

Module 4: Python Project

Import Libraries

```
In [95]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

In [41]: data = pd.read_csv('abc_company.csv')

In [51]: data.head

Out[51]: <bound method NDFrame.head of
0 Avery Bradley Boston Celtics 0 PG 25 06-Feb 180
1 Jae Crowder Boston Celtics 99 SF 25 06-Jun 235
2 John Holland Boston Celtics 30 SG 27 06-May 205
3 R.J. Hunter Boston Celtics 28 SG 22 06-May 185
4 Jonas Jerebko Boston Celtics 8 PF 29 06-Oct 231
.. ...
453 Shelvin Mack Utah Jazz 8 PG 26 06-Mar 203
454 Raul Neto Utah Jazz 25 PG 24 06-Jan 179
455 Tibor Pleiss Utah Jazz 21 C 26 07-Mar 256
456 Jeff Withey Utah Jazz 24 C 26 7-0 231
457 Priyanka Utah Jazz 34 C 25 07-Mar 231
College Salary
0 Texas 7730337.0
1 Marquette 6796117.0
2 Boston University NaN
3 Georgia State 1148640.0
4 NaN 5000000.0
.. ...
453 Butler 2433333.0
454 NaN 900000.0
455 NaN 2900000.0
456 Kansas 947276.0
457 Kansas 947276.0
[458 rows x 9 columns]>
```

```
In [61]: data.tail

Out[61]: <bound method NDFrame.tail of
0 Avery Bradley Boston Celtics 0 PG 25 06-Feb 180
1 Jae Crowder Boston Celtics 99 SF 25 06-Jun 235
2 John Holland Boston Celtics 30 SG 27 06-May 205
3 R.J. Hunter Boston Celtics 28 SG 22 06-May 185
4 Jonas Jerebko Boston Celtics 8 PF 29 06-Oct 231
.. ...
453 Shelvin Mack Utah Jazz 8 PG 26 06-Mar 203
454 Raul Neto Utah Jazz 25 PG 24 06-Jan 179
455 Tibor Pleiss Utah Jazz 21 C 26 07-Mar 256
456 Jeff Withey Utah Jazz 24 C 26 7-0 231
457 Priyanka Utah Jazz 34 C 25 07-Mar 231
College Salary
0 Texas 7730337.0
1 Marquette 6796117.0
2 Boston University NaN
3 Georgia State 1148640.0
4 NaN 5000000.0
.. ...
453 Butler 2433333.0
454 NaN 900000.0
455 NaN 2900000.0
456 Kansas 947276.0
457 Kansas 947276.0
[458 rows x 9 columns]>
```

```
In [64]: data.info

Out[64]: <bound method DataFrame.info of
0 Avery Bradley Boston Celtics 0 PG 25 06-Feb 180
1 Jae Crowder Boston Celtics 99 SF 25 06-Jun 235
2 John Holland Boston Celtics 30 SG 27 06-May 205
3 R.J. Hunter Boston Celtics 28 SG 22 06-May 185
4 Jonas Jerebko Boston Celtics 8 PF 29 06-Oct 231
.. ...
453 Shelvin Mack Utah Jazz 8 PG 26 06-Mar 203
454 Raul Neto Utah Jazz 25 PG 24 06-Jan 179
455 Tibor Pleiss Utah Jazz 21 C 26 07-Mar 256
456 Jeff Withey Utah Jazz 24 C 26 7-0 231
457 Priyanka Utah Jazz 34 C 25 07-Mar 231
College Salary height
0 Texas 7730337.0 171
1 Marquette 6796117.0 176
2 Boston University NaN 174
3 Georgia State 1148640.0 151
4 NaN 5000000.0 154
.. ...
453 Butler 2433333.0 179
454 NaN 900000.0 161
455 NaN 2900000.0 177
456 Kansas 947276.0 162
457 Kansas 947276.0 150
[458 rows x 10 columns]>
```

```
In [76]: data.isnull()

Out[76]:
Name Team Number Position Age Height Weight College Salary height
0 False False False False False False False False False
1 False False False False False False False False False
2 False False False False False False False False True False
3 False False False False False False False False False False
4 False False False False False False True False False
...
453 False False False False False False False False False
454 False False False False False False False True False False
455 False False False False False False True False False
456 False False False False False False False False False
457 False False False False False False False False False

458 rows x 10 columns
```

```
In [80]: # --- Null Value Handling ---
# 1. Check for Null Values
print("Null Values:\n", data.isnull().sum()) # Print the number of nulls in each column

Null Values:
Name 0
Team 0
Number 0
Position 0
Age 0
Height 0
Weight 0
College 84
Salary 11
height 0
dtype: int64

In [97]: data['Salary'].fillna(data['Salary'].median(), inplace=True) # Fill with median (or some other appropriate value)

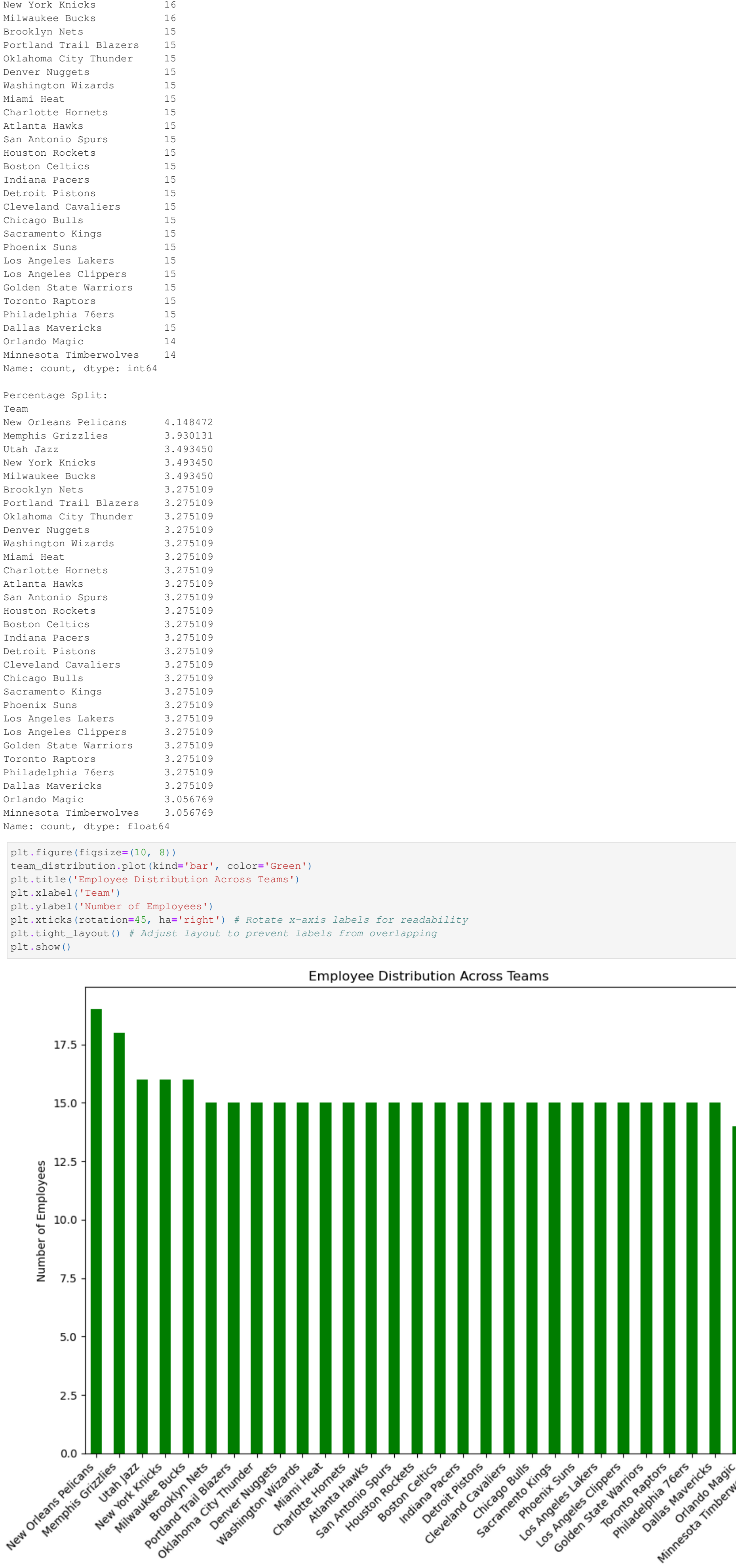
In [131]: data['College'].fillna(0, inplace=True) # Fill with 0 (or some other appropriate value)

In [99]: # Preprocessing: Correcting the 'height' column
data['height'] = np.random.randint(150, 181, size=len(data)) # Generates random integers between 150 and 180 (inclusive)
```

Analysis Tasks:

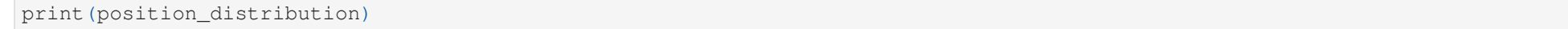
```
In [133]: # 1. Employee Distribution Across Teams
team_distribution = data['Team'].value_counts()
team_percentage = (team_distribution / len(data)) * 100
```

```
In [135]: print("\nEmployee Distribution Across Teams:")
print(team_distribution)
print("\nPercentage Split:")
print(team_percentage)
```



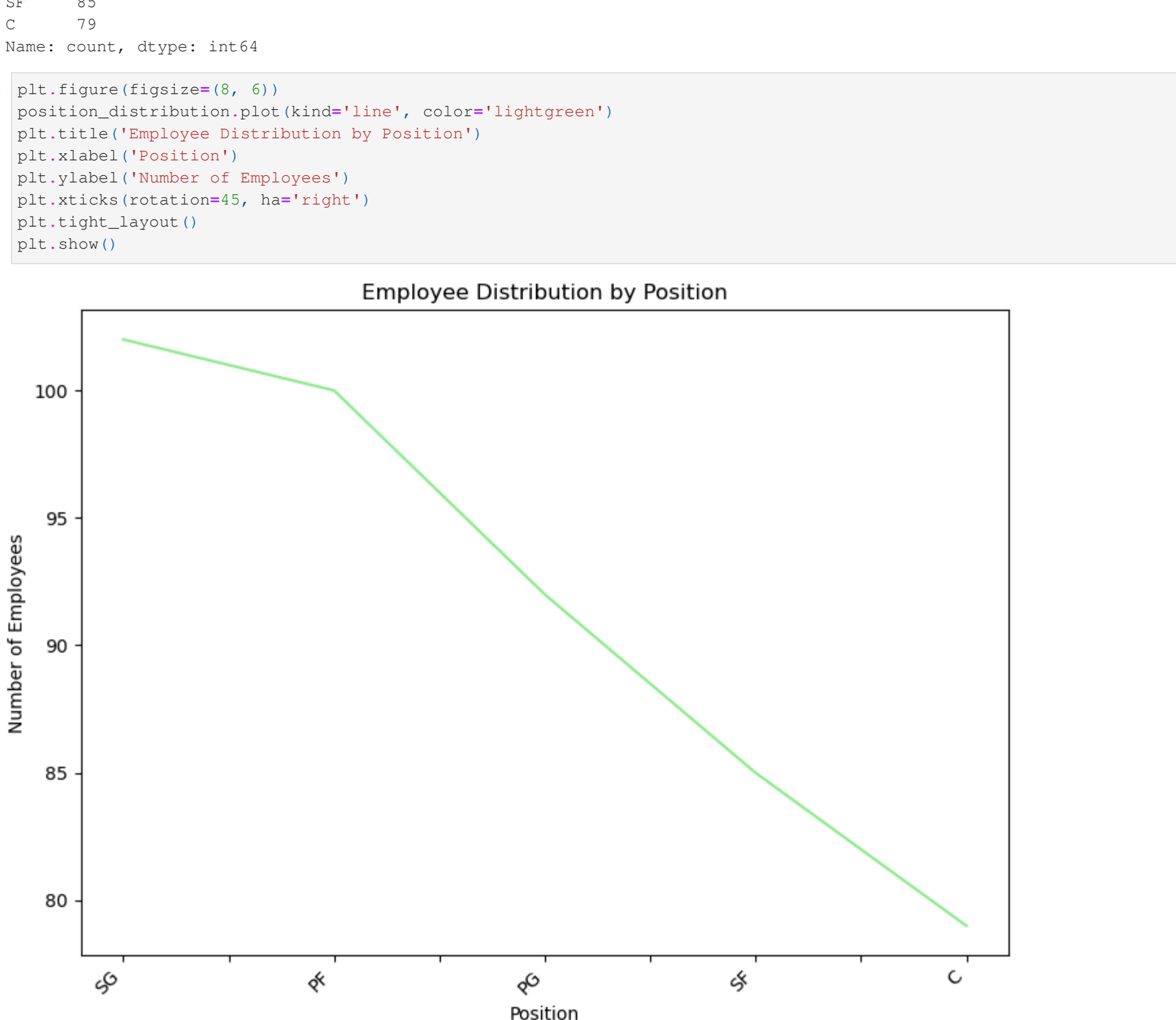
```
In [137]: plt.figure(figsize=(10, 8))
team_distribution.plot(kind='bar', color='Green')
plt.title('Employee Distribution Across Teams')
plt.xlabel('Team')
plt.ylabel('Number of Employees')
plt.xticks(rotation=5, ha='right') # Rotate x-axis labels for readability
plt.tight_layout() # Adjust layout to prevent labels from overlapping
plt.show()
```

Employee Distribution Across Teams



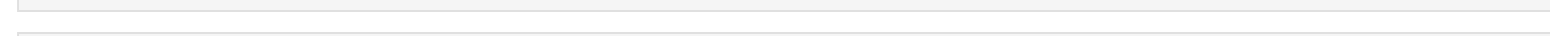
```
In [138]: # 2. Employee Segregation by Position
position_distribution = data['Position'].value_counts()
```

```
In [141]: print("\nEmployee Segregation by Position:")
print(position_distribution)
```



```
In [143]: plt.figure(figsize=(8, 6))
position_distribution.plot(kind='line', color='lightgreen')
plt.title('Employee Distribution by Position')
plt.xlabel('Position')
plt.ylabel('Number of Employees')
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.show()
```

Employee Distribution by Position



```
In [145]: # 3. Predominant Age Group
age_group = data['Age'].describe() # Use describe() for summary stats including median
median_age = data['Age'].median()
```

```
In [147]: print("\nAge Statistics:")
print(age_group)
print(f"Median Age: {median_age}") # Median is a better measure of central tendency than mean if data is skewed.
```



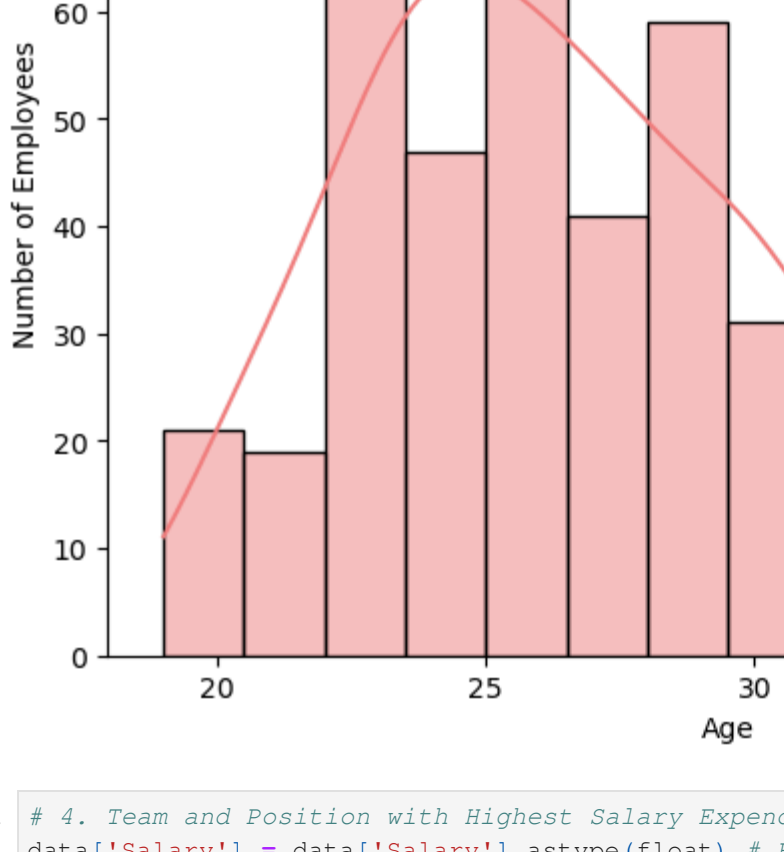
```
In [149]: plt.figure(figsize=(8, 6))
sns.histplot(data['Age'], kde=True, color='lightcoral') # Histogram with KDE for better visualization of distribution
plt.title('Age Distribution of Employees')
plt.xlabel('Age')
plt.ylabel('Number of Employees')
plt.show()
```

Age Distribution of Employees



```
In [151]: # 4. Team and Position with Highest Salary Expenditure
data['Salary'] = data['Salary'].astype(float) # Ensure salary is a float for calculations
salary_by_team_position = data.groupby(['Team', 'Position'])['Salary'].sum().reset_index()
highest_expenditure = salary_by_team_position.sort_values(by='Salary', ascending=False).iloc[0]
```

```
In [153]: print("\nTeam and Position with Highest Salary Expenditure:")
print(highest_expenditure)
```



```
In [155]: # 5. Correlation between Age and Salary
correlation = data['Age'].corr(data['Salary'])
print(f"Correlation between Age and Salary: {correlation}")
```

Correlation between Age and Salary: 0.20912419115196068

```
In [161]: plt.figure(figsize=(12, 10))
sns.scatterplot(data['Age'], y='Salary', data=data, color='red')
plt.title('Correlation between Age and Salary')
plt.xlabel('Age')
plt.ylabel('Salary')
plt.show()
```

Correlation between Age and Salary: 0.20912419115196068



