# PMI Technical Assessment
# People Analytics
# Data Science Lead

Applicant

Aldo Cantu

# Methodology

| Business Understanding | Design | Data Preparation | Feature Engineering & Modeling | Visualization & Presentation |
|---|---|---|---|---|
| • Understand and frame the problem<br>• Gather Requirements | • Data collection<br>• Initial EDA & Data assessment<br>• Bucketing Strategies<br>• Making assumptions | • Data Wrangling<br>• EDA | • Statistical methods<br>• Feature Engineering<br>• Model building<br>• Model assessment | • Graphs / Plot selection<br>• Make recommendations |

# Technical challenge:
# San Francisco Taxi Cabs 🚕

# San Francisco Taxi Cabs



**Problem Statement:**

Taxis that operate in San Francisco often roam without passengers. This yields unnecessary $CO_2$ emissions that affect our current global warming situation. The goal of this challenge is to analyse and calculate these emissions based on information collected of about 500 taxi cabs on a 30 day time frame.

**Specified Assumptions:**

Monthly taxi fleet change rate = 15%
Average $CO_2$ emission of a passenger vehicle = 400 gr./m

**Requirements:**

- Calculate potential yearly reduction in $CO_2$ emissions caused by taxi cabs without passengers.
- Data Visualization

**Bonus:**

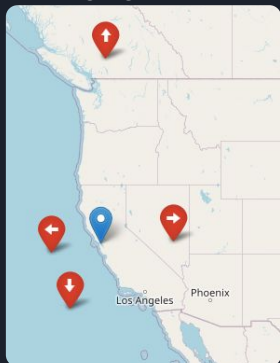- 📍 Predict the next place a passenger will hail a cab

🚕 San Francisco Taxi Cabs

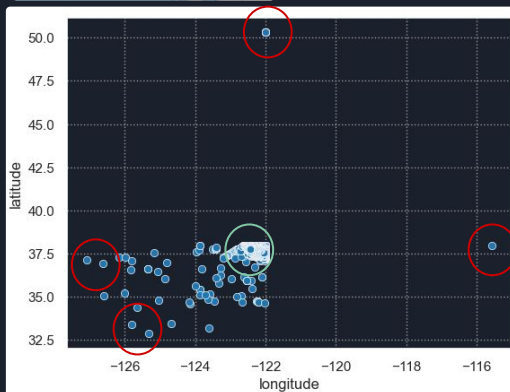# Initial EDA, Data assessment & Assumptions

📍 Location Analysis

💡 Assumptions          Correction

## As-is



| Farmost Cardinal Point | Limit |
|---|---|
| West Longitude | -122.51494 |
| East Longitude | -122.357034 |
| North Latitude | 37.811054 |
| South Latitude | 37.7080955 |

💡 Hard Limits (Physical)!

**OUTLIERS**
106275 Fares outside the assumed limits







6

San Francisco Taxi Cabs

# Initial EDA, Data assessment & Assumptions

Event Duration Analysis

**As-is**

💡Assumptions

Correction



| occupancy | min duration | max duration | Quantile |
|-----------|--------------|--------------|----------|
| empty | 179 sec | 5090 sec. | 5/75 |
| occupied | 242 sec | 2039 sec. | 5/95 |

Using Quantiles!



**OUTLIERS**
235206 Fares outside the assumed limits

🚕 San Francisco Taxi Cabs

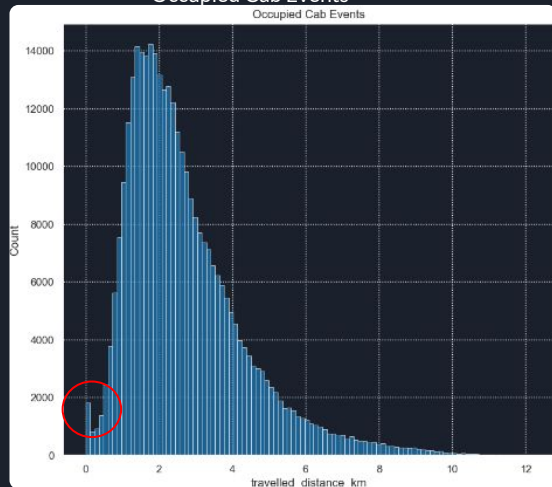# Initial EDA, Data assessment & Assumptions

### Distance Analysis

**As-is**

💡 **Assumptions**

**Correction**

Occupied Cab Events



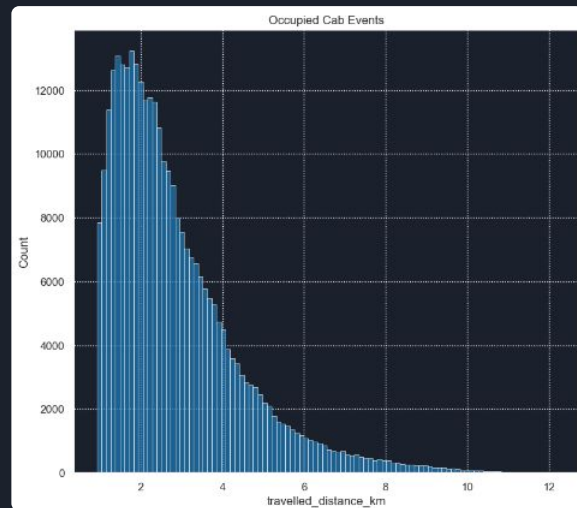Occupied cabs that travel less than 0.92 km are outliers (non-inclusive)
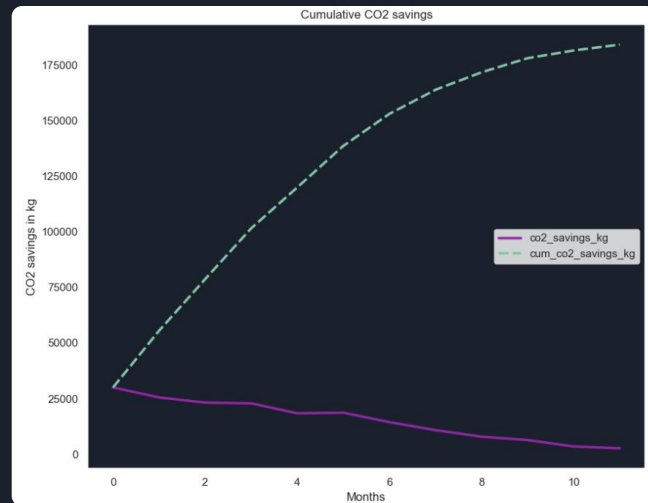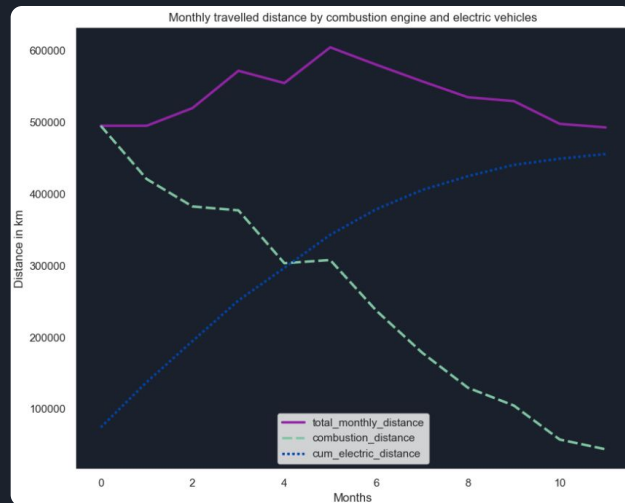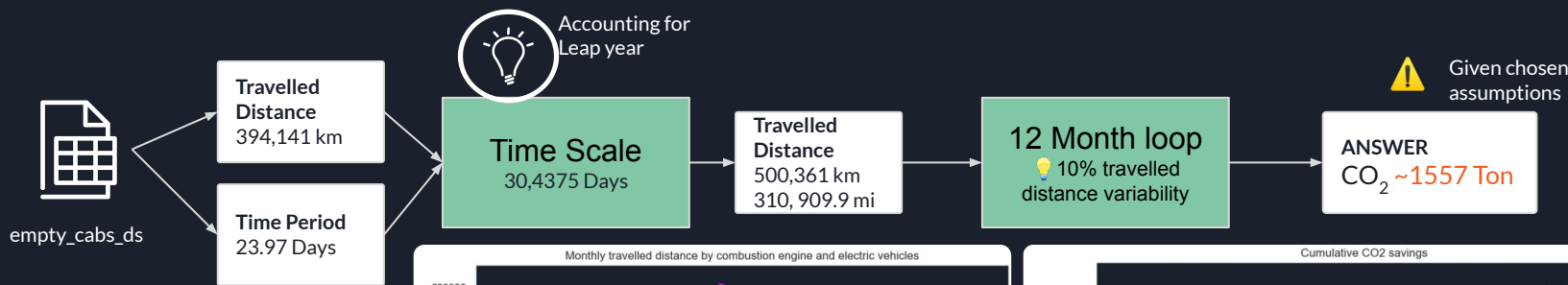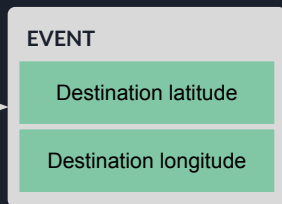
5% Quantile

Using Quantiles!



**OUTLIERS**
21046 Fares outside the assumed limits

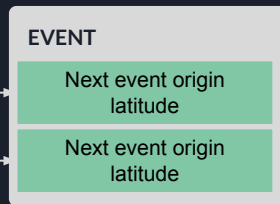# Next Fare Predictor

cabs_ds

**EVENT**
Destination latitude
Destination longitude

Shift

**EVENT**
Next event origin latitude
Next event origin latitude

**MSE**
0.6 mi
Distance

**EVENT**
Week day
hours

**Hour Bin**
Morning
Afternoon
Night
Late

Week Day

Start time of the Day per event

End time of Day per event

KMeans

Features

EVENT = state change ( free <-> occupied )

11

# Next Fare Predictor



**MSE**
1.64 mi

Distance

# Further Steps

- Adding a seasonality
- Taking weather into account
- Further time decomposition
- Implement a cross entropy Naive Bayes or KNN model with multivariate grid-system

Use case:
Turnover Analysis 👎

14

# Turnover Analysis



Problem Statement:

The CHRO has asked the Head of Talent Management to devise a retention strategy to make sure  top talent is retained.  The goal is to analyze employee turnover data that has been collected over the past twelve months. To answer the following questions:
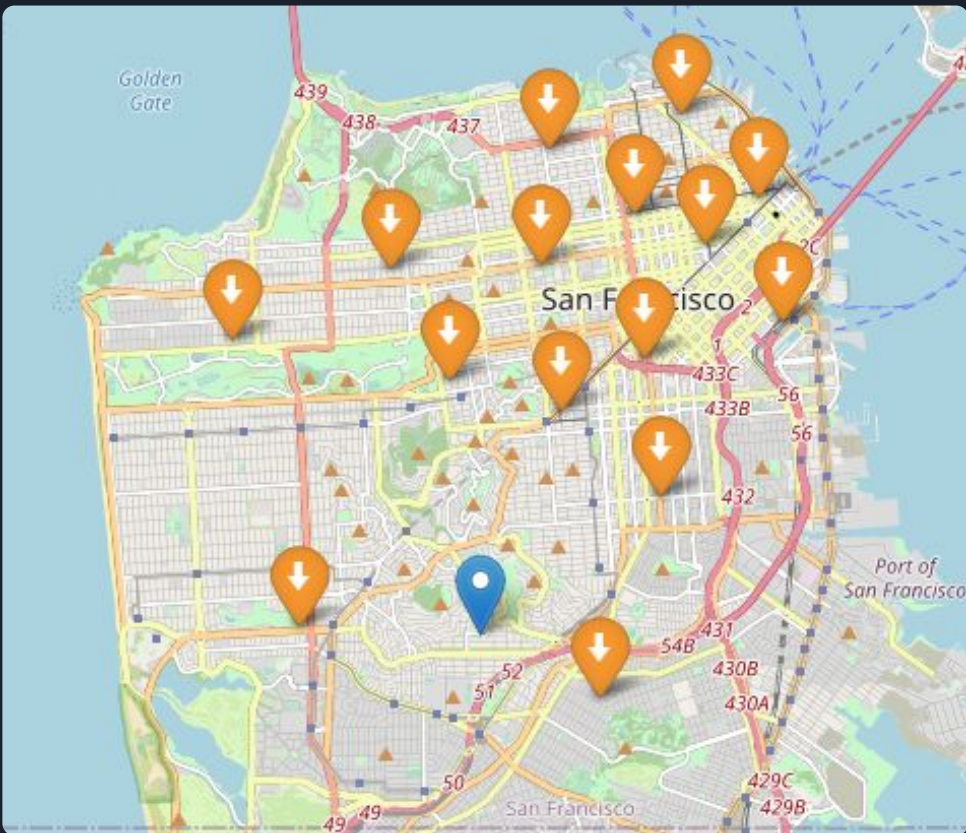
1. Is the company losing top talent?
2. If so, why and how can current top talent be retained

Requirements:

- Data Visualization
- Generate a strategy that will help lessen the turnover of top talented employees

15

# Reframing the problem

Since the question "*Is the company losing top talent?*" is not a testable hypothesis a reframing of the problem needs to be done.

On a study done by an employee retention platform [1], the overall average turnover rate in the US is about **47.2%**.

Moreover, another article written by Michael Mankins (2017)[2] on the Harvard Business Review site:
 "*On average, **15%** of a company's workforce - roughly one in seven employees - are A player or "stars".*"

So our hypothesis stands by the following:

"**Is the company losing more or equal to 6.7% of top talent?**"

# Is the company losing top talent?

|  | Company | Global |
|---|---|---|
| Overall company's turnover rate: | **11.2%** 👍 | 47.2% |
| High Potentials turnover rate: | **16.25%** 👎 | 6.7% |

$$Turnover\ Rate(\%) = \frac{\#Separations}{Avg.\#Employees} * 100$$

\* High potentials numbers used for High Potential turnover

**Yes**, the company is losing top talent with a significant number compared to what other companies have previously reported.

# Data assessment

From the initial data state, a few remarks have surfaced:

⚠️ "Tenure" and "Time in latest role" **features have negative values**

⚠️ Max values for "Tenure" and "Time in latest role" are **above 90 years** whereas realistically a maximum value of 50 years should be addressed

⚠️ "Actual monthly pay" feature seems to fit the description and values of a **week's worth pay**

⚠️ In "Country" and "Nationality" feature, *Monrovia* is declared as country/nationality whereas the correct value is *Liberia*

⚠️ In "Race" feature, "White" and "white" (upper- & lowercase) is treated as different category

⚠️ There is a typo on the "Sexuality" feature. Instead of *Non-hetersexual* it should be *Non-heterosexual*. (missing an "o")

⚠️ There is a typo on the "Career opportunities Survey score" feature. It should be "Career Opportunities Survey Score"
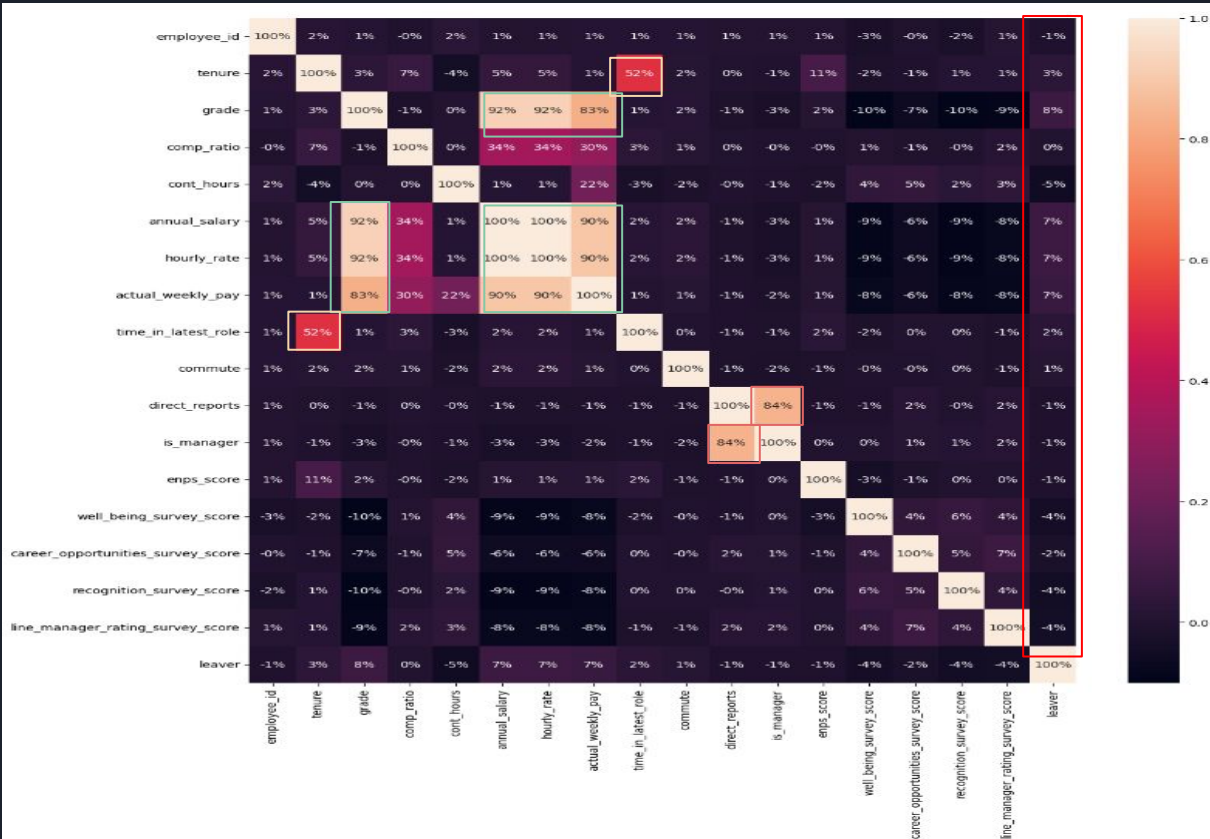
Business Understanding | Design | Data Preparation | Feature Engineering & Modeling | Visualization & Presentation

# Assumptions

💡 Negative values are filtered out

💡 Maximum tenure of 50 years

19

# Initial EDA

Pearson Correlation

Correlation matrix of Employee Turnover Data

# Initial EDA

A quick study of the
correlation matrix leads us to the following insights:

From the initial assessment of the information there are no significant correlations.
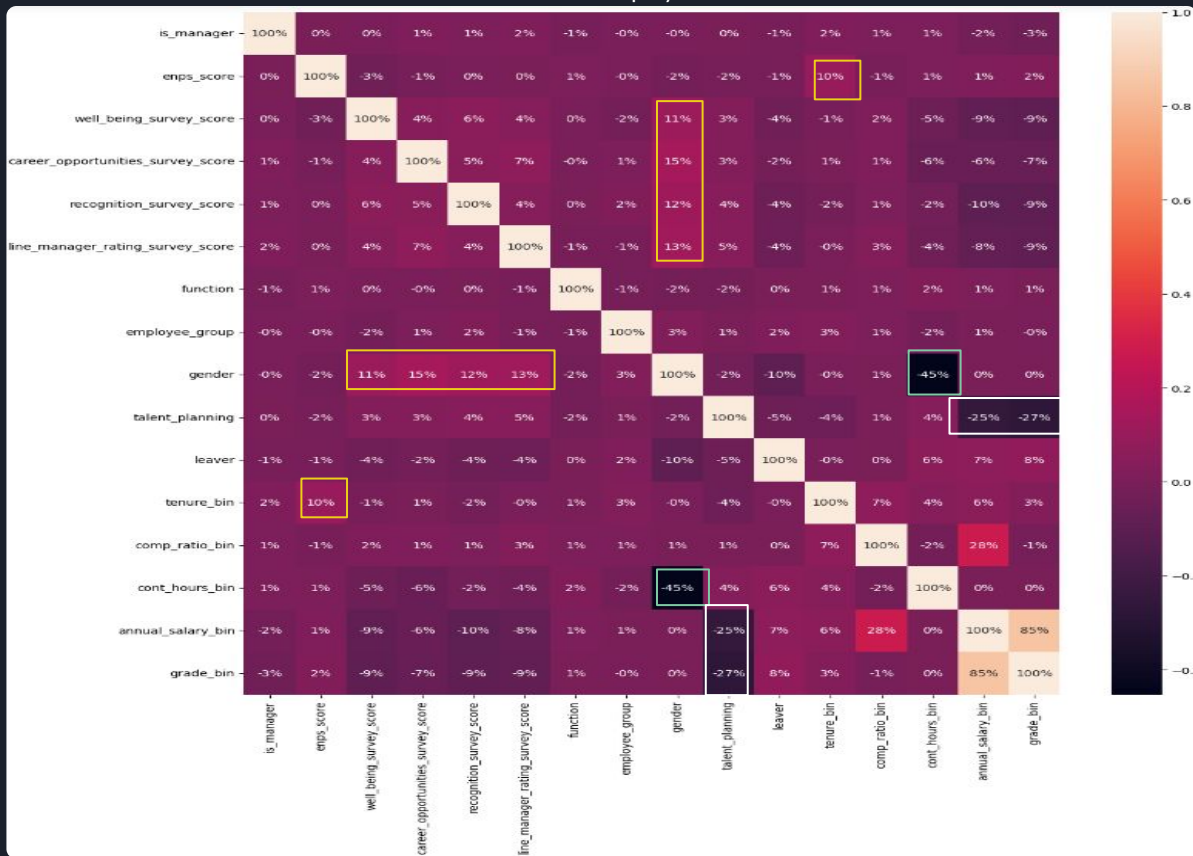
Some relationships are appreciated:
- Grade -> Annual Salary, Hourly Rate, Actual weekly pay
- Is Manager -> Direct Reports
- Tenure -> Time in latest Role

Further investigation should be done.

# Exploratory Data Analysis

Correlation matrix of cleaned Employee Turnover Data

# Exploratory Data Analysis

After some refactoring of the data the following information could be concluded:

Weak Negative Relationships:
-Talent Planning -> Annual Salary, Grade

Moderate Negative Relationship:
-Cont. Hours -> Gender

Weak positive Relationships:
-Survey Scores -> Gender
-eNPS Score -> Tenure

Not in scope:
-Comp. Ratio -> Annual Salary

Now a breakdown of personnel marked as "High Potential" from the Talent Planning Dimension

👎 Turnover Analysis

# Exploratory Data Analysis

Bucketing Strategies

💡 Since "Financial features" like Annual Salary, Hourly Rate and Actual week pay are fairly similar, only Annual Salary will be held on focus and further **bucketed in 10 thousand unit increments.**

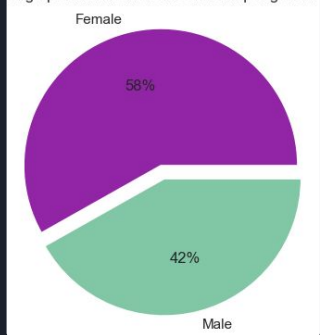💡 Grade is also going to be grouped in **increments of 5**

❌Ignored Features:
- Employee ID
- Hourly Rate
- Actual Week Pay
- Time in latest Role
- Commute
- Direct Reports
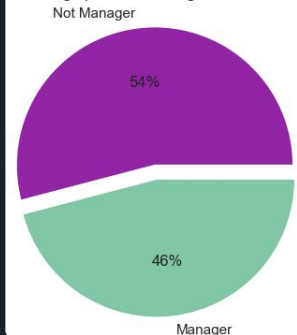- Country
- Nationality
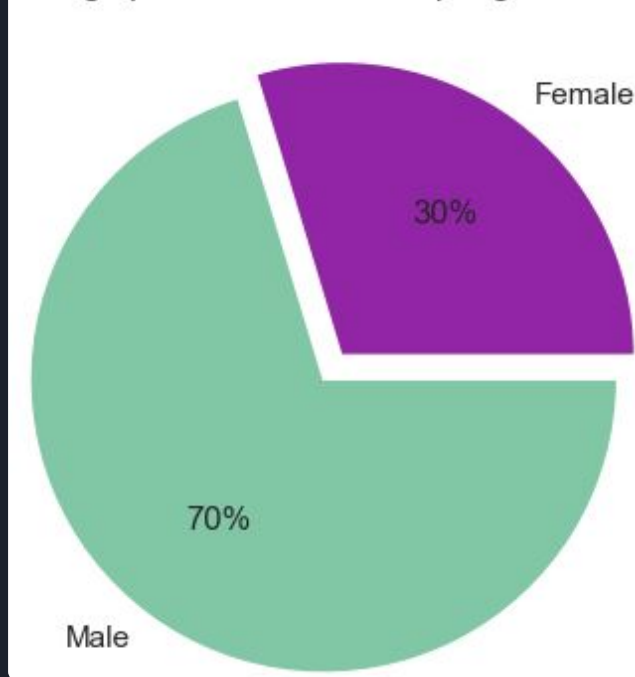- Race
- Sexuality

# Visualization

High potential leavers breakdown per gender
Female 58%
Male 42%

High potential Manager leavers
Not Manager 54%
Manager 46%

High potential leavers breakdown per Grade
Grade 10-15 52%
Grade 5-10 18%
Grade 20-25 5%
Grade 15-20 24%

High potential breakdown per gender
Female 30%
Male 70%

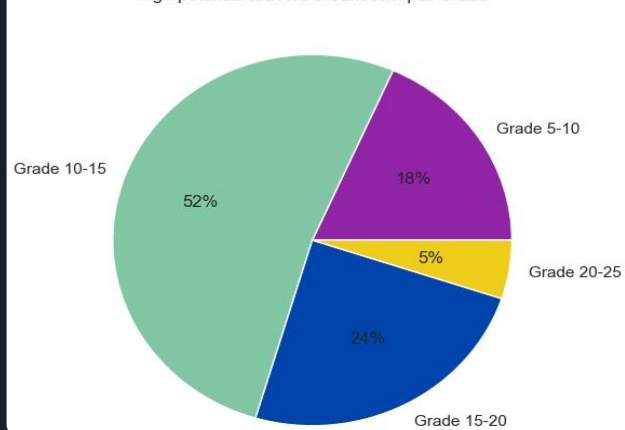⚠️A significant amount of people in management position are more likely to leave.

⚠️High potential females are more likely to leave than males

⚠️A significant amount of people in management position are more likely to leave.

⚠️Individuals within grade 10 to 15 tend to leave more frequently than the rest

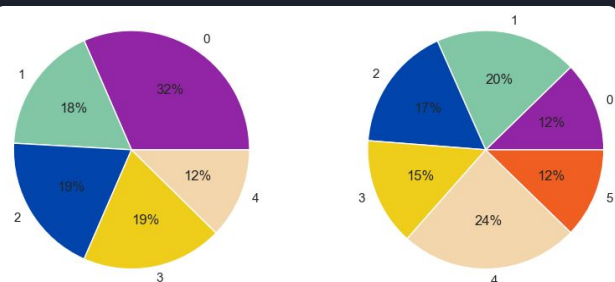⚠️Females are less prone to be classified as High Potential prospects.

# Visualization

Gender Comparison

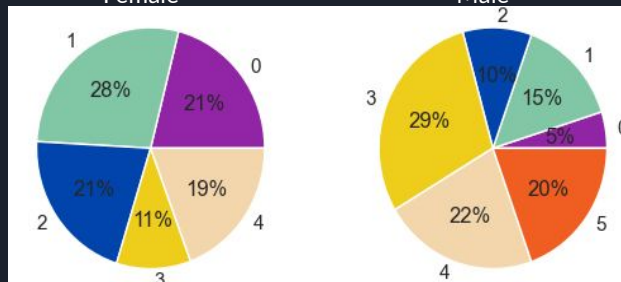High potential leavers breakdown per **Well-being**

Female / Male



High potential leavers breakdown per **Career Opportunities**
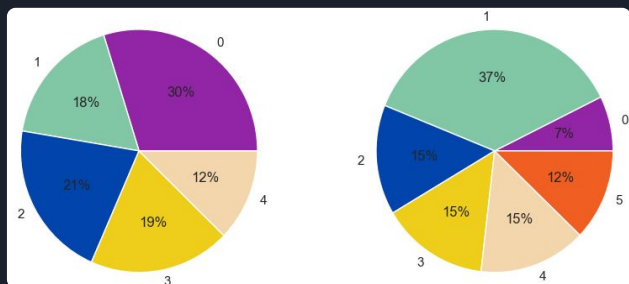
Female / Male



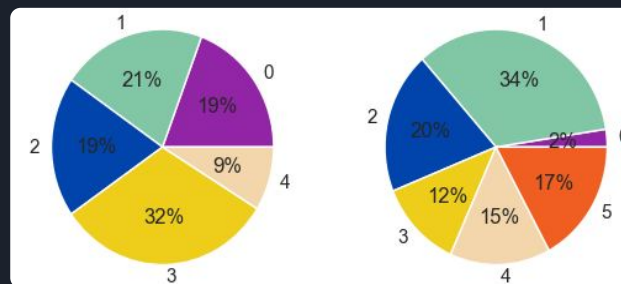High potential leavers breakdown per **Line Manager Rating**

Female / Male



High potential leavers breakdown per **Recognition**

Female / Male



⚠️ Females reported three times a value of Zero more than males

⚠️ Males are two times more likely to report a "Very Good" well being score

⚠️ Female have lower career opportunities compared to Males

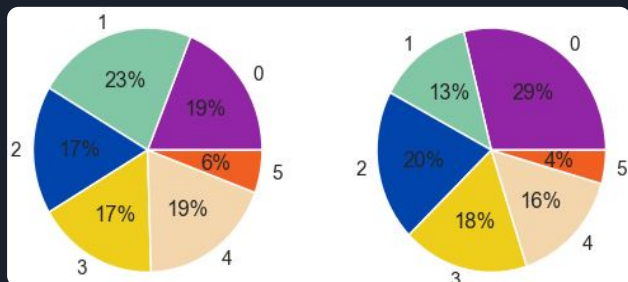⚠️ Overall Managers are not enabling personnel to be successful on their roles

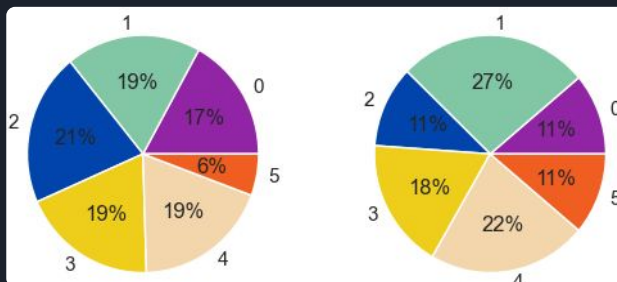⚠️ Males reported three times better recognition for their performance
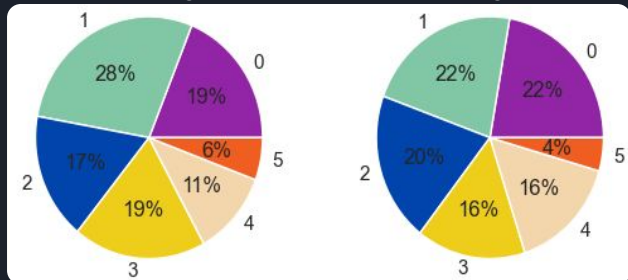
28

# Visualization

Manager Comparison

## High potential leavers breakdown per **Well-being**
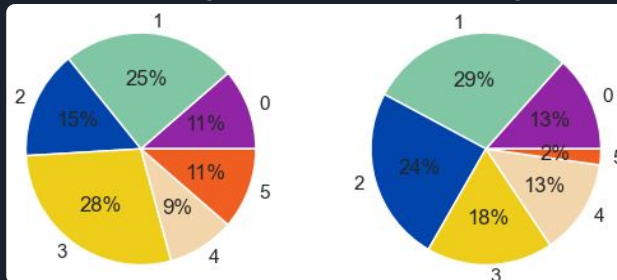
Not Manager — Manager



## High potential leavers breakdown per **Career Opportunities**

Not Manager — Manager



## High potential leavers breakdown per **Line Manager Rating**

Not Manager — Manager



## High potential leavers breakdown per **Recognition**

Not Manager — Manager



⚠️ Majority of manager reported less than poor "Well-being" scores

⚠️ Majority reported less than poor "Line Manager Ratings"

⚠️ Two thirds of overall High potential leavers reported less than poor "Career Opportunities" scores

⚠️ Two thirds of overall High potential leavers reported poor "Recognition" scores
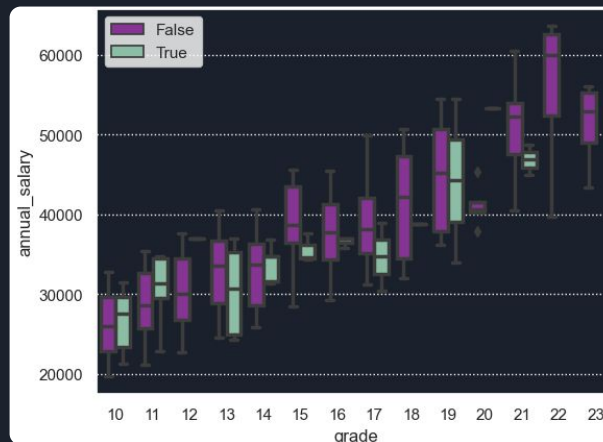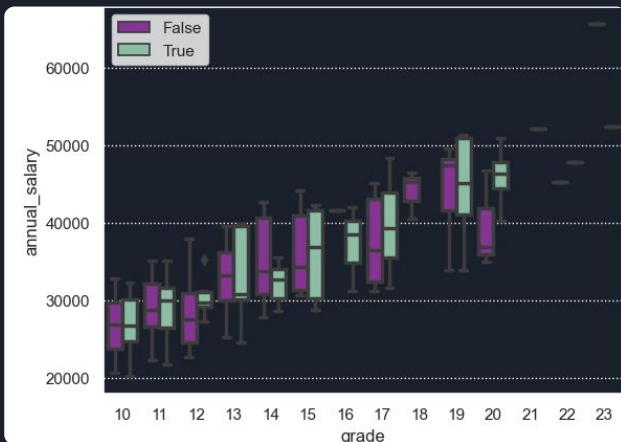
# Visualization

Comparison of grade and annual salary dimensioned by Leaver status of High Potential people
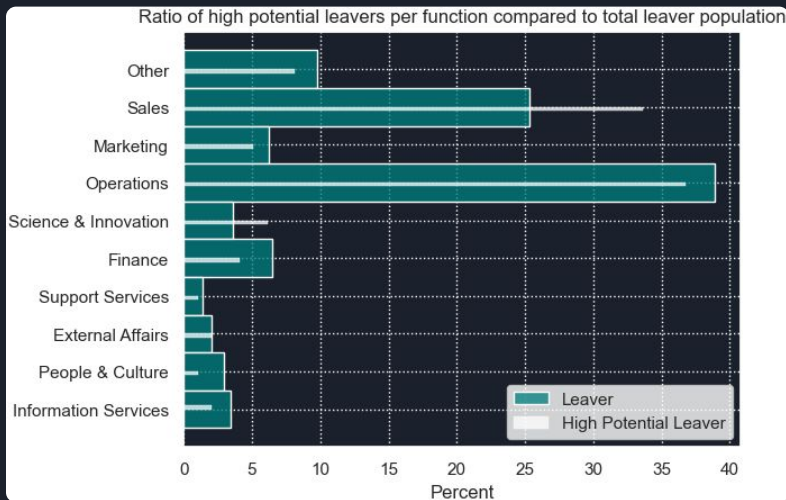


Female

Male

⚠️High potential people within grades 10, 13 and 19 are more likely to leave

⚠️Females within grade range from 10 to 19 tend to leave more often than males
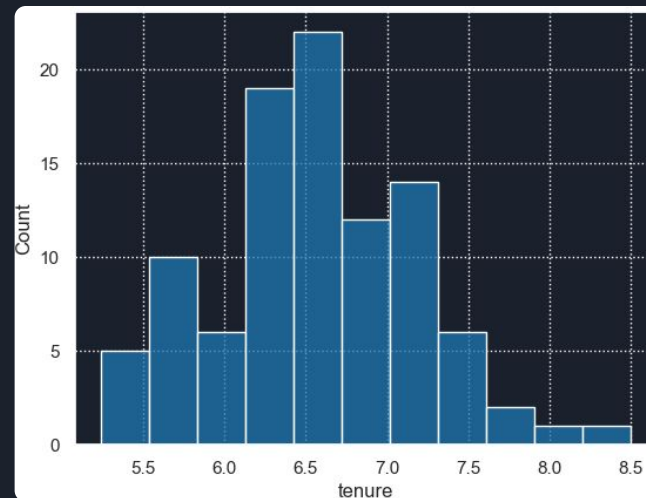
30

# Visualization

Ratio of high potential leavers per function compared to total leaver population

⚠️Sales and Science & Innovation functions report higher turnover rates compared to the overall turnover numbers

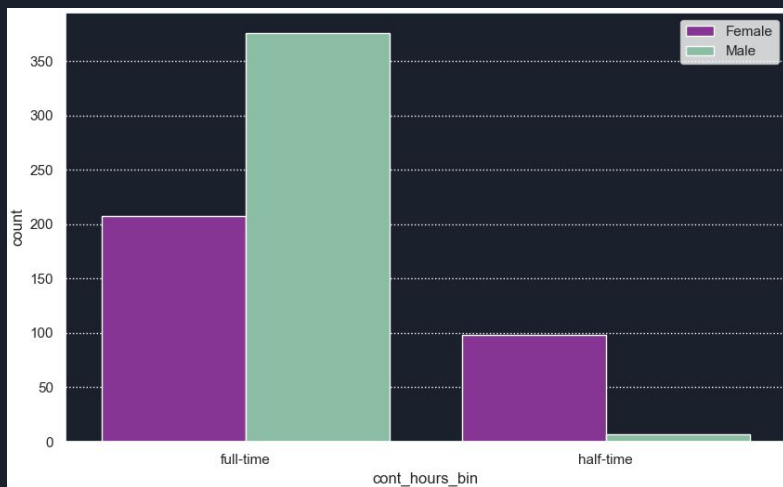Attrition of High Potential individuals by Tenure



⚠️People between 5.5 and 7.5 Years of Tenure are more likely to leave the company
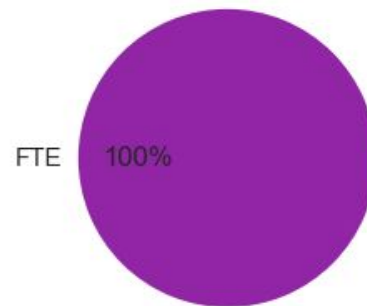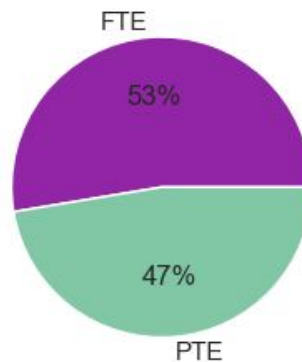
31

# Turnover Analysis

# Visualization

Leaver count of Contracted Hours Type per Gender



⚠️ Males tend to predominate the full-time employment contracts
⚠️ Females tend to leave the company a

High potential leavers by Contracted Hours breakdown per Gender

Female                                                         Male

# How can current top talent be retained

👍Since a significant part of people leaving the company (incl. managers) reside in the "mastered" years of tenure (between 5 and 10), a recommendation is to keep these resources **well tested** and look into **changing functions** to transfer their knowledge and skills exposing them to new challenges, teams and environments

👍Overall high potential people leaving the company reported poor career opportunities and recognition.  The company could offer **changing career paths** and **planning events more often** to give high talent the chance to be recognized by their peers

👍High potential managers tend to leave more often due to poor "Well-being". One action is to offer a range of **extracurricular activities and work-life balance plans.**

👍Offering **family-oriented advantages** could help mitigate high potential female employees. **Better parental leave plans, position retention after birth and daycare** are but a few ideas.

# Further Steps

Not all relationships were explored thoroughly, the following points would've been further analyzed:
- "Tenure" and "Time in latest role" correlations on turnover
- Overall data comparison (was heavily focused on high potential leavers)
- How hourly rate affects attrition instead of the used "Annual Salary"

Implemented a Random Forest Classifier (not shown in presentation)

**48%** of businesses say their **high-quality hires** come from **employee referrals**.[3]
(Linkedin)