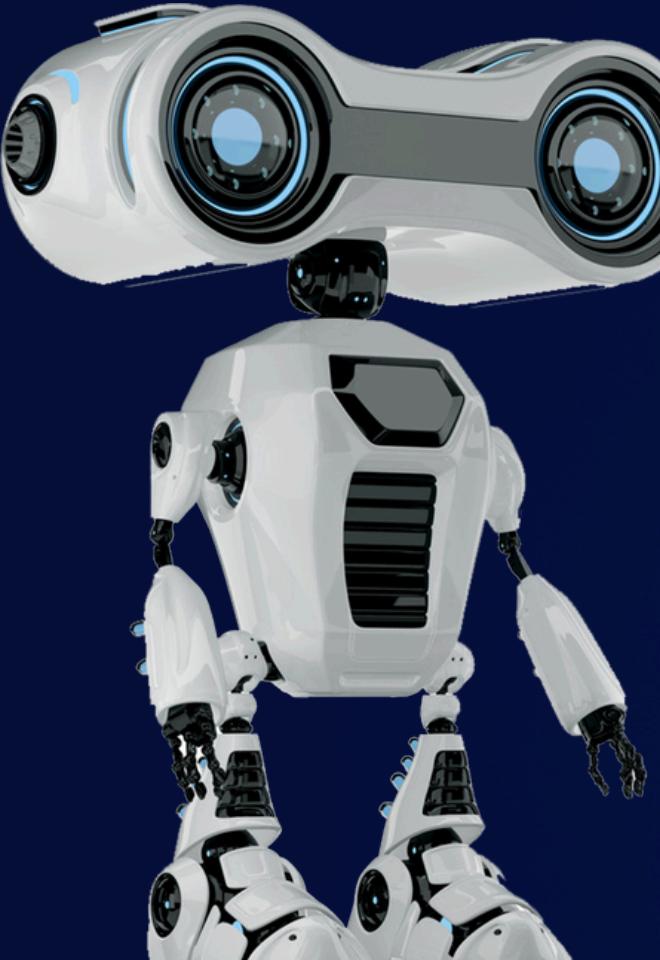
```
1 df.head()
```

	FIRST NAME	LAST NAME	SEX	DOJ	CURRENT DATE	DESIGNATION	AGE	SALARY	UNIT	LEAVES USED	LEAVES REMAINING	RATINGS	PAST EXP
0	TOMASA	ARMEN	F	5-18-2014	01-07-2016	Analyst	21.0	44570	Finance	24.0	6.0	2.0	0
1	ANNIE	Nan	F	Nan	01-07-2016	Associate	Nan	89207	Web	Nan	13.0	Nan	7
2	OLIVE	ANCY	F	7-28-2014	01-07-2016	Analyst	21.0	40955	Finance	23.0	7.0	3.0	0
3	CHERRY	AQUILAR	F	04-03-2013	01-07-2016	Analyst	22.0	45550	IT	22.0	8.0	3.0	0
4	LEON	ABOULAHoud	M	11-20-2014	01-07-2016	Analyst	Nan	43161	Operations	27.0	3.0	Nan	3



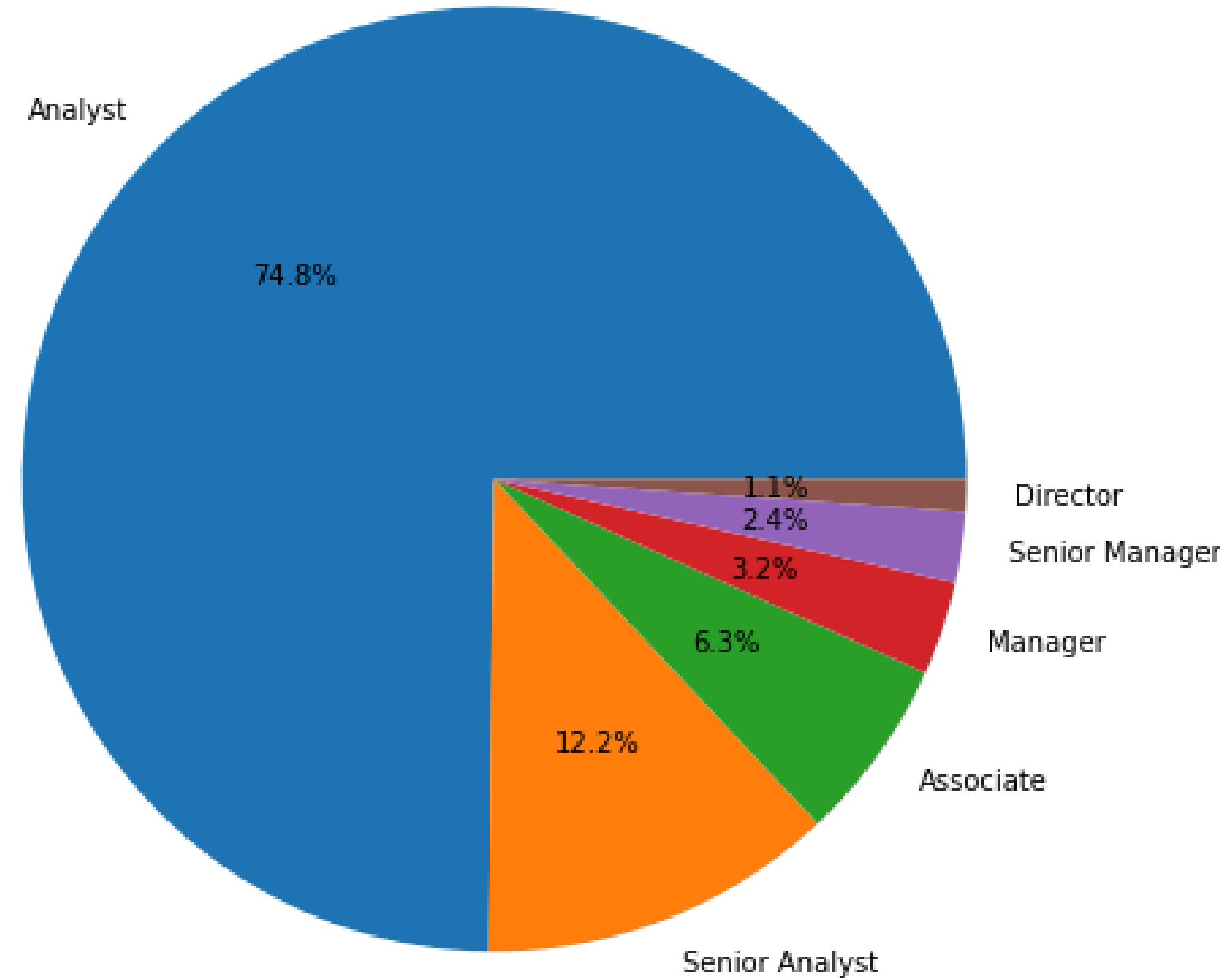
count	2636.000000	2639.000000	2636.000000	2637.000000	2637.000000	2639.000000
mean	24.756449	58136.678287	22.501517	7.503223	3.486159	1.566881
std	3.908228	36876.956944	4.604469	4.603193	1.114933	2.728416
min	21.000000	40001.000000	15.000000	0.000000	2.000000	0.000000
25%	22.000000	43418.000000	19.000000	4.000000	2.000000	0.000000
50%	24.000000	46781.000000	22.000000	8.000000	3.000000	1.000000
75%	25.000000	51401.500000	26.000000	11.000000	4.000000	2.000000
max	45.000000	388112.000000	30.000000	15.000000	5.000000	23.000000



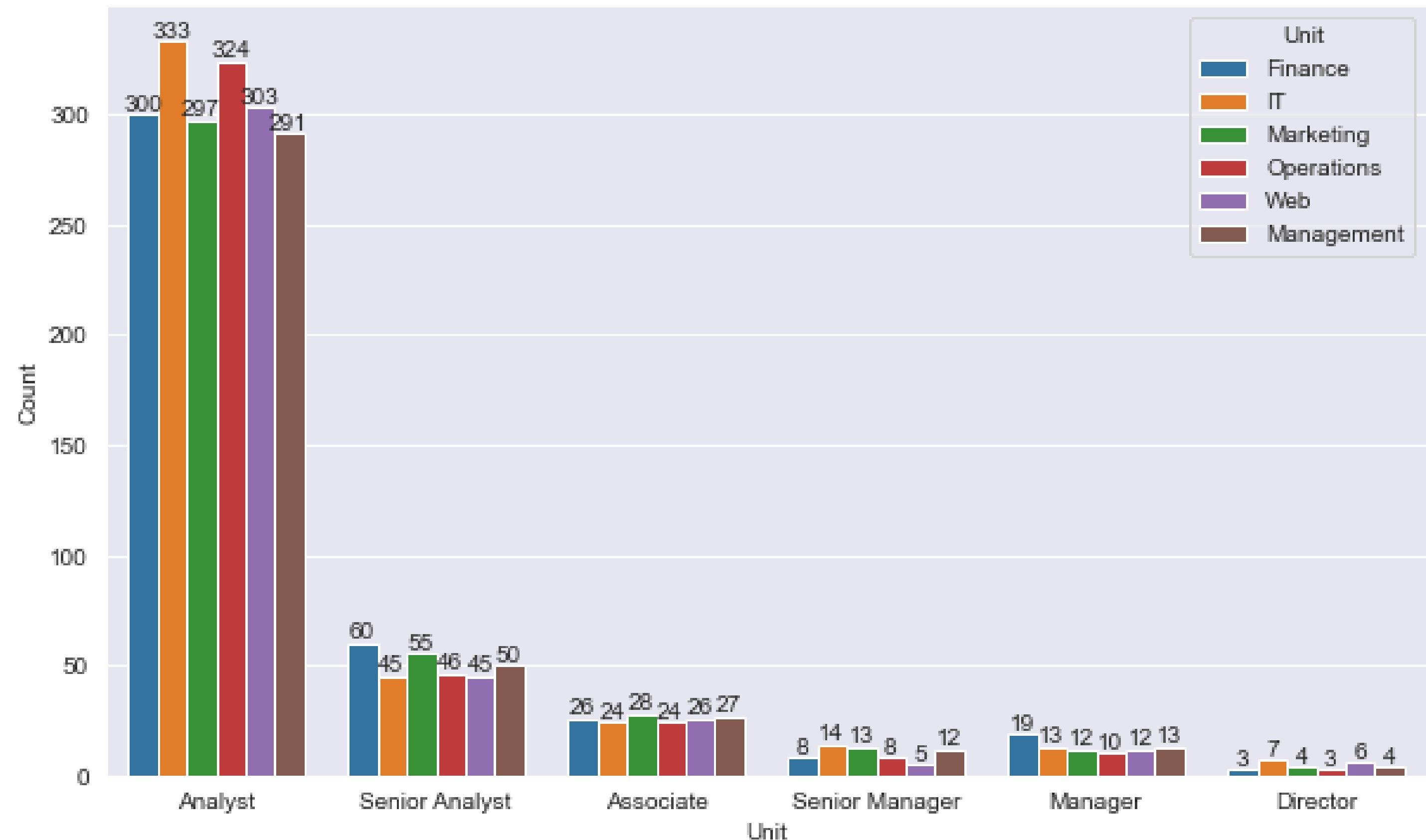


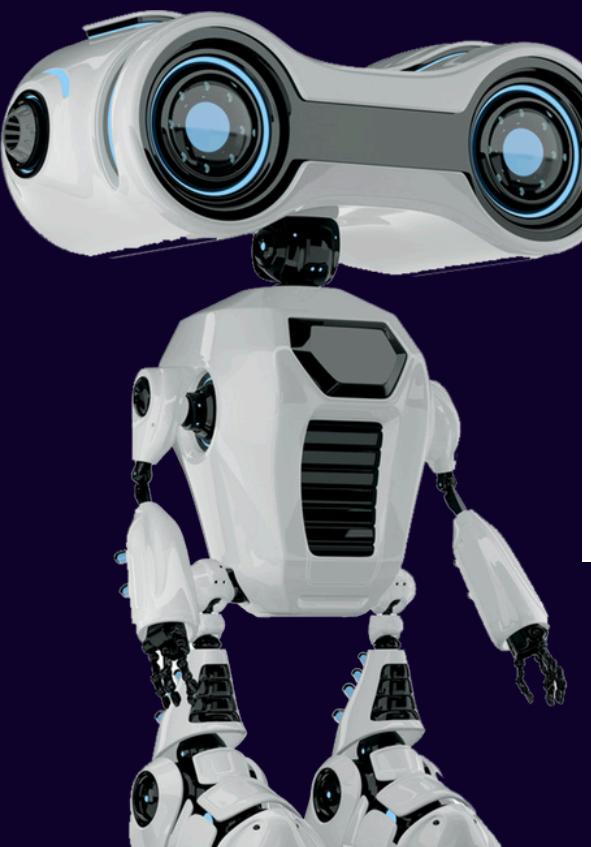
Designation

### Distribution of Designations

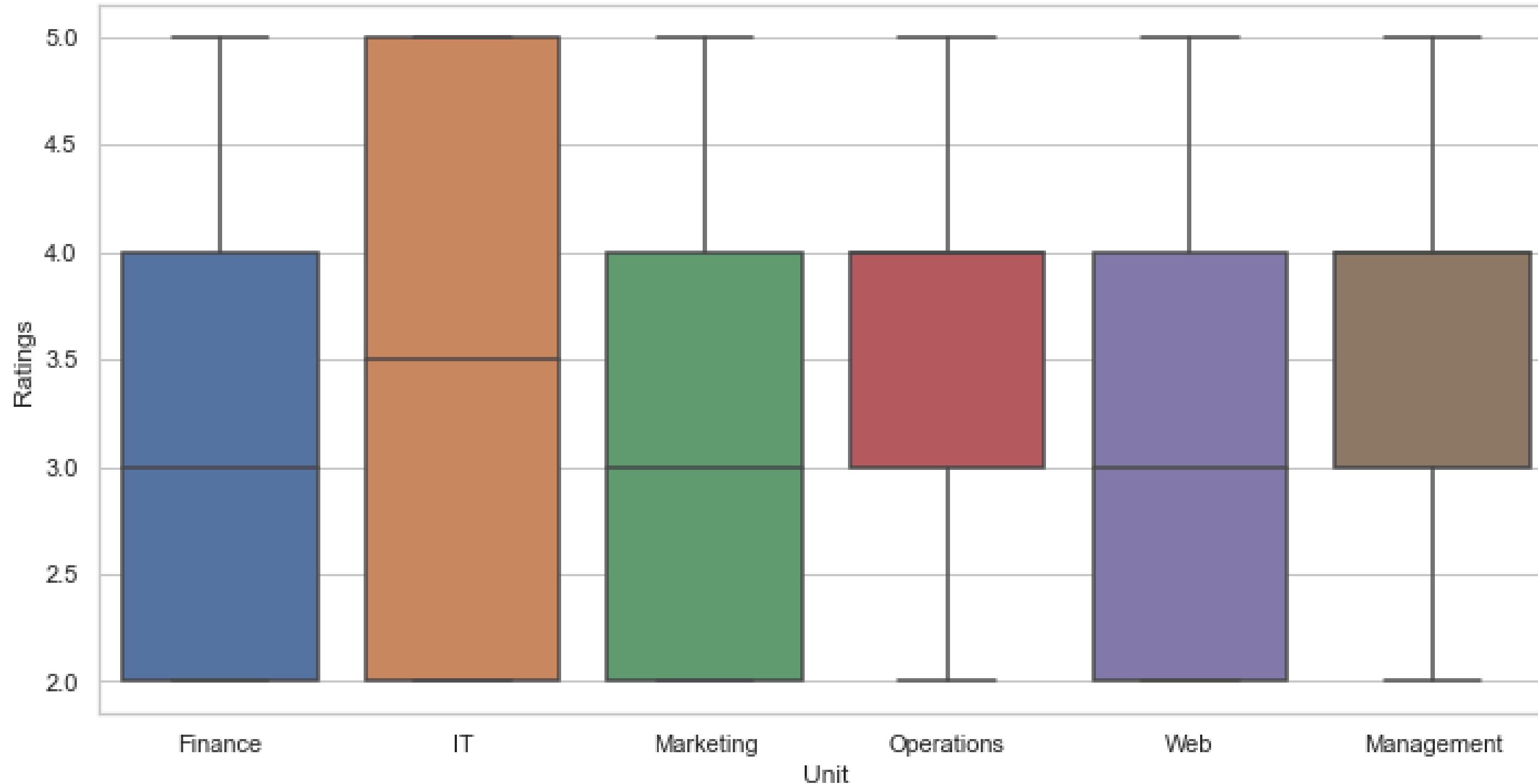


### Count Plot for Unit and Designation

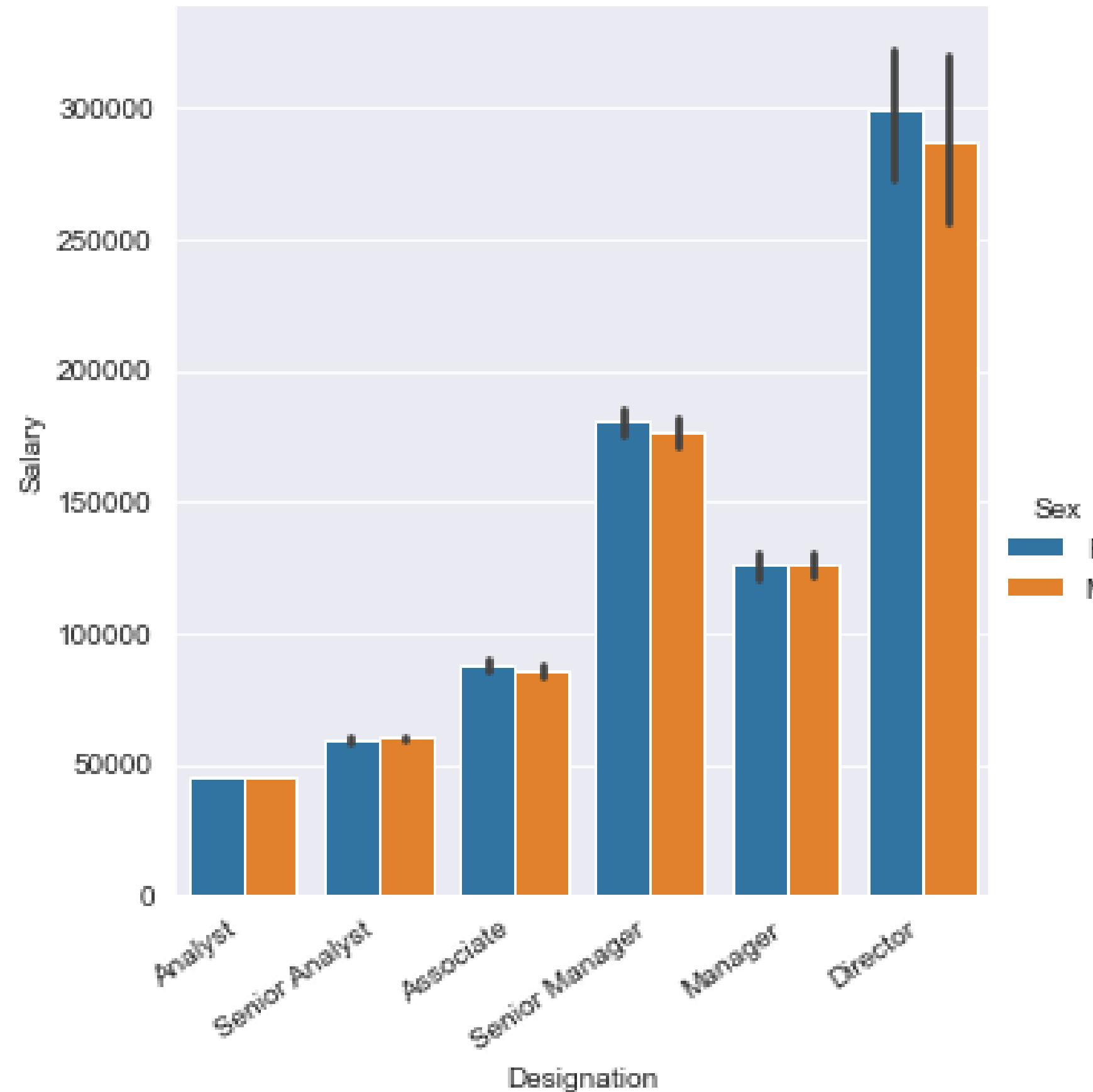


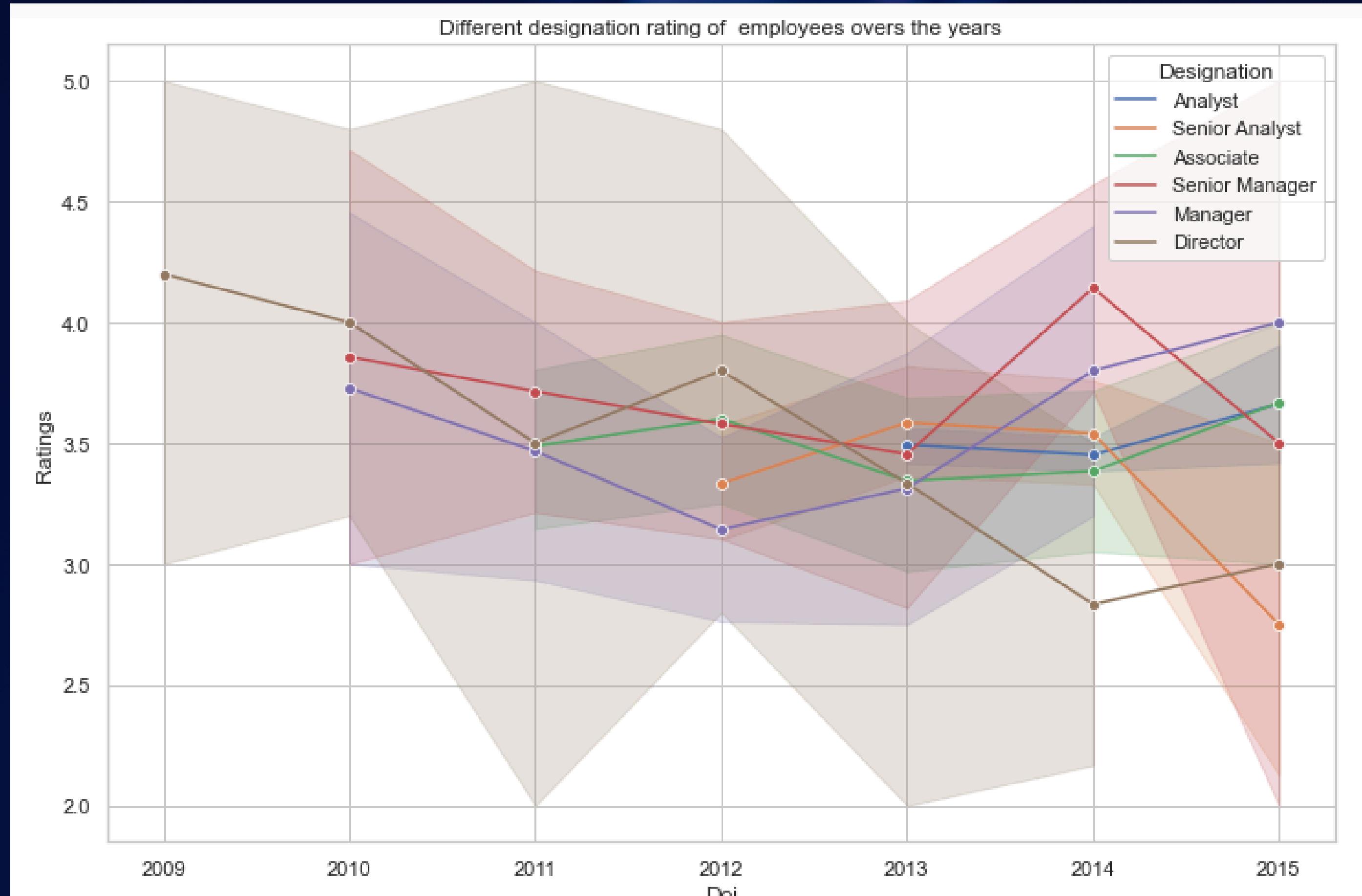
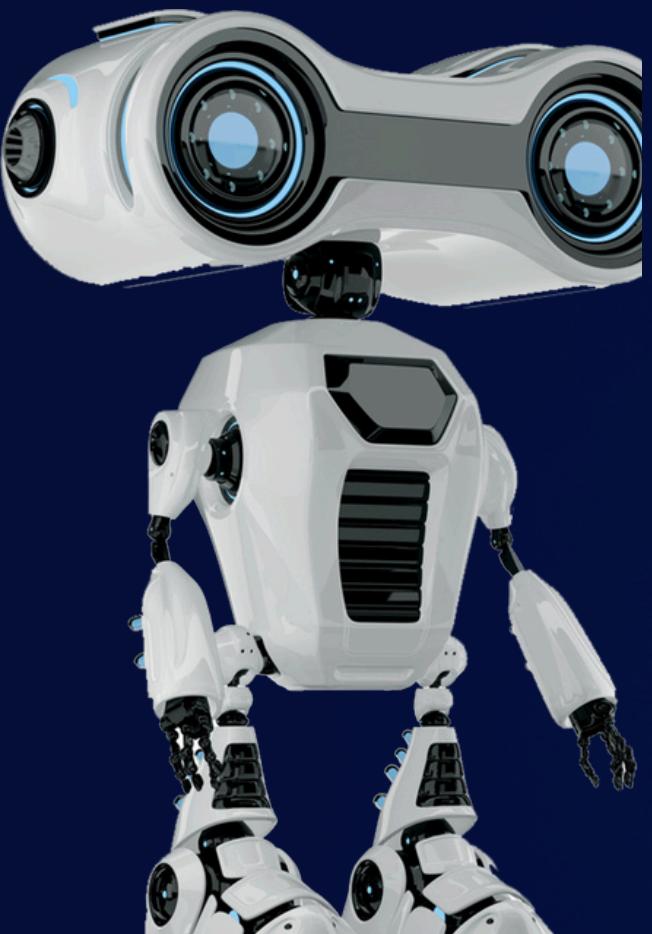


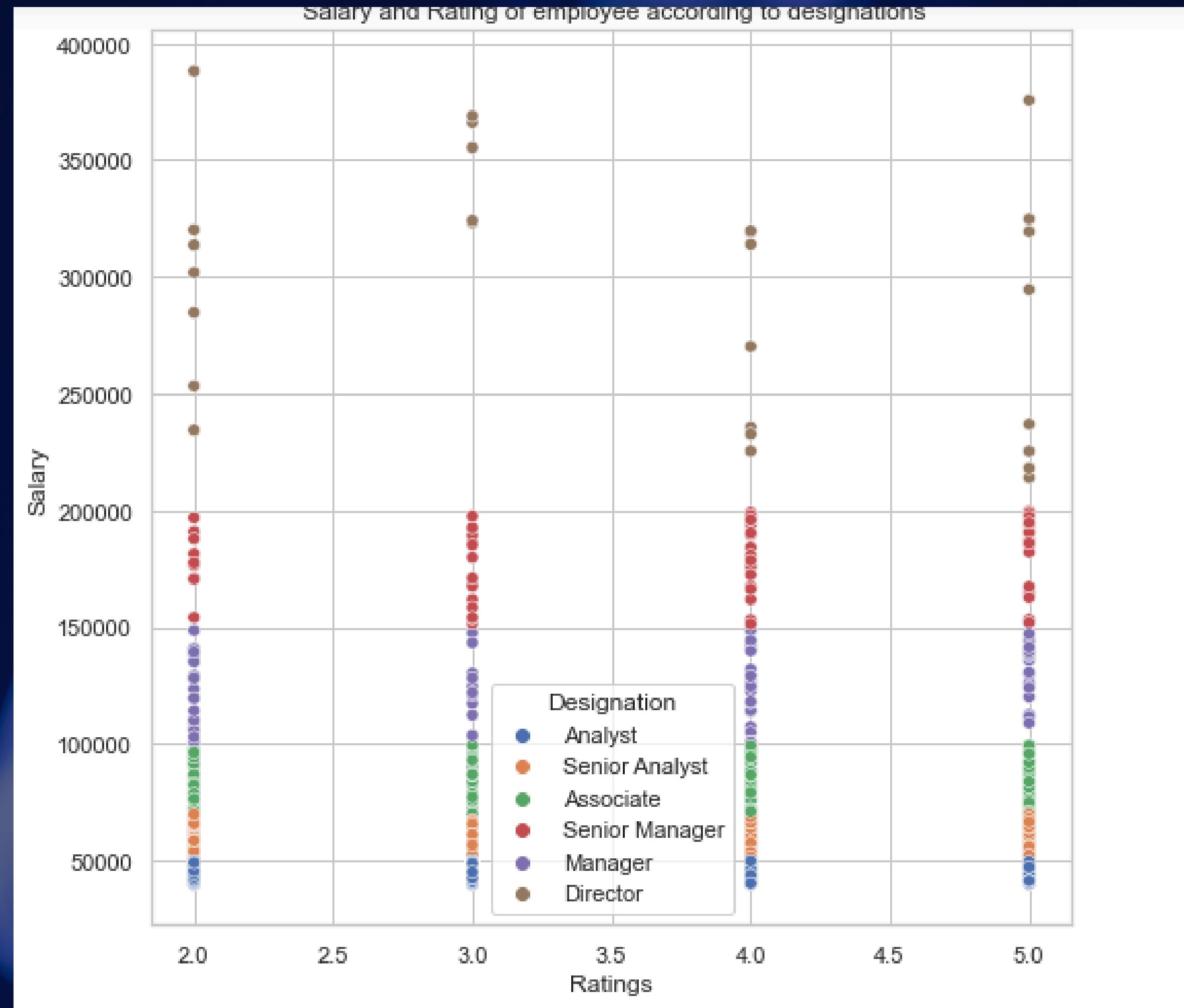
Comaprision of Rating of different Unit



<Figure size 864x432 with 0 Axes>







# Model Development

06

## ML Model Development

```
1 x = df.drop(columns = ['index','FirstName','LastName','CurrentDate','LeavesRemaining','Doj','Salary'])
2 y = df['Salary'].values
```

```
1 from sklearn.preprocessing import LabelEncoder
2 le = LabelEncoder()
```

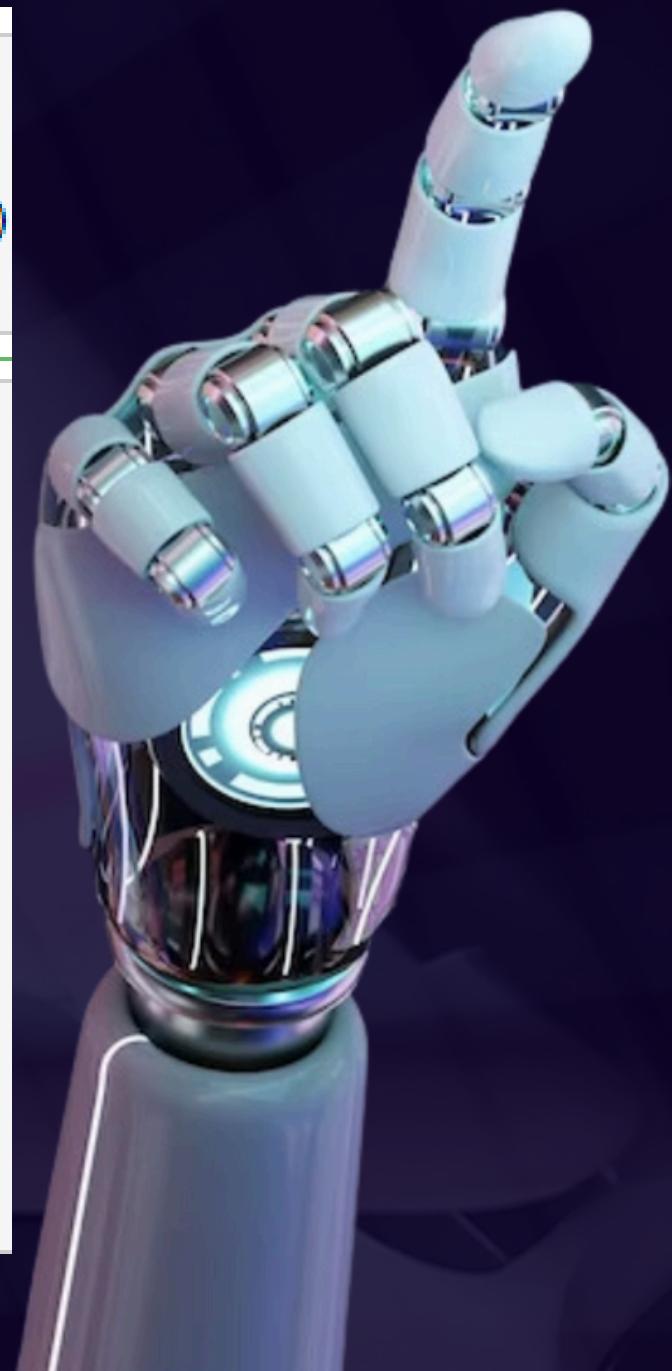
```
1 x['Sex'] = le.fit_transform(x['Sex'])
2 x['Designation'] = le.fit_transform(x['Designation'])
3 x['Unit'] = le.fit_transform(x['Unit'])
4 x = x.values
```

```
1 x
array([[ 0.,  0., 21., ...,  2.,  0.,  2.],
       [ 0.,  0., 21., ...,  3.,  0.,  1.],
       [ 0.,  0., 22., ...,  3.,  0.,  3.],
       ...,
       [ 0.,  0., 21., ...,  5.,  0.,  2.],
       [ 0.,  0., 24., ...,  3.,  1.,  2.],
       [ 1.,  0., 24., ...,  2.,  2.,  1.]])
```

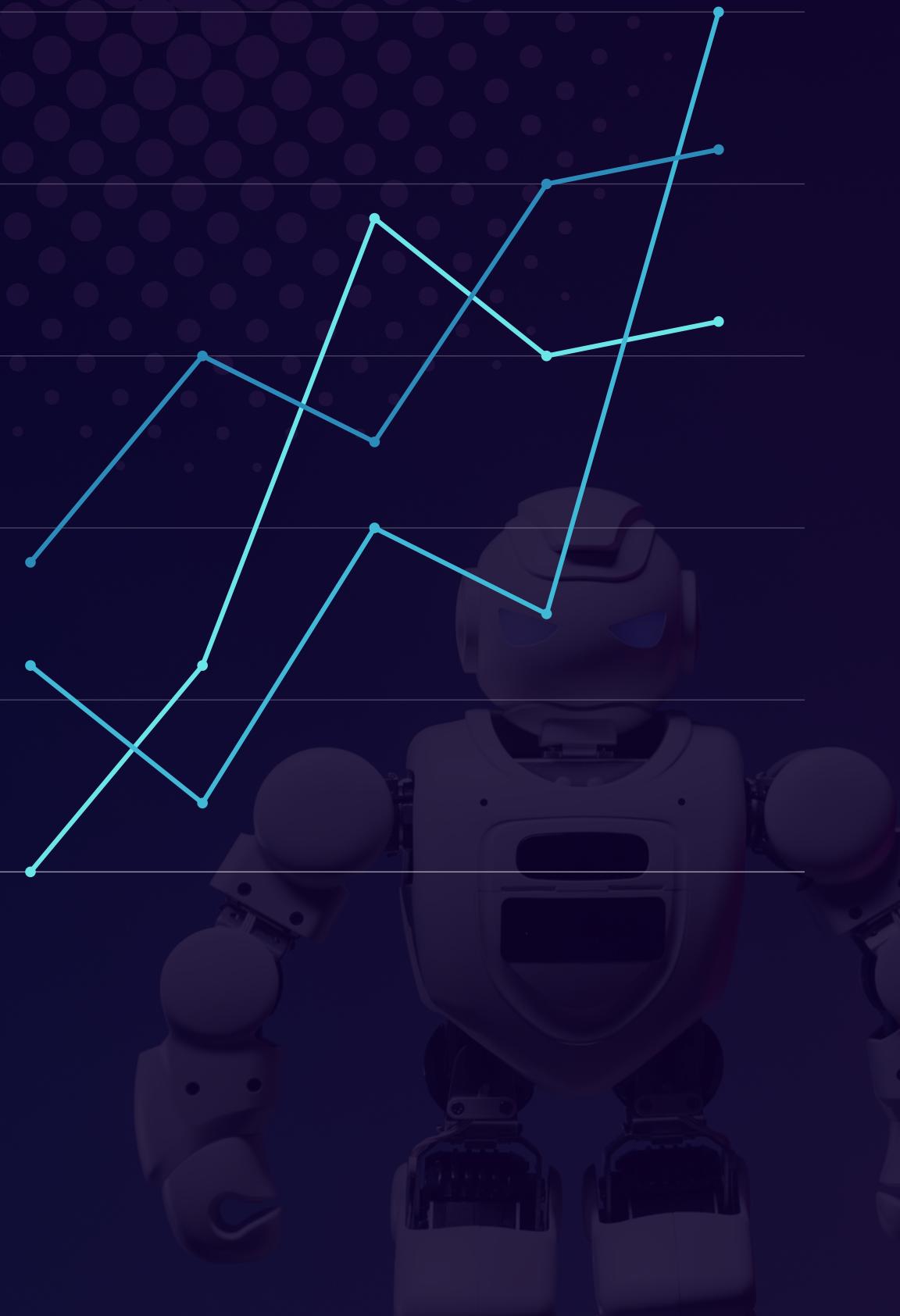


```
1 models = {'Random Forest': RandomForestRegressor(max_depth = 10 ,min_samples_leaf = 4,n_estimators = 10),
2           'Linear Regression': LinearRegression(),
3           'Decision Tree': DecisionTreeRegressor(max_depth = 10 ,min_samples_leaf = 2,min_samples_split = 10)
4           'GradeientBoost': GradientBoostingRegressor(min_samples_split = 10,n_estimators= 50)}
5
6
7
8
9
10
11
12
13
```

```
1 for name,model in models.items():
2     model.fit(X_train,y_train)
3     y_pred = model.predict(X_test)
4     mae = mean_absolute_error(y_test,y_pred)
5     mse = mean_squared_error(y_test,y_pred)
6     rmse = mean_squared_error(y_test,y_pred, squared = False)
7     r2 =r2_score(y_test,y_pred)
8     print(f' Evaluating {name} model')
9     print("Mean absolute Error (MAE) :", mae)
10    print("Mean squared Error (MSE) :" ,mse)
11    print(" Root Mean squared Error (RMSE) :" ,rmse)
12    print("R Squared error (R2) :" ,r2)
13    print('\n')
```



# Evaluation Metrics



Evaluating Random Forest model  
Mean absolute Error (MAE) : 3602.4431497168825  
Mean squared Error (MSE) : 31803874.223804925  
Root Mean squared Error (RMSE) : 5639.492372882946  
R Squared error (R2) : 0.9669763222597569

Evaluating Linear Regression model  
Mean absolute Error (MAE) : 10601.420801059045  
Mean squared Error (MSE) : 191750434.92097548  
Root Mean squared Error (RMSE) : 13847.39812820356  
R Squared error (R2) : 0.8008951826176555

Evaluating Decision Tree model  
Mean absolute Error (MAE) : 4116.304480419694  
Mean squared Error (MSE) : 42925278.649209544  
Root Mean squared Error (RMSE) : 6551.7385974418685  
R Squared error (R2) : 0.955428368284748

Evaluating GradientBoost model  
Mean absolute Error (MAE) : 3653.907678576465  
Mean squared Error (MSE) : 31995626.0611464  
Root Mean squared Error (RMSE) : 5656.467631052652  
R Squared error (R2) : 0.9667772159861658



# Thank You!