# Data Cleaning.

```
In [4]: import pandas as pd
        import matplotlib.pyplot as plt
        import seaborn as sns
```

```
In [5]: # Read the dataset
        df = pd.read_csv(r'F:\Technocolabs\WA_Fn-UseC_-HR-Employee-Attrition.csv')
```

```
In [6]: df
```

Out[6]:

|  | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education |
|---|---|---|---|---|---|---|---|
| 0 | 41 | Yes | Travel_Rarely | 1102 | Sales | 1 | 2 |
| 1 | 49 | No | Travel_Frequently | 279 | Research & Development | 8 | 1 |
| 2 | 37 | Yes | Travel_Rarely | 1373 | Research & Development | 2 | 2 |
| 3 | 33 | No | Travel_Frequently | 1392 | Research & Development | 3 | 4 |
| 4 | 27 | No | Travel_Rarely | 591 | Research & Development | 2 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 1465 | 36 | No | Travel_Frequently | 884 | Research & Development | 23 | 2 |
| 1466 | 39 | No | Travel_Rarely | 613 | Research & Development | 6 | 1 |
| 1467 | 27 | No | Travel_Rarely | 155 | Research & Development | 4 | 3 |
| 1468 | 49 | No | Travel_Frequently | 1023 | Sales | 2 | 3 |
| 1469 | 34 | No | Travel_Rarely | 628 | Research & Development | 8 | 3 |

1470 rows × 35 columns

```
In [7]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 35 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   Age                       1470 non-null   int64
 1   Attrition                 1470 non-null   object
 2   BusinessTravel            1470 non-null   object
 3   DailyRate                 1470 non-null   int64
 4   Department                1470 non-null   object
 5   DistanceFromHome          1470 non-null   int64
 6   Education                 1470 non-null   int64
 7   EducationField            1470 non-null   object
 8   EmployeeCount             1470 non-null   int64
 9   EmployeeNumber            1470 non-null   int64
 10  EnvironmentSatisfaction   1470 non-null   int64
 11  Gender                    1470 non-null   object
 12  HourlyRate                1470 non-null   int64
 13  JobInvolvement            1470 non-null   int64
 14  JobLevel                  1470 non-null   int64
 15  JobRole                   1470 non-null   object
 16  JobSatisfaction           1470 non-null   int64
 17  MaritalStatus             1470 non-null   object
 18  MonthlyIncome             1470 non-null   int64
 19  MonthlyRate               1470 non-null   int64
 20  NumCompaniesWorked        1470 non-null   int64
 21  Over18                    1470 non-null   object
 22  OverTime                  1470 non-null   object
 23  PercentSalaryHike         1470 non-null   int64
 24  PerformanceRating         1470 non-null   int64
 25  RelationshipSatisfaction  1470 non-null   int64
 26  StandardHours             1470 non-null   int64
 27  StockOptionLevel          1470 non-null   int64
 28  TotalWorkingYears         1470 non-null   int64
 29  TrainingTimesLastYear     1470 non-null   int64
 30  WorkLifeBalance           1470 non-null   int64
 31  YearsAtCompany            1470 non-null   int64
 32  YearsInCurrentRole        1470 non-null   int64
 33  YearsSinceLastPromotion   1470 non-null   int64
 34  YearsWithCurrManager      1470 non-null   int64
dtypes: int64(26), object(9)
memory usage: 402.1+ KB
```

```
In [8]: print(df.tail())
```

```
      Age Attrition      BusinessTravel  DailyRate              Department \
1465   36        No  Travel_Frequently        884  Research & Development
1466   39        No      Travel_Rarely        613  Research & Development
1467   27        No      Travel_Rarely        155  Research & Development
1468   49        No  Travel_Frequently       1023                   Sales
1469   34        No      Travel_Rarely        628  Research & Development

      DistanceFromHome  Education EducationField  EmployeeCount \
1465                23          2        Medical              1
1466                 6          1        Medical              1
1467                 4          3  Life Sciences              1
1468                 2          3        Medical              1
1469                 8          3        Medical              1

      EmployeeNumber  ...  RelationshipSatisfaction StandardHours \
1465            2061  ...                         3            80
1466            2062  ...                         1            80
1467            2064  ...                         2            80
1468            2065  ...                         4            80
1469            2068  ...                         1            80

      StockOptionLevel  TotalWorkingYears  TrainingTimesLastYear \
1465                 1                 17                      3
1466                 1                  9                      5
1467                 1                  6                      0
1468                 0                 17                      3
1469                 0                  6                      3

      WorkLifeBalance  YearsAtCompany YearsInCurrentRole \
1465                3               5                  2
1466                3               7                  7
1467                3               6                  2
1468                2               9                  6
1469                4               4                  3

      YearsSinceLastPromotion  YearsWithCurrManager
1465                        0                     3
1466                        1                     7
1467                        0                     3
1468                        0                     8
1469                        1                     2

[5 rows x 35 columns]
```

```python
In [9]: print(df.tail(10))  # Display the last 10 rows
```

```
      Age Attrition      BusinessTravel  DailyRate             Department
\
1460   29        No       Travel_Rarely        468  Research & Development
1461   50       Yes       Travel_Rarely        410                   Sales
1462   39        No       Travel_Rarely        722                   Sales
1463   31        No          Non-Travel        325  Research & Development
1464   26        No       Travel_Rarely       1167                   Sales
1465   36        No   Travel_Frequently        884  Research & Development
1466   39        No       Travel_Rarely        613  Research & Development
1467   27        No       Travel_Rarely        155  Research & Development
1468   49        No   Travel_Frequently       1023                   Sales
1469   34        No       Travel_Rarely        628  Research & Development

      DistanceFromHome  Education EducationField  EmployeeCount  \
1460                28          4        Medical              1
1461                28          3      Marketing              1
1462                24          1      Marketing              1
1463                 5          3        Medical              1
1464                 5          3          Other              1
1465                23          2        Medical              1
1466                 6          1        Medical              1
1467                 4          3  Life Sciences              1
1468                 2          3        Medical              1
1469                 8          3        Medical              1

      EmployeeNumber  ...  RelationshipSatisfaction StandardHours  \
1460            2054  ...                         2            80
1461            2055  ...                         2            80
1462            2056  ...                         1            80
1463            2057  ...                         2            80
1464            2060  ...                         4            80
1465            2061  ...                         3            80
1466            2062  ...                         1            80
1467            2064  ...                         2            80
1468            2065  ...                         4            80
1469            2068  ...                         1            80

      StockOptionLevel  TotalWorkingYears  TrainingTimesLastYear  \
1460                 0                  5                      3
1461                 1                 20                      3
1462                 1                 21                      2
1463                 0                 10                      2
1464                 0                  5                      2
1465                 1                 17                      3
1466                 1                  9                      5
1467                 1                  6                      0
1468                 0                 17                      3
1469                 0                  6                      3

      WorkLifeBalance  YearsAtCompany YearsInCurrentRole  \
1460                1               5                  4
1461                3               3                  2
1462                2              20                  9
1463                3               9                  4
1464                3               4                  2
1465                3               5                  2
1466                3               7                  7
1467                3               6                  2
1468                2               9                  6
1469                4               4                  3
```

|      | YearsSinceLastPromotion | YearsWithCurrManager |
|------|-------------------------|----------------------|
| 1460 | 0 | 4 |
| 1461 | 2 | 0 |
| 1462 | 9 | 6 |
| 1463 | 1 | 7 |
| 1464 | 0 | 0 |
| 1465 | 0 | 3 |
| 1466 | 1 | 7 |
| 1467 | 0 | 3 |
| 1468 | 0 | 8 |
| 1469 | 1 | 2 |

[10 rows x 35 columns]

```
In [10]: print(df.head(10))
```

```
   Age Attrition    BusinessTravel  DailyRate              Department  \
0   41       Yes     Travel_Rarely       1102                   Sales
1   49        No  Travel_Frequently      279  Research & Development
2   37       Yes     Travel_Rarely       1373  Research & Development
3   33        No  Travel_Frequently     1392  Research & Development
4   27        No     Travel_Rarely       591  Research & Development
5   32        No  Travel_Frequently     1005  Research & Development
6   59        No     Travel_Rarely       1324  Research & Development
7   30        No     Travel_Rarely       1358  Research & Development
8   38        No  Travel_Frequently      216  Research & Development
9   36        No     Travel_Rarely       1299  Research & Development

    DistanceFromHome  Education EducationField  EmployeeCount  EmployeeNumb
er  \
0                  1          2  Life Sciences              1
1
1                  8          1  Life Sciences              1
2
2                  2          2          Other              1
4
3                  3          4  Life Sciences              1
5
4                  2          1        Medical              1
7
5                  2          2  Life Sciences              1
8
6                  3          3        Medical              1
10
7                 24          1  Life Sciences              1
11
8                 23          3  Life Sciences              1
12
9                 27          3        Medical              1
13

   ...  RelationshipSatisfaction  StandardHours  StockOptionLevel  \
0  ...                         1             80                 0
1  ...                         4             80                 1
2  ...                         2             80                 0
3  ...                         3             80                 0
4  ...                         4             80                 1
5  ...                         3             80                 0
6  ...                         1             80                 3
7  ...                         2             80                 1
8  ...                         2             80                 0
9  ...                         2             80                 2

   TotalWorkingYears  TrainingTimesLastYear  WorkLifeBalance  YearsAtCompan
y  \
0                  8                      0                1
6
1                 10                      3                3             1
0
2                  7                      3                3
0
3                  8                      3                3
8
4                  6                      3                3
2
5                  8                      2                2
7
```

| | | | |
|---|---|---|---|
| 6 | 12 | 3 | 2 |
| 1 | | | |
| 7 | 1 | 2 | 3 |
| 1 | | | |
| 8 | 10 | 2 | 3 |
| 9 | | | |
| 9 | 17 | 3 | 2 |
| 7 | | | |

| | YearsInCurrentRole | YearsSinceLastPromotion | YearsWithCurrManager |
|---|---|---|---|
| 0 | 4 | 0 | 5 |
| 1 | 7 | 1 | 7 |
| 2 | 0 | 0 | 0 |
| 3 | 7 | 3 | 0 |
| 4 | 2 | 2 | 2 |
| 5 | 7 | 3 | 6 |
| 6 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 |
| 8 | 7 | 1 | 8 |
| 9 | 7 | 7 | 7 |

[10 rows x 35 columns]

```python
In [12]:  # Handle missing values
          df.dropna(inplace=True)
          print(df)
```

```
       Age Attrition     BusinessTravel  DailyRate              Department
\
0       41       Yes       Travel_Rarely       1102                   Sales
1       49        No   Travel_Frequently        279  Research & Development
2       37       Yes       Travel_Rarely       1373  Research & Development
3       33        No   Travel_Frequently       1392  Research & Development
4       27        No       Travel_Rarely        591  Research & Development
...    ...       ...                 ...        ...                     ...
1465    36        No   Travel_Frequently        884  Research & Development
1466    39        No       Travel_Rarely        613  Research & Development
1467    27        No       Travel_Rarely        155  Research & Development
1468    49        No   Travel_Frequently       1023                   Sales
1469    34        No       Travel_Rarely        628  Research & Development

      DistanceFromHome  Education EducationField  EmployeeCount  \
0                    1          2  Life Sciences              1
1                    8          1  Life Sciences              1
2                    2          2          Other              1
3                    3          4  Life Sciences              1
4                    2          1        Medical              1
...                ...        ...            ...            ...
1465                23          2        Medical              1
1466                 6          1        Medical              1
1467                 4          3  Life Sciences              1
1468                 2          3        Medical              1
1469                 8          3        Medical              1

      EmployeeNumber  ...  RelationshipSatisfaction StandardHours  \
0                  1  ...                         1            80
1                  2  ...                         4            80
2                  4  ...                         2            80
3                  5  ...                         3            80
4                  7  ...                         4            80
...              ...  ...                       ...           ...
1465            2061  ...                         3            80
1466            2062  ...                         1            80
1467            2064  ...                         2            80
1468            2065  ...                         4            80
1469            2068  ...                         1            80

      StockOptionLevel  TotalWorkingYears  TrainingTimesLastYear  \
0                    0                  8                      0
1                    1                 10                      3
2                    0                  7                      3
3                    0                  8                      3
4                    1                  6                      3
...                ...                ...                    ...
1465                 1                 17                      3
1466                 1                  9                      5
1467                 1                  6                      0
1468                 0                 17                      3
1469                 0                  6                      3

      WorkLifeBalance  YearsAtCompany YearsInCurrentRole  \
0                   1               6                  4
1                   3              10                  7
2                   3               0                  0
3                   3               8                  7
4                   3               2                  2
...               ...             ...                ...
1465                3               5                  2
```

```
1466             3              7              7
1467             3              6              2
1468             2              9              6
1469             4              4              3

      YearsSinceLastPromotion   YearsWithCurrManager
0                           0                      5
1                           1                      7
2                           0                      0
3                           3                      0
4                           2                      2
...                       ...                    ...
1465                        0                      3
1466                        1                      7
1467                        0                      3
1468                        0                      8
1469                        1                      2

[1470 rows x 35 columns]
```

```
In [31]: missing_values = df.isnull().sum()
         print("Missing values in each column:")
         print(missing_values)
```

```
Missing values in each column:
Age                         0
Attrition                   0
BusinessTravel              0
DailyRate                   0
Department                  0
DistanceFromHome            0
Education                   0
EducationField              0
EmployeeCount               0
EmployeeNumber              0
EnvironmentSatisfaction     0
Gender                      0
HourlyRate                  0
JobInvolvement              0
JobLevel                    0
JobRole                     0
JobSatisfaction             0
MaritalStatus               0
MonthlyIncome               0
MonthlyRate                 0
NumCompaniesWorked          0
Over18                      0
OverTime                    0
PercentSalaryHike           0
PerformanceRating           0
RelationshipSatisfaction    0
StandardHours               0
StockOptionLevel            0
TotalWorkingYears           0
TrainingTimesLastYear       0
WorkLifeBalance             0
YearsAtCompany              0
YearsInCurrentRole          0
YearsSinceLastPromotion     0
YearsWithCurrManager        0
dtype: int64
```

```python
In [37]:   # Check for missing values
           print(df.isnull().sum())
```

```
Age                          0
Attrition                    0
BusinessTravel               0
DailyRate                    0
Department                   0
DistanceFromHome             0
Education                    0
EducationField               0
EmployeeCount                0
EmployeeNumber               0
EnvironmentSatisfaction      0
Gender                       0
HourlyRate                   0
JobInvolvement               0
JobLevel                     0
JobRole                      0
JobSatisfaction              0
MaritalStatus                0
MonthlyIncome                0
MonthlyRate                  0
NumCompaniesWorked           0
Over18                       0
OverTime                     0
PercentSalaryHike            0
PerformanceRating            0
RelationshipSatisfaction     0
StandardHours                0
StockOptionLevel             0
TotalWorkingYears            0
TrainingTimesLastYear        0
WorkLifeBalance              0
YearsAtCompany               0
YearsInCurrentRole           0
YearsSinceLastPromotion      0
YearsWithCurrManager         0
dtype: int64
```

```python
In [38]:   # Drop rows with missing values
           df.dropna(inplace=True)
```

```python
In [39]:   # Remove duplicate rows
           df.drop_duplicates(inplace=True)
```

```
In [41]: # Check the data types after conversion
         print(df.dtypes)
```

```
Age                          int64
Attrition                    object
BusinessTravel               object
DailyRate                    int64
Department                   object
DistanceFromHome             int64
Education                    int64
EducationField               object
EmployeeCount                int64
EmployeeNumber               int64
EnvironmentSatisfaction      int64
Gender                       object
HourlyRate                   int64
JobInvolvement               int64
JobLevel                     int64
JobRole                      object
JobSatisfaction              int64
MaritalStatus                object
MonthlyIncome                int64
MonthlyRate                  int64
NumCompaniesWorked           int64
Over18                       object
OverTime                     object
PercentSalaryHike            int64
PerformanceRating            int64
RelationshipSatisfaction     int64
StandardHours                int64
StockOptionLevel             int64
TotalWorkingYears            int64
TrainingTimesLastYear        int64
WorkLifeBalance              int64
YearsAtCompany               int64
YearsInCurrentRole           int64
YearsSinceLastPromotion      int64
YearsWithCurrManager         int64
dtype: object
```

```
In [42]: # Save the cleaned dataset
         df.to_csv("cleaned_dataset.csv", index=False)
```

```
In [ ]:
```