

HR Attrition Analytics: Unlocking Insights on Employee Satisfaction & Performance

DEPI Initiative - Data Analyst Specialist Track

Group: ALX3_DAT1_G2

Under the Supervision of: Eng./Dina Ezzat



Team Members

(Group 3)

- Moataz Tammam Abdelaziz —> Data Cleaning & Analysis
- Yasmeen Farouq Saleh —> Data Cleaning & Analysis
- Muhammed Nasseim El-Saied —> Data Cleaning & Analysis
- Sara Ayed Wahba —> Data Visualization & Presentation Preparation
- Malak Tarek Muhammed —> Data Visualization & Presentation Preparation

Project Overview

This project focuses on uncovering the underlying factors that lead to ***Employee Attrition*** within the company. Our analysis examines three key evaluation dimensions:

- **Employees' Satisfaction and Engagement.**
- **Manager Assessments of Employees.**
- **Employees' Self-evaluation.**

These insights are derived from HR dataset across **three** core departments: **Sales**, **Technology**, and **Human Resources**.

Data Analysis Process

1

Data Exploration

2

Data Cleaning & Preparation

3

Data Analysis (EDA)

4

Data Visualization

5

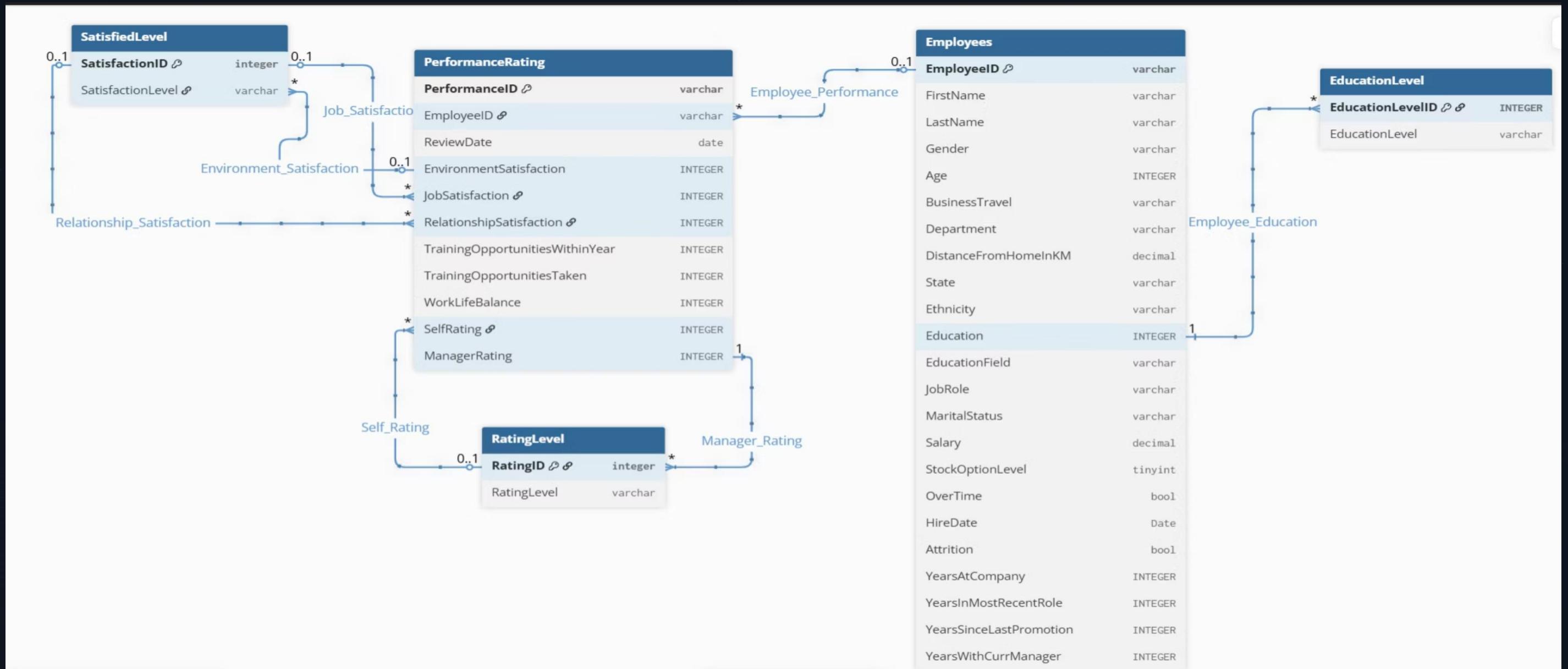
Data Insights & Reporting

Understanding Data

Our HR Dataset is built on a relational database consisting of **five interconnected tables**:

1. Employee table —> Master employee table (Demographics & Employment data, Employee Attrition & tenure fields). **Primary key: EmployeeID**
2. Performance Rating —> Employee Performance & Satisfaction scores per review (ReviewDate) and training participations.
Primary Key: PerformanceID, Foreign Key: EmployeeID.
3. Education Level —> Employees' Educational Background.
 - Lookup table: mapping **EducationLevelID “PK”** to Employees table “**FK: Education**”
4. Rating Level —> Categorical description for each rating.
 - Lookup table: **RatingID “PK”** → **Manager & Self-Rating “FK”** in PerformanceRating table.
5. Satisfied Level —> Categorical description for each Satisfaction Criteria.
 - Lookup table: **SatisfactionID “PK”** → **Environment, Job & Relationship Sats. “FK”** in PerformanceRating table.

Database Schema



Data Loading

```
1 -- Data Exploration PerformanceRating  
2 SELECT * FROM PerformanceRating;  
3 -- Data Exploration Employee  
4 SELECT * FROM Employees;  
5 -- Data Exploration SatisfiedLevel  
6 SELECT * FROM SatisfiedLevel;  
7 -- Data Exploration RatingLevel  
8 SELECT * FROM RatingLevel;  
9 -- Data Exploration EducationLevel  
10 SELECT * FROM EducationLevel;
```

Data Cleaning & Preparation



Handling Duplicates

```
16 -- Checking Duplicates in Employees table  
17 SELECT EmployeeID, COUNT (*) FROM Employees  
18 GROUP BY EmployeeID  
19 HAVING COUNT(*)> 1;  
20 -- Checking Duplicates in PerPerformanceRating table  
21 SELECT PerformanceID, COUNT (*) FROM PerformanceRating  
22 GROUP BY PerformanceID  
23 HAVING COUNT(*)> 1;
```

- It appears that there are no duplicates in our dataset.

Handling Null Values

```
24 -- Searching for Missing Values & Removing Trailing Whitespaces at key Fields in Employees table
25 SELECT
26   COUNT(IF(EmployeeID IS NULL OR TRIM(EmployeeID) = '', 1, NULL)) AS Missing_EmpID,
27   COUNT(IF(Gender IS NULL OR TRIM(Gender) = '', 1, NULL)) AS Missing_Gender,
28   COUNT(IF(Age IS NULL OR TRIM(Age) = '', 1, NULL)) AS Missing_Age,
29   COUNT(IF(Department IS NULL OR TRIM(Department) = '', 1, NULL)) AS Missing_Dept,
30   COUNT(IF(Salary IS NULL OR TRIM(Salary) = '', 1, NULL)) AS Missing_Salary,
31   COUNT(IF(YearsAtCompany IS NULL OR TRIM(YearsAtCompany) = '', 1, NULL)) AS Missing_Exp_Yrs
32 FROM Employees
33 -- Searching for Missing Values & Removing Trailing Whitespaces at key Fields in PerformanceRating table
34 SELECT
35   COUNT(IF(PerformanceID IS NULL OR TRIM(PerformanceID) = '', 1, NULL)) AS Missing_PerfID,
36   COUNT(IF(EmployeeID IS NULL OR TRIM(EmployeeID) = '', 1, NULL)) AS Missing_EmpID,
37   COUNT(IF(SelfRating IS NULL OR TRIM(SelfRating) = '', 1, NULL)) AS Missing_Self_Rate,
38   COUNT(IF(ManagerRating IS NULL OR TRIM(ManagerRating) = '', 1, NULL)) AS Missing_Mangager_Rate,
39   COUNT(IF(JobSatisfaction IS NULL OR TRIM(JobSatisfaction) = '', 1, NULL)) AS Missing_Job_Sat,
40   COUNT(IF(TrainingOpportunitiesWithinYear IS NULL OR TRIM(TrainingOpportunitiesWithinYear) = '', 1, NULL)) AS Miss_TrOppWithinYear,
41   COUNT(IF(RelationshipSatisfaction IS NULL OR TRIM(RelationshipSatisfaction) = '', 1, NULL)) AS Missing_Rel_Sat,
42   COUNT(IF(TrainingOpportunitiesTaken IS NULL OR TRIM(TrainingOpportunitiesTaken) = '', 1, NULL)) AS Missing_TrOppTaken
43 FROM PerformanceRating;
```

- It appears that there are no Null Values in our dataset.

Checking Invalid Data (Ensure Data Quality)

```
63 -- Checking Data Quality if there are any Invalid Values in Employees Table (Age & Salary)
64 SELECT COUNT ( * ) FROM Employee
65 WHERE Age < 0;
66 SELECT COUNT ( * ) FROM Employee
67 WHERE salary < 0;
68 -- Checking invalid data in PerformanceRating Table
69 SELECT COUNT(*) AS InvalidRatings
70 FROM PerformanceRating
71 WHERE SelfRating NOT BETWEEN 1 AND 5
72     OR ManagerRating NOT BETWEEN 1 AND 5
73     OR EnvironmentSatisfaction NOT BETWEEN 1 AND 5
74     OR JobSatisfaction NOT BETWEEN 1 AND 5
75     OR RelationshipSatisfaction NOT BETWEEN 1 AND 5
76     OR WorkLifeBalance NOT BETWEEN 1 AND 5;
77 -- Identifying Outliers in Employee Table
78 SELECT EmployeeID, Age
79 FROM Employees
80 WHERE Age < 18 OR Age > 60;
```

Exploratory Data Analysis (EDA)

- To examine & investigate datasets to discover patterns, trends & spot anomalies.

```
86 -- Summary Statistics
87 -- Salary & Employee Tenure by their Education Level
88 SELECT ed.EducationLevel,
89     COUNT(e.EmployeeID) AS cnt,
90     ROUND(AVG(COALESCE(e.Salary,0)), 2) AS Avg_Salary,
91     ROUND(AVG(COALESCE(e.YearsAtCompany,0)), 2) AS Avg_Years_At_Company
92 FROM Employees e
93 LEFT JOIN EducationLevel ed
94     ON e.Education = ed.EducationLevelID
95 GROUP BY ed.EducationLevel
96 ORDER BY Avg_Salary DESC;
97 -- Salary Distribution Statistics
98 SELECT MIN(Salary) AS MinSalary,
99     MAX(Salary) AS MaxSalary,
100    AVG(Salary) AS AvgSalary
101 FROM Employees;
102 -- Age Statistics
103 SELECT MIN(Age) AS MinAge,
104     MAX(Age) AS MaxAge,
105    AVG(Age) AS AvgAge
106 FROM Employees;
```

Focus Areas of Analysis

1

**Factors influencing
Employee Satisfaction and
Attrition**

2

**Overtime Impact on
Employee WLB &
Performance**

3

**Correlation between Salary
and Satisfaction or
Performance**

4

**Training Impact on Performance & Promotion
Potential**

5

Common Characteristics of Attritors

Performance Analysis (Performance Correlation with Age)

```
210 -- Correlation
211 -- Performance analysis and Age / Years At Company.
212 -- Employee Performance & Age
213 WITH perf AS (
214     SELECT
215         EmployeeID,
216         AVG(JobSatisfaction) AS Avg_JobSatisfaction,
217         AVG(SelfRating) AS Avg_SelfRating,
218         AVG(ManagerRating) AS Avg_ManagerRating
219     FROM PerformanceRating
220     GROUP BY EmployeeID
221 )
222 SELECT
223     CASE
224         WHEN e.Age < 25 THEN '<25'
225         WHEN e.Age BETWEEN 25 AND 35 THEN '25-35'
226         WHEN e.Age BETWEEN 36 AND 45 THEN '36-45'
227         ELSE '>45'
228     END AS Age_Range,
229     ROUND(AVG(p.Avg_JobSatisfaction),2) AS Avg_JobSatisfaction,
230     ROUND(AVG(p.Avg_SelfRating),2) AS Avg_SelfRating,
231     ROUND(AVG(p.Avg_ManagerRating),2) AS Avg_ManagerRating,
232     COUNT(*) AS Employee_Count
233 FROM Employee e
234 LEFT JOIN perf p
235     ON e.EmployeeID = p.EmployeeID
236 GROUP BY Age_Range
237 ORDER BY Age_Range;
238
```

Performance Analysis (Performance Correlation with tenure)

```
239 --Employee Performance & Years At Company
240 WITH perf AS (
241     SELECT
242         EmployeeID,
243         AVG(JobSatisfaction) AS JobSatisfaction,
244         AVG(SelfRating) AS SelfRating,
245         AVG(ManagerRating) AS ManagerRating
246     FROM PerformanceRating
247     GROUP BY EmployeeID
248 )
249 SELECT
250     CASE
251         WHEN YearsAtCompany < 2 THEN '<2'
252         WHEN YearsAtCompany BETWEEN 2 AND 5 THEN '2-5'
253         WHEN YearsAtCompany BETWEEN 6 AND 10 THEN '6-10'
254         ELSE '>10'
255     END AS YearsAtCompany_Range,
256     AVG(p.JobSatisfaction) AS Avg_JobSatisfaction,
257     AVG(p.SelfRating) AS Avg_SelfRating,
258     AVG(p.ManagerRating) AS Avg_ManagerRating,
259     COUNT(*) AS Employee_Count
260 FROM Employee e
261 JOIN perf p ON e.EmployeeID = p.EmployeeID
262 GROUP BY YearsAtCompany_Range
263 ORDER BY YearsAtCompany_Range;
264
```

Factors influencing Employee Satisfaction and Attrition

```
88 -- Identify the factors most associated with Low Employee Satisfaction or High Attrition
89 SELECT
90     e.Department,
91     ROUND(AVG(p.EnvironmentSatisfaction),2) AS Avg_Env_Satisfaction,
92     ROUND(AVG(p.JobSatisfaction),2) AS Avg_Job_Satisfaction,
93     ROUND(AVG(p.RelationshipSatisfaction),2) AS Avg_Relationship_Satisfaction,
94     ROUND(AVG(p.WorkLifeBalance),2) AS Avg_WorkLifeBalance,
95     ROUND(AVG(p.SelfRating),2) AS Avg_SelfRating,
96     ROUND(AVG(p.TrainingOpportunitiesWithinYear),2) AS Avg_Training_Opportunities
97 FROM Employees AS e
98 LEFT JOIN PerformanceRating AS p
99     ON e.EmployeeID = p.EmployeeID
100 GROUP BY e.Department
101 ORDER BY Avg_Job_Satisfaction ASC;
```

Overtime Impact on Work-Life Balance and Employee Performance:

```
122 -- Determine the Impact of Overtime on Work-Life balance and Overall Performance
123 SELECT
124     e.OverTime,
125     COUNT(DISTINCT(e.EmployeeID)) AS Num_Employees,
126     ROUND(AVG(p.WorkLifeBalance),2) AS Avg_WorkLife_Balance,
127     ROUND(AVG(p.JobSatisfaction),2) AS Avg_Job_Satisfaction,
128     ROUND(ROUND(AVG(p.ManagerRating),2) AS Avg_Manager_Rating,
129     ROUND(AVG(p.SelfRating),2) AS Avg_Self_Rating
130 FROM Employees AS e
131 LEFT JOIN PerformanceRating AS p
132     ON e.EmployeeID = p.EmployeeID
133 GROUP BY e.OverTime
134 ORDER BY Avg_WorkLife_Balance ASC;
```

Overtime Impact on Employee Performance



Overtime and Job Satisfaction

Overtime had a noticeable effect on job satisfaction. Employees who **worked extra hours** reported **higher satisfaction** compared to those who did not (3.51 vs 3.38).



Overtime and Performance

Employees who **did not work overtime** have **higher manager rating** average of 3.47 compared to 3.45 for those who do.

Relation between Salary and Employee Satisfaction & Performance

```
135 -- Relation between Salary and Employee Satisfaction & Performance
136 SELECT
137     e.Department,
138     ROUND(AVG(e.Salary), 2) AS Avg_Salary,
139     ROUND(AVG(p.ManagerRating), 2) AS Avg_Manager_Rating,
140     ROUND(AVG(p.SelfRating), 2) AS Avg_Self_Rating,
141     ROUND(AVG(p.JobSatisfaction), 2) AS Avg_Job_Satisfaction
142 FROM Employees AS e
143 LEFT JOIN PerformanceRating AS p
144     ON e.EmployeeID = p.EmployeeID
145 GROUP BY e.Department
146 ORDER BY Avg_Salary DESC;
147 SELECT
148     CASE
149         WHEN Salary < 80000 THEN '< 80k'
150         WHEN Salary BETWEEN 80000 AND 120000 THEN '80k-120k'
151         WHEN Salary BETWEEN 120001 AND 180000 THEN '120k-180k'
152         WHEN Salary BETWEEN 180001 AND 300000 THEN '180k-300k'
153         ELSE '> 300k'
154     END AS Salary_Band,
155     COUNT(*) AS Employee_Count
156 FROM Employees
157 GROUP BY Salary_Band
158 ORDER BY Employee_Count DESC;
```

Training Opportunities Impact on Performance & Promotion Potential

```
207 -- The Impact of Training Opportunities on Performance and Promotion Rate
208 SELECT
209     CASE
210         WHEN p.TrainingOpportunitiesWithinYear > 0 THEN 'Received Training'
211         ELSE 'No Training'
212     END AS Training_Status,
213     COUNT(DISTINCT e.EmployeeID) AS Num_Employees,
214     ROUND(AVG(p.ManagerRating), 2) AS Avg_Manager_Rating,
215     ROUND(AVG(p.SelfRating), 2) AS Avg_Self_Rating,
216     ROUND(AVG(e.YearsAtCompany), 2) AS Avg_Years_At_Company
217 FROM Employees e
218 LEFT JOIN PerformanceRating p
219     ON e.EmployeeID = p.EmployeeID
220 GROUP BY Training_Status
221 ORDER BY Avg_Years_At_Company DESC;
222 SELECT
223     e.Department,
224     ROUND(AVG(p.TrainingOpportunitiesWithinYear), 2) AS Avg_Training_Opportunities_WithinYear,
225     ROUND(AVG(p.TrainingOpportunitiesTaken), 2) AS Avg_Training_Opportunities_Taken,
226     ROUND(AVG(p.ManagerRating), 2) AS Avg_ManagerRating,
227     ROUND(AVG(p.SelfRating), 2) AS Avg_SelfRating,
228     ROUND(AVG(e.YearsAtCompany), 2) AS Avg_Years_At_Company,
229     ROUND(AVG((ManagerRating + SelfRating) / 2.0), 2) AS PromotionPotentialIndex
230 FROM Employees AS e
231 LEFT JOIN PerformanceRating AS p
232     ON e.EmployeeID = p.EmployeeID
233 GROUP BY e.Department
234 ORDER BY Avg_Training_Opportunities_Taken DESC;
```

- **HR Department** has the **lowest training opportunities** (Average: 0.96), but the **highest Performance "Manager Rating"** (Average: 3.49).
- **Sales Department** has **slightly higher training opportunities** (Average: 0.99), but the **lowest Performance** (Average: 3.43).
- **Technology Department** has the **highest training opportunities** (Average: 1.04) and **Manager rating of (Average: 3.47)**.

Common Characteristics of Employees who Left

```
323 -- Common Characteristics of Employees who left (Attrition Rate = Yes)
324 SELECT
325   UPPER(TRIM(e.Department)) AS Department,
326   ROUND(AVG(e.Age), 2) AS Avg_Age,
327   ROUND(AVG(e.Salary), 2) AS Avg_Salary,
328   ROUND(AVG(e.YearsAtCompany), 2) AS Avg_YearsAtCompany,
329   ROUND(AVG(p.JobSatisfaction), 2) AS Avg_JobSatisfaction,
330   ROUND(AVG(p.ManagerRating), 2) AS Avg_ManagerRating,
331   ROUND(AVG(p.WorkLifeBalance), 2) AS Avg_WorkLifeBalance,
332   ROUND(AVG(p.EnvironmentSatisfaction), 2) AS Avg_EnvironmentSatisfaction,
333   ROUND(AVG(p.RelationshipSatisfaction), 2) AS Avg_RelationshipSatisfaction,
334   COUNT(DISTINCT e.EmployeeID) AS Attrition_Count,
335   ROUND(100.0 * COUNT(DISTINCT e.EmployeeID) /
336     (SELECT COUNT(DISTINCT EmployeeID) FROM Employees WHERE LOWER(TRIM(Attrition)) = 'yes'), 2) AS Attrition_Rate
337 FROM Employees e
338 LEFT JOIN PerformanceRating p ON e.EmployeeID = p.EmployeeID
339 WHERE LOWER(TRIM(e.Attrition)) = 'yes'
340 GROUP BY e.Department
341 ORDER BY Attrition_Count DESC;
```

Employees who are **younger, less incentivized**, and have **shorter tenure** tend to **leave at higher rates**.

HR Project Jupyter Notebook "Python"

HR Analytics: Employees' Performance & Satisfaction

Project Overview:

This project explores employees' data to uncover insights into workforce performance & satisfaction by analyzing key factors such as demographics, education levels & other factors. This project aims to support HR Decision-making and provide recommendations to improve employee engagement and overall organizational performance.



Importing libraries & Loading our Dataset

1. Importing Libraries

```
[2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

2. Loading Datasets

```
[36]: # Loading datasets
Employees = pd.read_csv ('Employees.csv')
Educational_Level = pd.read_csv('EducationLevel.csv')
Performance = pd.read_csv('PerformanceRating.csv')
Rating_Level = pd.read_csv('RatingLevel.csv')
Satisfied_Level = pd.read_csv('SatisfiedLevel.csv')

# Merging tables
# Merge employees with educational level (EducationID)
Employees = Employees.merge(Educational_Level, left_on="Education", right_on="EducationLevelID", how="left")

# Merge employees with performance rating (EmployeeID)
Employees = Employees.merge(Performance, on="EmployeeID", how="left")

# Merge rating Level with both (Manager & Self-Rating)
Employees = Employees.merge(Rating_Level, left_on="ManagerRating", right_on="RatingID", how="left"
                           ).rename(columns={"RatingLevel": "Mananger_RatingLevel"}).drop(columns=["RatingID"])

Employees = Employees.merge(Rating_Level, left_on="SelfRating", right_on="RatingID", how="left"
                           ).rename(columns={"RatingLevel": "Self_RatingLevel"}).drop(columns=["RatingID"])

# Merge Satisfied Level with (Job, Environment, Relationship Satisfaction & Work-Life Balance)
# Environment Satisfaction
Employees = Employees.merge(Satisfied_Level, left_on="EnvironmentSatisfaction", right_on="SatisfactionID", how="left"
                           ).rename(columns={"SatisfactionLevel": "Environment_SatisfactionLevel"}).drop(columns=["SatisfactionID"])

# Job Satisfaction
Employees = Employees.merge(Satisfied_Level, left_on="JobSatisfaction", right_on="SatisfactionID", how="left"
                           ).rename(columns={"SatisfactionLevel": "Job_SatisfactionLevel"}).drop(columns=["SatisfactionID"])

# Relationship Satisfaction
Employees = Employees.merge(Satisfied_Level, left_on="RelationshipSatisfaction", right_on="SatisfactionID", how="left"
                           ).rename(columns={"SatisfactionLevel": "Relationship_SatisfactionLevel"}).drop(columns=["SatisfactionID"])

# Work-Life Balance
Employees = Employees.merge(Satisfied_Level, left_on="WorkLifeBalance", right_on="SatisfactionID", how="left"
                           ).rename(columns={"SatisfactionLevel": "WLB_Level"}).drop(columns=[ "SatisfactionID"])
```

```
[37]: #Showing first few rows  
Employees.head(10)
```

	EmployeeID	FirstName	LastName	Gender	Age	BusinessTravel	Department	DistanceFromHomeInKM	State	Ethnicity	...	TrainingOpportunitiesTaken	WorkLifeBalance
0	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	0.0	
1	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	1.0	
2	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	0.0	
3	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	1.0	
4	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	0.0	
5	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	0.0	
6	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	0.0	
7	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	0.0	
8	3012-1A41	Leonelle	Simco	Female	30	Some Travel	Sales	27	IL	White	...	1.0	
9	CBCB-9C9D	Leonerd	Aland	Male	38	Some Travel	Sales	23	CA	White	...	0.0	

10 rows × 41 columns

Exploring our Dataset (EDA)

3. Exploratory Data Analysis

In this section, we will explore the datasets to understand its structure and get Summary Statistics on it before performing data analysis

```
[38]: # Checking shape of the data  
Employees.shape
```

```
[38]: (6899, 41)
```

```
[39]: # Checking info of the data  
Employees.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 6899 entries, 0 to 6898  
Data columns (total 41 columns):  
 #   Column           Non-Null Count  Dtype     
---  --     
 0   EmployeeID      6899 non-null   object    
 1   FirstName       6899 non-null   object    
 2   LastName        6899 non-null   object    
 3   Gender          6899 non-null   object    
 4   Age              6899 non-null   int64     
 5   BusinessTravel  6899 non-null   object    
 6   Department      6899 non-null   object    
 7   DistanceFromHomeInKM 6899 non-null   int64     
 8   State            6899 non-null   object    
 9   Ethnicity        6899 non-null   object    
 10  Education        6899 non-null   int64     
 11  EducationField   6899 non-null   object    
 12  JobRole          6899 non-null   object    
 13  MaritalStatus    6899 non-null   object    
 14  Salary            6899 non-null   int64     
 15  StockOptionLevel 6899 non-null   int64     
 16  Overtime          6899 non-null   object    
 17  HireDate         6899 non-null   object    
 18  Attrition        6899 non-null   object    
 19  YearsAtCompany   6899 non-null   int64     
 20  YearsInMostRecentRole 6899 non-null   int64     
 21  YearsSinceLastPromotion 6899 non-null   int64     
 22  YearsWithCurrManager 6899 non-null   int64     
 23  EducationLevelID 6899 non-null   int64     
 24  EducationLevel    6899 non-null   object    
 25  PerformanceID    6709 non-null   object    
 26  ReviewDate        6709 non-null   object    
 27  EnvironmentSatisfaction 6709 non-null   float64   
 28  JobSatisfaction   6709 non-null   float64   
 29  RelationshipSatisfaction 6709 non-null   float64   
 30  TrainingOpportunitiesWithinYear 6709 non-null   float64   
 31  TrainingOpportunitiesTaken    6709 non-null   float64   
 32  WorkLifeBalance    6709 non-null   float64   
 33  SelfRating         6709 non-null   float64   
 34  ManagerRating      6709 non-null   float64
```

Descriptive/ Summary Statistics

```
[40]: #Showing Summary Statistics
```

```
Employees.describe()
```

	Age	DistanceFromHomeInKM	Education	Salary	StockOptionLevel	YearsAtCompany	YearsInMostRecentRole	YearsSinceLastPromotion	YearsW
count	6899.000000	6899.000000	6899.000000	6899.000000	6899.000000	6899.000000	6899.000000	6899.000000	6899.000000
mean	30.604146	22.327874	2.866647	110898.374112	0.725467	5.578055	2.778953	4.143934	
std	7.986542	12.899799	1.028396	98427.862382	0.839724	3.410087	2.810170	3.203770	
min	18.000000	1.000000	1.000000	20387.000000	0.000000	0.000000	0.000000	0.000000	
25%	25.000000	12.000000	2.000000	44646.000000	0.000000	3.000000	0.000000	1.000000	
50%	28.000000	22.000000	3.000000	74458.000000	1.000000	6.000000	2.000000	4.000000	
75%	36.000000	33.000000	4.000000	137219.500000	1.000000	9.000000	5.000000	7.000000	
max	51.000000	45.000000	5.000000	547204.000000	3.000000	10.000000	10.000000	10.000000	

Data Cleaning (Checking Nulls & Duplicates)

4. Data Cleaning

Before analyzing the dataset, it is important to perform data cleaning to ensure accuracy, consistency, and reliability of results. In this process, we can handle (Missing values, Duplicate Records, Inconsistencies and any errors/ Outliers detected in the datasets).

This process helps us improve data quality and make it ready for exploration, visualization, and modeling which ensures that insights and conclusions drawn from the analysis are trustworthy.

A. Checking Duplicates

```
[1]: Employees.duplicated().sum()  
[1]: np.int64(0)
```

It appears that there are no duplicate values within this dataset.

B. Checking Null Values

```
[2]: Employees.isna().sum()  
[2]: EmployeeID          0  
FirstName           0  
LastName            0  
Gender              0  
Age                 0  
BusinessTravel      0  
Department          0  
DistanceFromHomeInKM 0  
State               0  
Ethnicity            0  
Education            0  
EducationField       0  
JobRole              0  
MaritalStatus        0  
Salary               0  
StockOptionLevel     0  
OverTime             0  
HireDate             0  
Attrition            0  
YearsAtCompany       0  
YearsInMostRecentRole 0  
YearsSinceLastPromotion 0  
YearsWithCurrManager 0  
EducationLevelID     0  
EducationLevel       0  
PerformanceID       190  
ReviewDate           190  
EnvironmentSatisfaction 190  
JobSatisfaction     190  
RelationshipSatisfaction 190
```

Data Cleaning (Checking Nulls & Duplicates)

```
TrainingOpportunitiesWithinYear    190
TrainingOpportunitiesTaken        190
WorkLifeBalance                  190
SelfRating                       190
ManagerRating                     190
Mananger_RatingLevel             190
Self_RatingLevel                 190
Environment_SatisfactionLevel   190
Job_SatisfactionLevel            190
Relationship_SatisfactionLevel  190
WLB_Level                         190
dtype: int64

43]: missing_rows = Employees[Employees.isnull().any(axis=1)]
print(missing_rows)

EmployeeID FirstName LastName Gender Age BusinessTravel \
6421 3BAB-C57E Hashim Payn Male 19 Some Travel
6422 BC20-368E Flem Simenon Male 39 Some Travel
6423 F3AE-DEF2 Vern Foord Male 20 Frequent Traveller
6424 F5B1-A266 Arly Fleetham Female 22 Some Travel
6425 2EB9-9D31 Cybill Mitkin Female 27 No Travel
... ...
6894 467E-977A Jud Melanaphy Male 20 Some Travel
6895 6FB9-A624 Marc Calver Non-Binary 27 Some Travel
6896 EBF4-5928 Rudolph MacDearmont Male 21 Some Travel
6897 60E6-B1D9 Merill Agg Male 21 Some Travel
6898 84D4-D4C3 Naoma Hebbard Female 20 No Travel

Department DistanceFromHomeInKM State           Ethnicity \
6421 Technology 28 CA White
6422 Sales      37 NY White
6423 Technology 39 CA White
6424 Technology 34 CA White
6425 Human Resources 4 CA White
... ...
6894 Technology 28 CA Black or African American
6895 Technology 8 CA Black or African American
6896 Sales      4 CA Black or African American
6897 Technology 7 CA Black or African American
6898 Technology 28 CA Black or African American

... TrainingOpportunitiesTaken WorkLifeBalance SelfRating \
6421 ... NaN NaN NaN
6422 ... NaN NaN NaN
6423 ... NaN NaN NaN
6424 ... NaN NaN NaN
6425 ... NaN NaN NaN
... ...
6894 ... NaN NaN NaN
6895 ... NaN NaN NaN
6896 ... NaN NaN NaN
6897 ... NaN NaN NaN
6898 ... NaN NaN NaN

[190 rows x 41 columns]
```

It appears that there are 190 Employees that have no Performance Records

Checking Unique Values

```
[54]: # Checking Unique Values
for col in Employees:
    print('-'*30)
    print(col)
    print(Employees[col].unique())
-----
EmployeeID
['3012-1A41' 'CBCB-9C9D' '95D7-1CE9' ... 'EBF4-5928' '60E6-B1D9'
 '84D4-D4C3']
-----
FirstName
['Leonelle' 'Leonerd' 'Ahmed' ... 'Rudolph' 'Merill' 'Naoma']
-----
LastName
['Simco' 'Aland' 'Sykes' ... 'MacDearmont' 'Agg' 'Hebbard']
-----
Gender
['Female' 'Male' 'Non-Binary' 'Prefer Not To Say']
-----
Age
[30 38 43 39 29 34 42 40 31 32 35 37 33 45 36 48 47 44 28 41 46 27 26 25
 24 23 22 21 49 50 20 19 51 18]
-----
BusinessTravel
['Some Travel' 'No Travel' 'Frequent Traveller']
-----
Department
['Sales' 'Human Resources' 'Technology']
-----
DistanceFromHomeInKM
[27 23 29 12 30 45 3 20 4 42 8 35 21 34 19 1 17 36 37 41 25 44 22 40
 14 15 7 5 31 9 13 38 16 39 28 32 2 6 11 33 26 43 10 18 24]
-----
State
['IL' 'CA' 'NY']
-----
Ethnicity
['White' 'Asian or Asian American' 'Mixed or multiple ethnic groups'
 'Black or African American' 'Native Hawaiian' 'Other'
 'American Indian or Alaska Native']
-----
Education
['Doctorate' 'Masters' 'Bachelors' 'High School'
 'No Formal Qualifications']
-----
EducationField
['Marketing' 'Marketing' 'Computer Science' 'Technical Degree'
 'Information Systems' 'Other' 'Economics' 'Human Resources'
 'Business Studies']
```

Trimming Extra Spaces & Removing Inconsistencies

In Education Field: there are 'Marketing' and 'Marketing ' (with an empty space)

```
[62]: # Replacing Marketing with white space with (Marketing)
Employees['EducationField'] = Employees['EducationField'].replace('Marketing ', 'Marketing')

[66]: # Ensure all data in Marketing Department are standardized
Employees['EducationField'].unique()

[66]: array(['Marketing', 'Computer Science', 'Technical Degree',
           'Information Systems', 'Other', 'Economics', 'Human Resources',
           'Business Studies'], dtype=object)

[60]: # Converting the data type into datetime
Employees['HireDate'] = pd.to_datetime(Employees['HireDate'], errors='coerce')
Employees['ReviewDate'] = pd.to_datetime(Employees['ReviewDate'], errors='coerce')

[65]: # Removing trailing white Extra Spaces from Columns and Text Fields

Employees.columns = Employees.columns.str.strip()
text_columns = Employees.select_dtypes(include='object').columns
Employees[text_columns] = Employees[text_columns].apply(lambda x: x.str.strip())

[67]: # Removing Logically Inconsistent Records
Employees = Employees[
    (Employees['YearsWithCurrManager'] <= Employees['YearsAtCompany']) &
    (Employees['YearsSinceLastPromotion'] <= Employees['YearsAtCompany'])]

Employees = Employees[Employees['YearsAtCompany'] <= Employees['Age']]

[68]: # Extracting Year from Review Date
Employees['ReviewYear'] = Employees['ReviewDate'].dt.year
```

Checking if there are any Outliers

```
[70]: # Check for Outliers
numeric_cols = Employees.select_dtypes(include=['int64', 'float64']).columns

for col in numeric_cols:
    Q1 = Employees[col].quantile(0.25)
    Q3 = Employees[col].quantile(0.75)
    IQR = Q3 - Q1

    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR

    outliers = Employees[(Employees[col] < lower_bound) | (Employees[col] > upper_bound)]

    if len(outliers) > 0:
        print(f"\n♦ Column: {col}")
        print(f"    Number of Outliers: {len(outliers)}")
    else :
        print("No Outliers")

No Outliers
No Outliers

    ♦ Column: Salary
    Number of Outliers: 561

    ♦ Column: StockOptionLevel
    Number of Outliers: 361
No Outliers
```

Experience Level in relation to Years in Recent Role

Define and Apply Experience Level Based on Years in Most Recent Role

```
[71]: def experience_level(x):
    if x < 2:
        return 'Junior'
    elif x < 5:
        return 'Mid'
    else:
        return 'Senior'

Employees['ExperienceLevel'] = Employees['YearsInMostRecentRole'].apply(experience_level)

[72]: # Detect and Count Salary Outliers Before Capping (Grouped by Year, Job Role, and Experience Level)

outliers_before = []

for (year, role, exp), group in Employees.groupby(['ReviewYear', 'JobRole', 'ExperienceLevel']):
    Q1 = group['Salary'].quantile(0.25)
    Q3 = group['Salary'].quantile(0.75)
    IQR = Q3 - Q1
    lower = Q1 - 1.5 * IQR
    upper = Q3 + 1.5 * IQR

    group_outliers = group[(group['Salary'] < lower) | (group['Salary'] > upper)]
    outliers_before.append(group_outliers)

outliers_before_df = pd.concat(outliers_before)
num_outliers_before = len(outliers_before_df)
print(f"Outliers_Num_Before_Preprocessing: {num_outliers_before}")

Outliers_Num_Before_Preprocessing: 277
```

```
[73]: # Apply Contextual Capping to Adjust Salary Outliers (Grouped by Year, Job Role, and Experience Level)

data_capped = Employees.copy()

for (year, role, exp), group in Employees.groupby(['ReviewYear', 'JobRole', 'ExperienceLevel']):
    Q1 = group['Salary'].quantile(0.25)
    Q3 = group['Salary'].quantile(0.75)
    IQR = Q3 - Q1
    lower = Q1 - 1.5 * IQR
    upper = Q3 + 1.5 * IQR

    data_capped.loc[Employees.index, 'New_Salary'] = Employees['Salary'].clip(lower, upper)
```

```
[74]: # Detect and Count Salary Outliers After Capping (Grouped by Year, Job Role, and Experience Level)
```

```
outliers_after = []

for (year, role, exp), group in data_capped.groupby(['ReviewYear', 'JobRole', 'ExperienceLevel']):
    Q1 = group['New_Salary'].quantile(0.25)
    Q3 = group['New_Salary'].quantile(0.75)
    IQR = Q3 - Q1
    lower = Q1 - 1.5 * IQR
    upper = Q3 + 1.5 * IQR

    group_outliers = group[(group['New_Salary'] < lower) | (group['New_Salary'] > upper)]
    outliers_after.append(group_outliers)

remaining_outliers = pd.concat(outliers_after)
num_outliers_after = len(remaining_outliers)
print(f"Outliers_Num_After_Preprocessing: {num_outliers_after}")

Outliers_Num_After_Preprocessing: 198
```

```
[77]: improvement = ((num_outliers_before - num_outliers_after) / num_outliers_before) * 100
```

```
print(f"Outliers_Num_Before_Preprocessing: {num_outliers_before}")
print(f"Outliers_Num_After_Preprocessing: {num_outliers_after}")
print(f"Improvement: {num_outliers_before - num_outliers_after} less")
print(f"Improvement_Rate: {improvement:.2f}%")
```

```
Outliers_Num_Before_Preprocessing: 277
Outliers_Num_After_Preprocessing: 198
Improvement: 79 less
Improvement_Rate: 28.52%
```

```
[76]: remaining_outliers[['EmployeeID', 'ReviewYear', 'JobRole', 'ExperienceLevel', 'YearsInMostRecentRole', 'New_Salary']].sort_values('New_Salary', ascending=False)
```

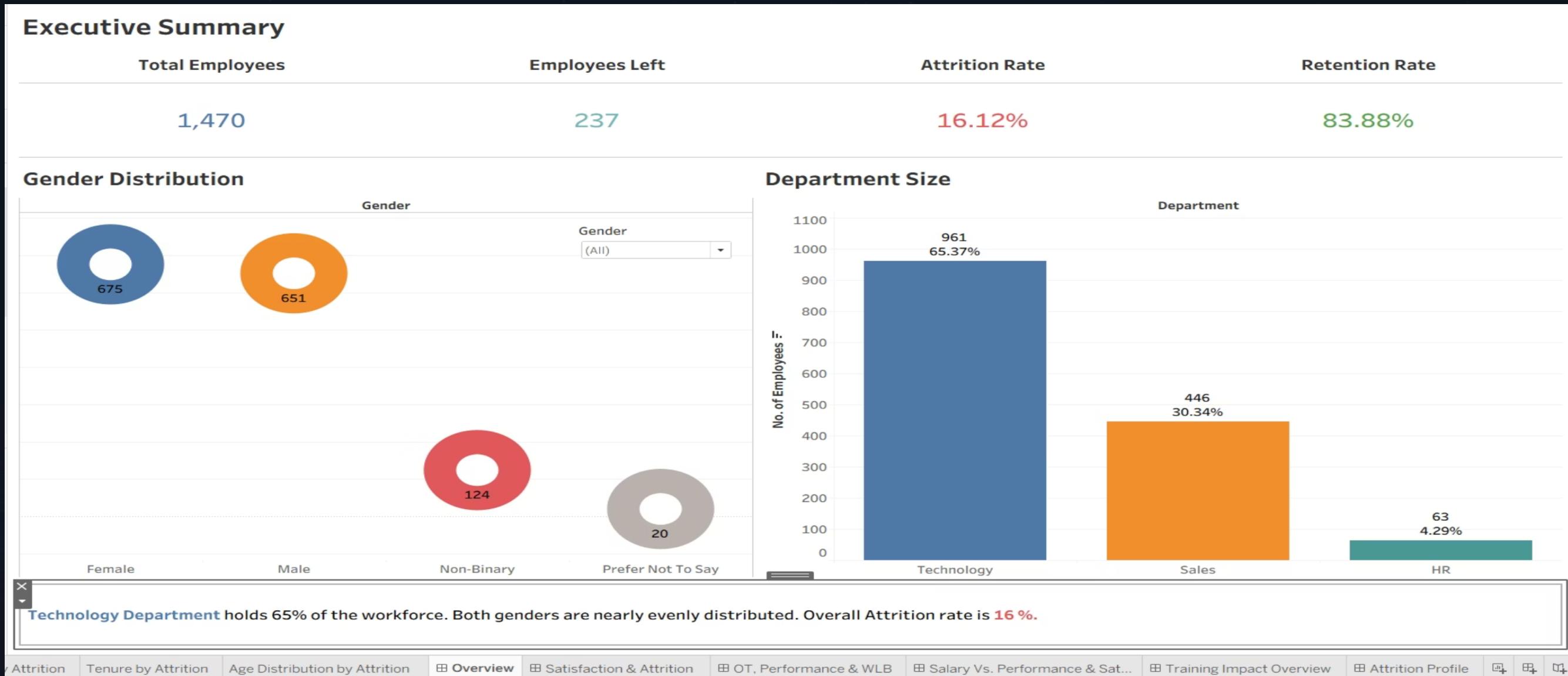
	EmployeeID	ReviewYear	JobRole	ExperienceLevel	YearsInMostRecentRole	New_Salary
2234	79F7-78EC	2013.0	Sales Executive	Senior	5	133668.0
1445	DF11-9C46	2016.0	Software Engineer	Mid	4	133668.0
156	ED73-F078	2016.0	Software Engineer	Mid	3	133668.0
1822	CD4D-7BBA	2016.0	Software Engineer	Junior	0	133668.0
163	ED73-F078	2015.0	Software Engineer	Mid	3	133668.0
1860	9811-F345	2016.0	Data Scientist	Senior	9	133668.0
1861	9811-F345	2017.0	Data Scientist	Senior	9	133668.0
160	ED73-F078	2020.0	Software Engineer	Mid	3	133668.0
1449	DF11-9C46	2020.0	Software Engineer	Mid	4	133668.0
3020	DEB9-1BBF	2020.0	Software Engineer	Mid	3	133668.0
580	4969-9A35	2020.0	Software Engineer	Senior	9	133668.0
920	3994-472A	2020.0	Software Engineer	Senior	10	133668.0
6367	AFC3-E23F	2021.0	Data Scientist	Junior	1	133668.0
5865	66D5-3142	2021.0	Data Scientist	Junior	0	133668.0
1858	819A-2C9C	2021.0	Data Scientist	Mid	2	133668.0
1865	9811-F345	2021.0	Data Scientist	Senior	9	133668.0
1827	CD4D-7BBA	2021.0	Software Engineer	Junior	0	133668.0
921	3994-472A	2021.0	Software Engineer	Senior	10	133668.0
3021	DEB9-1BBF	2021.0	Software Engineer	Mid	3	133668.0
581	4969-9A35	2021.0	Software Engineer	Senior	9	133668.0

```
[78]: data_capped.drop(
    data_capped[
        (data_capped["EmployeeID"] == "1A7C-19DB") &
        (data_capped["ReviewYear"] == 2015)
    ].index,
    inplace=True
)
```

```
[79]: data_capped.to_excel('HR_Cleaned.xlsx', index=False)
data_capped.to_csv('HR_Cleaned.csv', index=False)
```

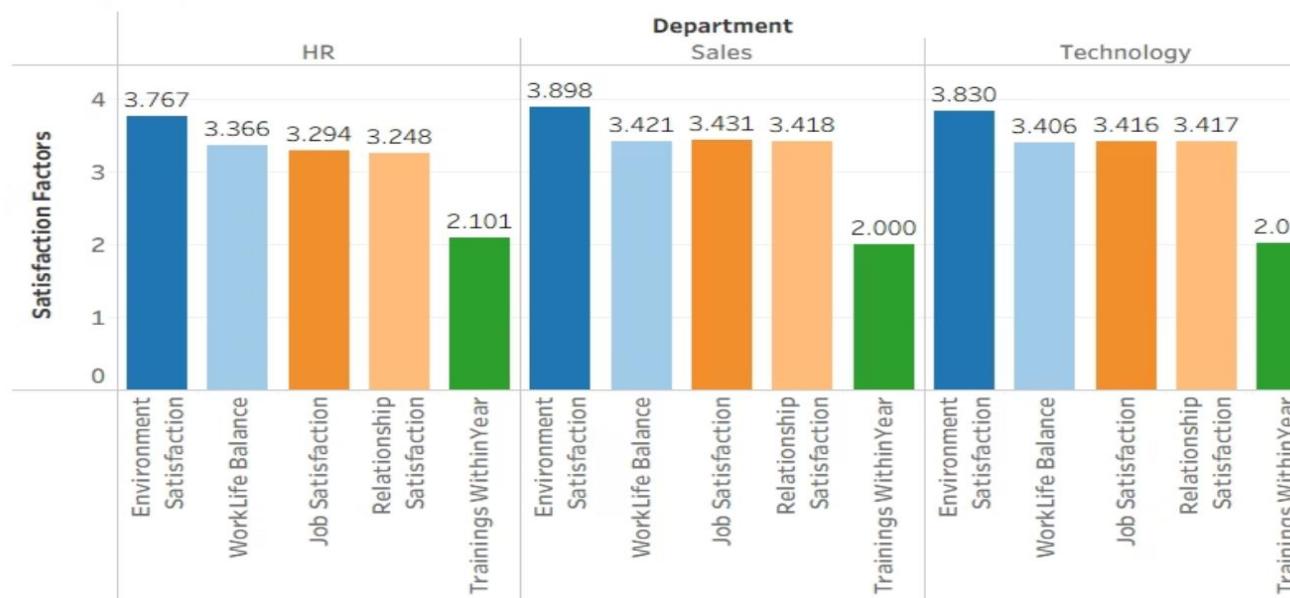
Data Visualization (Tableau Dashboard)

An Overview for HR Attrition Project

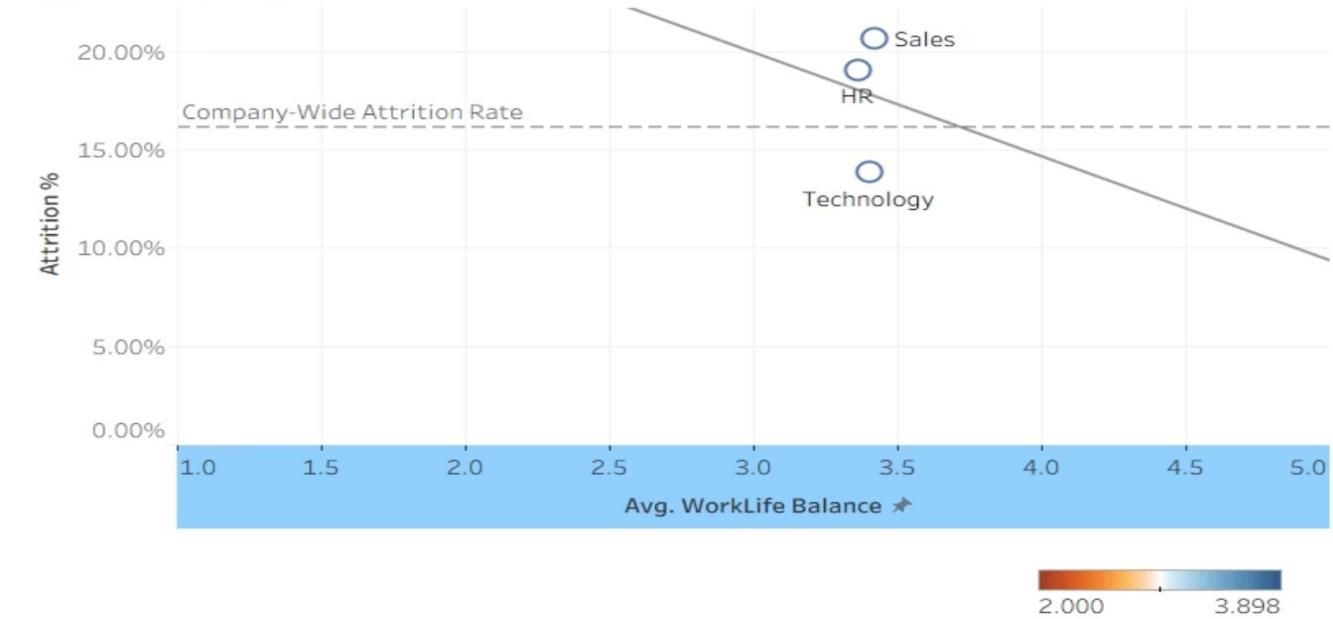


Satisfaction Drivers Dashboard

Employee Satisfaction & Attrition Drivers



WLB Vs. Attrition



Department Satisfaction Metrics

Department	Environment Satisfaction	Job Satisfaction	Relationship Satisfaction	WorkLife Balance	Trainings WithinYear
HR	3.767	3.294	3.248	3.366	2.101
Sales	3.898	3.431	3.418	3.421	2.000
Technology	3.830	3.416	3.417	3.406	2.016



Insights:

Satisfaction Scores across departments vary only slightly, indicating that Attrition is not driven by satisfaction alone. Other factors such as WLB, Salary, and Overtime appear more impactful.

Overtime Impact Dashboard

Impact of Overtime on Attrition & Performance

Total Employees	Employees Left	Attrition Rate	WorkLife Balance
1,470	237	16.12%	3

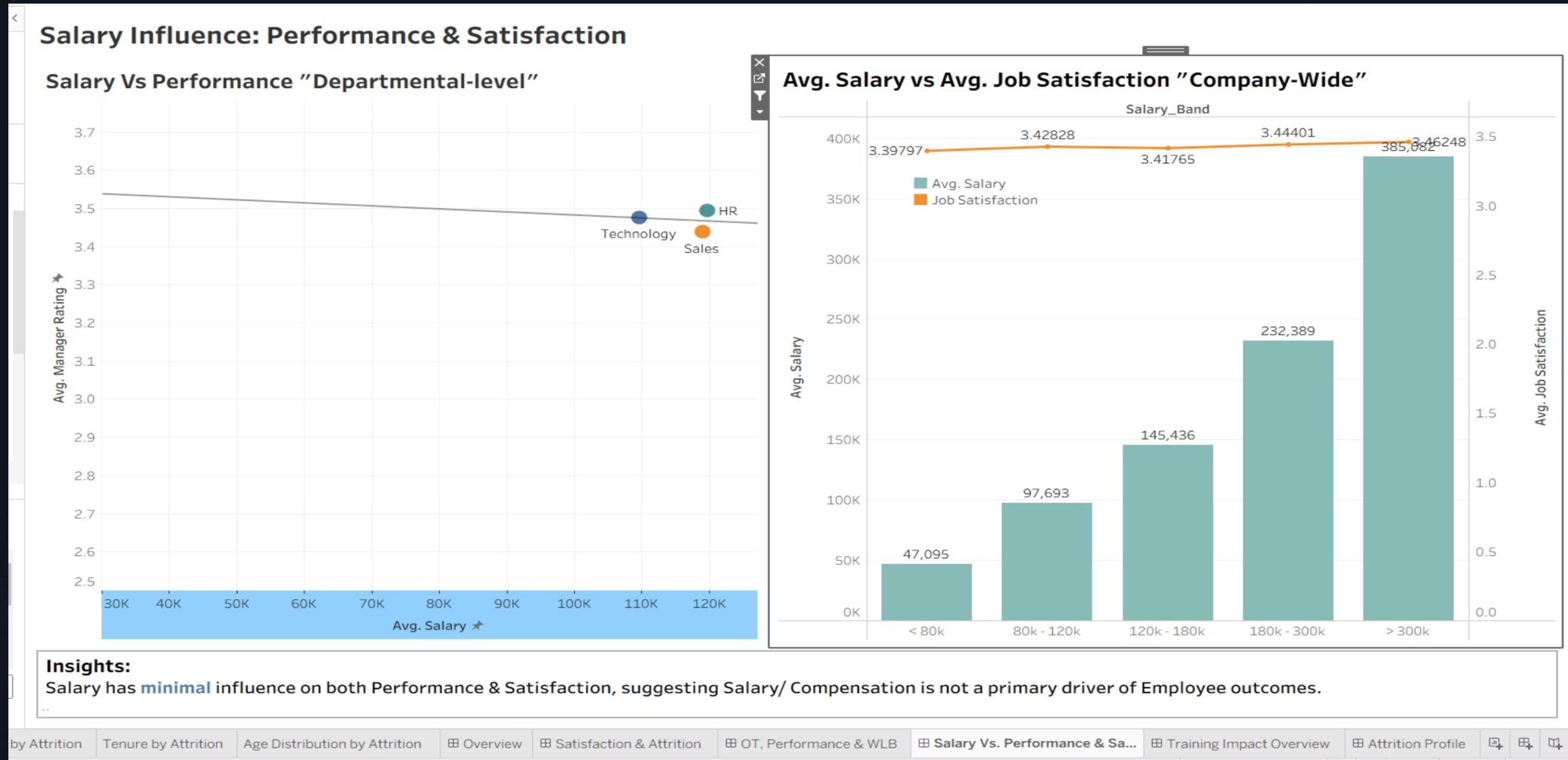
Overtime Vs. WLB & Performance



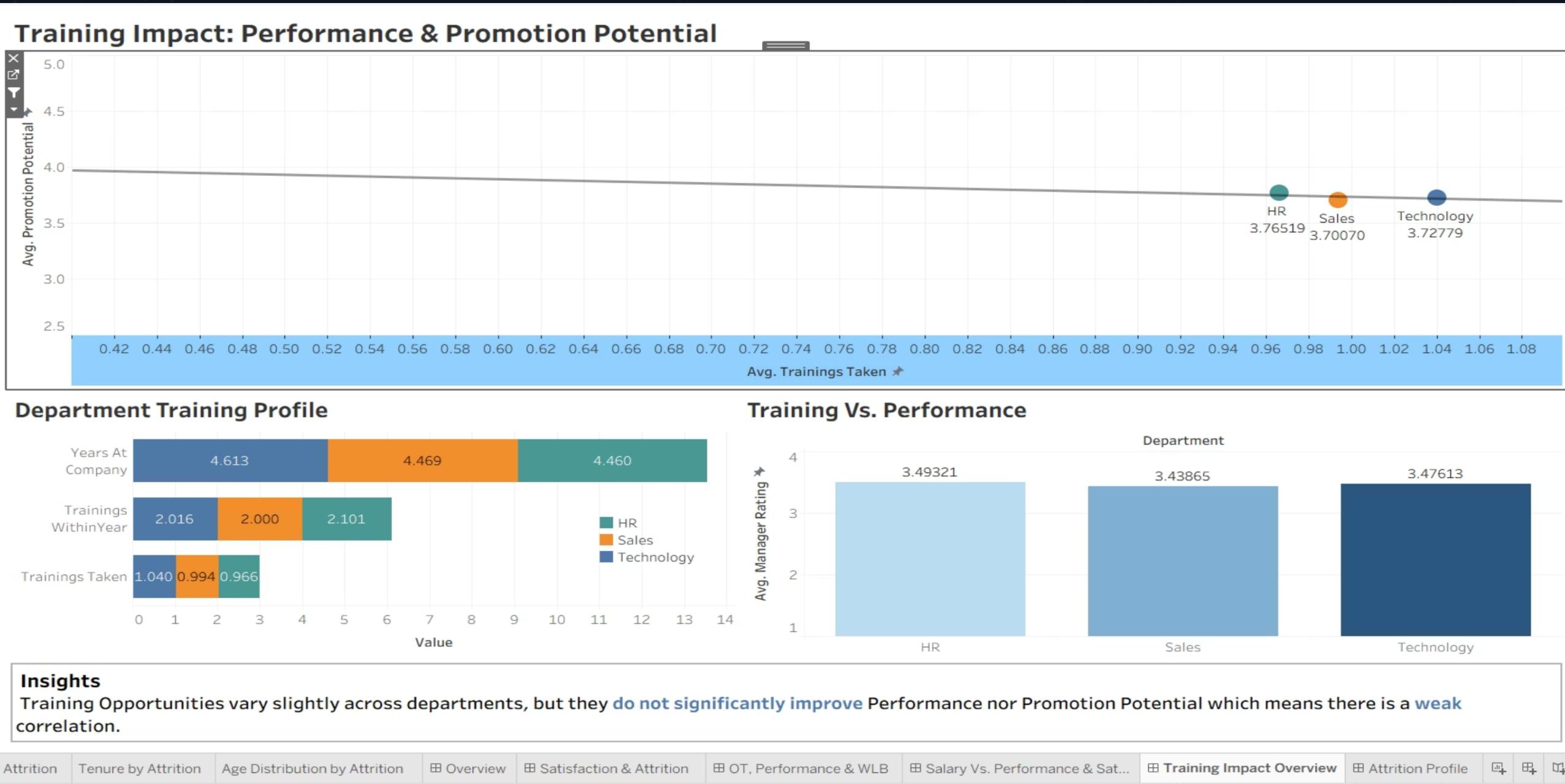
Insights:

Employees working overtime have higher Job Satisfaction, however, lower Performance indicating overtime affects employee performance.

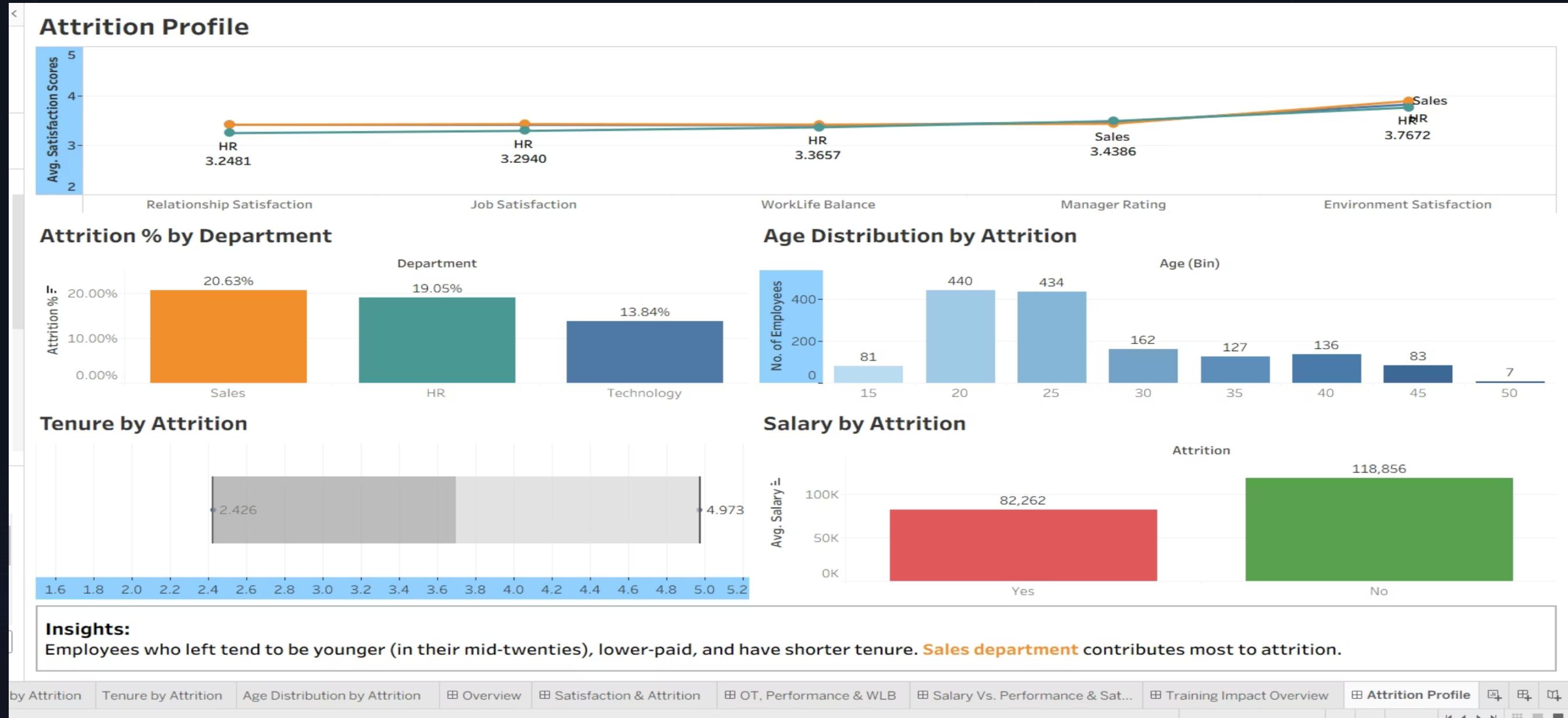
Salary Influence Dashboard



Training Impact Dashboard



Attrition Profile Dashboard



Conclusion

Our HR Analytics project revealed several key patterns shaping Employee Satisfaction and turnover:

- The **Technology department dominates the workforce (65%)**, and gender distribution is nearly **balanced** across the organization.
- Attrition is **moderately high at 16%**, but satisfaction scores across departments remain **relatively consistent**, suggesting that **Attrition is not driven by Satisfaction alone**, but rather by other operational and behavioral factors.
- **Salary level shows Limited influence** on Satisfaction or Performance, suggesting **compensation is not the main concern** for most employees.
- **Training opportunities show Weak impact** on Performance or Promotion readiness, indicating current programs **may lack effectiveness or alignment** with role needs.
- Employees who leave are generally **younger, earlier in their careers, lower-paid, and have shorter tenures**. The **Sales department contributes most** to the **overall attrition rate**.



Recommendations to effectively reduce Employee turnover & strengthen Employee Experience

Improve Early-Career Retention

- Create structured career paths, mentorship programs, and faster development cycles targeting younger, low-tenure employees who are most at risk.
- Create growth pathways to reduce early exits.



Enhance WLB Company-Wide

- Introduce flexible scheduling, remote-work options, or mandatory time-off initiatives.
- Regularly monitor WLB metrics for early warning signs.

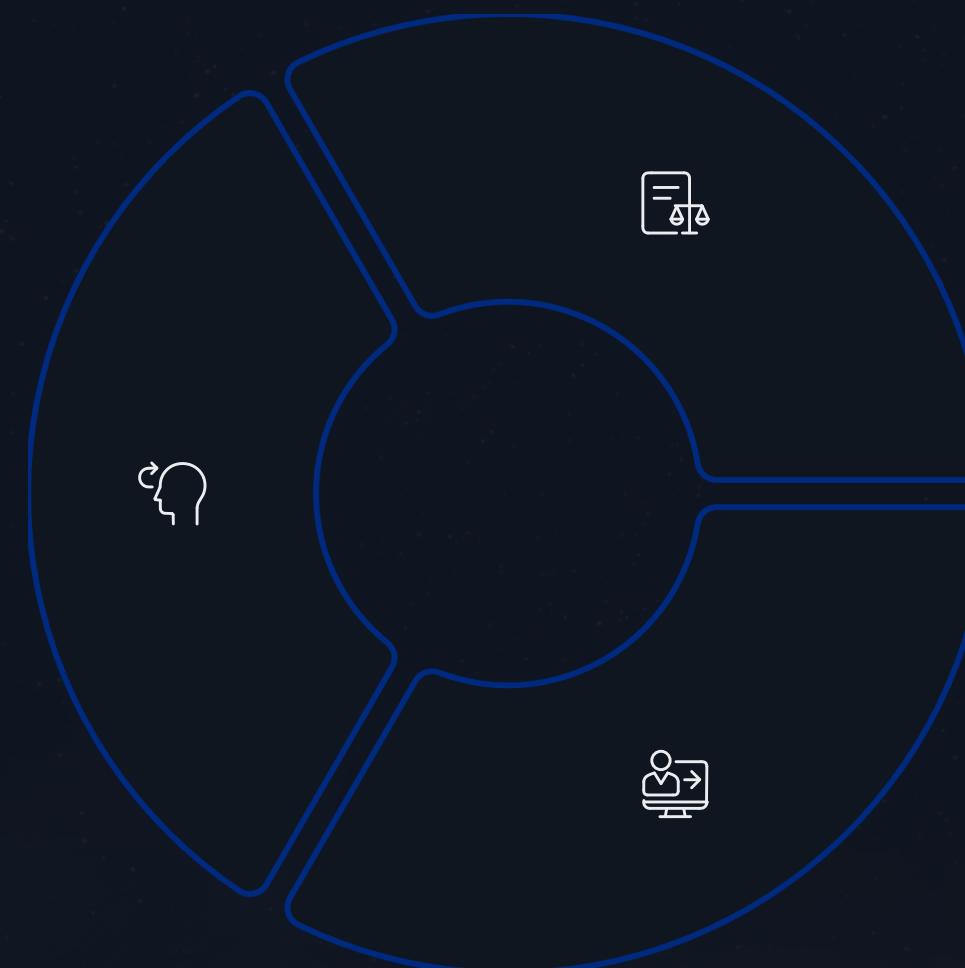
Address High Attrition in Sales Department

- Review workload, targets, incentives, and managerial support practices.
- Conduct focused interviews or targeted pulse surveys to uncover department-specific stressors.

Recommendations to effectively reduce Employee turnover & strengthen Employee Experience

Reassess Training Strategies

- Shift from quantity to **quality and relevance**: design role-specific, skill-based, and progression-aligned development programs.
- Evaluate training effectiveness and ensure it contributes to skill growth and promotion readiness.



Maintain Transparent Compensation Strategy

- While salary is not a primary turnover driver, ensure fairness across levels and roles to prevent future dissatisfaction.

Adopt Predictive Analytics (Strategic Add-On)

- Build attrition prediction models to proactively identify high-risk employees and intervene early.

Thank You!

Feel free to ask any questions!

