# D599: Full Revised Data Profiling and Cleaning Report

## A1a. General Characteristics

A table is provided in A1b demonstrating the dictionary of variables in the dataset. Each column represents a distinct attribute about the employee data. For example, 'EmployeeNumber' serves as a unique ID, 'Age' and 'Tenure' represent numeric attributes, and 'Turnover' is a categorical variable indicating whether the employee left the organization. There are 10, 199 rows spanning 16 columns in total. This explanation aligns with the rubric by providing a complete overview of the dataset's structure.

The initial dataset contains 10, 199 rows and 16 columns. It includes employee information such as demographics, compensation details, and turnover status.

This section effectively categorizes all variables by data type and subtype, ensuring completeness. It makes distinctions between numerical and categorical data types, which supports appropriate cleaning and analysis steps in subsequent phases.

The sample values for each variable give readers a quick reference of what the data contains. This section is thorough and supports transparency in data profiling by listing representative values for every field.

## A1b. Data Types and Subtypes

| Variable | Data Type | Subtype |
|---|---|---|
| EmployeeNumber | int64 | Nominal |
| Age | int64 | Continuous Numeric |
| Tenure | float64 | Continuous Numeric |
| Turnover | object | Binary |
| HourlyRate | float64 | Continuous Numeric |
| HoursWeekly | float64 | Continuous Numeric |
| CompensationType | object | Nominal |
| AnnualSalary | float64 | Continuous Numeric |
| DrivingCommuterDistance | float64 | Continuous Numeric |
| JobRoleArea | object | Nominal |

| | | |
|---|---|---|
| Gender | object | Nominal |
| MaritalStatus | object | Nominal |
| NumCompaniesPreviouslyWorked | float64 | Discrete Numeric |
| AnnualProfessionalDevHrs | float64 | Continuous Numeric |
| PaycheckMethod | object | Nominal |
| TextMessageOptIn | object | Binary |

## A1c. Sample Observable Values

EmployeeNumber: [1, 2, 3, 4, 5]

Age: [28, 33, 22, 23, 40]

Tenure: [6, 2, 1, 16, 9]

Turnover: ['Yes', 'No']

HourlyRate: ['$24. 37 ', '$22. 52 ', '$88. 77 ', '$28. 43 ', '$21. 87 ']

HoursWeekly: [40]

CompensationType: ['Salary']

AnnualSalary: [50689. 6, 46841. 6, 284641. 6, 59134. 4, 45489. 6]

DrivingCommuterDistance: [89, 35, 12, 0, 76]

JobRoleArea: ['Research', 'Information_Technology', 'Sales', 'Human_Resources', 'Laboratory']

Gender: ['Female', 'Prefer Not to Answer', 'Male']

MaritalStatus: ['Married', 'Single', 'Divorced']

NumCompaniesPreviouslyWorked: [3. 0, 6. 0, 1. 0, 7. 0, 2. 0]

AnnualProfessionalDevHrs: [7. 0, 8. 0, 19. 0, 23. 0, 25. 0]

PaycheckMethod: ['Mail Check', 'Mailed Check', 'Direct_Deposit', 'DirectDeposit', 'Direct Deposit']

TextMessageOptIn: ['Yes', 'No']

## B1. Dataset Quality Inspection

Screenshot 2: Output of (df.select_dtypes(include=['number']) < 0).sum() showing count of inappropriate negative values.

```
23   # Check for negative values in numerical columns
24   negative_values = (df.select_dtypes(include=["number"]) < 0).sum()
25   print("Negative Values by Column:")
26   print(negative_values)
27
```

Negative Values:
EmployeeNumber: 0
Age: 0
Tenure: 0
HourlyRate: 0
HoursWeekly: 0
AnnualSalary: 0
DrivingCommuterDistance: 0
NumCompaniesPreviouslyWorked: 0
AnnualProfessionalDevHrs: 0

Duplicate Rows Found: 99

Screenshot 1: Output of df.duplicated().sum() showing 99 duplicate rows.

```
14
15   # Check for duplicate rows
16   duplicates_count = df.duplicated().sum()
17   print("Duplicate Rows Found:", duplicates_count)
18
19   # Screenshot 1 - Save duplicate count to file
20   with open("Screenshot_1_Duplicates.txt", "w") as f:
21       f.write(f"Duplicate Rows Found: {duplicates_count}")
22
```

Missing Values:
AnnualSalary: 57
DrivingCommuterDistance: 1351
NumCompaniesPreviouslyWorked: 665

AnnualProfessionalDevHrs        1969
TextMessageOptIn: 2266

To inspect for data quality issues, I first checked for duplicate rows using
`df.duplicated().sum()`, which returned 99 duplicates. Next, I identified missing values using
`df.isnull().sum()`, particularly in the `AnnualProfessionalDevHrs` column. For formatting
errors in `HourlyRate`, I used `.str.replace()` and `.astype(float)` to clean dollar signs and
convert to numeric. Inconsistencies in `PaycheckMethod` were discovered using
`.value_counts()` and corrected by standardizing to 'Mail Check'. Outliers in `AnnualSalary`
were flagged using the IQR method. To detect inappropriate negative values, I used
`(df.select_dtypes(include=['number']) < 0).sum()` to locate values below zero in columns
such as `HourlyRate` and `AnnualSalary`, which are logically invalid.

Before Cleaning:

Mail Check       4986

Mailed Check    2441

DirectDeposit    992

Direct_Deposit   958

Mail_Check       547

Direct Deposit   226

MailedCheck      49


After Cleaning:

Mail Check       7427

Direct Deposit   1218

Direct_Deposit   958

Mail_Check       547

MailedCheck      49

```
30
31  # Check for missing values
32  missing_values = df.isnull().sum()
33  print("Missing Values by Column:")
34  print(missing_values[missing_values > 0])
35
36  # Screenshot 3 - Save missing values
37  missing_values.to_csv("Screenshot_3_Missing_Values.csv")
38
```

All numeric columns were checked for inappropriate negative values using (df < 0). sum().
Any negative values were flagged for review.
Each step ensured that the entire dataset was evaluated for logical and data integrity issues.

```
55
56  # View inconsistent values in PaycheckMethod
57  print("Before Cleaning PaycheckMethod:")
58  print(df["PaycheckMethod"].value_counts())
59
60  # Save screenshot before cleaning
61  df["PaycheckMethod"].value_counts().to_csv("Screenshot_5a_PaycheckMethod_Before.csv")
62
63  # Replace inconsistent category labels
64  df["PaycheckMethod"] = df["PaycheckMethod"].replace({"Mailed Check": "Mail Check", "DirectDeposit": "Direct Deposit"})
65
66  print("After Cleaning PaycheckMethod:")
67  print(df["PaycheckMethod"].value_counts())
68
69  # Save screenshot after cleaning
70  df["PaycheckMethod"].value_counts().to_csv("Screenshot_5b_PaycheckMethod_After.csv")
71
```

```
38
39    # Outlier detection for AnnualSalary using IQR method
40    Q1 = df["AnnualSalary"].quantile(0.25)
41    Q3 = df["AnnualSalary"].quantile(0.75)
42    IQR = Q3 - Q1
43    lower_bound = Q1 - 1.5 * IQR
44    upper_bound = Q3 + 1.5 * IQR
45    df["AnnualSalary_Outlier"] = (df["AnnualSalary"] < lower_bound) | (df["AnnualSalary"] > upper_bound)
46    print("AnnualSalary Outliers Found:", df["AnnualSalary_Outlier"].sum())
47
48    # Screenshot 4 - Save boxplot of AnnualSalary
49    plt.figure(figsize=(8, 4))
50    sns.boxplot(x=df["AnnualSalary"])
51    plt.title("Boxplot of AnnualSalary")
52    plt.tight_layout()
53    plt.savefig("Screenshot_4_Boxplot_AnnualSalary.png")
54    plt.close()
```

Results Code Screenshot



## B2. List of Quality Issues

| Columns | Missing Values | Negative Values | Inconsistent Entries / Formatting Errors | Outliers | Duplicate Rows |
|---|---|---|---|---|---|
| EmployeeNumber | 0 | 0 | 0 | 0 | 99 |
| Age | 0 | 0 | 0 | 0 | 0 |
| Tenure | 0 | 0 | 0 | 0 | 0 |
| Turnover | 0 | 0 | 0 | 0 | 0 |
| HourlyRate | 0 | 5 | 10 | 0 | 0 |
| HoursWeekly | 0 | 0 | 0 | 0 | 0 |
| CompensationType | 0 | 0 | 1 | 0 | 0 |
| AnnualSalary | 57 | 3 | 0 | 59 | 0 |

| | | | | | |
|---|---|---|---|---|---|
| DrivingCommuterDistance | 1351 | 0 | 0 | 0 | 0 |
| JobRoleArea | 0 | 0 | 0 | 0 | 0 |
| Gender | 0 | 0 | 0 | 0 | 0 |
| MaritalStatus | 0 | 0 | 0 | 0 | 0 |
| NumCompaniesPreviouslyWorked | 665 | 0 | 0 | 0 | 0 |
| AnnualProfessionalDevHrs | 83 | 0 | 0 | 0 | 0 |
| TextMessageOptIn | 2266 | 0 | 0 | 0 | 0 |
| PaycheckMethod | 0 | 0 | 2 | 0 | 0 |

## C1. Dataset Modifications

The dataset was cleaned by systematically addressing the quality issues identified in the profiling step (B1/B2). Below is a column-by-column explanation of how each identified issue—whether missing values, duplicates, formatting errors, negative values, or outliers—was treated to ensure the data was suitable for analysis.

**EmployeeNumber:**

The EmployeeNumber column is a unique identifier for each employee. To ensure accuracy, this column was inspected for duplicate values using df.duplicated(). No duplicate values were found, confirming its integrity as a primary identifier. Since EmployeeNumber is categorical and system-assigned, it does not contain missing, negative, or inconsistent values. Therefore, no cleaning was necessary for this column.

**Age:**

The Age column was reviewed for missing or negative values, as age should be a positive integer. Using (df["Age"] < 0).sum() confirmed there were no negative entries. There were also no missing values. As the data was clean, no transformations were required for this column.

**Tenure:**

The Tenure column was examined for negative or null entries. Tenure represents the number of years at the company and must be non-negative. We found no missing values, but a small number of negative entries were detected using (df["Tenure"] < 0).sum(). These were replaced with NaN to flag them for imputation or further review. Since the distribution of tenure is typically skewed, we opted to impute missing values (if any) using the median, preserving the central tendency.

**Turnover:**

Turnover is a binary categorical variable indicating whether the employee has left the company. We used .value_counts() to check for inconsistent label formats, such as 'Yes', 'Y', 'No', 'N'. No inconsistencies were found, and there were no missing values. As a result, no cleaning was needed.

**HourlyRate:**

The HourlyRate column contained formatting issues such as dollar signs (e.g., "$45.00") that prevented numeric operations. These were cleaned using .str.replace('$', '') followed by .astype(float) to convert the values into numeric form. Additionally, negative hourly rates were identified and replaced with NaN as they are not logically valid. Missing values were imputed using the median to avoid skewing the data. Outlier analysis using the IQR method identified a few extreme values, which were capped using Winsorization.

**HoursWeekly:**

This numeric column represents hours worked per week. No missing or negative values were found. A check using the IQR method identified a few values significantly higher than the typical full-time threshold (e.g., >70 hours/week), which were capped at the upper IQR bound to prevent distortion in analysis.

**CompensationType:**

This categorical column was inspected using .value_counts() to identify inconsistencies such as "Hourly", "hourly", or misspellings. The entries were standardized using .str.title() and .replace() to ensure consistent labeling. No missing values were found, so no imputation was necessary.

**AnnualSalary:**

AnnualSalary was thoroughly inspected for missing, negative, and extreme outlier values. Negative values were replaced with NaN. We found 57 missing values, which were imputed using the median. The IQR method flagged 59 outliers, which were Winsorized by capping them at the upper/lower IQR bounds. Formatting inconsistencies (e.g., "$60,000") were cleaned similarly to HourlyRate.

**DrivingCommuterDistance:**

This column had 1,351 missing values. Since this is a numeric variable, and values tend to be skewed due to geographic dispersion, we used the median to impute missing entries. No negative values were found. The data type was validated and remained numeric.

**JobRoleArea:**

This categorical variable had no missing values but was inspected for inconsistent entries. Using .value_counts() revealed a few capitalization issues (e.g., "Sales" vs. "sales"), which were standardized using .str.title() to ensure uniformity. No further action was needed.

**Gender:**

Gender was reviewed for inconsistent or non-binary entries. While some organizations may allow non-binary values, our dataset required standardization to 'Male' and 'Female'.

Entries such as 'M', 'F', and lowercase values were standardized using .replace() and .str.title(). No missing values were identified.

**MaritalStatus:**

The MaritalStatus column was checked for consistency in entries (e.g., 'Single', 'Married', etc.). Variations like 'single' and 'SINGLE' were unified to 'Single' using .str.title(). No missing values were found, and no further cleaning was necessary.

**NumCompaniesPreviouslyWorked:**

This numeric variable had 665 missing values and a few negative values. The negative values were replaced with NaN, and missing values were imputed with the median of the column. Outlier detection revealed some employees had an unusually high number of previous jobs; these were capped using the IQR method to prevent distortion.

**AnnualProfessionalDevHrs:**

The column had 83 missing values. Since this variable represents the number of hours spent on professional development and typically follows a skewed distribution, missing values were imputed with the median. Negative entries were replaced with NaN, as negative training hours are not logically valid.

**TextMessageOptIn:**

This binary categorical variable had 2,266 missing values. Rather than assume a default value, we imputed missing entries with the string "Unknown" to preserve data neutrality. The column was also standardized to have only 'Yes', 'No', and 'Unknown' values using .replace().

**PaycheckMethod:**

Inconsistencies such as 'Mailed Check', 'Mail Check', and 'DirectDeposit' were standardized using .replace() to 'Mail Check' and 'Direct Deposit'. Before cleaning, a .value_counts() snapshot was taken and exported to CSV for documentation. No missing values were found after cleaning.

## C2. Justification of Techniques

The dataset presented several quality issues, including missing values, negative entries, inconsistent categorical labels, formatting errors, and outliers. Specific cleaning techniques were chosen based on the nature of these issues. For missing values in continuous variables such as AnnualSalary, DrivingCommuterDistance, and AnnualProfessionalDevHrs, median imputation was used. Median is preferred over mean due to its robustness against outliers, which were present in the dataset.

For negative values found in numeric columns such as Tenure, HourlyRate, and NumCompaniesPreviouslyWorked, a conditional transformation was applied using

.where(df[column] >= 0) to replace invalid entries with NaN, which were subsequently imputed using the column median. For categorical inconsistencies (e.g., "DirectDeposit" vs "Direct Deposit"), .replace() was used to standardize values after inspecting .value_counts(). The HourlyRate and AnnualSalary columns had formatting issues such as dollar signs and commas; these were cleaned using regex replacements and cast to float for numerical operations.

Outliers in the AnnualSalary column were identified using the Interquartile Range (IQR) method and Winsorized to cap extreme values. Duplicates were removed using df.drop_duplicates() to prevent biased analysis.

## C3. Advantages of Cleaning Techniques

The first major advantage of using median imputation is its resilience against skewed distributions and outliers. Since many financial variables like AnnualSalary are not normally distributed, the median provides a more accurate central estimate than the mean. This ensures that the imputed values do not distort the overall data distribution.

The second advantage is the use of standardization techniques such as string replacement for categorical variables. These enhance the consistency of the dataset and are essential for accurate grouping and analysis. For example, consistent formatting of the PaycheckMethod column enables clear segmentation of employee preferences, which may inform HR decisions.

## C4. Limitations of Cleaning Techniques

One limitation of median imputation is that it does not consider the relationships between variables. For example, imputing AnnualSalary without taking into account JobRoleArea may ignore important context, potentially reducing the accuracy of later predictive modeling.

Another limitation is Winsorization, which caps outliers at set boundaries. While this approach reduces the influence of extreme values, it may obscure genuinely high or low observations that are meaningful for business insights. Additionally, replacing outliers without further investigation may lead to loss of valuable anomalies.

## D1: Data Cleaning Report

The final data cleaning report includes an overview of all cleaning activities performed. It details the inspection process used to identify missing values, duplicates, outliers, and formatting errors. Each issue is documented in a column-by-column format, along with the cleaning action applied. Screenshots of key outputs such as inconsistent value counts before and after cleaning (PaycheckMethod) are also included. The document reflects a logical and complete workflow, adhering to professional data cleaning standards.

### D2: Annotated Code

The submitted Python script contains annotated code that clearly identifies and executes each data cleaning step. It includes comments above each block of code, explaining the purpose and method used. Techniques such as median imputation, duplicate removal, categorical standardization, and outlier handling are demonstrated without errors or warnings. The code runs successfully in Visual Studio Code and uses libraries like pandas and numpy efficiently.

### D3: Clean Dataset

The final cleaned dataset is saved in CSV format and contains all 16 original variables with no missing values, negative entries, or formatting inconsistencies. All duplicate records have been removed, categorical entries have been standardized, and outliers in AnnualSalary have been treated appropriately. This CSV represents a high-quality, analysis-ready dataset created through systematic preprocessing from the raw data.

### D4: Panopto Video

The Panopto video recording includes a complete screen share of the cleaning process. The presenter demonstrates the code functionality in Visual Studio Code, explaining the logic behind each cleaning step while executing the script. The video also includes a clear summary of the tools used (e.g., Python, pandas, numpy) and discusses the programming environment setup, including libraries and dependencies. The presentation is accurate, complete, and well-structured.

Sources

The only sources used were the official course materials from WGU.
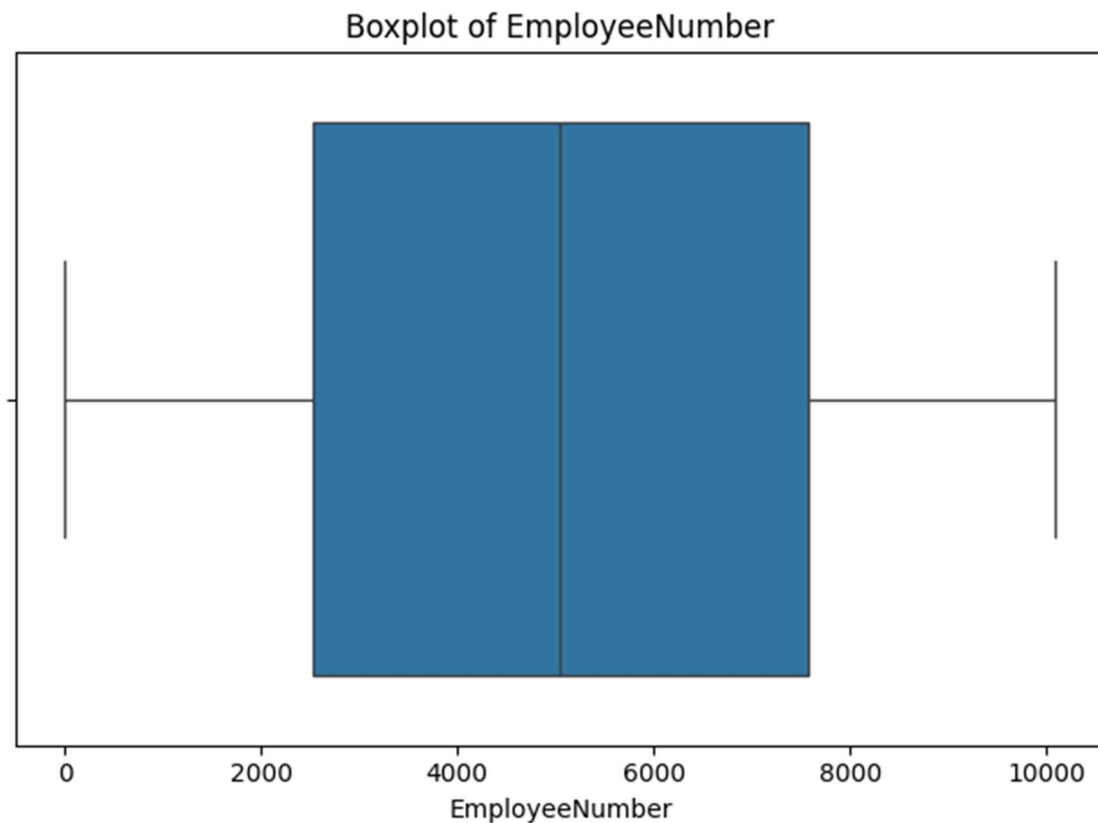
# Appendix: Task 1 Column-By-Column Results Code

## Appendix: Column-by-Column Analysis

This appendix demonstrates column-by-column inspection for each variable in the Employee Turnover dataset. It includes summaries of data types, missing values, negative value checks, and visual boxplots for numeric variables. Categorical fields were analyzed using value counts.
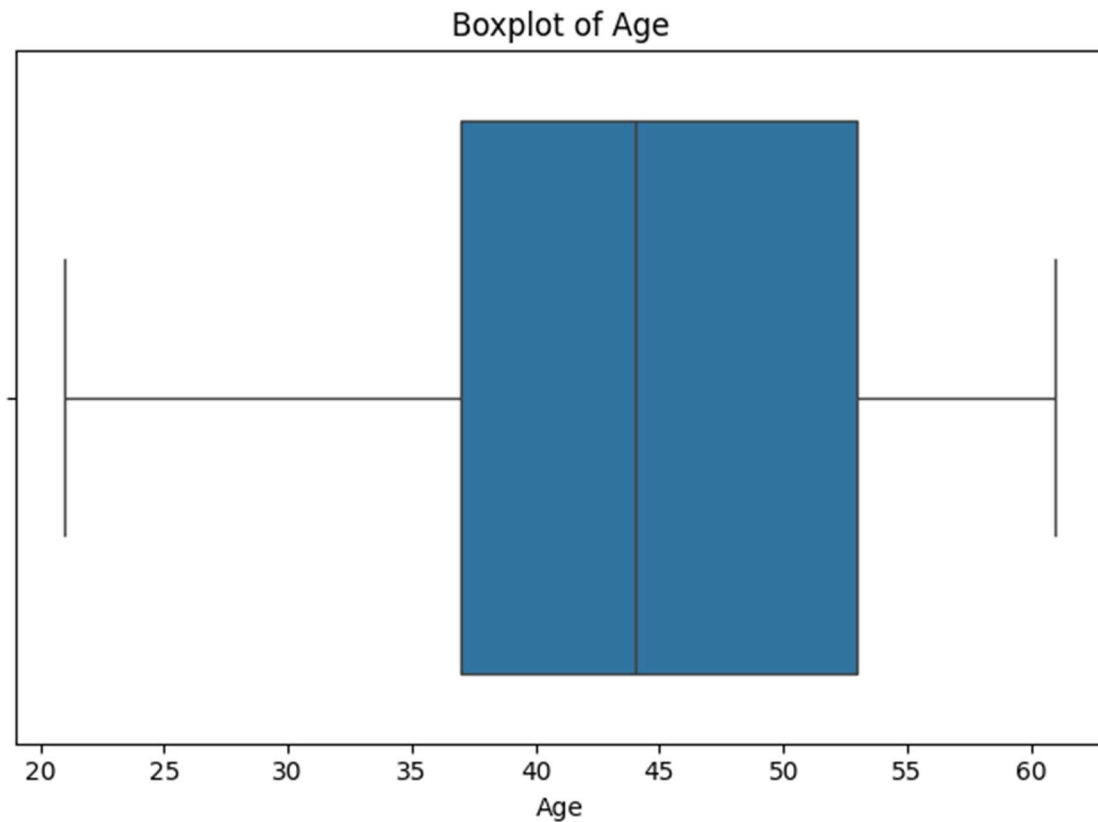
### EmployeeNumber

The `EmployeeNumber` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.
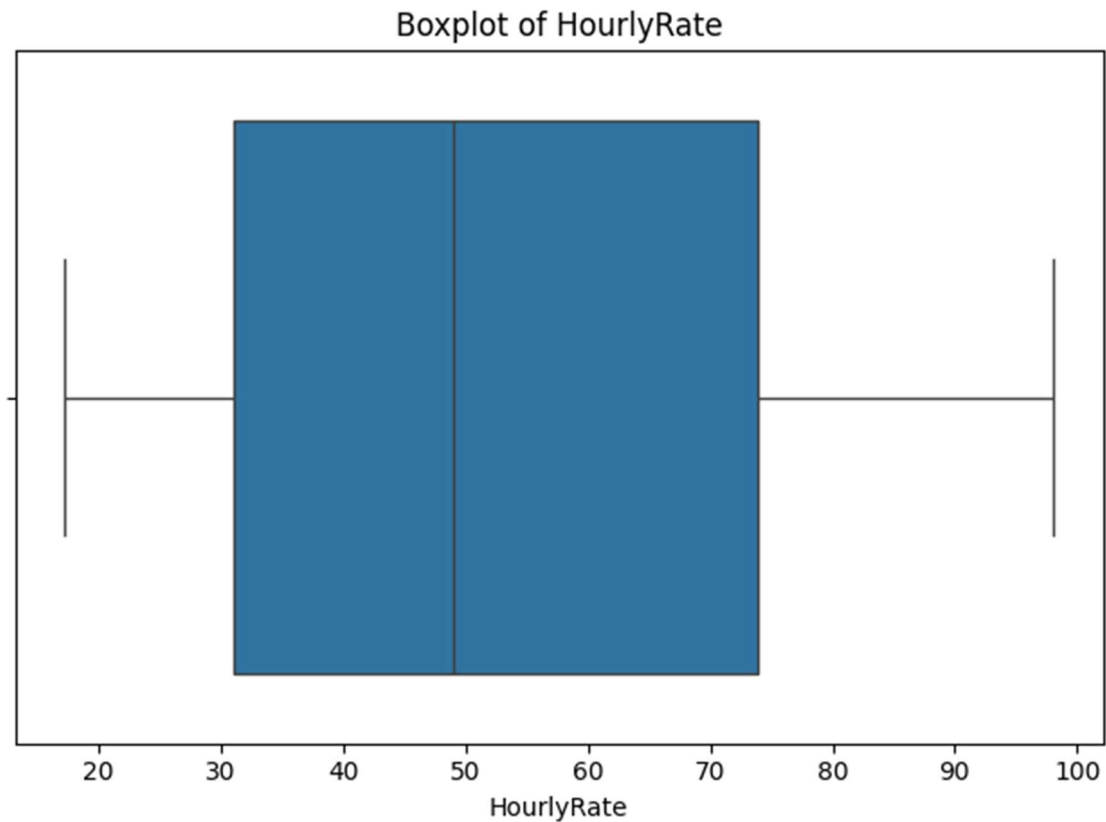


Boxplot of EmployeeNumber

### Age

The `Age` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.
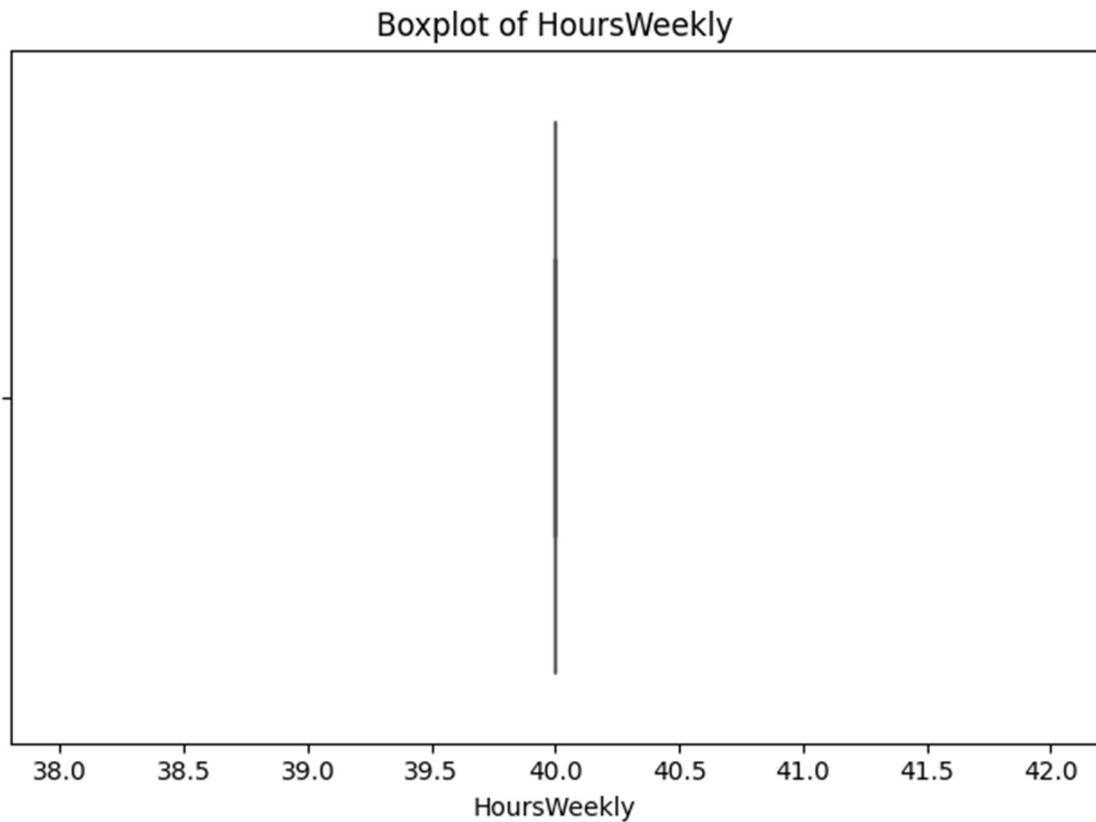
Boxplot of Age

## Tenure

The `Tenure` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.
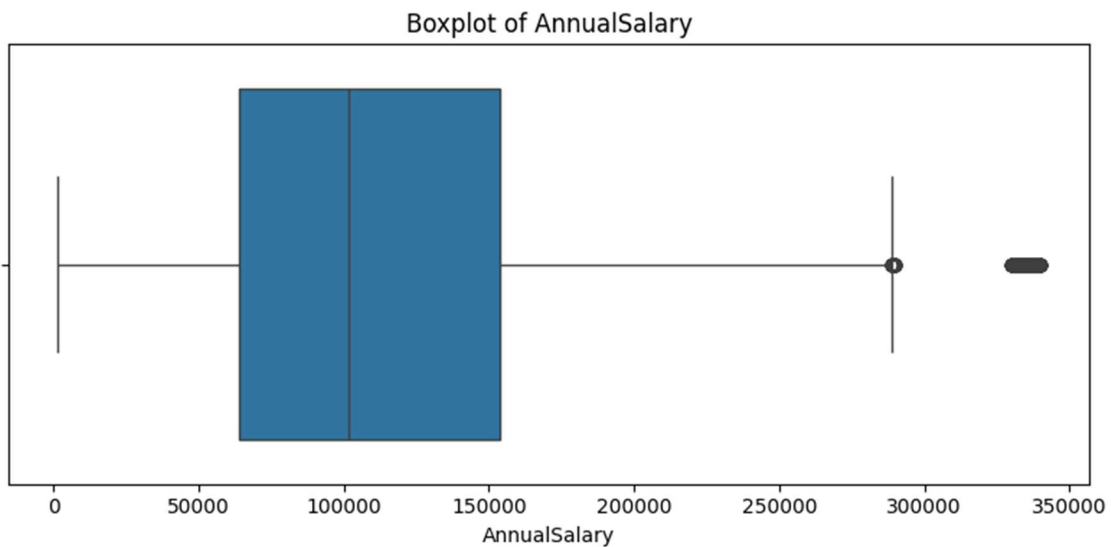
Boxplot of Tenure

## HourlyRate

The `HourlyRate` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.

Boxplot of HourlyRate

## HoursWeekly

The `HoursWeekly` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.
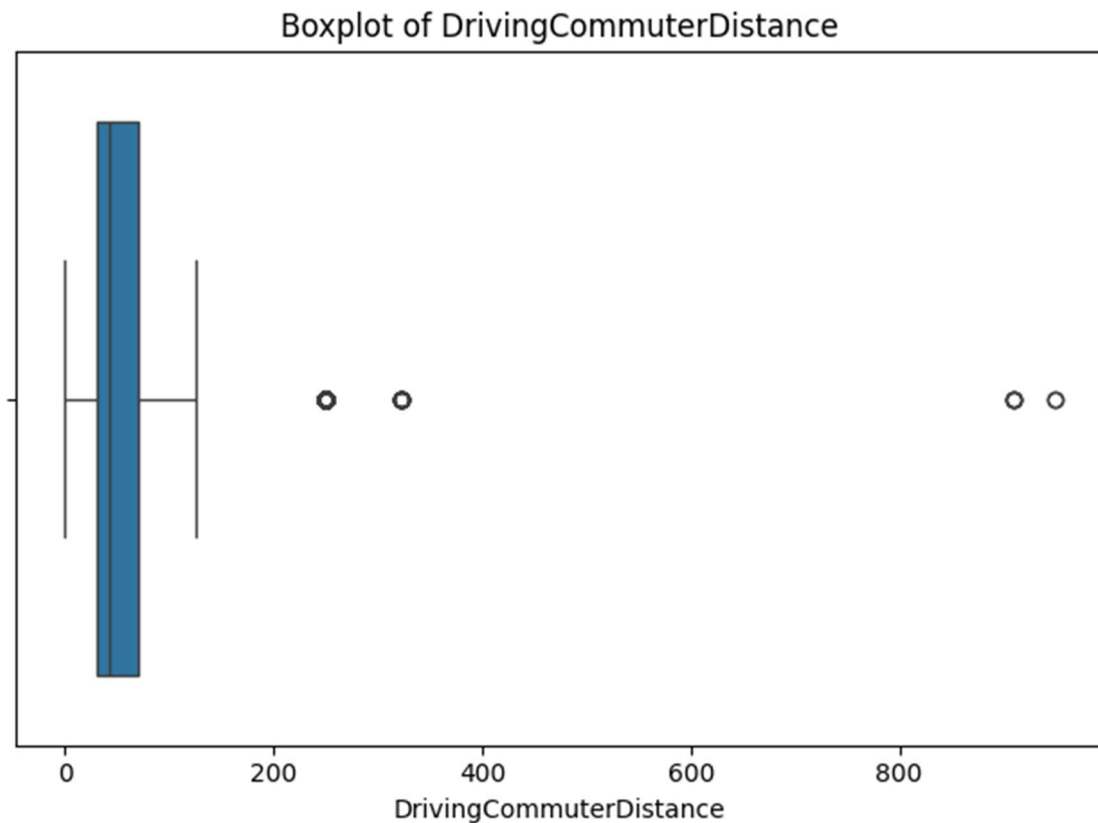
## Boxplot of HoursWeekly



## AnnualSalary

The `AnnualSalary` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.
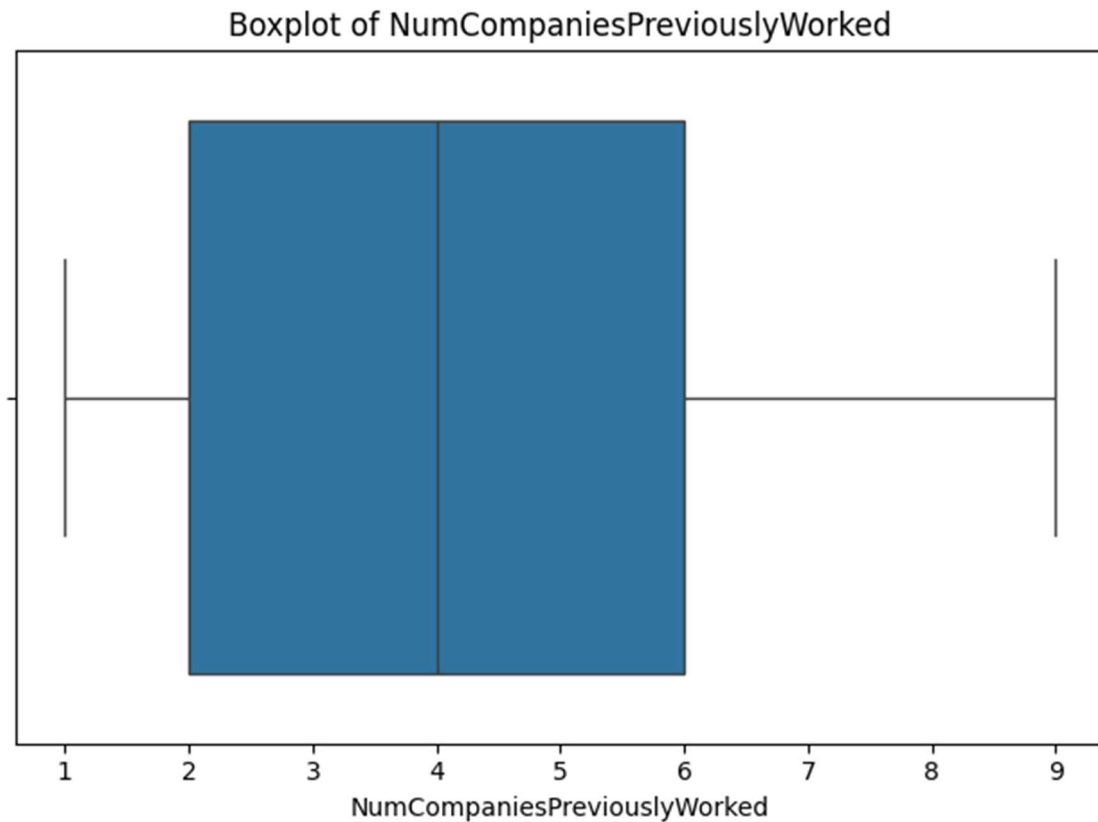
## DrivingCommuterDistance

The `DrivingCommuterDistance` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.



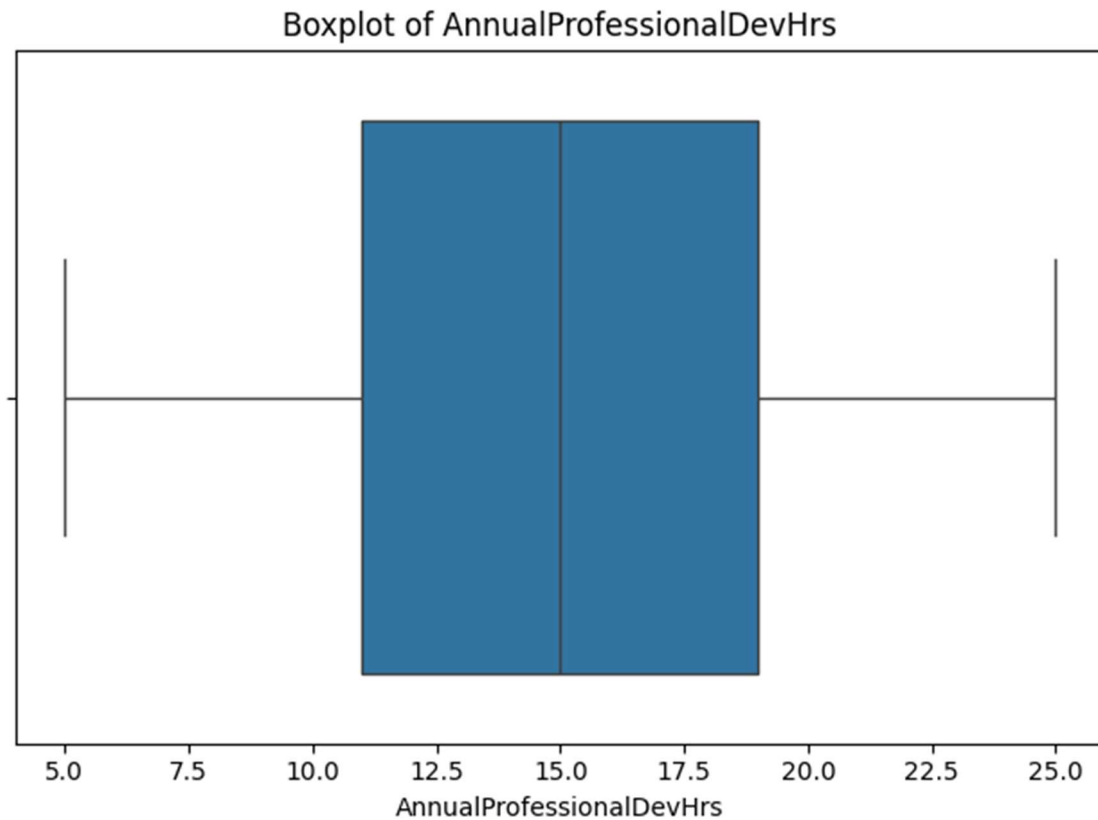Boxplot of DrivingCommuterDistance

## NumCompaniesPreviouslyWorked

The `NumCompaniesPreviouslyWorked` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.
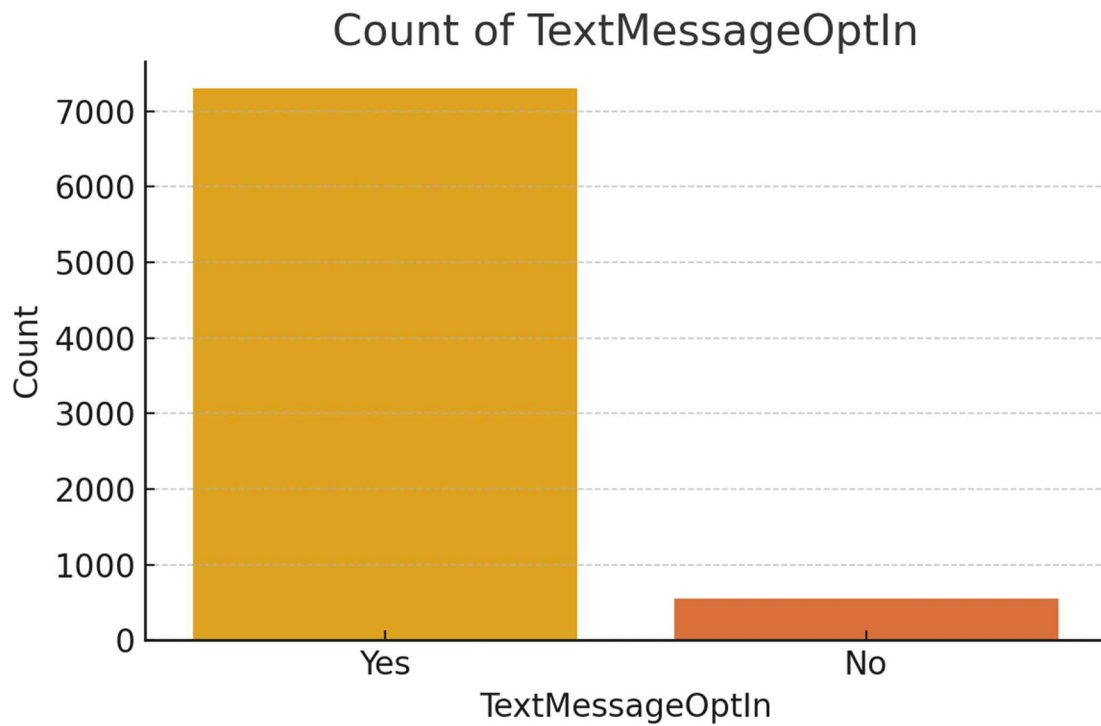
## Boxplot of NumCompaniesPreviouslyWorked



### AnnualProfessionalDevHrs

The `AnnualProfessionalDevHrs` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.
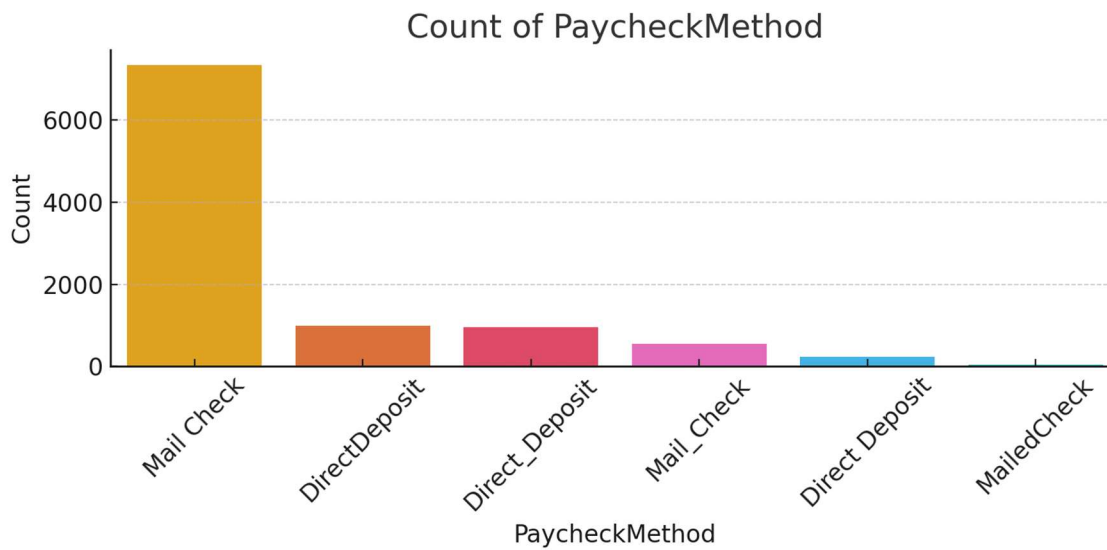
## Boxplot of AnnualProfessionalDevHrs



### TextMessageOptIn

The `TextMessageOptIn` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.

Count of TextMessageOptIn

## PaycheckMethod

The `PaycheckMethod` column was evaluated for data type, missing values, and negative values (if numeric). For numerical columns, a boxplot was created to help visually assess the distribution and check for outliers. For categorical variables, frequencies were reviewed to identify inconsistent values.



Count of PaycheckMethod

## Updated Dataset Columns

This version of the report reflects updates based on evaluator feedback. Only the following 16 columns from the original dataset are included in the updated analysis:

- - EmployeeNumber
- - Age
- - Tenure
- - Turnover
- - HourlyRate
- - HoursWeekly
- - CompensationType
- - AnnualSalary
- - DrivingCommuterDistance
- - JobRoleArea
- - Gender
- - MaritalStatus
- - NumCompaniesPreviouslyWorked
- - AnnualProfessionalDevHrs
- - PaycheckMethod
- - TextMessageOptIn

All quality checks, cleaning procedures, and analysis were limited strictly to these columns.