

A photograph of two women in an office environment. One woman in the foreground, wearing a dark blue sweater over a white collared shirt, is smiling broadly. Another woman in the background, wearing a grey turtleneck, is also smiling. They are positioned in front of a white wall decorated with several framed photographs of landscapes.

Human Resources

Dataset Analysis

Group C

Under Supervision Of:

Prof. Dr/ Wael Mahrous

Presented By:

Rashad Hazem Aly	1112542203	Leader
Mariem Korashi Khalaf	1112760217	Presenter
Basem Barakat Allah	1113635223	Member
Rasha Gaafar Osman	1123640976	Member
Mohamed AbdElAlim	1112535391	Member
Esraa Magdy Elsaid	1110329116	Member

What is HR analytics?

*HR analytics (usually referred to as people analytics or workforce analytics) involves gathering, analysing, and reporting HR data to drive business results.

*HR analytics enables organizations to better understand of their workforce, make decisions based on data, test the effectiveness of HR policies and interventions and measure the impact of a range of HR metrics, ultimately improving overall business performance to as much as a 25% rise in business productivity, a 50% decrease in attrition rates, and 80% increase in recruiting efficiency.

*In other words, HR analytics is a data-driven approach to Human Resources Management.



HR Database:



HR database is a system where you store and manage data on your company's employees. HR databases can be used to track a variety of information, including HR metrics, which give the HR team insights for better decision-making.

It may include:

Employee's name, Address, Contact info., Job title, Who they report to, Annual leaves, Benefits, Absenteeism, Skills, Salaries, Training and Development Programs

HR Dataset:

Employees Table

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X
1	Emplo	FirstN	LastNa	Gender	Age	Busine	Depart	DistanceFromHome	State	Ethnic	Educati	EducationFiel	JobRol	MaritalSt	Salary	StockOptionLe	Overtime	HireDate	Attritio	YearsAtCompany	YearsInMostRecentRole	YearsSinceLastPromotio	YearsWithCurrManager	
2	3012-1A4	Leonelle	Simco	Female	30	Some Tr Sales		27	IL	White	5	Marketing	Sales Ex	Divorced	102059	1 No	03/01/2012	No	10	4	9	7		
3	CBCB-9	Leonerd	Aland	Male	38	Some Tr Sales		23	CA	White	4	Marketing	Sales Ex	Single	157718	0 Yes	04/01/2012	No	10	6	10	0		
4	95D7-1C	Ahmed	Sykes	Male	43	Some Tr Human F		29	CA	Asian or	4	Marketing	HR Busi	Married	309964	1 No	04/01/2012	No	10	6	10	8		
5	47A0-55	Ermentri	Berrie	Non-Bin	39	Some Tr Technol		12	IL	White	3	Computer Scien	Engineer	Married	293132	0 No	05/01/2012	No	10	10	10	0		
6	42CC-04	Stace	Savage	Female	29	Some Tr Human F		29	CA	White	2	Technical Degre	Recruite	Single	49606	0 No	05/01/2012	Yes	6	1	1	6		
7	C219-6C	Clerkolar	Hinkins	Male	34	Some Tr Sales		30	NY	Mixed or	2	Marketing	Sales Ex	Divorced	133468	1 No	05/01/2012	No	10	3	7	9		
8	D906-B6	Uta	Melmar	Female	42	No Trav Technol		45	NY	Black or	3	Information Syst	Engineer	Married	259284	1 No	09/01/2012	No	10	2	6	6		
9	3C7D-86	Joyan	Brason	Female	40	Some Tr Sales		3	CA	Native H	2	Other	Sales Ex	Divorced	104426	1 No	11/01/2012	No	10	3	4	6		
10	3D71-8D	Ali	Blazejew	Male	38	Some Tr Sales		20	IL	Black or	4	Marketing	Sales Ex	Married	147098	1 No	11/01/2012	No	10	5	8	2		
11	5476-CA	Kagley	Snoad	Female	31	Frequent Technol		4	NY	Native H	2	Information Syst	Data Sci	Single	69747	0 No	12/01/2012	Yes	6	5	5	1		
12	73CF-49	Hannis	Waslin	Female	32	Some Tr Technol		42	CA	White	4	Computer Scien	Data Sci	Single	102022	0 No	12/01/2012	No	10	4	5	8		
13	277A-AE	Annabel	Pablos	Female	35	Some Tr Technol		8	NY	Other	4	Information Syst	Machine	Single	272175	0 No	12/01/2012	No	10	7	8	2		
14	8BAB-B	Torey	Abram	Male	38	Some Tr Sales		35	NY	Asian or	3	Marketing	Manager	Single	340229	0 Yes	13/01/2012	No	10	7	9	3		
15	11ID-E5E	Edna	Alison	Non-Bin	37	Some Tr Technol		3	IL	White	3	Computer Scien	Software	Divorced	48395	1 No	15/01/2012	No	10	2	2	4		
16	97F4-OB	Vernen	Powne	Male	33	Some Tr Technol		4	NY	Mixed or	4	Other	Senior S	Single	97126	0 No	15/01/2012	No	10	8	10	2		
17	5C03-10C	Willetta	Lurriman	Female	42	Some Tr Technol		21	IL	Black or	3	Information Syst	Engineer	Married	316208	1 No	17/01/2012	No	10	8	8	7		
18	BD1B-53	Wendall	Dryden	Male	43	Some Tr Sales		27	CA	Mixed or	3	Marketing	Sales Ex	Single	128885	0 Yes	17/01/2012	No	10	1	10	2		
19	DFA9-9C	Cale	Holston	Male	43	No Trav Sales		34	NY	White	4	Marketing	Sales Ex	Married	108315	2 No	17/01/2012	No	10	9	10	10		
20	ED73-FC	Ermaline	Napolior	Female	45	Frequent Technol		19	NY	Asian or	1	Computer Scien	Software	Married	136521	1 No	18/01/2012	No	10	3	6	1		
21	C6EC-F1	Charlene	Severwri	Female	38	Some Tr Sales		1	CA	White	1	Economics	Sales Ex	Single	151141	0 No	19/01/2012	No	10	3	6	9		
22	D7EE-5E	Zsazsa	Evered	Female	39	Frequent Sales		17	CA	White	2	Technical Degre	Sales Ex	Married	107863	1 Yes	19/01/2012	Yes	8	8	8	5		
23	C395-8C	Curcio	Franek	Male	33	Some Tr Human F		3	CA	Mixed or	1	Technical Degre	Recruite	Divorced	53616	1 Yes	19/01/2012	No	10	2	3	7		
24	E348-E1	Burnaby	Guillet	Male	36	No Trav Technol		36	IL	White	2	Information Syst	Software	Divorced	61298	1 Yes	19/01/2012	No	10	3	9	3		
25	B3AF-7E	Elvira	Ianelli	Female	45	Some Tr Human F		34	NY	Black or	2	Human Resourc	Recruit	Divorced	54132	1 Yes	20/01/2012	No	10	10	10	10		
26	469A-81	Baxie	Rising	Male	30	Some Tr Technol		36	CA	America	4	Information Syst	Engineer	Married	328415	0 No	20/01/2012	No	10	1	10	1		
27	9E22-62	Gifford	Poynter	Non-Bin	48	Some Tr Technol		37	CA	White	3	Computer Scien	Machine	Single	145337	0 No	21/01/2012	No	10	10	10	0		
28	D5DA-3I	Rickey	Shere	Male	33	Some Tr Sales		41	CA	Asian or	2	Marketing	Sales Ex	Married	71201	2 No	24/01/2012	No	10	8	10	4		
29	BFF3-A	Collen	Sedman	Female	31	No Trav Human F		25	NY	America	2	Marketing	Recruite	Divorced	55682	1 No	25/01/2012	No	10	0	2	0		
30	C0BE-F1	Bertram	Dolemar	Male	40	Some Tr Sales		35	NY	White	2	Marketing	Sales Re	Married	63455	1 Yes	26/01/2012	No	10	4	8	7		
31	D565-28	Bessie	Bellson	Female	47	Frequent Sales		21	CA	White	3	Economics	Sales Ex	Single	65626	0 No	27/01/2012	No	10	2	5	1		
32	00D4-DC	Joyce	Goor	Female	30	Frequent Technol		44	CA	Black or	1	Computer Scien	Software	Single	68508	0 Yes	28/01/2012	Yes	5	4	4	4		
33	7FFD-C1	Claresta	Impy	Female	32	Some Tr Technol		22	IL	White	2	Other	Data Sci	Married	75821	1 No	30/01/2012	No	10	1	5	1		
34	07B2-D6	Rossie	Everleigh	Male	44	Some Tr Sales		40	NY	White	2	Economics	Manager	Divorced	285620	1 No	30/01/2012	No	10	6	8	4		
35	1749-8IA	Elora	Bentjens	Female	48	Some Tr Technol		17	CA	Black or	3	Information Syst	Software	Single	65880	0 No	01/02/2012	No	10	6	7	1		
36	FF14-A4	Koenraa	Nannizzi	Male	47	Some Tr Technol		41	NY	America	1	Computer Scien	Data Sci	Married	40786	0 No	01/02/2012	No	10	3	5	5		
37	3CD6-55	Dorise	Klishin	Female	31	Some Tr Technol		20	CA	White	2	Computer Scien	Software	Single	59697	0 Yes	02/02/2012	Yes	7	5	6	0		
38	4FC2-A4	Gagle	Riseley	Female	32	No Trav Technol		12	NY	Black or	3	Computer Scien	Engineer	Married	316725	1 No	02/02/2012	No	10	6	10	2		
39	9AAB-D	Staci	Leith	Female	28	Some Tr Technol		14	IL	White	2	Business Studie	Software	Divorced	105984	2 Yes	04/02/2012	No	10	7	9	10		
40	427A-8C	Pasquali	Abrehea	Male	32	Some Tr Technol		15	CA	Mixed or	3	Other	Analytics	Married	393294	1 No	04/02/2012	No	10	2	2	2		
41	A923-9E	Alaine	Hinrichs	Female	41	Frequent Technol		7	CA	Black or	1	Business Studie	Machine	Single	107008	0 No	05/02/2012	No	10	9	10	9		
42	5C61-8F	Dyana	Gallie	Female	34	Some Tr Technol		5	CA	White	1	Information Syst	Data Sci	Single	127432	0 Yes	06/02/2012	No	10	8	9	9		
43	CF2F-8C	Lindy	Rawstor	Male	45	Some Tr Sales		20	CA	White	2	Marketing	Sales Ex	Married	76693	1 No	06/02/2012	No	10	2	2	10		
44	40A9-8E	Joannes	McFadd	Female	39	Some Tr Technol		31	CA	White	3	Information Syst	Software	Single	90017	0 Yes	07/02/2012	No	10	7	7	3		
45	2E72-4B	Grace	Gohier	Male	44	Some Tr Sales		21	CA	White	3	Marketing	Manager	Married	426142	0 No	07/02/2012	No	10	4	5	4		

Education Level Table

A	B	C	D
EducationID		EducationLevel	
1		No Formal Qualifications	
2		High School	
3		Bachelors	
4		Masters	
5		Doctorate	

Rating Level Table

A	B	C	D
RatingID		RatingLevel	
1		Unacceptable	
2		Needs Improvement	
3		Meets Expectation	
4		Exceeds Expectation	
5		Above and Beyond	

Performance Level Table

A	B	C	D	E	F	G	H	I	J	K	L
PerformanceID	EmployeeID	ReviewDate	EnvironmentSatisfaction	JobSatisfaction	RelationshipScore	TrainingOpportunitiesWithinLastYear	TrainingOpportunitiesTaken	WorkLifeBalance	OverallSelfRating	ManagerRating	
PR01	79F7-78EC	01/02/2013	5	4	5	1	0	4	4	4	4
PR02	B61E-0F26	01/03/2013	5	4	4	1	3	4	4	4	3
PR03	F5E3-48BB	01/03/2013	3	4	5	3	2	3	5	4	4
PR04	0678-748A	01/04/2013	5	3	2	2	0	2	3	3	2
PR05	541F-3E19	01/04/2013	5	2	3	1	0	4	4	3	3
PR06	F93E-BDEF	01/04/2013	3	3	2	2	0	4	4	4	4
PR07	9E7A-1F70	01/08/2013	3	4	5	2	1	5	4	3	3
PR08	05ED-92F1	01/10/2013	4	5	4	1	1	3	3	3	2
PR09	F72D-261D	01/10/2013	4	5	2	1	1	4	5	4	4
PR10	774E-685D	01/11/2013	5	4	3	2	3	4	5	4	4
PR100	B013-7D0C	04/10/2013	4	3	3	2	0	4	3	3	3
PR1000	528C-3E0D	3/16/2016	4	4	2	2	2	4	5	5	5
PR1001	D077-169C	3/17/2016	3	5	3	2	2	3	5	5	5
PR1002	9727-BC84	3/18/2016	4	3	3	2	2	2	4	3	3
PR1003	DA8E-9496	3/18/2016	3	5	4	1	0	5	5	5	5
PR1004	DEC5-9319	3/18/2016	3	4	3	2	3	2	4	4	4
PR1005	888B-EB84	3/19/2016	3	4	2	3	1	4	5	5	5
PR1006	9C57-828C	3/19/2016	5	4	2	1	1	2	3	3	3
PR1007	E1B4-9AA1	3/22/2016	5	4	3	3	2	3	4	3	3
PR1008	3CD6-5587	3/23/2016	5	4	2	2	0	4	4	4	3
PR1009	BAFA-86DF	3/23/2016	3	3	4	2	1	2	3	3	3
PR101	152E-8DB1	04/12/2013	5	2	5	1	0	5	5	4	4

Satisfied Level Table

A	B	C	D
SatisfactionID		SatisfactionLevel	
1		Very Dissatisfied	
2		Dissatisfied	
3		Neutral	
4		Satisfied	
5		Very Satisfied	

Analyzing the Data:



DATA CLEANING

DATA CLEANING STEPS

Removing unwanted observations

- Duplicate/ redundant or irrelevant values deletion .

Missing Data handling

- Fixing issue of unknown missing values

Structural error solving

- Fixing problems with mislabeled classes, types in names of features, same attribute with different name etc.

Outliers Management

- Unwanted values which are not fitting in datasets.

Importing Necessary Libraries For Data Cleaning

```
import pandas as pd
import plotly.express as px
import numpy as np
from matplotlib import pyplot as plt
%matplotlib inline
import seaborn as sns
pd.options.display.max_columns=None
pd.options.display.max_rows=None
pd.options.display.float_format='{:,.2f}'.format
```

Python

Recalling Data From Its Previous Format (CSV) And Storing It In A Custom Data Frame For Each File

```
df_employee=pd.read_csv('Employee.csv')
df_RatingLevel=pd.read_csv('RatingLevel.csv')
df_satisfiedLevel=pd.read_csv('SatisfiedLevel.csv')
df_EducationLevel=pd.read_csv('EducationLevel.csv')
df_PerformanceRating=pd.read_csv('PerformanceRating.csv')
```

Python

Displaying Information About Employees Table

```
df_employee.info()
```

Python

Data columns (total 23 columns):			
#	Column	Non-Null Count	Dtype
0	EmployeeID	1470	non-null
1	FirstName	1470	non-null
2	LastName	1470	non-null
3	Gender	1470	non-null
4	Age	1470	non-null
5	BusinessTravel	1470	non-null
6	Department	1470	non-null
7	DistanceFromHome (KM)	1470	non-null
8	State	1470	non-null
9	Ethnicity	1470	non-null
10	Education	1470	non-null
15	StockOptionLevel	1470	non-null
16	Overtime	1470	non-null
17	HireDate	1470	non-null
18	Attrition	1470	non-null
19	YearsAtCompany	1470	non-null
...			
21	YearsSinceLastPromotion	1470	non-null
22	YearsWithCurrManager	1470	non-null
dtypes: int64(9), object(14)			

(No empty values, but HireDate is object Not DateTime)

Ensuring No Empty Data Or Nulls In Employees Table

```
df_employee.isnull().sum()
```

Python

EmployeeID	0	JobRole	0
FirstName	0	MaritalStatus	0
LastName	0	Salary	0
Gender	0	StockOptionLevel	0
Age	0	Overtime	0
BusinessTravel	0	HireDate	0
Department	0	Attrition	0
DistanceFromHome (KM)	0	YearsAtCompany	0
State	0	YearsInMostRecentRole	0
Ethnicity	0	YearsSinceLastPromotion	0
Education	0	YearsWithCurrManager	0
EducationField	0	dtype: int64	

No nulls found

Solving All Problems In Employees Table

```
df_employee=df_employee.drop_duplicates(subset="EmployeeID")
df_employee.HireDate=pd.to_datetime(df_employee.HireDate)
df_employee['Full Name']=df_employee.FirstName+' '+df_employee.LastName
```

Python

```
df_employee.drop("FirstName",axis=1,inplace=True)
df_employee.drop("LastName",axis=1,inplace=True)
```

Python

- Ensured deleting all duplicates in EmployeeID.
- Converted the hireDate type from object to datetime.
- Collected the first and last name in a table named full name.
- Removed fristName and LastName columns.

Displaying Employees Table After Cleaning

```
df_employee.info()
```

Python

```
Data columns (total 22 columns):
 #   Column           Non-Null Count  Dtype  
 ---  --  
 0   EmployeeID      1470 non-null    object 
 1   Gender          1470 non-null    object 
 2   Age              1470 non-null    int64  
 3   BusinessTravel   1470 non-null    object 
 4   Department       1470 non-null    object 
 5   DistanceFromHome (KM) 1470 non-null    int64  
 6   State            1470 non-null    object 
 7   Ethnicity        1470 non-null    object 
 8   Education        1470 non-null    int64  
 9   EducationField   1470 non-null    object 
 10  JobRole          1470 non-null    object 
 11  MaritalStatus    1470 non-null    object 
 12  Salary           1470 non-null    int64  
 13  StockOptionLevel 1470 non-null    int64  
 14  OverTime         1470 non-null    object 
 15  HireDate         1470 non-null    datetime64[ns]
 16  Attrition        1470 non-null    object 
 17  YearsAtCompany   1470 non-null    int64  
 18  YearsInMostRecentRole 1470 non-null    int64  
 19  YearsSinceLastPromotion 1470 non-null    int64  
 20  YearsWithCurrManager 1470 non-null    int64  
 21  Full Name        1470 non-null    object 

dtypes: datetime64[ns](1), int64(9), object(12)
memory usage: 252.8+ KB
```

Ensured everything is well.

Looking For Any Outlayers In All Tables

```
df_employee.describe()
```

Python

	Age	DistanceFromHome (KM)	Education	Salary	StockOptionLevel	HireDate	YearsAtCompany	YearsInMgmt
count	1,470.00	1,470.00	1,470.00	1,470.00	1,470.00	1470	1,470.00	1,470.00
mean	28.99	22.50	2.91	112,956.50	0.79	2017-07-05 14:50:26.938775296		4.56
min	18.00	1.00	1.00	20,387.00	0.00	2012-01-03 00:00:00		0.00
25%	23.00	12.00	2.00	43,580.50	0.00	2014-07-10 06:00:00		2.00
50%	26.00	22.00	3.00	71,199.50	1.00	2017-09-02 00:00:00		4.00
75%	34.00	33.00	4.00	142,055.75	1.00	2020-05-14 18:00:00		7.00
max	51.00	45.00	5.00	547,204.00	3.00	2022-12-31 00:00:00		10.00
std	7.99	12.81	1.02	103,342.89	0.85	NaN		3.29

Looking For Any Outlayers In All Tables (Autoclean Library)

```
1 from AutoClean import AutoClean
2 import pandas as pd
3 resultant = pd.read_csv("./Employee.csv")
4 pipeline = AutoClean(resultant)
5 x = pipeline.output
6 m = resultant.isnull().sum()
7 print(m)
8 mr = resultant[resultant.isnull().any(axis=1)]
9 print(mr)
10 |
```

Those outlayers had explanation
and did not affect the results.

So, they were not deleted.

```
1 13-10-2024 21:31:10.00 - INFO - Started validation of input parameters...
2 13-10-2024 21:31:10.00 - INFO - Completed validation of input parameters
3 13-10-2024 21:31:10.00 - INFO - Started handling of duplicates... Method: "AUTO"
4 13-10-2024 21:31:10.02 - DEBUG - 0 missing values found
5 13-10-2024 21:31:10.02 - INFO - Completed handling of duplicates in 0.014074 seconds
6 13-10-2024 21:31:10.02 - INFO - Started handling of missing values...
7 13-10-2024 21:31:10.02 - DEBUG - 0 missing values found
8 13-10-2024 21:31:10.02 - INFO - Completed handling of missing values in 0.002094 seconds
9 13-10-2024 21:31:10.02 - INFO - Started handling of outliers... Method: "WIND"
10 13-10-2024 21:31:10.02 - DEBUG - Outlier imputation of 1 value(s) succeeded for feature "Age"
11 13-10-2024 21:31:11.03 - DEBUG - Outlier imputation of 124 value(s) succeeded for feature "Salary"
12 13-10-2024 21:31:11.19 - DEBUG - Outlier imputation of 85 value(s) succeeded for feature "StockOptionLevel"
13 13-10-2024 21:31:11.19 - INFO - Completed handling of outliers in 0.376335 seconds
14 13-10-2024 21:31:11.19 - INFO - Started conversion of DATETIME features... Granularity: s
15 13-10-2024 21:31:11.27 - DEBUG - Conversion to DATETIME succeeded for feature "HireDate"
16 13-10-2024 21:31:11.28 - INFO - Completed conversion of DATETIME features in 0.0769 seconds
17 13-10-2024 21:31:11.28 - INFO - Started encoding categorical features... Method: "AUTO"
18 13-10-2024 21:31:11.28 - DEBUG - Encoding to ONEHOT succeeded for feature "Gender"
19 13-10-2024 21:31:11.28 - DEBUG - Encoding skipped for feature "LastName"
20 13-10-2024 21:31:11.28 - DEBUG - Encoding skipped for feature "EmployeeID"
21 13-10-2024 21:31:11.28 - DEBUG - Encoding skipped for feature "FirstName"
22 13-10-2024 21:31:11.30 - DEBUG - Encoding to ONEHOT succeeded for feature "OverTime"
23 13-10-2024 21:31:11.30 - DEBUG - Encoding to ONEHOT succeeded for feature "Department"
24 13-10-2024 21:31:11.30 - DEBUG - Encoding to ONEHOT succeeded for feature "Ethnicity"
25 13-10-2024 21:31:11.30 - DEBUG - Encoding to ONEHOT succeeded for feature "Attrition"
26 13-10-2024 21:31:11.30 - DEBUG - Encoding to ONEHOT succeeded for feature "BusinessTravel"
```

Ensuring No Empty Data Or Nulls In Performance Rating Table

```
df_PerformanceRating.isnull().sum()
```

Python

```
PerformanceID          0  
EmployeeID             0  
ReviewDate              0  
EnvironmentSatisfaction 0  
JobSatisfaction         0  
RelationshipSatisfaction 0  
TrainingOpportunitiesWithinYear 0  
TrainingOpportunitiesTaken    0  
WorkLifeBalance          0  
SelfRating               0  
ManagerRating             0  
dtype: int64
```

No nulls found

Solving All Problems In Performance Rating Table

```
df_PerformanceRating=df_PerformanceRating.drop_duplicates(subset="PerformanceID")
df_PerformanceRating. ReviewDate=pd.to_datetime(df_PerformanceRating. ReviewDate)
df_PerformanceRating.info()
```

Python

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6709 entries, 0 to 6708
Data columns (total 11 columns):
 #   Column           Non-Null Count  Dtype  
 ---  --  
 0   PerformanceID   6709 non-null    object  
 1   EmployeeID      6709 non-null    object  
 2   ReviewDate       6709 non-null    datetime64[ns]
 3   EnvironmentSatisfaction  6709 non-null    int64  
 4   JobSatisfaction  6709 non-null    int64  
 5   RelationshipSatisfaction 6709 non-null    int64  
 6   TrainingOpportunitiesWithinYear 6709 non-null    int64  
 7   TrainingOpportunitiesTaken     6709 non-null    int64  
 8   WorkLifeBalance    6709 non-null    int64  
 9   SelfRating        6709 non-null    int64  
 10  ManagerRating     6709 non-null    int64  
dtypes: datetime64[ns](1), int64(8), object(2)
memory usage: 576.7+ KB
```

- Ensured deleting all duplicates in PerformanceID.
- Converted the reviewDate type from object to datetime.

Ensuring No Empty Data Or Nulls In **SatisfiedLevel** Table And Changing The Columns Names

```
df_satisfiedLevel.info()
```

Python

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5 entries, 0 to 4
Data columns (total 2 columns):
 #   Column           Non-Null Count  Dtype  
---  --  
 0   SatisfactionID  5 non-null      int64  
 1   SatisfactionLevel 5 non-null    object  
dtypes: int64(1), object(1)
memory usage: 212.0+ bytes
```

No nulls found

```
df_satisfiedLevel=df_satisfiedLevel.rename(columns={'SatisfactionID':'SatisfactionLevelId'})
df_satisfiedLevel=df_satisfiedLevel.rename(columns={'SatisfactionLevel':'SatisfactionLevelName'})
```

Python

Ensuring No Empty Data Or Nulls In EducationLevel And RatingLevel Table

```
df_EducationLevel.info()
```

Python

```
df_RatingLevel.info()
```

Python

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5 entries, 0 to 4
Data columns (total 2 columns):
 #   Column           Non-Null Count  Dtype  
---  --  
 0   EducationLevelID 5 non-null      int64  
 1   EducationLevel    5 non-null      object 
dtypes: int64(1), object(1)
memory usage: 212.0+ bytes
```

No nulls found

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5 entries, 0 to 4
Data columns (total 2 columns):
 #   Column           Non-Null Count  Dtype  
---  --  
 0   RatingID         5 non-null      int64  
 1   RatingLevel      5 non-null      object 
dtypes: int64(1), object(1)
memory usage: 212.0+ bytes
```

Show Sample From Dataframe

```
df_employee.sample()
```

Python

```
df_employee.describe()
```

Python

	Age	DistanceFromHome (KM)	Education	Salary	StockOptionLevel	HireDate	YearsAtCompany	YearsInMgmt
count	1,470.00	1,470.00	1,470.00	1,470.00	1,470.00	1470	1,470.00	1,470.00
mean	28.99	22.50	2.91	112,956.50	0.79	2017-07-05 14:50:26.938775296	4.56	4.56
min	18.00	1.00	1.00	20,387.00	0.00	2012-01-03 00:00:00	0.00	0.00
25%	23.00	12.00	2.00	43,580.50	0.00	2014-07-10 06:00:00	2.00	2.00
50%	26.00	22.00	3.00	71,199.50	1.00	2017-09-02 00:00:00	4.00	4.00
75%	34.00	33.00	4.00	142,055.75	1.00	2020-05-14 18:00:00	7.00	7.00
max	51.00	45.00	5.00	547,204.00	3.00	2022-12-31 00:00:00	10.00	10.00
std	7.99	12.81	1.02	103,342.89	0.85	NaN	3.29	3.29

DATA QUALITY

```
import pandas as pd
# Merge the two datasets on EmployeeID if not already merged
merged_data = df_employee.merge(df_PerformanceRating, on='EmployeeID')

# Convert the 'HireDate' and 'ReviewDate' to datetime format
merged_data['HireDate'] = pd.to_datetime(merged_data['HireDate'])
merged_data['ReviewDate'] = pd.to_datetime(merged_data['ReviewDate'])

# Calculate the number of employees where ReviewDate is after HireDate
employees_review_after_hire = merged_data[merged_data['ReviewDate'] < merged_data['HireDate']]
distinct_count = employees_review_after_hire['EmployeeID'].nunique()
# Count the number of such employees
num_employees_review_after_hire = distinct_count

print(f"Number of employees where the review date is after the hire date: {num_employees_review_after_hire}")
```

✓ 0.1s

Number of employees where the review date is after the hire date: 289

Some employees were found to have review dates in dates before their hiring dates ----- Did No Change

```
education_values = df_employee['EducationField'].unique()

# Display the unique values
print("Unique values in the 'Education' field:")
for value in education_values:
    print(value)

```

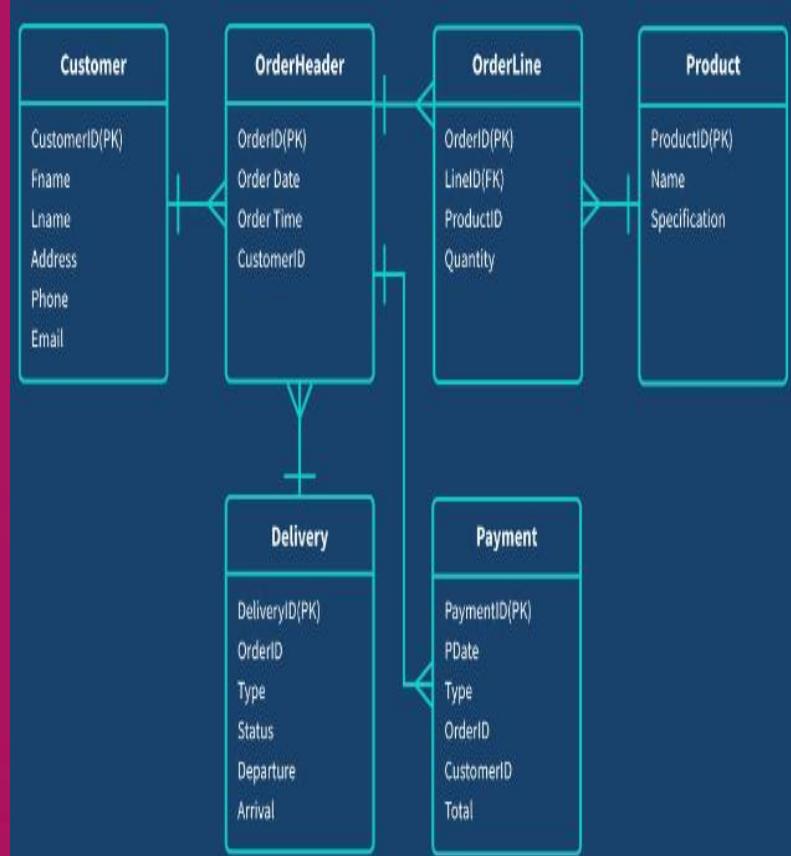
✓ 0.0s

Unique values in the 'Education' field:

Marketing
Marketing
Computer Science
Technical Degree
Information Systems
Other
Economics
Human Resources
Business Studies

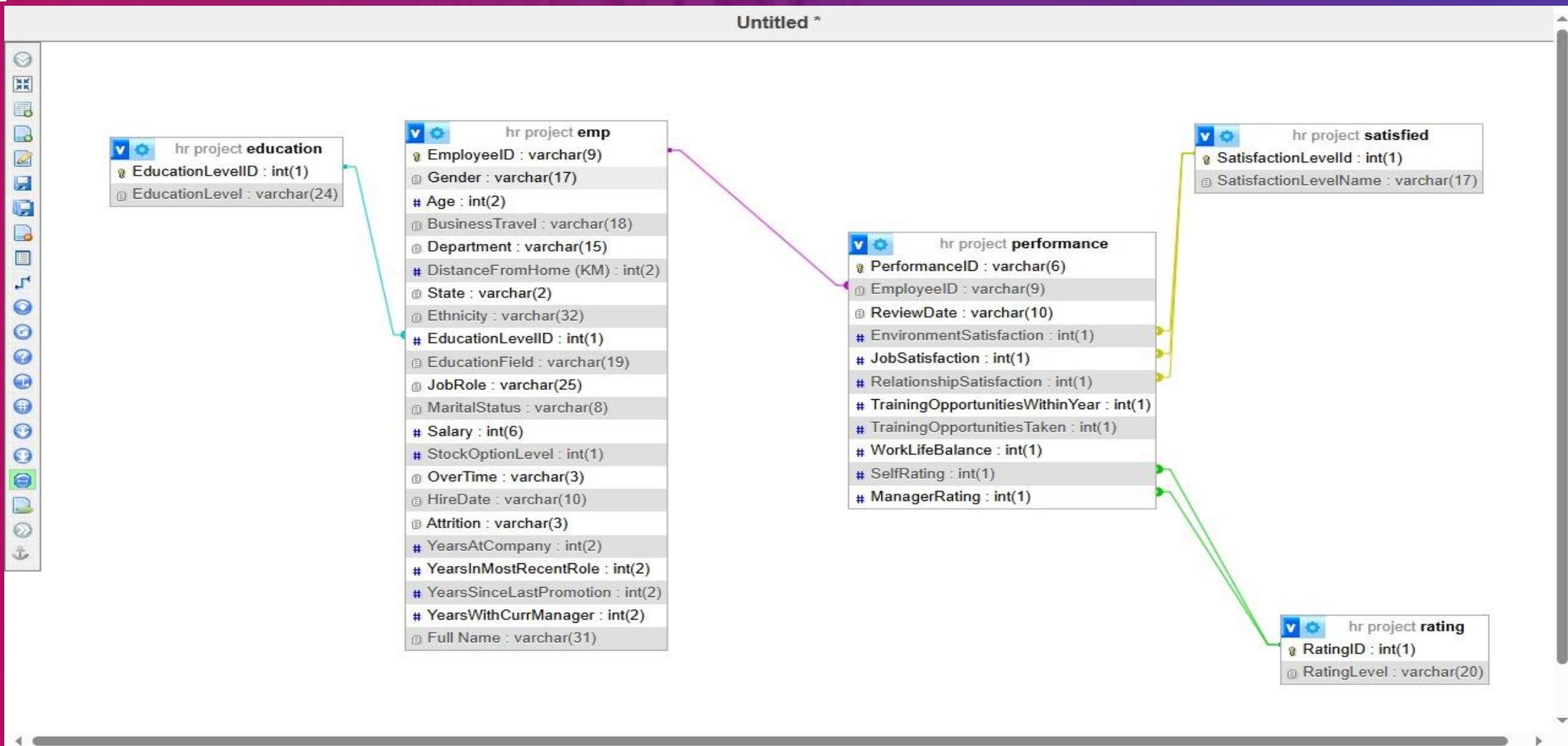
- Two education value fields were found to be marketing without any significant difference ----- Did no change.
- The Gender column in Employee Table had 2 items beside male and female and both of them could affect the total no. of males and females if known---- No change was done.
- It was unethical to put a column for Ethnicity as it was of no value for analytical purposes.
- There were many outliers between employees' salaries----- No change was done----- as we suggested it may be due to difference in levels of employees, as between general manager and senior staff and normal employees
- More sources of data were needed for better understanding, analysing and consequent decision making.

DATA MODELING

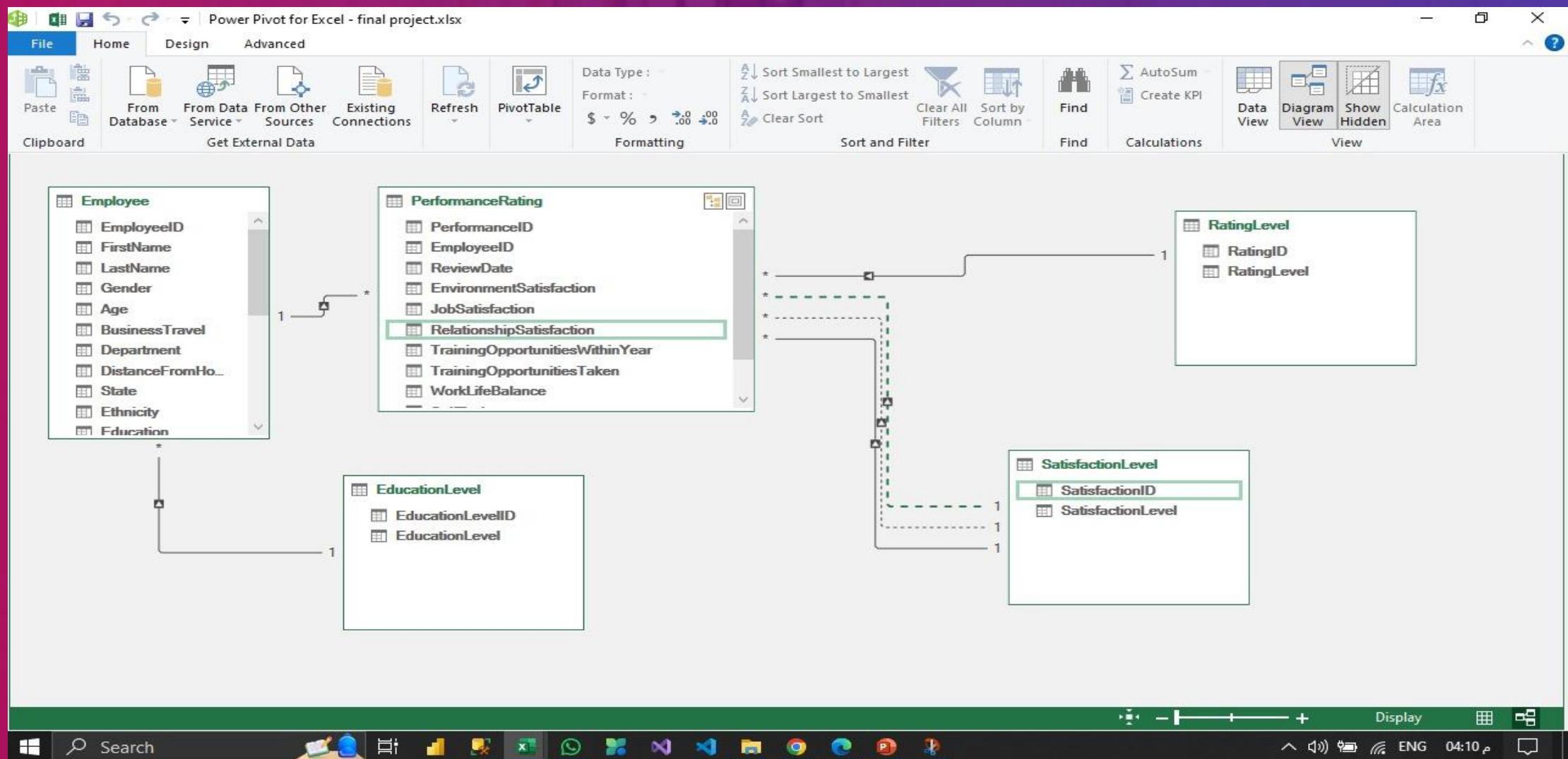


- Data modelling is the process of creating a diagram that represents your data system and defines the structure, attributes, and relationships of your data entities.
- Data modelling organizes and simplifies data in a way that makes it easy to understand, manage, and query, while also ensuring data integrity and consistency. Data models inform about data architecture, database design, and restructuring legacy systems.

Data Modeling Using SQL



Data Modeling Using Excel



DETERMINING DATA ANALYSIS QUESTIONS AND ANSWERING THEM



1. How many employees are in the company?

```
SELECT COUNT(*) AS "Total Employees" FROM employee;
```

[Extra options](#)

Total Employees

1470

2. What is the average age of employees?

```
SELECT ROUND(AVG(employee.Age),0) AS "Avarge of age " FROM employee;
```

Profiling | [Edit inline](#) | [Edit](#) | [Explain SQL](#) | [Create PHP code](#) | [Refresh](#)

Show all | Number of rows: 25

[Extra options](#)

Avarge of age

29

3. What is the average salary of employees by department?

```
SELECT Department, AVG(Salary) AS AvgSalary FROM Employee GROUP BY Department;
```

Profiling | [Ed](#)

Show all | Number of rows: 25

Filter rows:

Search this table

[Extra options](#)

		Department	Avg Salary	
<input type="checkbox"/>	Edit	Copy	Delete	Human Resources 119698.8095
<input type="checkbox"/>	Edit	Copy	Delete	Sales 119117.6099
<input type="checkbox"/>	Edit	Copy	Delete	Technology 109655.1228

4. What is the gender distribution among employees?

```
SELECT Gender, COUNT(*) AS TotalEmployees FROM Employee GROUP BY Gender;
```

Profiling

Show all

Number of rows:

25

Filter rows:

Search this table

Extra options

Gender	TotalEmployees
Female	675
Male	651
Non-Binary	124
Prefer Not To Say	20

5. How many employees have been with the company for more than 5 years?

```
SELECT COUNT(*) AS EmployeesOver5Years FROM Employee WHERE YearsAtCompany > 5;
```

Profiling [Edit]

Extra options

EmployeesOver5Years

587

6. What is the distribution of employees by job title?

```
SELECT JobRole, COUNT(*) AS TotalEmployees FROM Employee GROUP BY JobRole;
```

Profiling [Edit inline] [Edit] [Explain SQL] [Create PHP code] [Refresh]

Show all

Number of rows:

25 ▾

Filter rows:

Search this table

Extra options

JobRole	TotalEmployees
Analytics Manager	52
Data Scientist	261
Engineering Manager	75
HR Business Partner	7
HR Executive	28
HR Manager	4
Machine Learning Engineer	146
Manager	37
Recruiter	24
Sales Executive	327
Sales Representative	83
Senior Software Engineer	132
Software Engineer	294

```

import matplotlib.pyplot as plt
import seaborn as sns

# Your existing code
plt.figure(figsize=(12, 6))

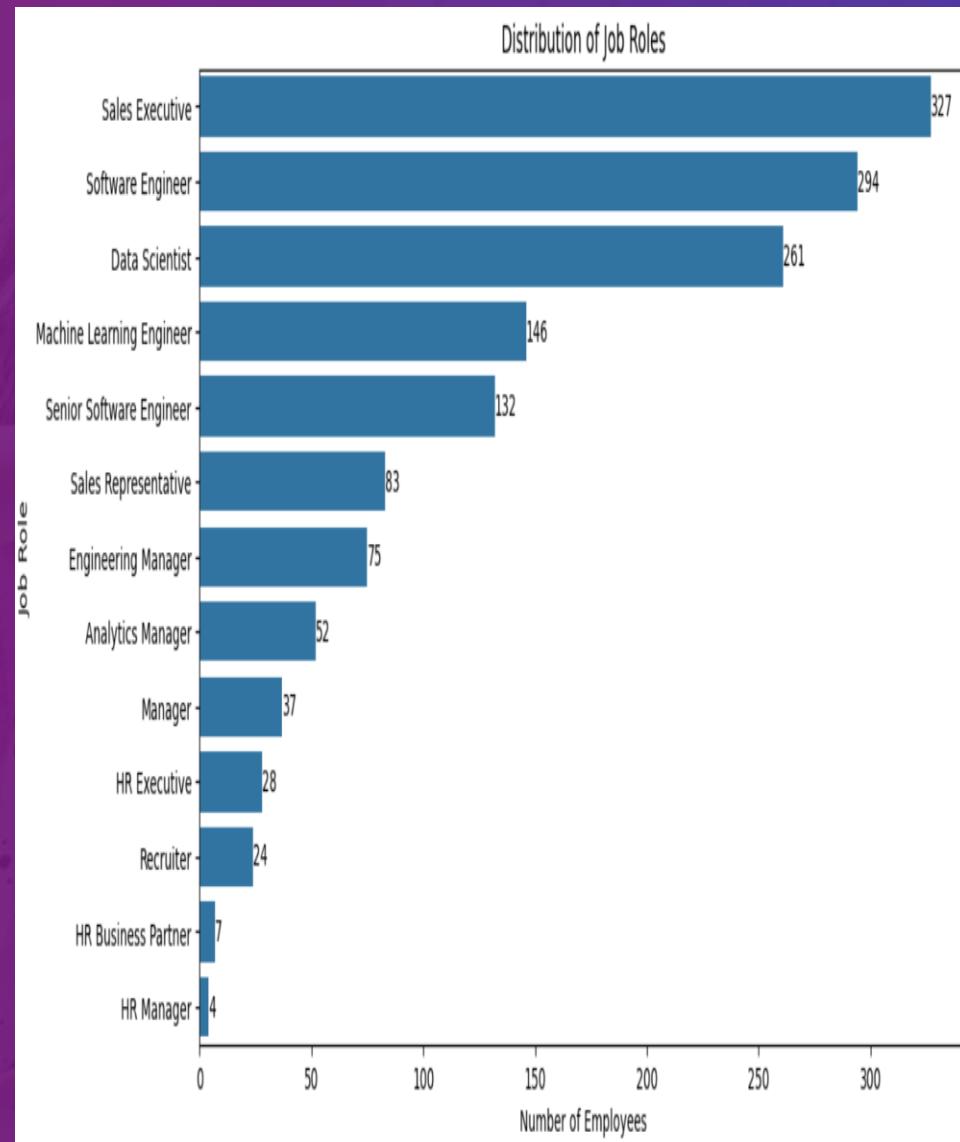
ax = sns.countplot(data=df_employee, y='JobRole', order=df_employee['JobRole'].value_counts().index)
plt.title('Distribution of Job Roles')
plt.xlabel('Number of Employees')
plt.ylabel('Job Role')

# Adding the numbers at the end of each bar
for container in ax.containers:
    ax.bar_label(container, fmt='%d') # '%d' for integer formatting

plt.show()

```

Python



7. What is the distribution of employees by age group?

```
SELECT CASE WHEN Age BETWEEN 20 AND 29 THEN '20-29' WHEN Age BETWEEN 30 AND 39 THEN '30-39' WHEN Age BETWEEN 40 AND 49 THEN '40-49' ELSE '50+'  
END AS AgeGroup, COUNT(*) AS TotalEmployees FROM Employee GROUP BY CASE WHEN Age BETWEEN 20 AND 29 THEN '20-29' WHEN Age BETWEEN 30 AND 39 THEN  
'30-39' WHEN Age BETWEEN 40 AND 49 THEN '40-49' ELSE '50+' END;
```

Profiling | [Edit inline](#) | [Edit](#) | [Explain SQL](#) | [Create PHP code](#) | [Refresh](#)

Show all | Number of rows: 25 ▾ Filter rows: Search this table

[Extra options](#)

AgeGroup	TotalEmployees
20-29	874
30-39	289
40-49	219
50+	88

8. What is the number of employees in each education level?

```
SELECT education.EducationLevel,COUNT(emp.EmployeeID)AS "Number of Employee"  
FROM emp INNER JOIN education ON emp.EducationLevelID =  
education.EducationLevelID GROUP BY education.EducationLevel;
```

Profiling | [Edit inline](#) | [Edit](#) | [Explain SQL](#) | [Create PHP code](#) | [Refresh](#)

Show all | Number of rows: 25 ▾

[Extra options](#)

EDUCATIONLEVEL	NUMBER OF EMPLOYEE
Bachelors	572
Doctorate	48
High School	282
Masters	398
No Formal Qualifications	170

9. What is the distribution of education levels across departments?

```
SELECT employee.Department, educationlevel.EducationLevel, COUNT(employee.EmployeeID) AS "number of emp" FROM employee INNER JOIN educationlevel ON employee.Education=educationlevel.EducationLevelID GROUP BY employee.Department, educationlevel.EducationLevel;
```

Profiling | [Edit inline](#) | [Edit](#) | [Explain SQL](#) | [Create PHP code](#) | [Refresh](#)

Show all | Number of rows: 25 ▾ Filter rows: Search this table

Extra options

Department	EducationLevel	number of emp
Human Resources	Bachelors	27
Human Resources	Doctorate	3
Human Resources	High School	13
Human Resources	Masters	15
Human Resources	No Formal Qualifications	5
Sales	Bachelors	166
Sales	Doctorate	15
Sales	High School	87
Sales	Masters	126
Sales	No Formal Qualifications	50
Technology	Bachelors	379
Technology	Doctorate	30
Technology	High School	182
Technology	Masters	256
Technology	No Formal Qualifications	115

10. What is the average performance rating of employees based on their education level?

```
SELECT educationlevel.EducationLevel, AVG(performancerating.ManagerRating) AS AvgManagerRating,  
AVG(performancerating.SelfRating) AS AvgSelfRating FROM performancerating JOIN employee ON  
performancerating.EmployeeID = employee.EmployeeID JOIN educationlevel ON employee.Education =  
educationlevel.EducationLevelID GROUP BY educationlevel.EducationLevel;
```

Profiling | [Edit inline](#) | [Edit](#) | [Explain SQL](#) | [Create PHP](#)

Show all | Number of rows: 25 ▾ Filter rows: Search this table

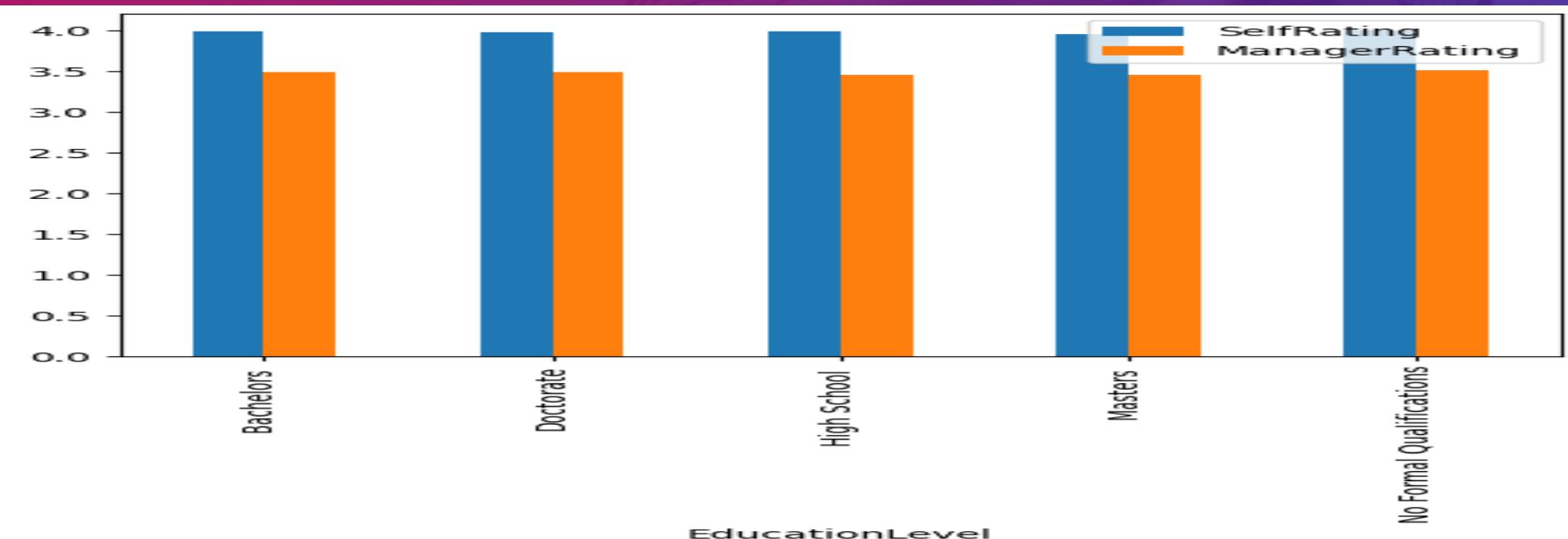
Extra options

EducationLevel	AvgManagerRating	Avg SelfRating
Bachelors	3.4833	3.9924
Doctorate	3.4882	3.9716
High School	3.4567	3.9904
Masters	3.4531	3.9582
No Formal Qualifications	3.5060	4.0024

```
merged_df = df_PerformanceRating.merge(df_employee, on='EmployeeID').merge(df_EducationLevel, on='EducationLevelID')
education_performance = merged_df.groupby('EducationLevel')[['SelfRating', 'ManagerRating']].mean().reset_index()
education_performance.plot(x='EducationLevel', kind='bar')
```

Python

```
<Axes: xlabel='EducationLevel'>
```



11. in Each Job role what is the number of Employee with attrition

```
SELECT employee.JobRole,COUNT(employee.EmployeeID) AS "number of Emp" FROM employee  
WHERE employeeAttrition="Yes" GROUP by employee.JobRole ORDER by "number of Emp";
```

Profiling [Edit inline] [Edit] [Explain SQL] [Create PHP code] [Re-

Show all

Number of rows:

25

Filter rows:

Search this table

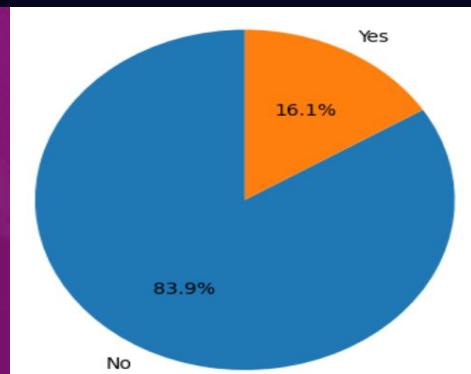
Extra options

JobRole	number of Emp
Engineering Manager	2
Manager	2
Analytics Manager	3
HR Executive	3
Recruiter	9
Senior Software Engineer	9
Machine Learning Engineer	10
Sales Representative	33
Software Engineer	47
Sales Executive	57
Data Scientist	62

What is the percentage of attrition among employees?

```
Attrition_counts = df_employee['Attrition'].value_counts()  
plt.pie(x=Attrition_counts.values, labels=Attrition_counts.index, autopct='%.1f%%', startangle=90)  
plt.show()
```

Python



12. What is the avg salary grouped by the education field?

```
SELECT EducationField, SUM(Salary) AS "sum of salary" FROM employee GROUP BY EducationField ORDER BY "sum of salary" DESC;
```

Profiling [Edit inline] [Edit] [Explain SQL] [Create PHP code]

Show all

Number of rows:

25



Filter rows:

Search this table

Extra options



EducationField

sum of salary

v 1

	Edit	Copy	Delete	EducationField	sum of salary
<input type="checkbox"/>				Computer Science	48115757
<input type="checkbox"/>				Information Systems	41520135
<input type="checkbox"/>				Marketing	40390223
<input type="checkbox"/>				Economics	11334205
<input type="checkbox"/>				Business Studies	9250227
<input type="checkbox"/>				Other	7900032
<input type="checkbox"/>				Human Resources	3930278
<input type="checkbox"/>				Technical Degree	3605195

13. What is the distribution of employees' education levels by gender?

```
SELECT educationlevel.EducationLevel,Gender,COUNT(employee.EmployeeID) AS 'number of emp' FROM employee INNER JOIN educationlevel on employee.Education=educationlevel.EducationLevelID GROUP BY educationlevel.EducationLevel,Gender;
```

Profiling [Edit inline] [Edit] [Explain SQL] [Create PHP code] [Refresh]

Show all

Number of rows:

25



Filter rows:

Search this table

Extra options

EducationLevel

Gender

number of emp

Bachelors	Female	269
Bachelors	Male	254
Bachelors	Non-Binary	43
Bachelors	Prefer Not To Say	6
Doctorate	Female	21
Doctorate	Male	23
Doctorate	Non-Binary	4
High School	Female	136
High School	Male	119
High School	Non-Binary	23
High School	Prefer Not To Say	4
Masters	Female	164
Masters	Male	190
Masters	Non-Binary	37
Masters	Prefer Not To Say	7
No Formal Qualifications	Female	85
No Formal Qualifications	Male	65
No Formal Qualifications	Non-Binary	17
No Formal Qualifications	Prefer Not To Say	3

14. Which department has the highest average manager performance rating?

```
SELECT Department, AVG(ManagerRating) AS AvgPerformanceRating FROM employee e JOIN performanceRating p ON e.EmployeeID = p.EmployeeID GROUP BY Department ORDER BY AvgPerformanceRating DESC;
```

Profiling [Edit inline] [Edit] [Explain SQL] [Create PHP code] [R

Show all

Number of rows:

25 ▾

Filter rows:

Search this table

Extra options

Department	AvgPerformanceRating
Technology	3.4874
Sales	3.4500
Human Resources	3.4422

15. Is there a significant difference in performance ratings between Genders of employees?

```
SELECT e.Gender, AVG(p.ManagerRating) AS AvgPerformanceRating FROM Employee e JOIN PerformanceRating p ON e.EmployeeID = p.EmployeeID GROUP BY e.Gender;
```

Profiling [Edit inline] [Edit] [Explain SQL] [Create PHP code] [R

Show all

Number of rows:

25 ▾

Filter rows:

Search this table

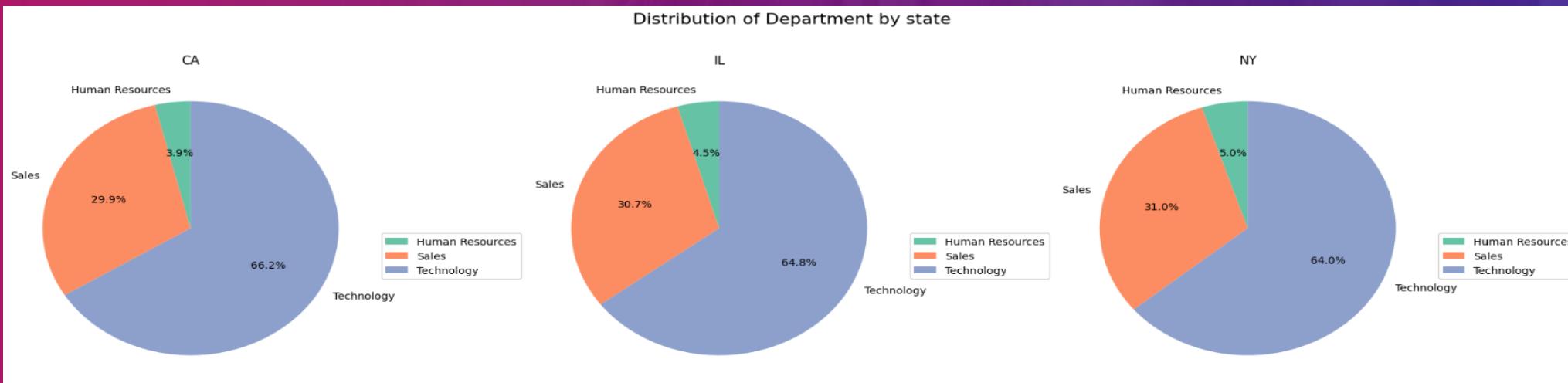
Extra options

Gender	AvgPerformanceRating
Female	3.4866
Male	3.4609
Non-Binary	3.4822
Prefer Not To Say	3.3279

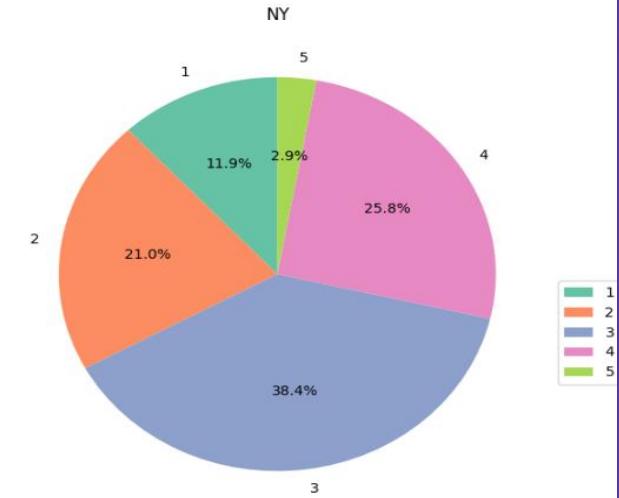
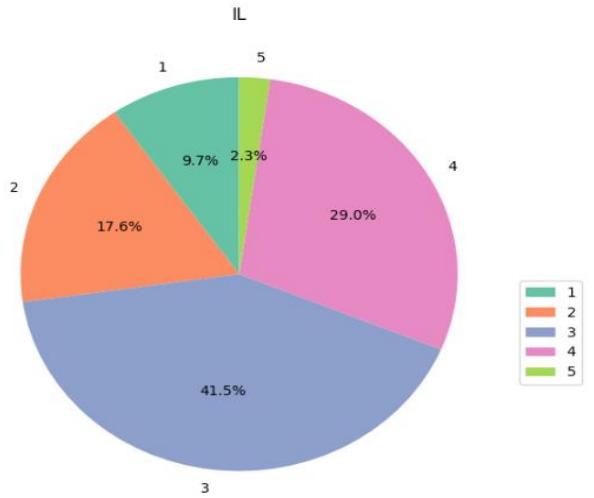
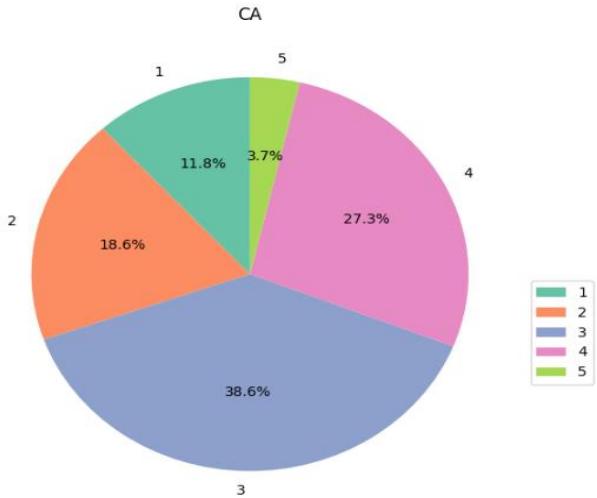
What Is The Effect Of The State On Other Parameters?

```
def plot_pie_chart_by_state(feature):
    colors = ['#66c2a5', '#fc8d62', '#8da0cb', '#e78ac3', '#a6d854']
    fig, axes = plt.subplots(1, 3, figsize=(20, 6), subplot_kw=dict(aspect="equal"))
    grouped_data = df_employee.groupby([feature, 'State']).size().unstack().T
    for ax, (gender, data) in zip(axes, grouped_data.iterrows()):
        data.plot(kind='pie', ax=ax, autopct='%.1f%%', startangle=90, colors=colors)
        ax.set_title(gender)
        ax.set_ylabel('') # Hide the y-label
        ax.legend(loc='best', bbox_to_anchor=(1, 0.5))
    plt.suptitle(f'Distribution of {feature} by state', fontsize=16)
    plt.tight_layout()
    plt.show()

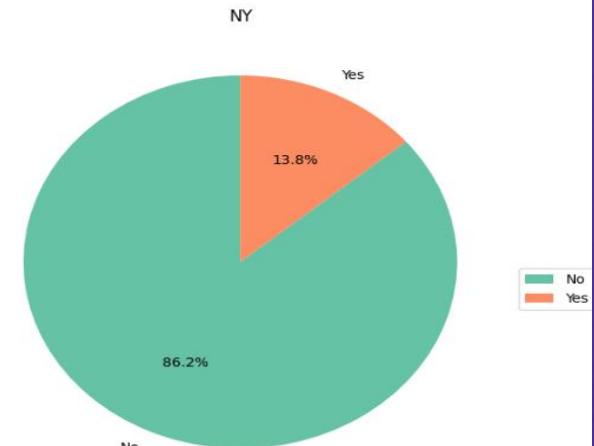
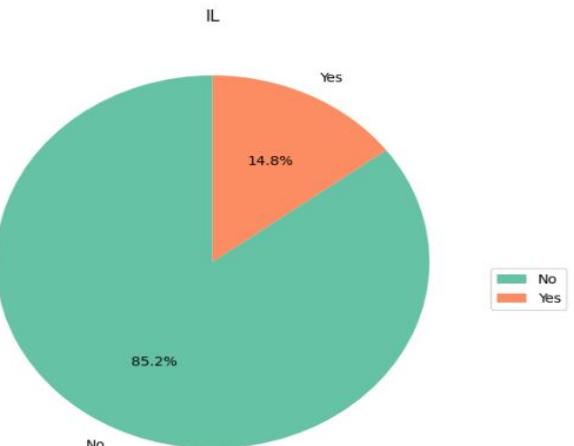
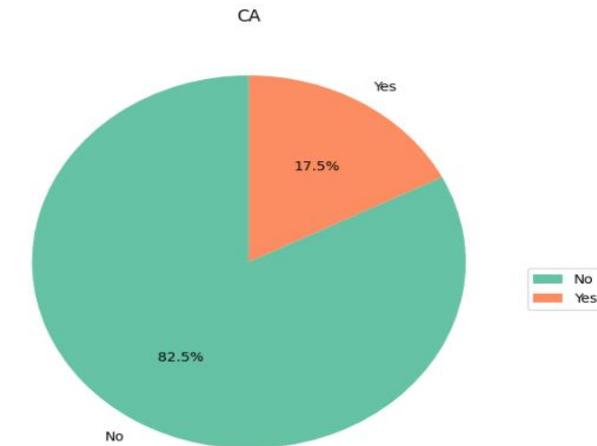
for feature in ['Department', 'EducationLevelID', 'Attrition', 'BusinessTravel']:
    plot_pie_chart_by_state(feature)
```



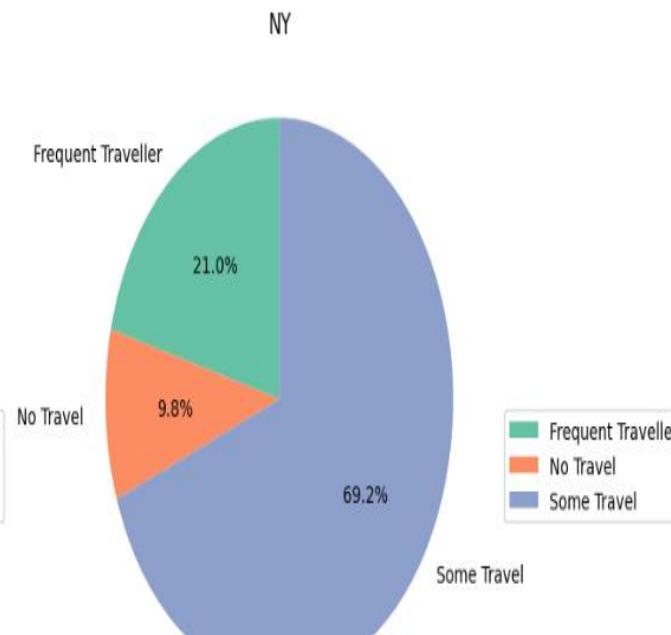
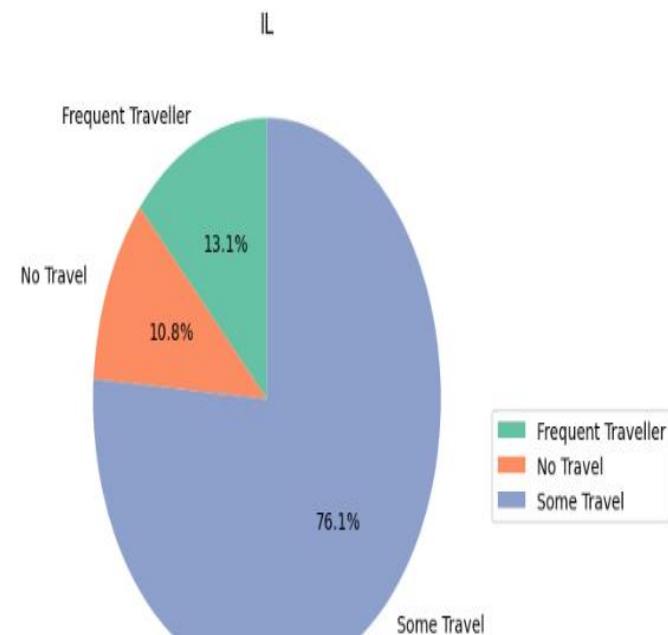
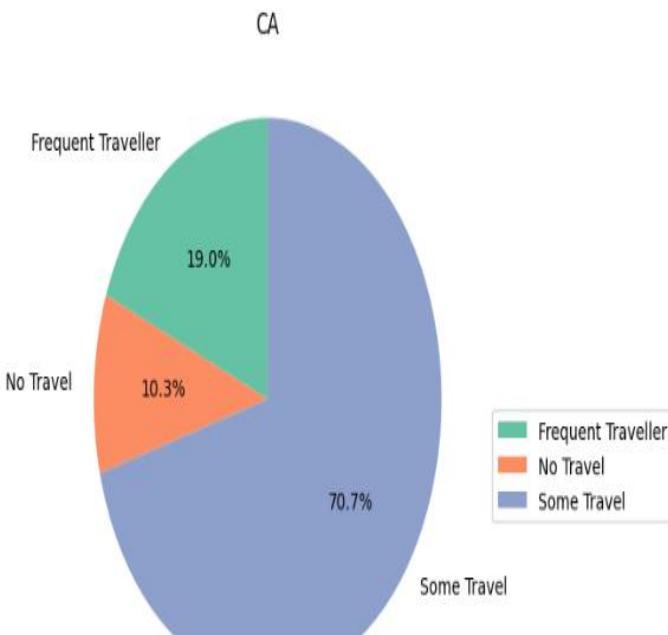
Distribution of EducationLevelID by state



Distribution of Attrition by state



Distribution of BusinessTravel by state



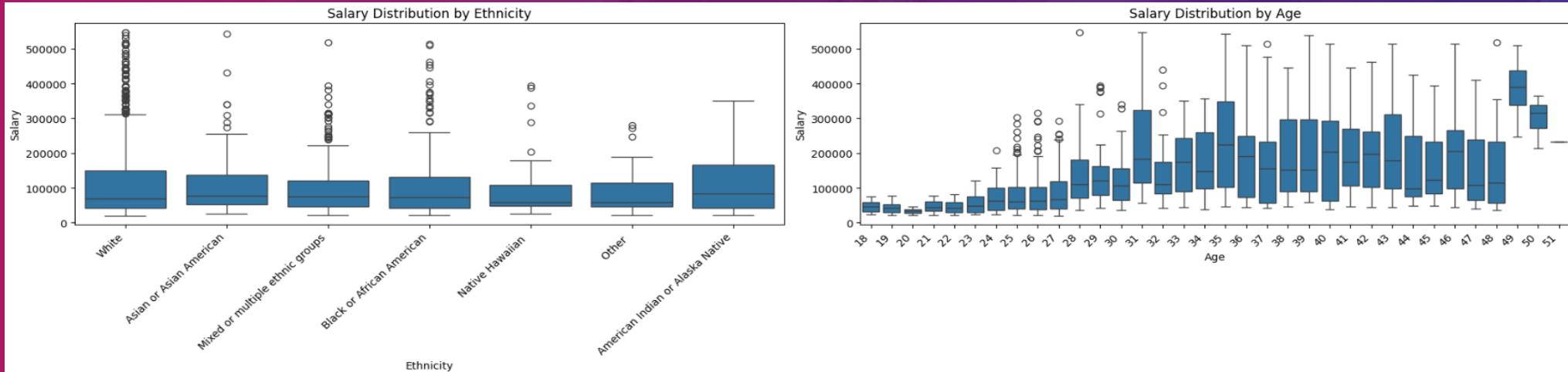
What Affects Employees Salaries?

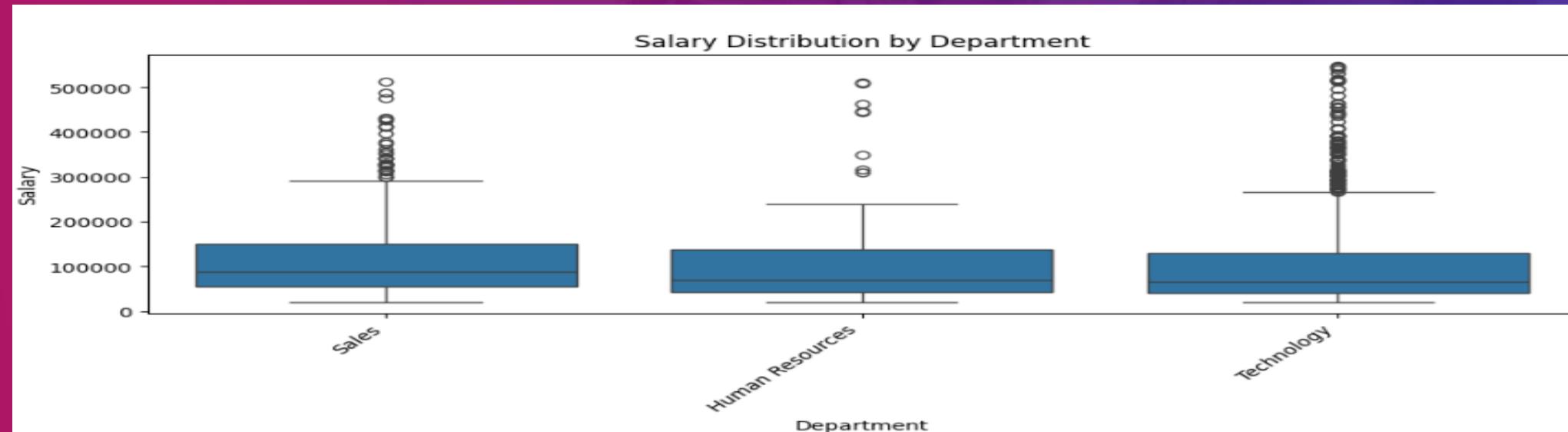
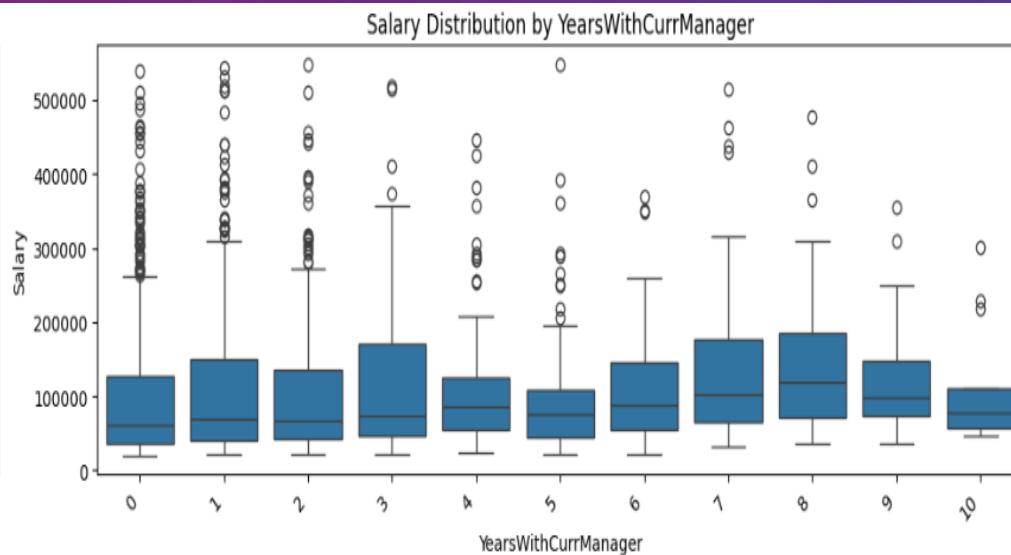
```
features_to_compare = ['Ethnicity', 'Age', 'YearsInMostRecentRole', 'YearsWithCurrManager', 'Department']
plt.figure(figsize=(20, 15))

for i, feature in enumerate(features_to_compare, 1):
    plt.subplot(3, 2, i) # Adjusted to 3 rows, 2 columns to fit 5 plots
    sns.boxplot(x=feature, y='Salary', data=df_employee)
    plt.title(f'Salary Distribution by {feature}')
    plt.xticks(rotation=45, ha='right')
    plt.xlabel(feature)
    plt.ylabel('Salary')

plt.tight_layout()
plt.show()
```

Python





12. What is the avg salary grouped by the education field?

```
SELECT emp.EducationField,AVG(emp.Salary) AS"Avg of salary" FROM emp GROUP BY  
emp.EducationField ORDER BY `Avg of salary` DESC;
```

Profiling [Edit inline] [Edit] [Explain SQL] [Create PHP code] [Refresh]

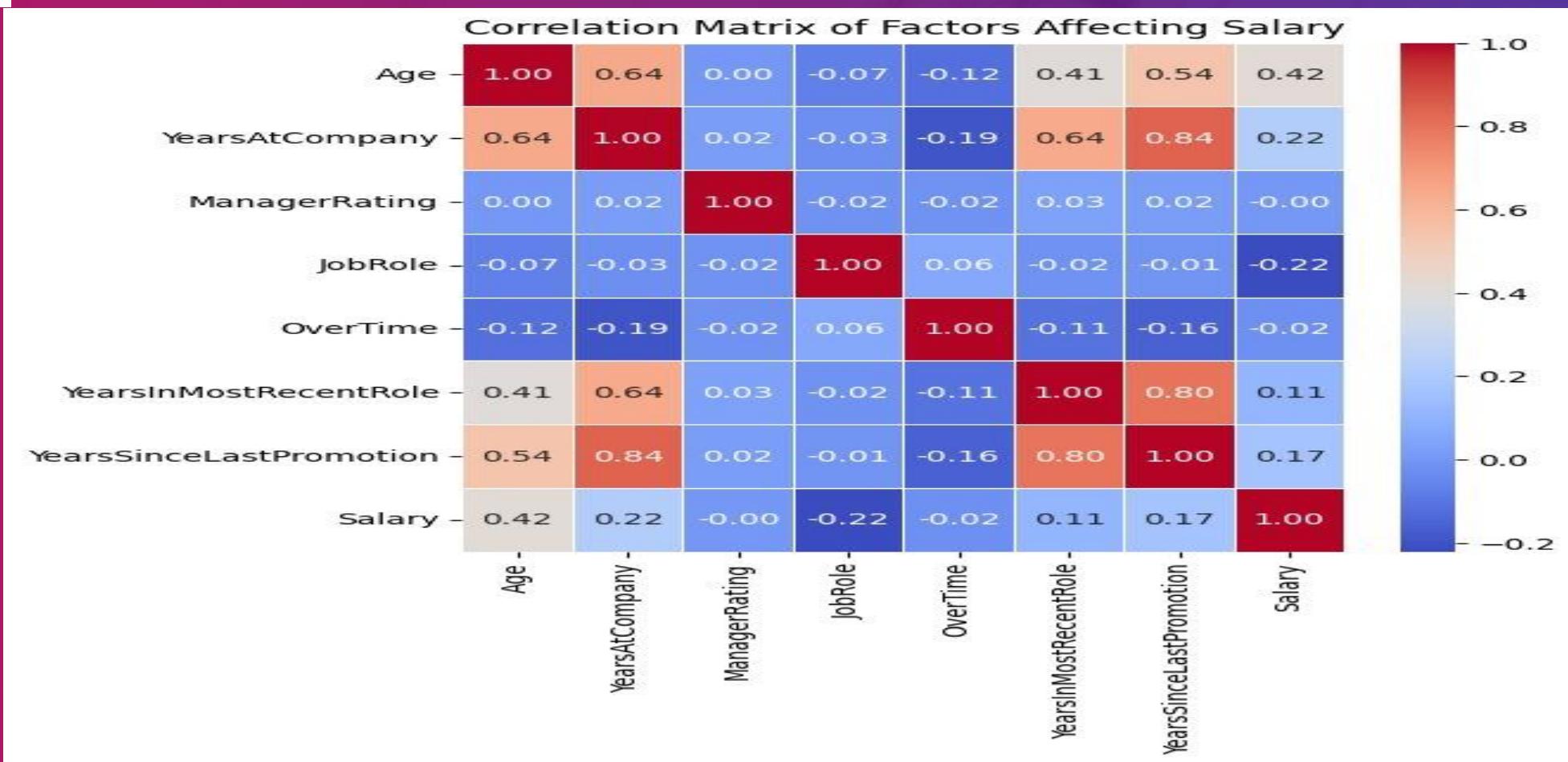
Show all

Number of rows:

25 ✓

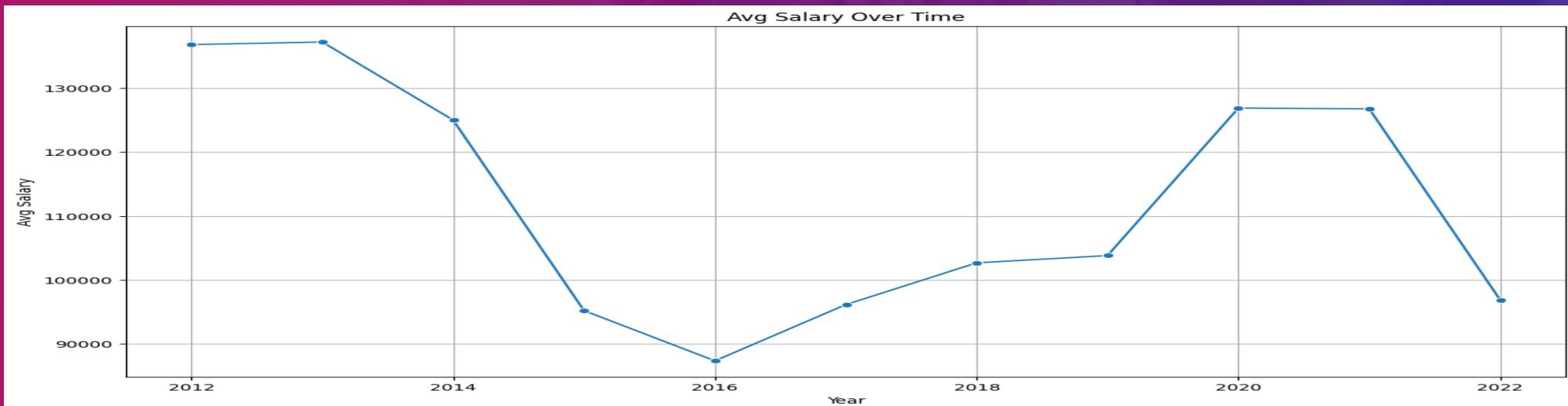
Extra options

				EducationField	Avg of salary
<input type="checkbox"/>	 Edit	 Copy	 Delete	Human Resources	145565.8519
<input type="checkbox"/>	 Edit	 Copy	 Delete	Marketing	124277.6092
<input type="checkbox"/>	 Edit	 Copy	 Delete	Information Systems	114380.5372
<input type="checkbox"/>	 Edit	 Copy	 Delete	Economics	112219.8515
<input type="checkbox"/>	 Edit	 Copy	 Delete	Computer Science	109353.9932
<input type="checkbox"/>	 Edit	 Copy	 Delete	Business Studies	98406.6702
<input type="checkbox"/>	 Edit	 Copy	 Delete	Other	96341.8537
<input type="checkbox"/>	 Edit	 Copy	 Delete	Technical Degree	94873.5526



What Is The Average Of Salaries Overtime?

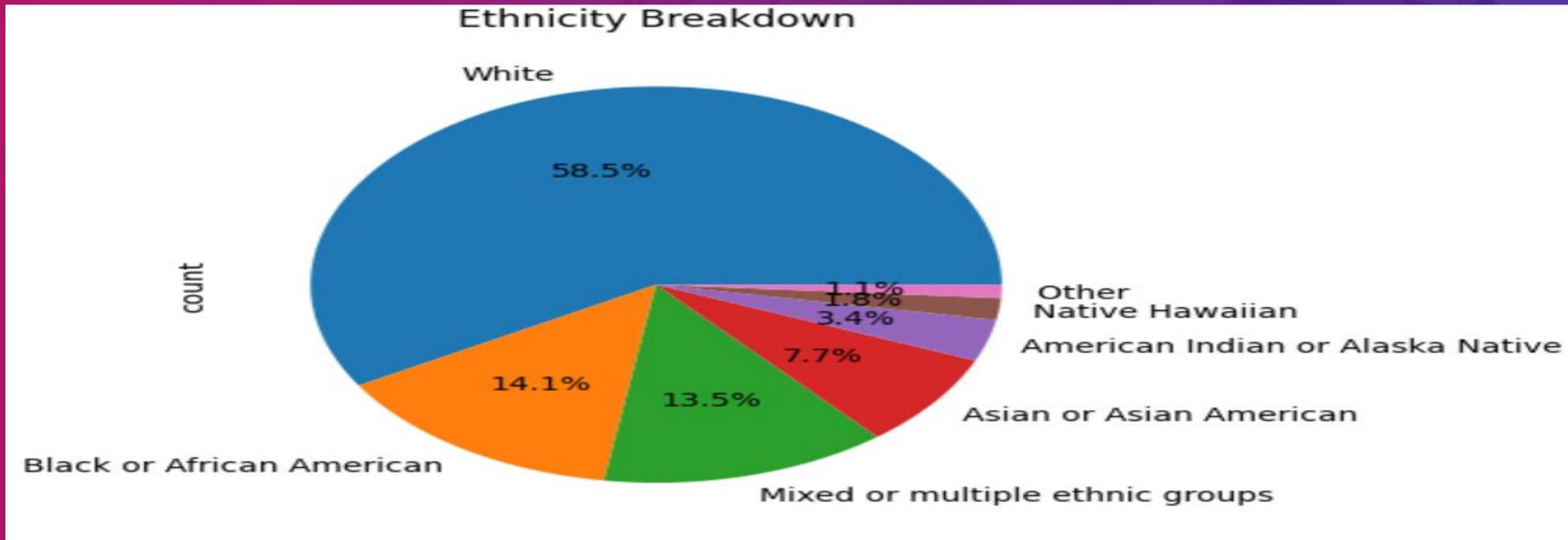
```
salary_over_time = df_employee.groupby(df_employee['HireDate'].dt.year).agg({
    'Salary': 'mean'
}).reset_index()
plt.figure(figsize=(14, 7))
sns.lineplot(x='HireDate', y='Salary', data=salary_over_time, marker='o')
plt.title('Avg Salary Over Time')
plt.xlabel('Year')
plt.ylabel('Avg Salary')
plt.grid(True)
plt.show()
```



How Are Employees Distributed According To Ethnicity?

```
df_employee['Ethnicity'].value_counts().plot(kind='pie', autopct='%1.1f%%', title='Ethnicity Breakdown')  
plt.show()
```

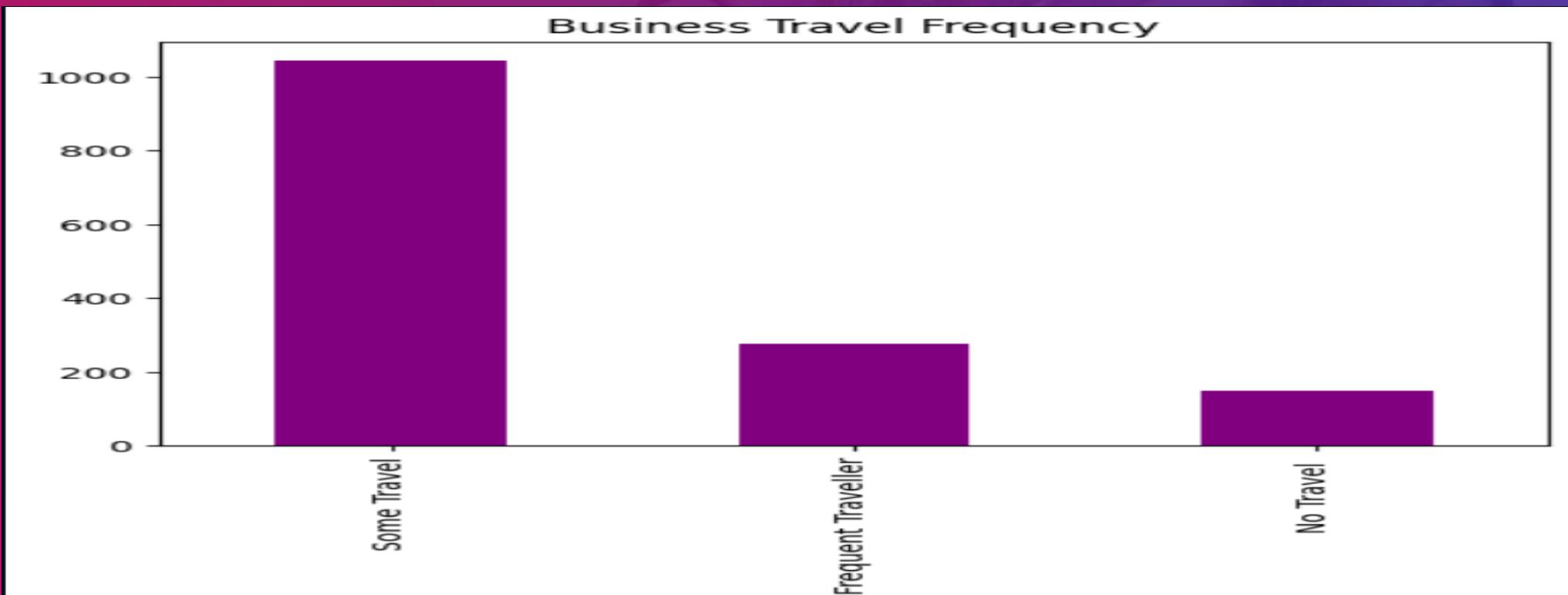
Python



Which Departments Have More Business Travels

```
df_employee['BusinessTravel'].value_counts().plot(kind='bar', title='Business Travel Frequency', color='purple');
```

Python



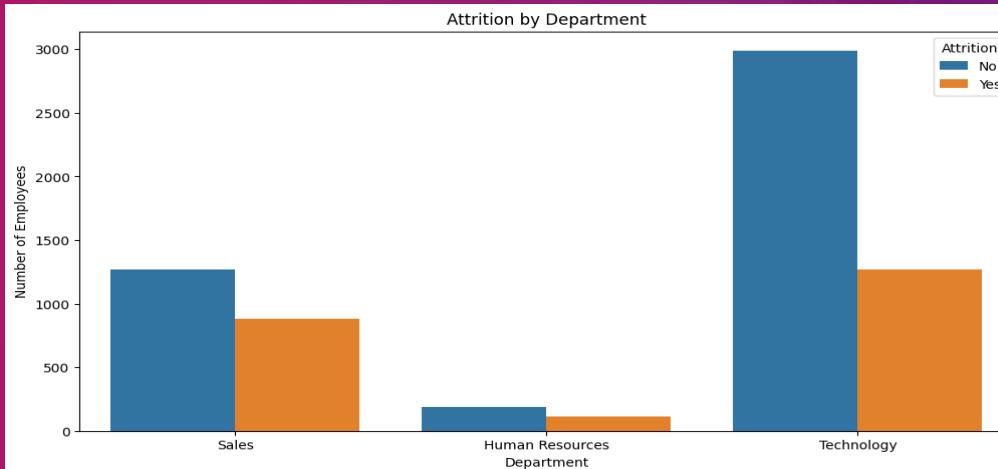
What Are The Main Causes Of Employees Attrition

```
plt.figure(figsize=(12, 6))
sns.boxplot(data=merged_data, x='Attrition', y='Salary')
plt.title('Salary Distribution by Attrition')
plt.xlabel('Attrition')
plt.ylabel('Salary')
plt.show()
```

Python

```
plt.figure(figsize=(12, 6))
sns.countplot(data= merged_data, x='Department', hue='Attrition')
plt.title('Attrition by Department')
plt.xlabel('Department')
plt.ylabel('Number of Employees')
plt.legend(title='Attrition', loc='upper right')
plt.show()
```

Python

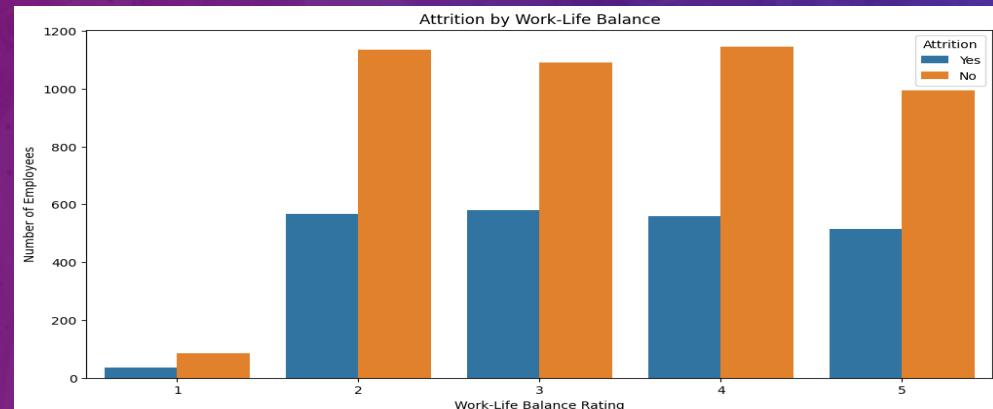


```
# Plotting training opportunities vs attrition
plt.figure(figsize=(12, 6))
sns.countplot(data=merged_data, x='TrainingOpportunitiesTaken', hue='Attrition')
plt.title('Attrition by Training Opportunities Taken')
plt.xlabel('Training Opportunities Taken')
plt.ylabel('Number of Employees')
plt.legend(title='Attrition', loc='upper right')
plt.show()
```

Python

```
plt.figure(figsize=(12, 6))
sns.countplot(data=merged_data, x='WorkLifeBalance', hue='Attrition')
plt.title('Attrition by Work-Life Balance')
plt.xlabel('Work-Life Balance Rating')
plt.ylabel('Number of Employees')
plt.legend(title='Attrition', loc='upper right')
plt.show()
```

Python



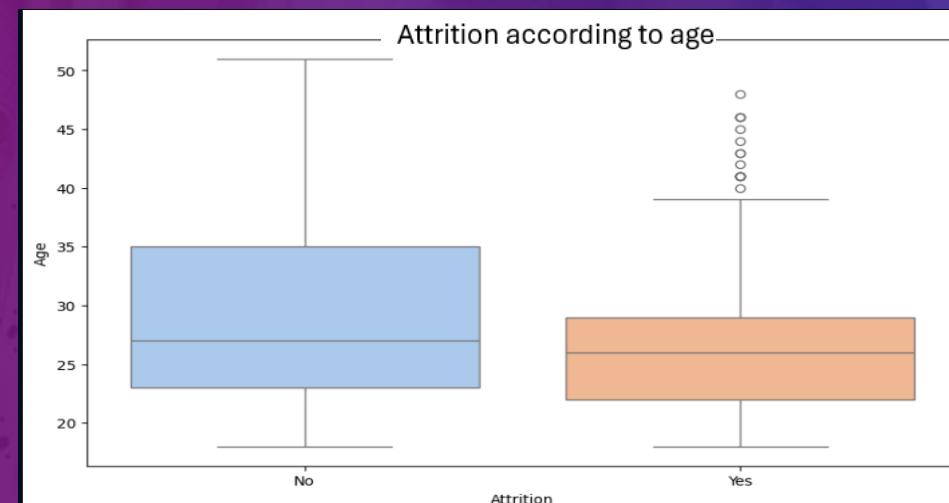
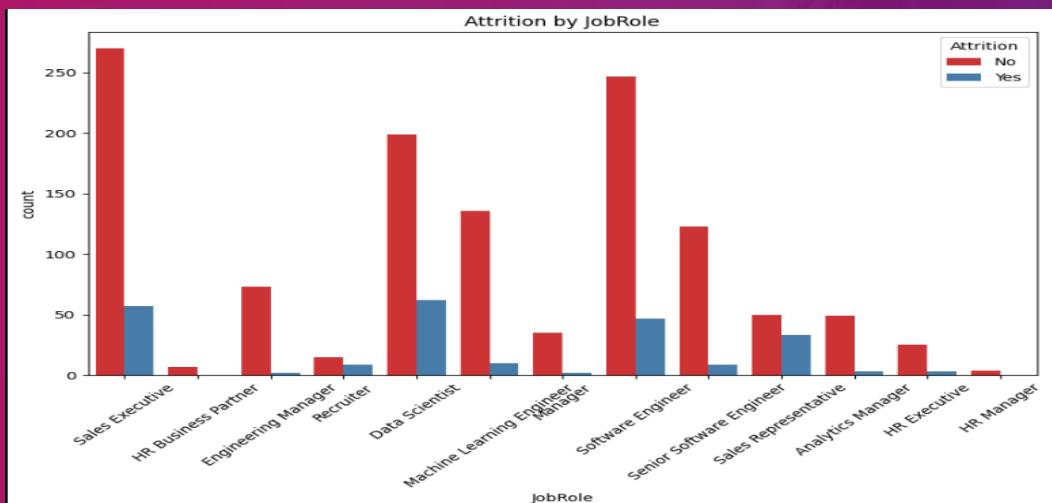
What Is the effect of JobRole & Age On Attrition?

```
plt.figure(figsize=(10, 6))
sns.countplot(x='JobRole', hue='Attrition', data=df_employee, palette='Set1')
plt.title('Attrition by JobRole')
plt.xticks(rotation=45)
plt.show()
```

Python

```
plt.figure(figsize=(10, 6))
sns.boxplot(x='Attrition', y='Age', data=df_employee, palette='pastel')
plt.title('Attrition according to age')
plt.show()
```

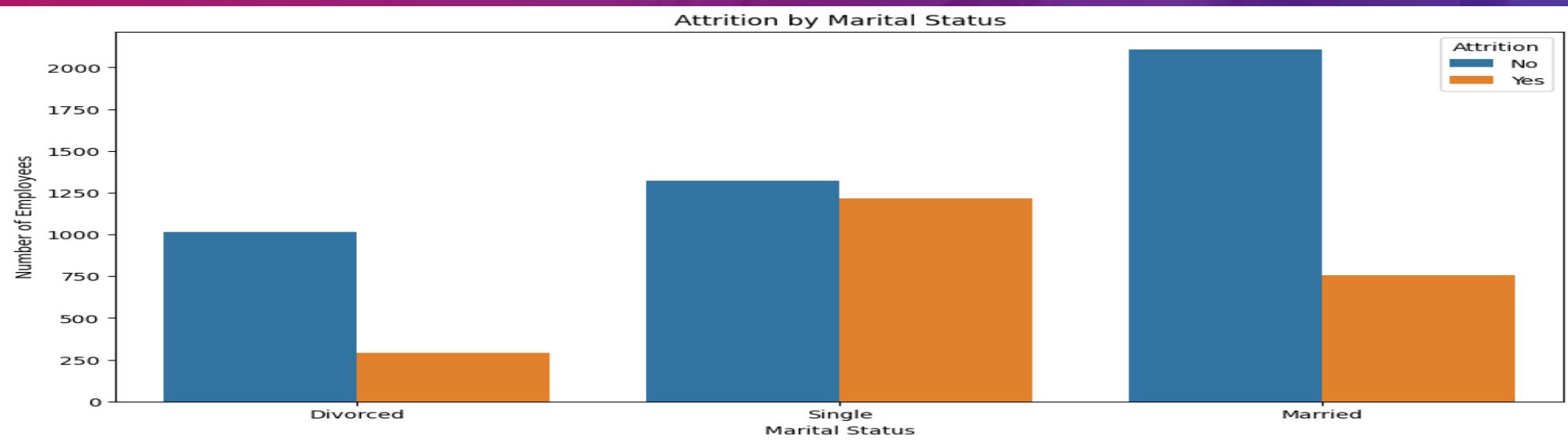
Python



What Is The Effect Of Marital Status On Attrition?

```
# Plotting marital status vs. attrition
plt.figure(figsize=(12, 6))
sns.countplot(data=merged_data, x='MaritalStatus', hue='Attrition')
plt.title('Attrition by Marital Status')
plt.xlabel('Marital Status')
plt.ylabel('Number of Employees')
plt.legend(title='Attrition', loc='upper right')
plt.show()
```

Python



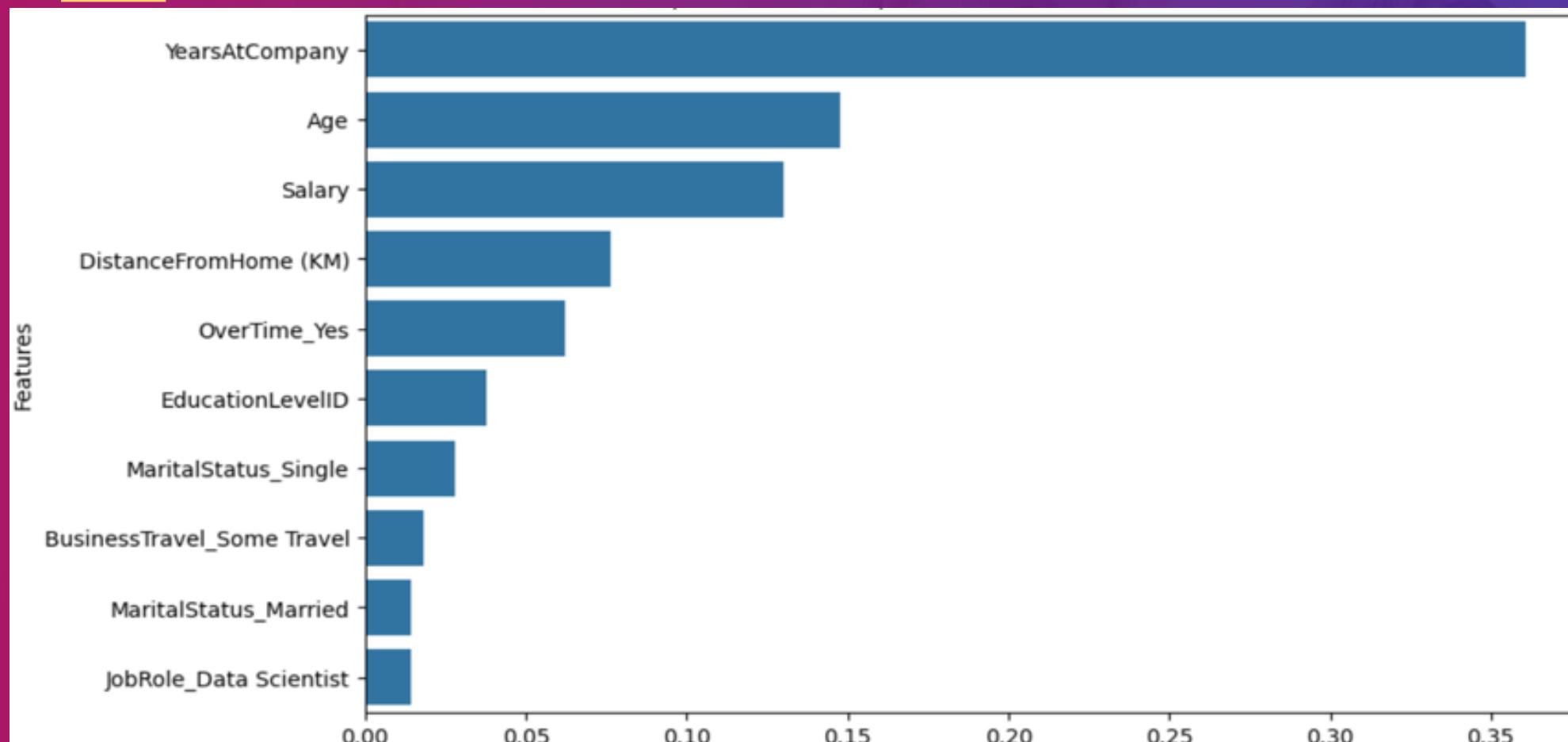
```
average_age_leavers = merged_data[merged_data['Attrition'] == 'Yes']['Age'].mean()  
print(f"Average age of employees who leave: {average_age_leavers:.2f} years")
```

Python

```
Average age of employees who leave: 26.63 years
```

Comment:

The Main Causes of Attrition are:



What Is The Percentage Of Employees With Stock Options?

```
# Calculate percentage of employees with stock options  
stock_option_percentage = (merged_data['StockOptionLevel'].value_counts(normalize=True) * 100).get(1, 0)  
print(f"Percentage of employees with stock options: {stock_option_percentage:.2f}%")
```

Python

Percentage of employees with stock options: 37.50%

What Is The Average Job Satisfaction By Department?

```
plt.figure(figsize=(12, 6))

avg_job_satisfaction = df_PerformanceRating.groupby('EmployeeID')['JobSatisfaction'].mean().reset_index()

merged_df = pd.merge(df_employee, avg_job_satisfaction, on='EmployeeID')

ax=sns.barplot(data=merged_df, x='Department', y='JobSatisfaction')

plt.title('Average Job Satisfaction by Department')

plt.xlabel('Department')

plt.ylabel('Average Job Satisfaction')

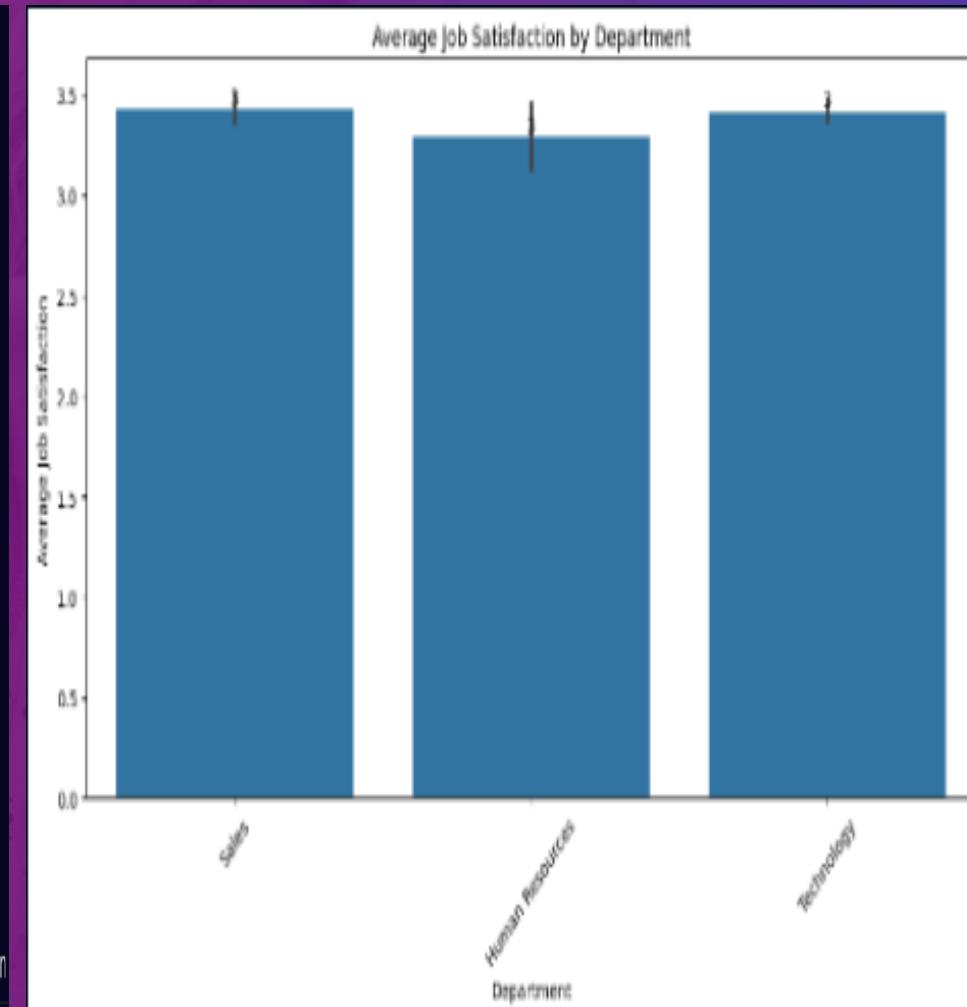
plt.xticks(rotation=45)

for container in ax.containers:

    ax.bar_label(container, fmt='%d')

plt.show()
```

Python



Which Ages Have More Job Satisfaction

```
plt.figure(figsize=(10, 6))
sns.boxplot(data=merged_data, x=pd.cut(merged_data['Age'], bins=5), y='JobSatisfaction')
plt.title('Job Satisfaction across Age Groups')
plt.xlabel('Age Groups')
plt.ylabel('Job Satisfaction')
plt.show()
```

Python

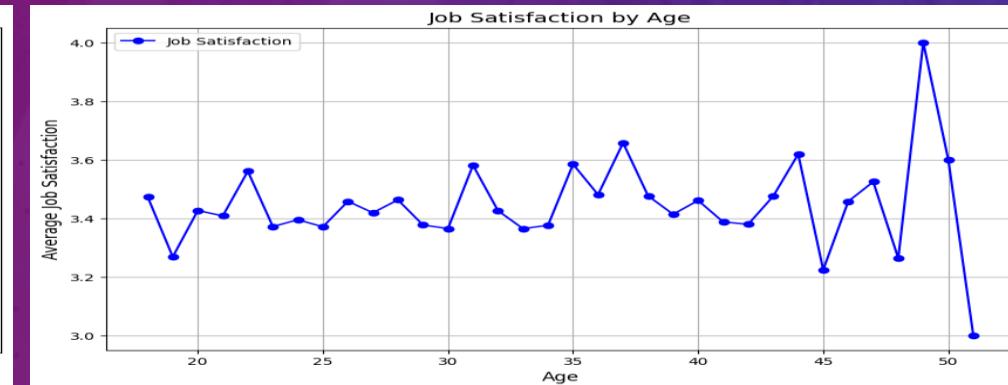
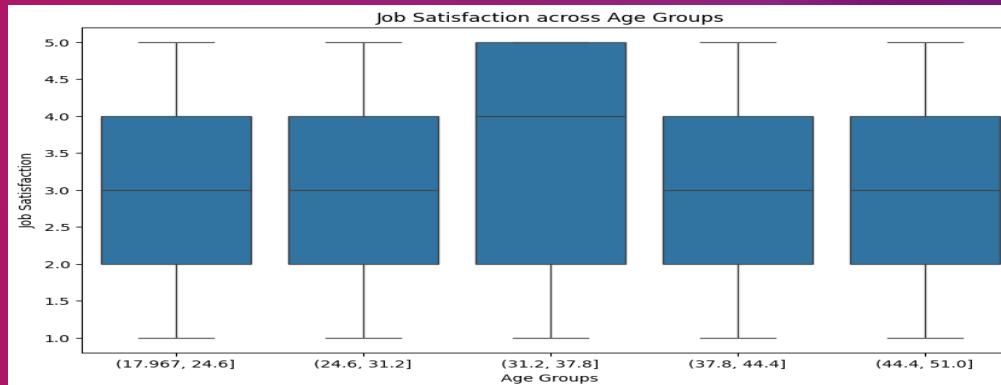
```
# Group by Age and calculate the average Job Satisfaction
age_job_satisfaction = merged_data.groupby('Age')['JobSatisfaction'].mean().reset_index()

# Plot a line chart for Age vs Average Job Satisfaction
plt.figure(figsize=(10, 6))
plt.plot(age_job_satisfaction['Age'], age_job_satisfaction['JobSatisfaction'], marker='o', color='b', label='Job Satisfaction')

# Add labels and title
plt.title('Job Satisfaction by Age', fontsize=14)
plt.xlabel('Age', fontsize=12)
plt.ylabel('Average Job Satisfaction', fontsize=12)
plt.grid(True)
plt.legend()

# Show the plot
plt.show()
```

Python



16. Avarage job Sataisfaction for each age group for all departments

```
SELECT emp.Department, emp.AgeGroup, AVG(pr.JobSatisfaction) AS "Avg JobSatisfaction" FROM employee AS emp JOIN performanceRating AS pr ON emp.EmployeeID = pr.EmployeeID GROUP BY Department, AgeGroup;
```

Profiling | [Edit mine](#) | [Edit](#) | [Explain SQL](#) | [Create PHP code](#) | [Refresh](#)

Show all | Number of rows: 25 ▾ Filter rows: Search this table

Extra options

Department	AgeGroup	Avg Job Satisfaction
Human Resources	[18-24]	3.1639
Human Resources	[25-31]	3.5292
Human Resources	[32-38]	3.3824
Human Resources	[39-45]	3.4067
Human Resources	[46-52]	3.5000
Sales	[18-24]	3.4021
Sales	[25-31]	3.4306
Sales	[32-38]	3.4864
Sales	[39-45]	3.3624
Sales	[46-52]	3.4545
Technology	[18-24]	3.4526
Technology	[25-31]	3.4088
Technology	[32-38]	3.4752
Technology	[39-45]	3.4625
Technology	[46-52]	3.3556

17. Avarage performance rating for all education levels

```
SELECT el.EducationLevel, (AVG(pr.SelfRating) + AVG(pr.ManagerRating)) / 2 AS "Avg PerformanceRating" FROM employee AS emp JOIN educationLevel AS el ON emp.Education = el.EducationLevelID JOIN performanceRating AS pr ON emp.EmployeeID = pr.EmployeeID GROUP BY el.EducationLevel ORDER BY "Avg PerformanceRating" DESC;
```

Profiling | [Edit mine](#) | [Edit](#) | [Explain SQL](#) | [Create PHP code](#) | [Refresh](#)

Show all | Number of rows: 25 ▾ Filter rows: Search this table

Extra options

EducationLevel	Avg PerformanceRating
No Formal Qualifications	3.75416667
Bachelors	3.73795118
Doctorate	3.72985712
High School	3.72353812
Masters	3.70561865

18. Average salary by Ethnicity

```
SELECT emp.Ethnicity ,AVG(emp.Salary) AS"Avg Salary by Ethnicity" FROM emp GROUP BY emp.Ethnicity ORDER by `Avg Salary by Ethnicity` DESC;
```

Profiling [Edit inline] [Edit] [Explain SQL] [Create PHP code] [Refresh]

Show all

Number of rows:

25

Extra options

		Ethnicity	Avg Salary by Ethnicity	
<input type="checkbox"/>	 Edit	 Copy	 Delete	White 115316.9349
<input type="checkbox"/>	 Edit	 Copy	 Delete	Native Hawaiian 115274.1923
<input type="checkbox"/>	 Edit	 Copy	 Delete	Black or African American 112176.8019
<input type="checkbox"/>	 Edit	 Copy	 Delete	American Indian or Alaska Native 112037.2800
<input type="checkbox"/>	 Edit	 Copy	 Delete	Asian or Asian American 109850.9646
<input type="checkbox"/>	 Edit	 Copy	 Delete	Mixed or multiple ethnic groups 106132.8485
<input type="checkbox"/>	 Edit	 Copy	 Delete	Other 101652.1250

19. Average job satisfaction by Ethnicity

Showing rows 0 - 6 (/ total, query took 0.1285 seconds.)

```
SELECT emp.Ethnicity ,AVG(performance.JobSatisfaction) AS"Avg job satisfaction by Ethnicity" FROM emp INNER JOIN performance ON emp.EmployeeID=performance.EmployeeID GROUP BY emp.Ethnicity ORDER by `Avg job satisfaction by Ethnicity` DESC;
```

Profiling [Edit inline] [Edit] [Explain SQL] [Create PHP code] [Refresh]

Show all

Number of rows:

25

Extra options

Ethnicity	Avg job satisfaction by Ethnicity
American Indian or Alaska Native	3.5211
Black or African American	3.4631
Asian or Asian American	3.4275
Other	3.4253
White	3.4240
Mixed or multiple ethnic groups	3.4086
Native Hawaiian	3.3567

Monthly Average Manager Rating Trend

```
merged_df['ReviewDate'] = pd.to_datetime(merged_df['ReviewDate'])
trend = merged_df.groupby(merged_df['ReviewDate'].dt.to_period('M'))['ManagerRating'].mean()
plt.figure(figsize=(12, 6))
trend.plot(marker='o')
plt.title('Monthly Average Manager Rating Trend', fontsize=16)
plt.xlabel('Month', fontsize=14)
plt.ylabel('Average Self Rating', fontsize=14)
plt.xticks(rotation=45)
plt.grid()
plt.tight_layout()
plt.show()
```

Python



Plotting Count Plots For Categorical Columns

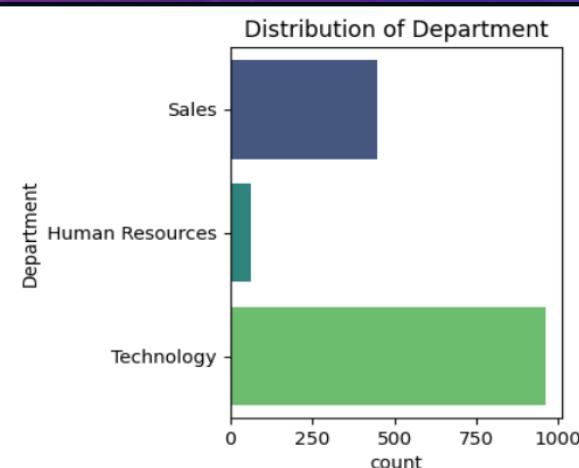
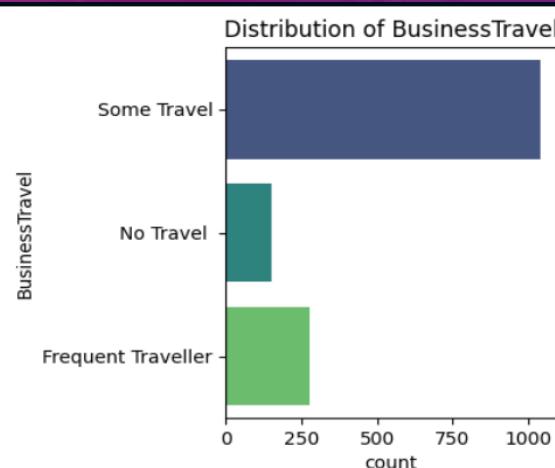
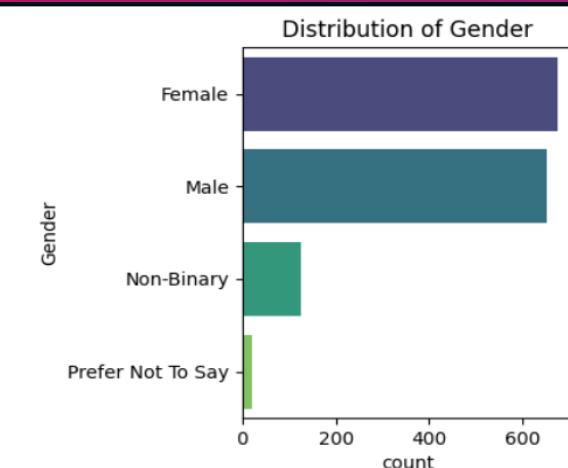
```
categorical_cols = ['Gender', 'BusinessTravel', 'Department', 'State', 'Ethnicity', 'EducationLevelID', 'EducationField', 'JobRole']

plt.figure(figsize=(15, 15))

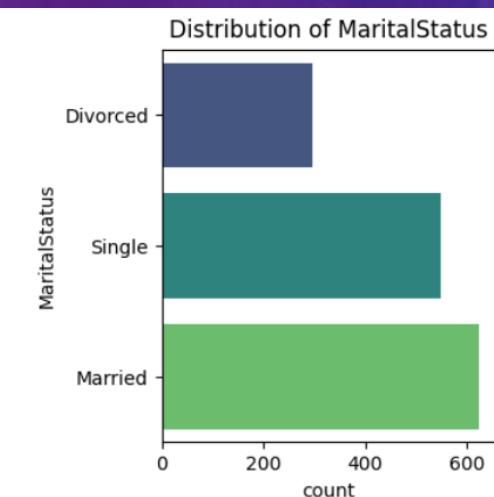
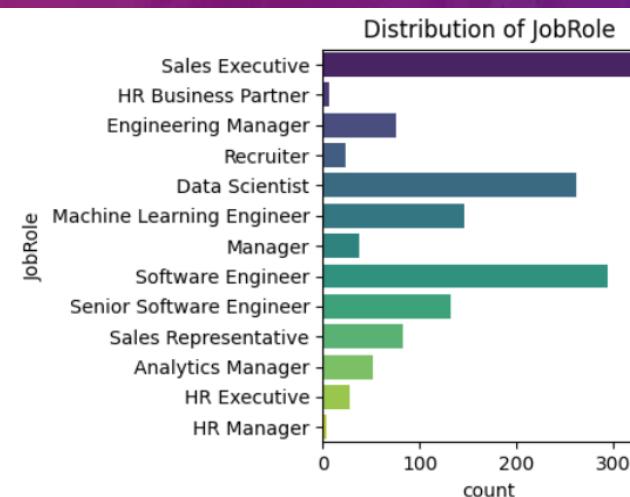
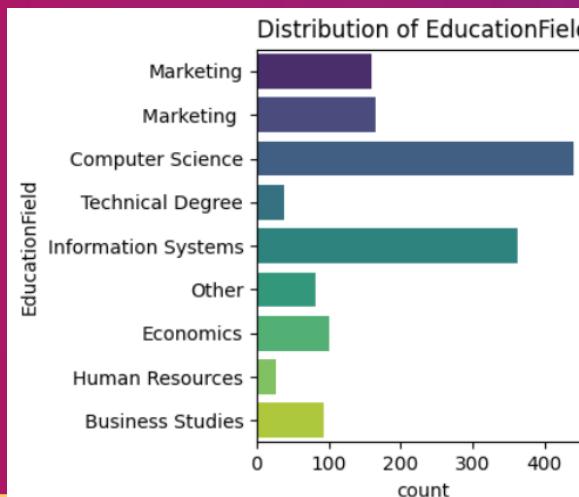
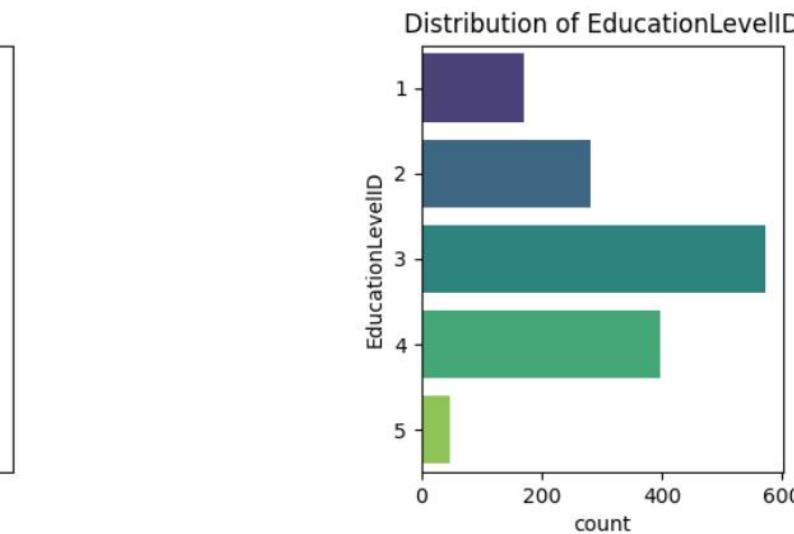
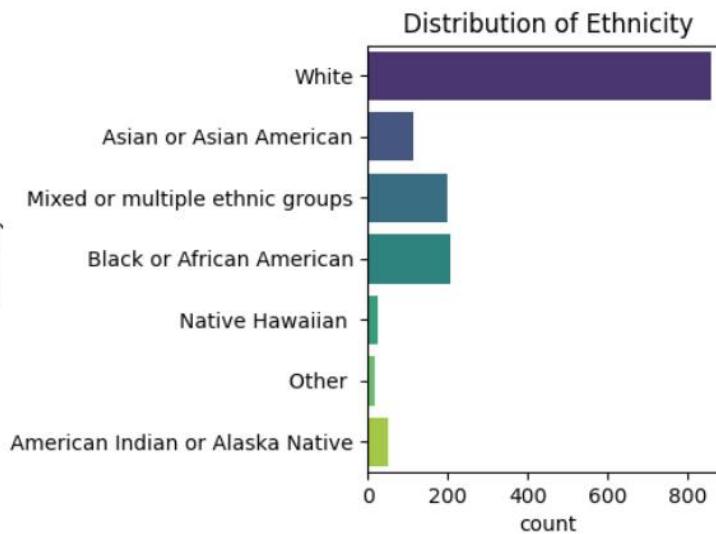
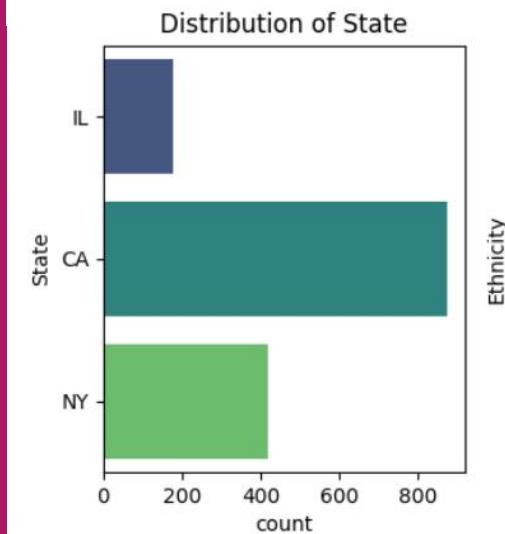
for i, col in enumerate(categorical_cols, 1):
    plt.subplot(4, 3, i)
    sns.countplot(y = df_employee[col], palette='viridis')
    plt.title(f'Distribution of {col}')

plt.tight_layout()
plt.show()
```

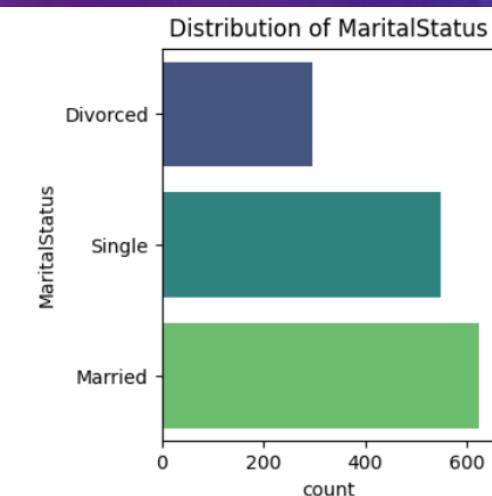
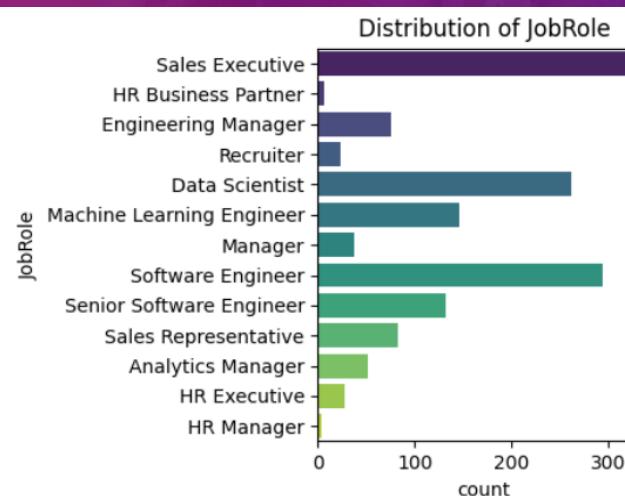
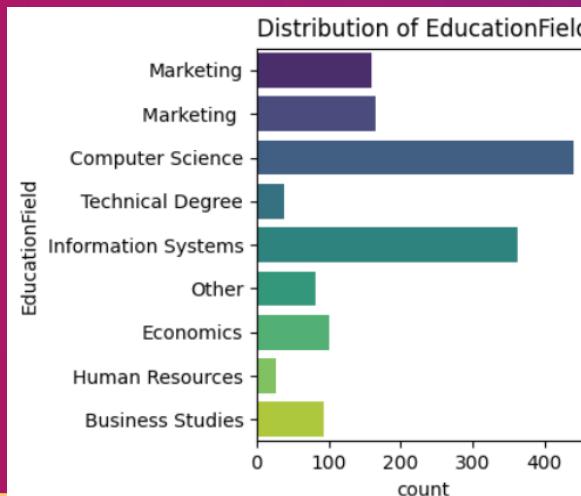
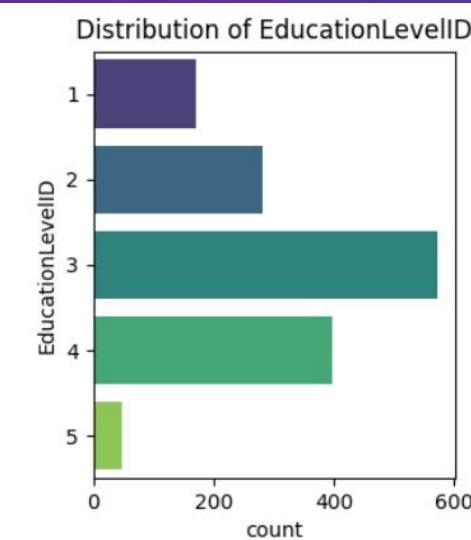
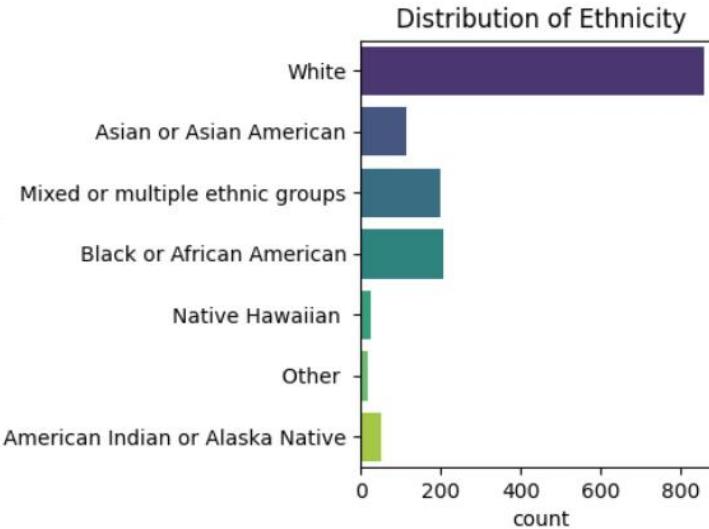
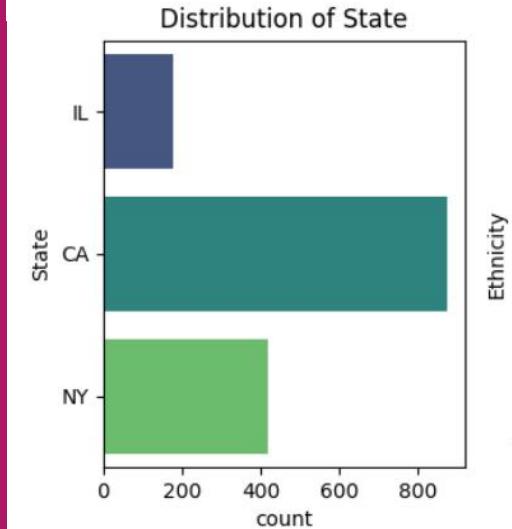
Python



Plotting Count Plots For Categorical Columns

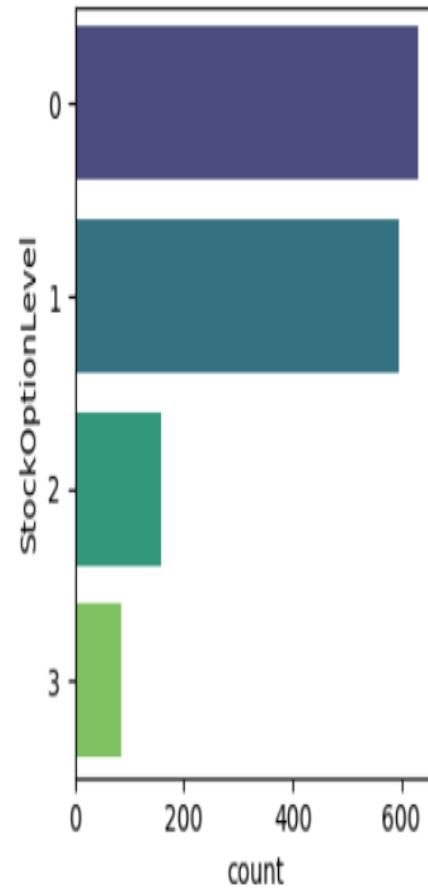


Plotting Count Plots For Categorical Columns

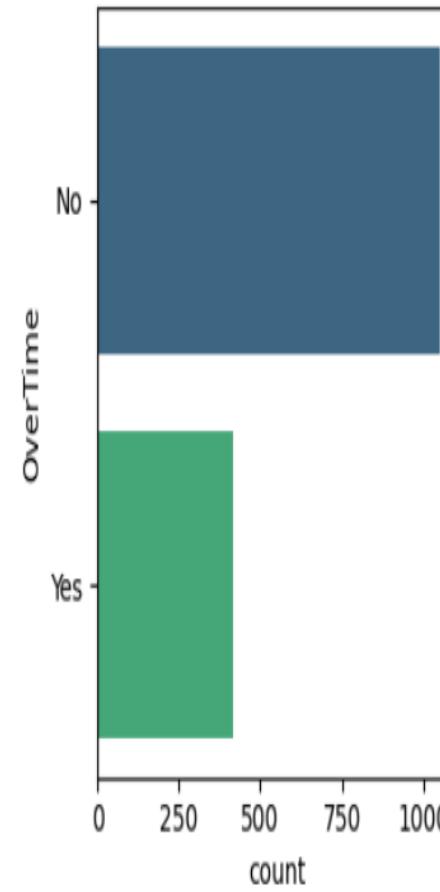


Plotting Count Plots For Categorical Columns

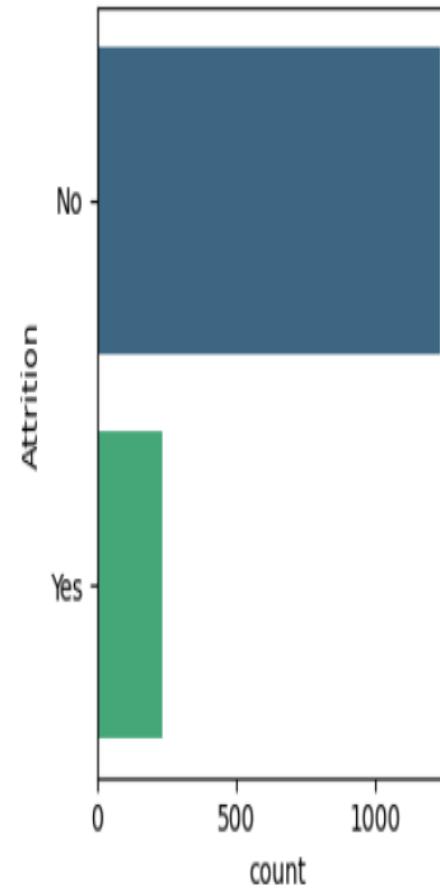
Distribution of StockOptionLevel



Distribution of OverTime



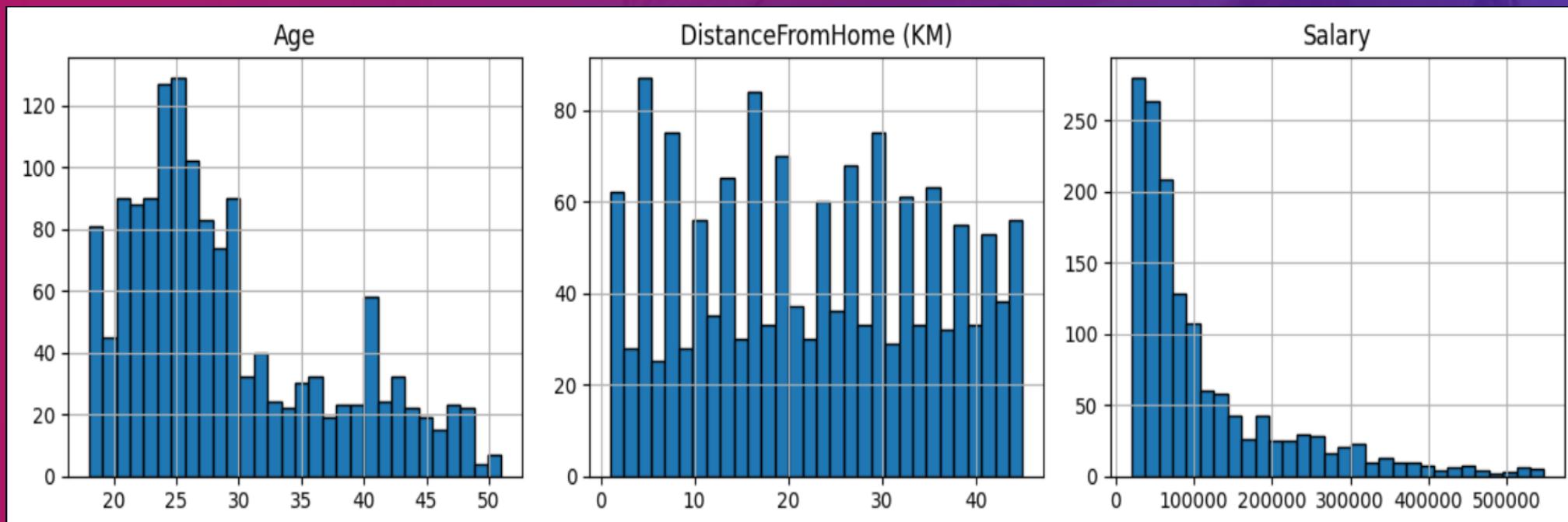
Distribution of Attrition

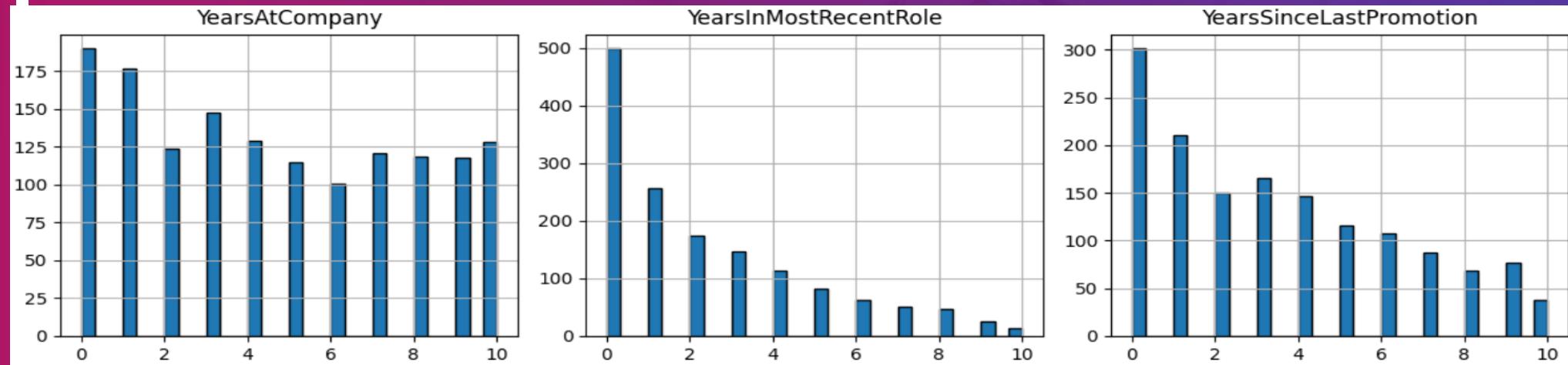


Plotting Count Plots For Numerical Columns

```
numerical_cols = ['Age', 'DistanceFromHome (KM)', 'Salary', 'YearsAtCompany', 'YearsInMostRecentRole', 'YearsSinceLastPromotion',  
df_employee[numerical_cols].hist(figsize=(12, 10), bins=30, edgecolor='black')  
plt.tight_layout()  
plt.show()
```

Python



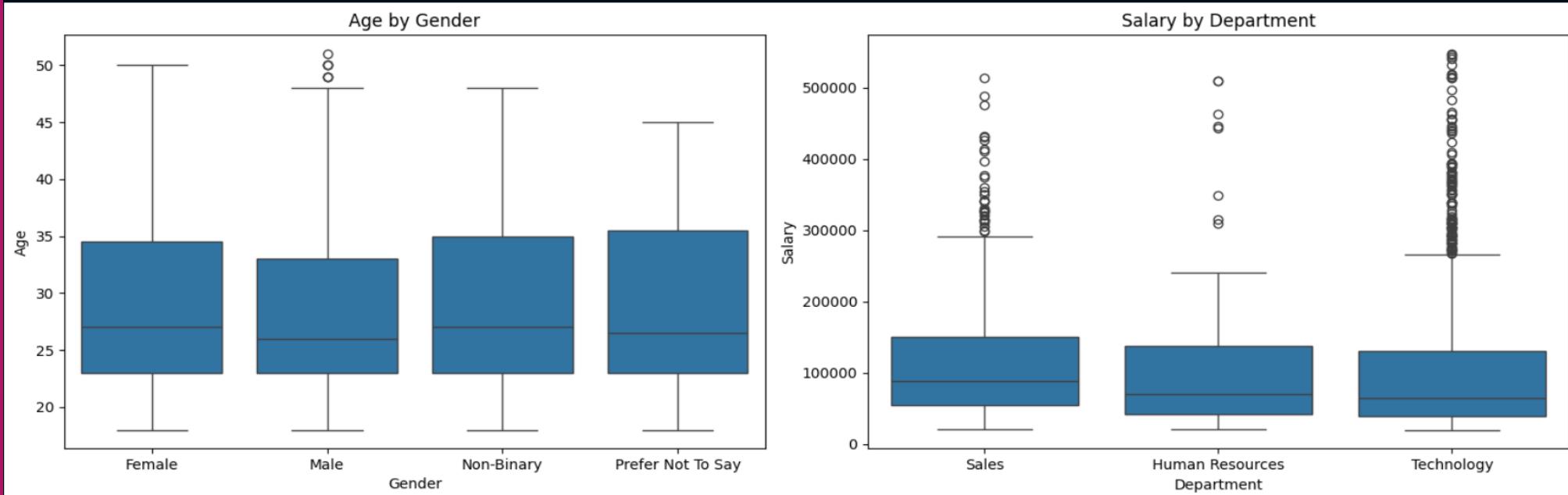


```

plt.figure(figsize=(15, 10))
# Age by Gender
plt.subplot(2, 2, 1)
sns.boxplot(x='Gender', y='Age', data=df_employee)
plt.title('Age by Gender')

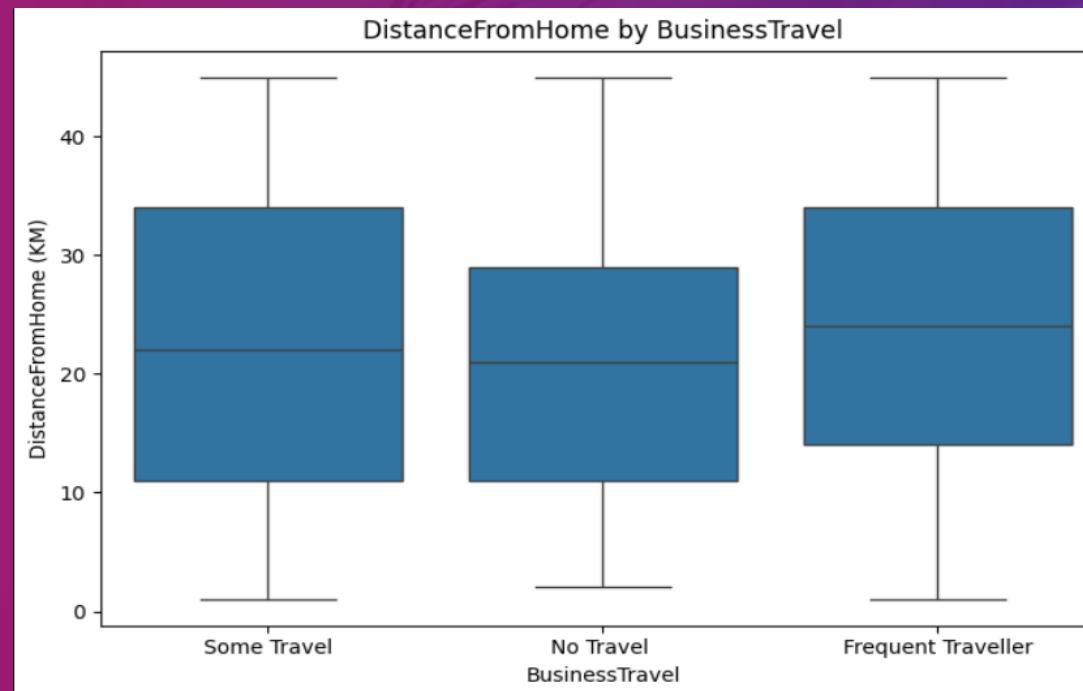
# Salary by Department
plt.subplot(2, 2, 2)
sns.boxplot(x='Department', y='Salary', data=df_employee)
plt.title('Salary by Department')

```



```
# DistanceFromHome by BusinessTravel  
plt.subplot(2, 2, 3)  
sns.boxplot(x='BusinessTravel', y='DistanceFromHome (KM)', data=df_employee)  
plt.title('DistanceFromHome by BusinessTravel')  
plt.tight_layout()
```

Python



FORECASTING

1-Predicting The Number of Employees that may have attrition in the company by the next year

Import libraries and merge 2 tables in one new table

```
import pandas as pd
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
from sklearn.preprocessing import OneHotEncoder

# Assuming df_employee and df_PerformanceRating are your original DataFrames

# Merge the two DataFrames
merged_data = df_employee.merge(df_PerformanceRating, on='EmployeeID', how='inner')
```

```

features = ['Age', 'JobRole', 'Salary', 'YearsAtCompany', 'WorkLifeBalance', 'JobSatisfaction']
X = merged_data[features]
y = merged_data['Attrition'].apply(lambda x: 1 if x == 'Yes' else 0) # Convert Attrition Yes/No to binary
X = pd.get_dummies(X, columns=['JobRole'], drop_first=True) # Convert JobRole to dummy variables
assert len(X) == len(y), "X and y must have the same number of samples."
X_train, X_test, y_train, y_test, employee_ids_train, employee_ids_test = train_test_split(
    X, y, merged_data['EmployeeID'], test_size=0.3, random_state=42
)
model = RandomForestClassifier(random_state=42) # Set random_state for reproducibility
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
print(classification_report(y_test, y_pred))
results = pd.DataFrame({
    'EmployeeID': employee_ids_test, # The correct EmployeeID for each test prediction
    'Attrition_Prediction': y_pred # The predictions (0 = stay, 1 = leave)
})
employees_to_leave = results[results['Attrition_Prediction'] == 1]
distinct_count = employees_to_leave['EmployeeID'].nunique()
print("Employees likely to leave:")
print(distinct_count)

```

	precision	recall	f1-score	support
0	0.99	0.97	0.98	1356
1	0.94	0.97	0.96	657
accuracy			0.97	2013
macro avg	0.96	0.97	0.97	2013
weighted avg	0.97	0.97	0.97	2013
Employees likely to leave:				
265				

2- Predict the Avg Salary to Each Job role:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
features = ['Age', 'YearsAtCompany', 'JobSatisfaction', 'ManagerRating', 'JobRole']
X = merged_data[features]
y = merged_data['Salary']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
model = LinearRegression()
model.fit(X_train, y_train)
job_roles_columns = [col for col in X_test.columns if 'JobRole' in col]
job_roles_encoded = pd.get_dummies(job_roles, columns=['JobRole'], drop_first=True)
job_role_names = job_roles_encoded.columns
average_salary_by_job_role = X_test.groupby('JobRole')['PredictedSalary'].mean()
Print("متوسط الرواتب المتوقعة لكل دور وظيفي في السنة القادمة")
print(average_salary_by_job_role)
```

Python

:متوسط الرواتب المتوقعة لكل دور وظيفي في السنة القادمة

JobRole	PredictedSalary
Data Scientist	97,731.42
Engineering Manager	294,319.15
HR Business Partner	311,464.38
HR Executive	96,212.39
HR Manager	424,422.60
Machine Learning Engineer	131,475.11
Manager	340,470.42
Recruiter	40,073.01
Sales Executive	128,787.84
Sales Representative	40,446.79
Senior Software Engineer	131,410.38
Software Engineer	54,650.61
Name: PredictedSalary, dtype: float64	

3- Which departments are likely to have the highest attrition rates next year?

```
# Group data by Department and calculate attrition rate  
attrition_by_department = merged_data.groupby('Department')['Attrition'].apply(lambda x: (x == 'Yes').mean())  
  
# Sort by highest attrition rate  
attrition_by_department = attrition_by_department.sort_values(ascending=False)  
  
# Display results  
print(attrition_by_department)
```

Python

```
Department  
Sales          0.41  
Human Resources 0.38  
Technology      0.30  
Name: Attrition, dtype: float64
```

```
Department  
Sales          0.41  
Human Resources 0.38  
Technology      0.30  
Name: Attrition, dtype: float64
```

DATA VISUALIZATION



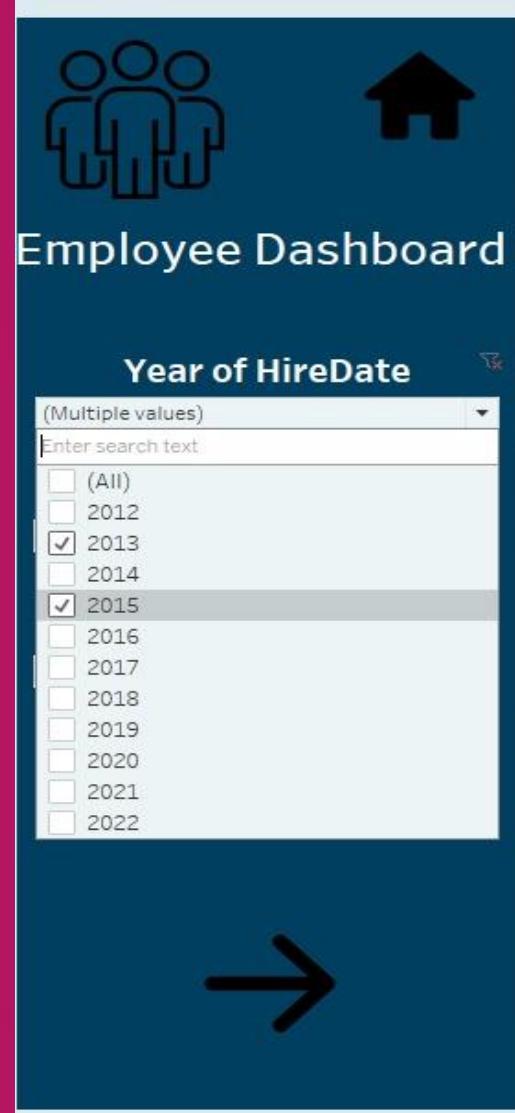
Human Resource Dashboard

Go to Employee Dashboard

Go to Performance Rating Dashboard

Members Team:

- 1- Rashad Hazem Ali
 - 2- Rasha Gaafar
 - 3- Mariam Korashy
 - 4- Mohamed Abd El Aleem mohamed
 - 5- Esraa Eleraky
 - 6- Basem Barakat-Allah Abdulqader Abdulkarim



Total Employee

263

Total Salary

30,746,522

Avg Salary

116,907

Hire Employee per years



Hiring by Each Department and Date



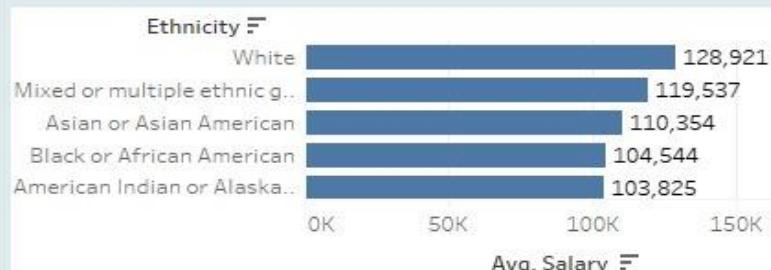
Attrition by Department

Department	Attrition	
	No	Yes
Human Resources	8	1
Sales	54	19
Technology	163	18

Avg Salary by Education Level

EducationLevel	Count of Employee	Avg. Salary
Bachelors	95	124,074
Doctorate	12	168,634
High School	63	109,701
Masters	61	115,546
No Formal Qualifications	32	93,011

Average Salary by Ethnicity



Over Time by Department



The dashboard features a sidebar on the left with filters for Year of HireDate, Department, and Select Gender, each with dropdown menus for All. The main area contains six primary sections: 1) Total Employee (1,470), 2) Total Salary (\$166,046,052), 3) Avg Salary (\$112,956), 4) Attrition by Department (table showing counts for Human Resources, Sales, and Technology), 5) Avarage Salary by Ethncity (horizontal bar chart showing salaries for Black or African American, American Indian or Alaska Native, Asian or Asian American, Mixed or multiple ethnic groups, and Other), and 6) Over Time by Department (three pie charts showing Yes vs No responses for Human Resources, Sales, and Technology).

Employee Dashboard

Year of HireDate
(All)

Department
(All)

Select Gender
All

Total Employee

1,470

Hire Employee per years

HireDate

Number of ...

Year	Hires
2012	150
2013	140
2014	135
2015	125
2016	110
2017	100
2018	130
2019	140
2020	125
2021	135
2022	150

Total Salary

166,046,052

Hiring by Each Department and Date

Count of Depa...

Year of HireDate

Avg Salary by Education Level

EducationLevel	Count of Employee	Avg. Salary
Bachelors	572	\$115,405
Doctorate	48	\$154,269
High School	282	\$105,181
Masters	398	\$117,641
No Formal Qualifications	170	\$94,983

Attrition by Department

Department	Attrition No	Yes
Human Resources	51	12
Sales	354	92
Technology	828	133

Avarage Salary by Ethncity

Ethnicity

Ethnicity	Avg. Salary
Black or African American	\$112,177
American Indian or Alaska Native	\$112,037
Asian or Asian American	\$109,851
Mixed or multiple ethnic groups	\$106,133
Other	\$101,652

Over Time by Department

Department

Human Resources

Sales

Technology

Yes

No



Total Employee

564



Total Salary

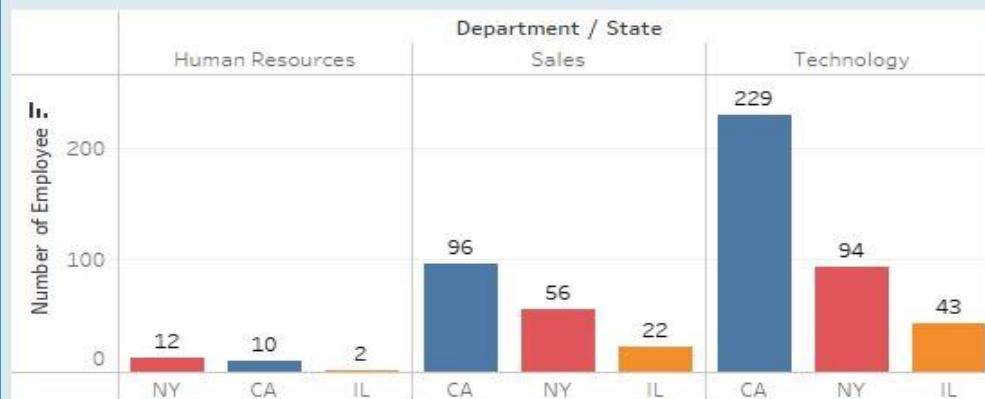
63,539,050



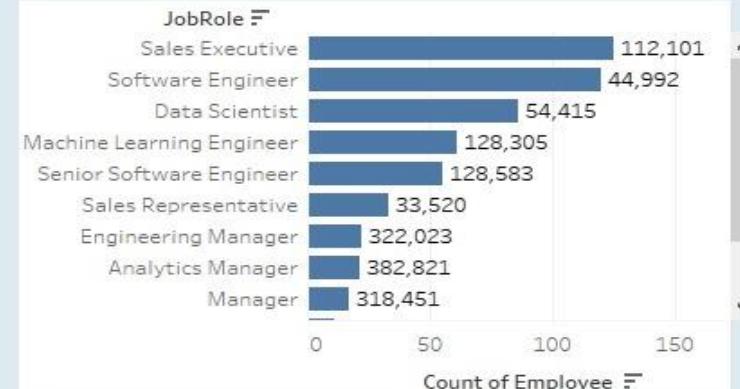
Avg Salary

112,658

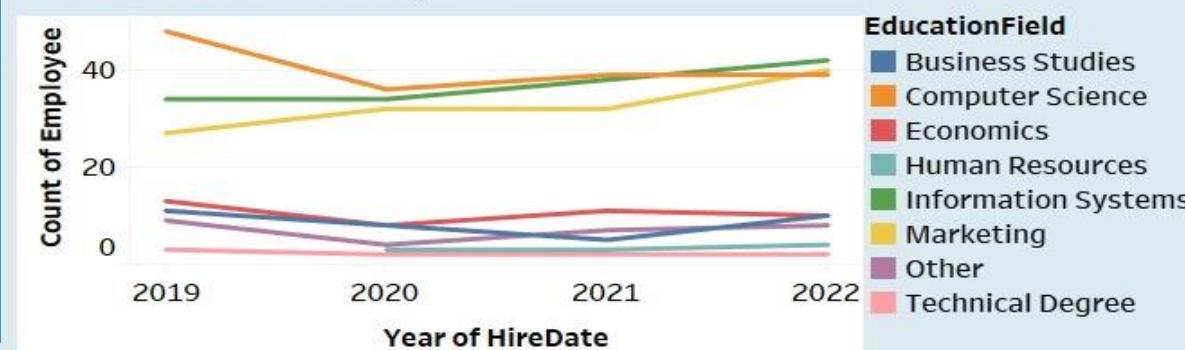
Num of Emp in Each state and department



Jop Role count



Education Field Hired by time



Salary by Each Education field

Education Field	Count
Business Studies	3,541,368
Computer Science	16,775,966
Economics	3,900,357
Human Resources	1,708,988
Information Systems	17,819,650
Marketing	16,357,520
Other	2,578,752
Technical Degree	856,449

The dashboard provides a comprehensive overview of employee data across different dimensions.

Total Employee: 1,470

Total Salary: 166,046,052

Avg Salary: 112,956

Job Role count:

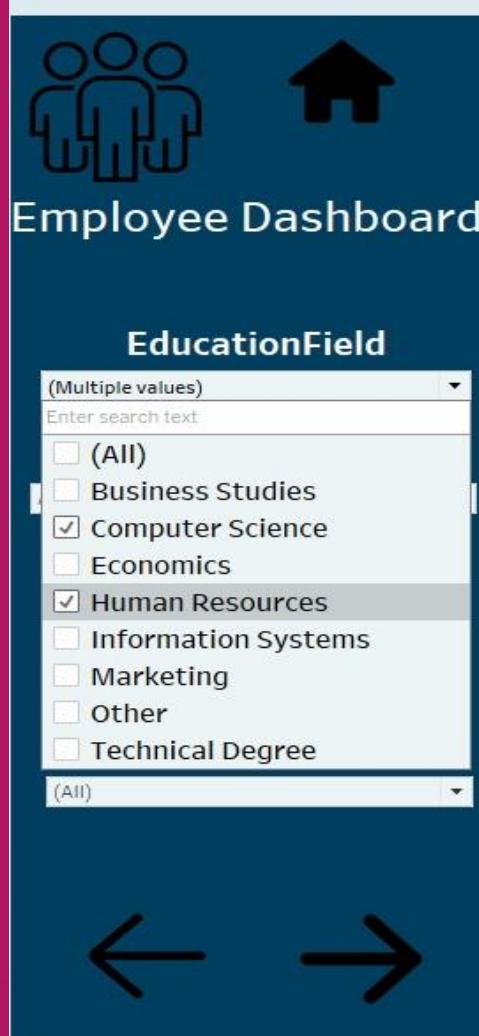
JobRole	Count of Employee
Sales Executive	117,196
Software Engineer	51,967
Data Scientist	56,079
Machine Learning Engineer	130,165
Senior Software Engineer	126,161
Sales Representative	40,656
Engineering Manager	286,259
Analytics Manager	346,484
Manager	317,531

Education Field Hired by time:

Year of HireDate	Business Studies	Computer Science	Economics	Human Resources	Information Systems	Marketing	Other	Technical Degree
2013	10	45	5	30	35	30	5	10
2015	15	45	5	25	30	25	5	10
2017	10	25	10	30	35	30	5	10
2019	10	45	10	30	35	30	5	10
2021	10	40	10	35	40	35	5	10

Salary by Each Education field:

EducationField	Salary
Business Studies	9,250,227
Computer Science	48,115,757
Economics	11,334,205
Human Resources	3,930,278
Information Systems	41,520,135
Marketing	40,390,223
Other	7,900,032
Technical Degree	3,605,195



Total Employee

467



Total Salary

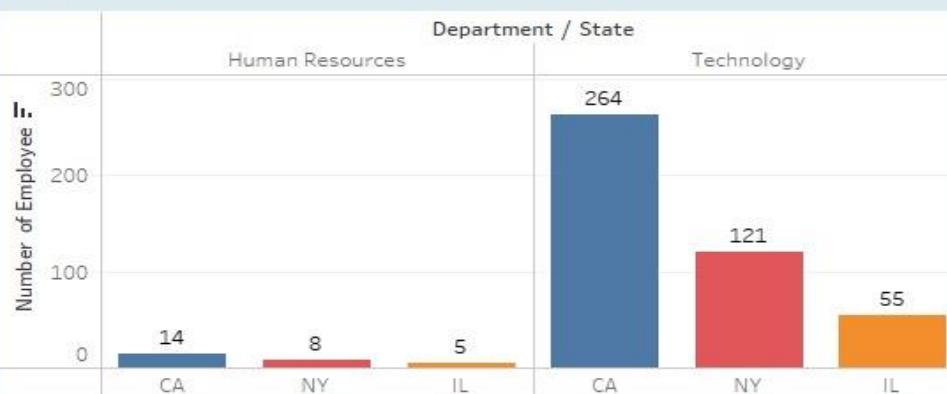
52,046,035



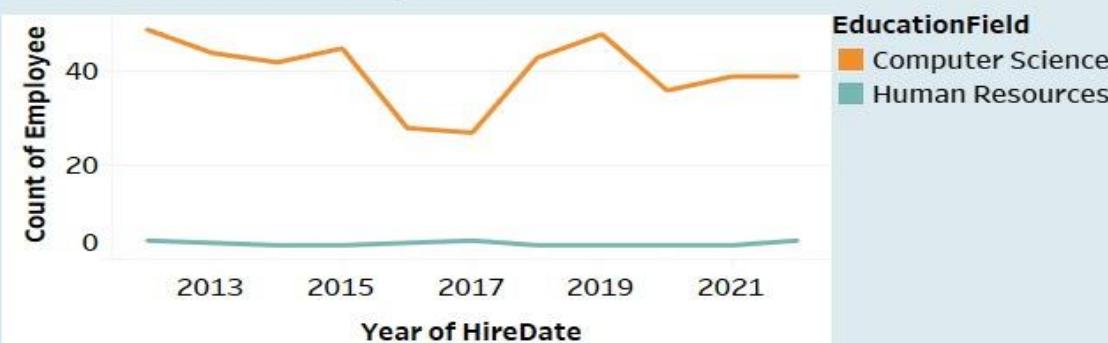
Avg Salary

111,448

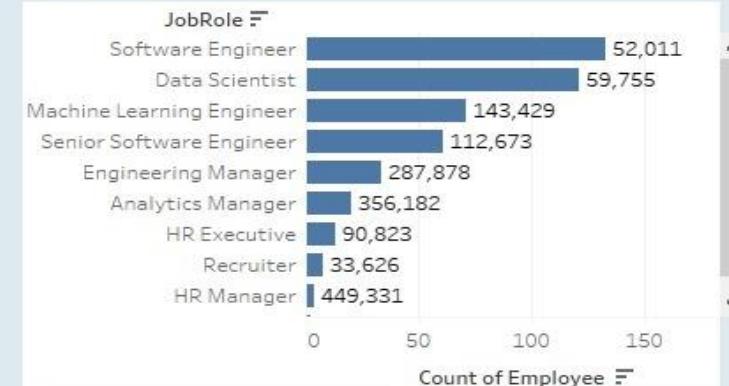
Num of Emp in Each state and department



Education Field Hired by time

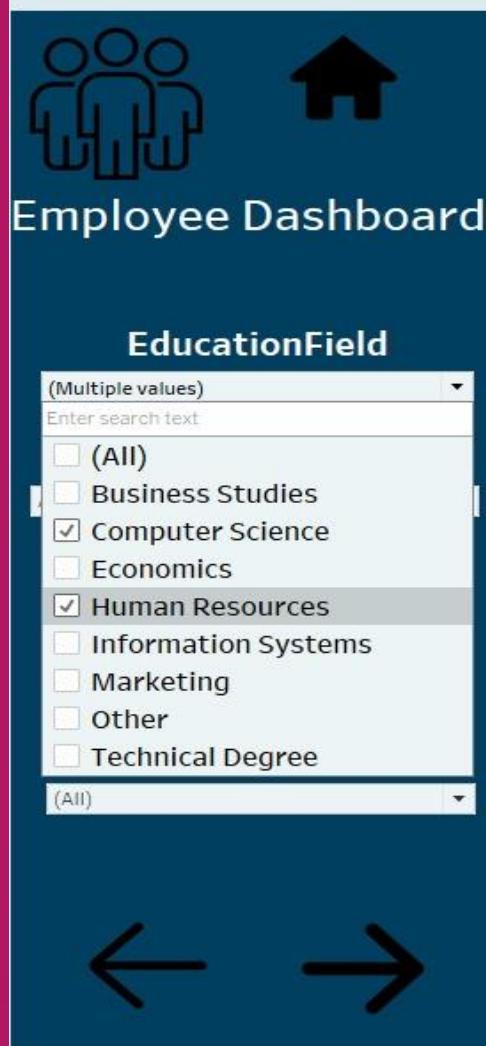


Jop Role count



Salary by Each Education field

Education Field	Number of Students
Computer Science	48,115,757
Human Resources	3,930,278



Total Employee

467



Total Salary

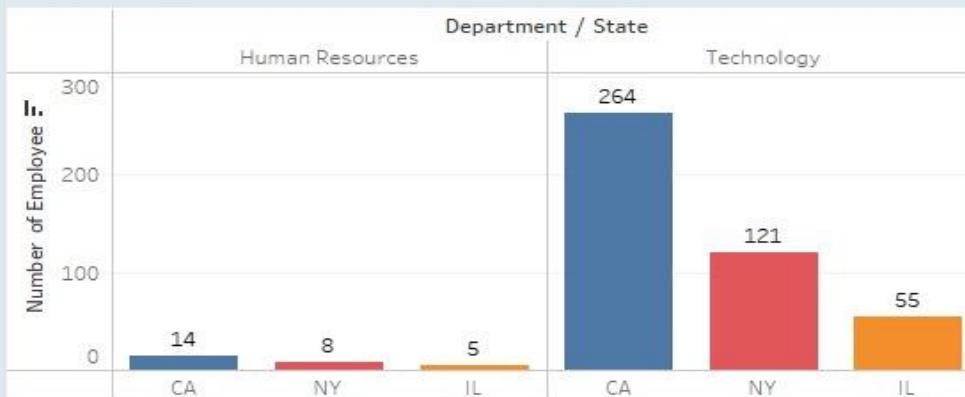
52,046,035



Avg Salary

111,448

Num of Emp in Each state and department



Education Field Hired by time



Jop Role count

JobRole

A horizontal bar chart titled "Count of Employee" showing the number of employees for ten different job titles. The x-axis represents the count of employees, ranging from 0 to 150. The y-axis lists the job titles. The bars are blue.

Job Title	Count of Employees
Software Engineer	52,011
Data Scientist	59,755
Machine Learning Engineer	143,429
Senior Software Engineer	112,673
Engineering Manager	287,878
Analytics Manager	356,182
HR Executive	90,823
Recruiter	33,626
HR Manager	449,331

Salary by Each Education field

Education Field	Number of Jobs
Computer Science	48,115,757
Human Resources	3,930,278



Total Employee

1,470



Total Salary

166,046,052



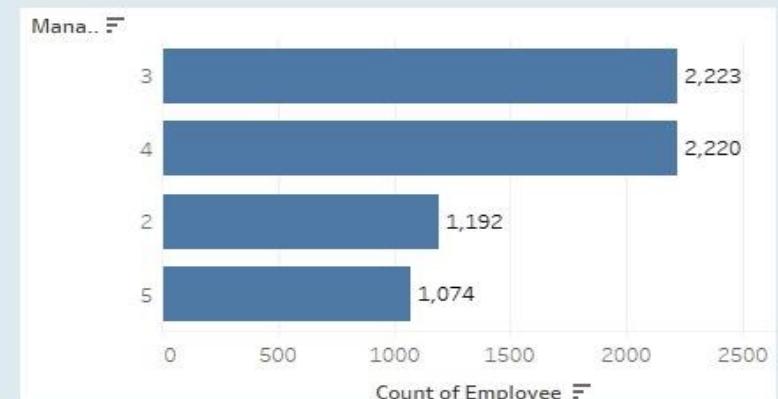
Avg Salary

112,956

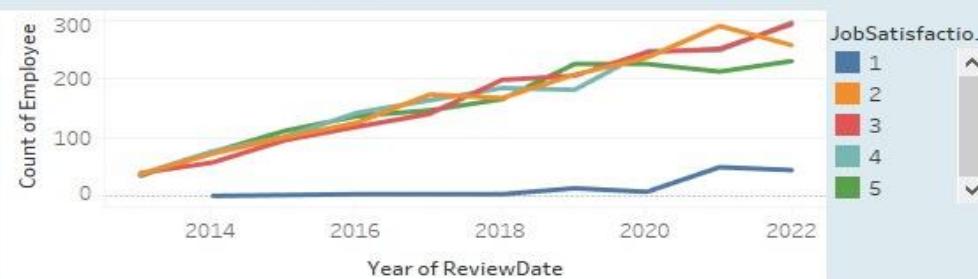
Manager Rating For Each Department

Manager	Role	Department	Count
2	Human Resources		62
	Sales		417
	Technology		708
3	Human Resources		87
	Sales		683
	Technology		1,457
4	Human Resources		97
	Sales		714
	Technology		1,409
5	Human Resources		52
	Sales		335
	Technology		687

Manager rating to all Employee



Job Satisfaction over years



Environment satisfaction by Ethnicity

	Ethnicity						
Environment Satisfaction	American Indian or Alaskan Native	Asian or Pacific Islander	Black or African American	Mixed Race	Native Hawaiian/Pacific Islander	Other	White
1	4	11	18	20	2	1	80
2	5	14	21	24	6	1	70
3	94	220	356	350	54	30	1,107
4	90	224	350	373	48	29	1,061
5	91	214	326	332	47	26	1,010



Total Employee

1,121



Total Salary

125,226,438



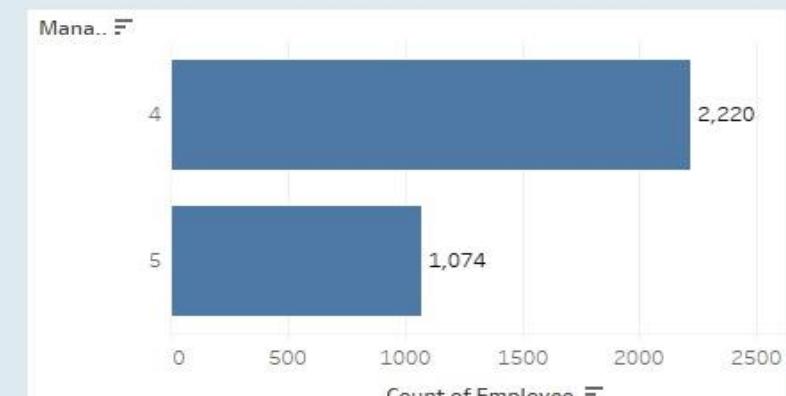
Avg Salary

111,710

Manager Rating For Each Department

ManagerRan.	Department	
4	Human Resources	914
	Sales	714
	Technology	1,409
5	Human Resources	521
	Sales	335
	Technology	687

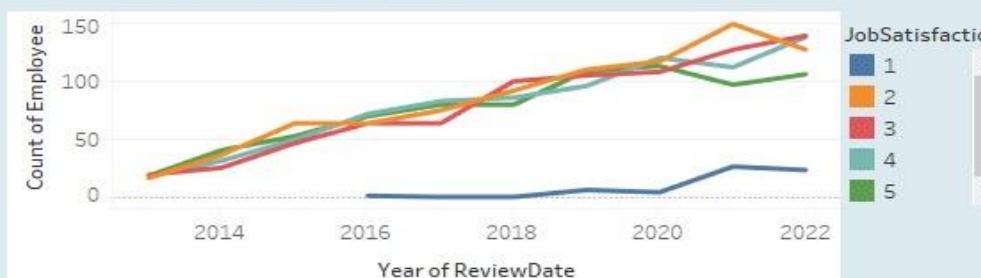
Manager rating to all Employee



Year of ReviewDate

←

Job Satisfaction over years



Environment satisfaction by Ethnicity

	Ethnicity						
	Environment	America	Asian	Black	Mixed	Native	
Satisfaction	Asian	Asian or African	African	Middle Eastern	Hawaiian	Other	White
1			7	5	11	1	1
2		5	5	10	15	3	1
3	37	110	177	158	24	17	560
4	42	130	169	177	25	14	515
5	43	100	156	159	19	16	505



Total Employee

1,470



Total Salary

166,046,052



Avg Salary

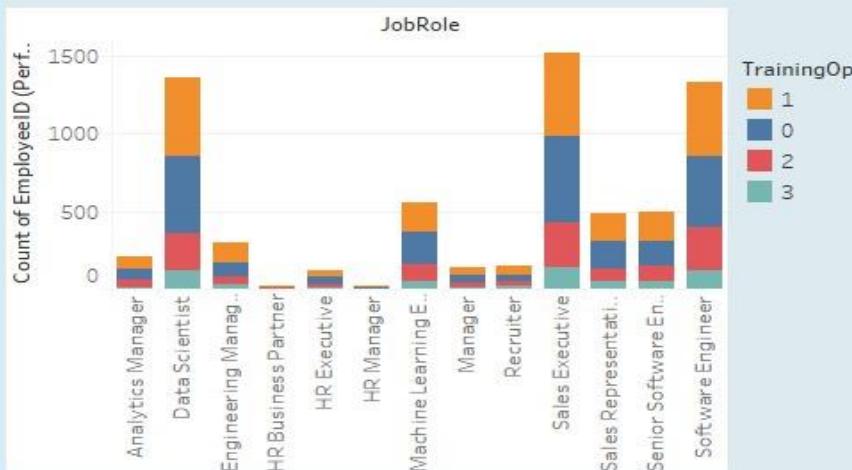
112,956

Work life balance for Each Gender

	WorkLifeBalance (PerformanceRating U.csv)				
Gender	1	2	3	4	5
Female	52	783	762	831	663
Male	61	750	744	710	703
Non-Binary	8	151	153	145	132
Prefer Not To Say		18	11	20	12

Training opportunities token for Each job

role



Manager Rating for Each Education level

Manager Rating	Bachelors	Doctorate	High School	Masters	No Formal Qualifications
2	450	50	250	300	150
3	850	50	450	580	280
4	880	50	480	550	280
5	420	50	220	280	150



Total Employee

898



Total Salary

100.034.146



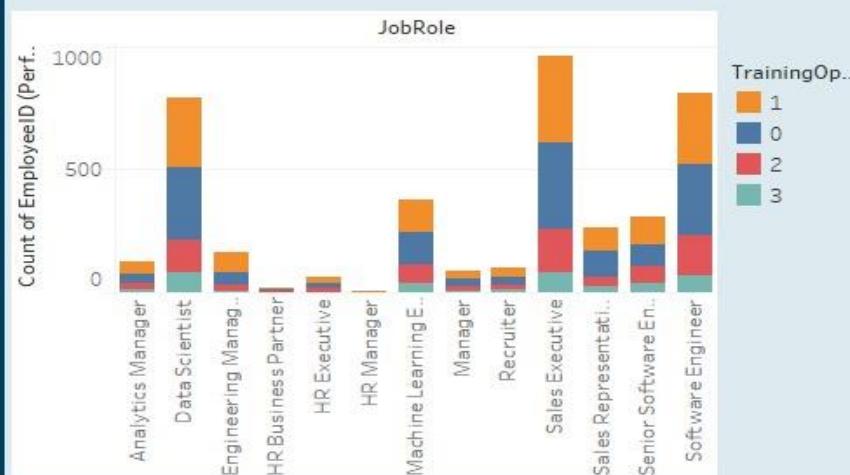
Avg Salary

111,397

Work life balance for Each Gender

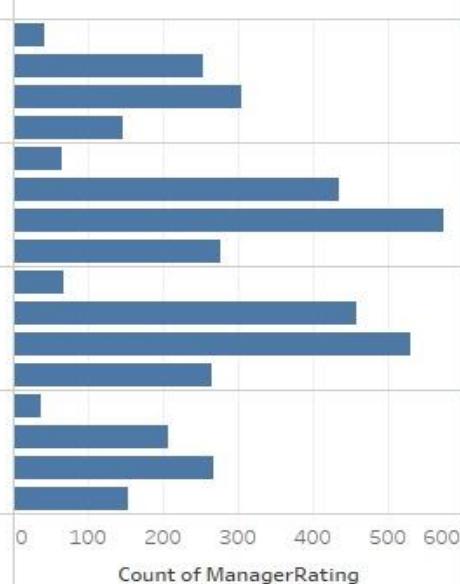
	WorkLifeBalance (PerformanceRating U.csv)				
Gender	1	2	3	4	5
Female	27	477	453	507	402
Male	36	462	446	446	418
Non-Binary	6	93	94	91	70
Prefer Not To Say		15	8	15	9

Training opportunities token for Each job role



Manager Rating for Each Education level

ManagerRef.	EducationLevel
2	Doctorate High School Masters No Formal Qualifications
3	Doctorate High School Masters No Formal Qualifications
4	Doctorate High School Masters No Formal Qualifications
5	Doctorate High School Masters No Formal Qualifications



CONCLUSION

ereby grants
use of any and all pictures or
videotapes including, but not limited to,

Please Listen to the Final Project Conclusion



A photograph of a young woman with dark, curly hair and glasses, smiling warmly at the camera. She is wearing a dark turtleneck sweater. The background is a bright, modern interior space with large windows and a painting on the wall.

THANK YOU!
