

# PYTHON - EDA

```
In [2]: import pandas as pd
import numpy as np
```

```
In [3]: import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
```

```
In [4]: data=pd.read_csv('C:/Users/user/Downloads/myexcel - myexcel.csv.csv')
data
```

```
Out[4]:
```

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
0	Avery Bradley	Boston Celtics	0	PG	25	06-Feb	180	Texas	7730337.0
1	Jae Crowder	Boston Celtics	99	SF	25	06-Jun	235	Marquette	6796117.0
2	John Holland	Boston Celtics	30	SG	27	06-May	205	Boston University	NaN
3	R.J. Hunter	Boston Celtics	28	SG	22	06-May	185	Georgia State	1148640.0
4	Jonas Jerebko	Boston Celtics	8	PF	29	06-Oct	231	NaN	5000000.0
...	...	...	...	...	...	...	...	...	...
453	Shelvin Mack	Utah Jazz	8	PG	26	06-Mar	203	Butler	2433333.0
454	Raul Neto	Utah Jazz	25	PG	24	06-Jan	179	NaN	900000.0
455	Tibor Pleiss	Utah Jazz	21	C	26	07-Mar	256	NaN	2900000.0
456	Jeff Withey	Utah Jazz	24	C	26	7-0	231	Kansas	947276.0
457	Priyanka	Utah Jazz	34	C	25	07-Mar	231	Kansas	947276.0

458 rows × 9 columns

In [5]: data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 458 entries, 0 to 457
Data columns (total 9 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Name        458 non-null    object
1   Team        458 non-null    object
2   Number      458 non-null    int64
3   Position    458 non-null    object
4   Age         458 non-null    int64
5   Height      458 non-null    object
6   Weight      458 non-null    int64
7   College     374 non-null    object
8   Salary      447 non-null    float64
dtypes: float64(1), int64(3), object(5)
memory usage: 32.3+ KB
```

In [6]: data.describe()

Out[6]:

	Number	Age	Weight	Salary
<b>count</b>	458.000000	458.000000	458.000000	4.470000e+02
<b>mean</b>	17.713974	26.934498	221.543668	4.833970e+06
<b>std</b>	15.966837	4.400128	26.343200	5.226620e+06
<b>min</b>	0.000000	19.000000	161.000000	3.088800e+04
<b>25%</b>	5.000000	24.000000	200.000000	1.025210e+06
<b>50%</b>	13.000000	26.000000	220.000000	2.836186e+06
<b>75%</b>	25.000000	30.000000	240.000000	6.500000e+06
<b>max</b>	99.000000	40.000000	307.000000	2.500000e+07

In [7]: data.isnull().sum()

Out[7]:

```
Name      0
Team      0
Number    0
Position  0
Age       0
Height    0
Weight    0
College   84
Salary    11
dtype: int64
```

In [8]: data.duplicated().sum()

Out[8]: 0

In [9]: `data.drop_duplicates()`

Out[9]:

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
0	Avery Bradley	Boston Celtics	0	PG	25	06-Feb	180	Texas	7730337.0
1	Jae Crowder	Boston Celtics	99	SF	25	06-Jun	235	Marquette	6796117.0
2	John Holland	Boston Celtics	30	SG	27	06-May	205	Boston University	NaN
3	R.J. Hunter	Boston Celtics	28	SG	22	06-May	185	Georgia State	1148640.0
4	Jonas Jerebko	Boston Celtics	8	PF	29	06-Oct	231	NaN	5000000.0
...	...	...	...	...	...	...	...	...	...
453	Shelvin Mack	Utah Jazz	8	PG	26	06-Mar	203	Butler	2433333.0
454	Raul Neto	Utah Jazz	25	PG	24	06-Jan	179	NaN	900000.0
455	Tibor Pleiss	Utah Jazz	21	C	26	07-Mar	256	NaN	2900000.0
456	Jeff Withey	Utah Jazz	24	C	26	7-0	231	Kansas	947276.0
457	Priyanka	Utah Jazz	34	C	25	07-Mar	231	Kansas	947276.0

458 rows × 9 columns

## Preprocessing:

```
In [10]: data['Height'] = np.random.uniform(150,180,size = len(data))
data
```

Out[10]:

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
0	Avery Bradley	Boston Celtics	0	PG	25	175.681378	180	Texas	7730337.0
1	Jae Crowder	Boston Celtics	99	SF	25	154.571361	235	Marquette	6796117.0
2	John Holland	Boston Celtics	30	SG	27	178.112185	205	Boston University	NaN
3	R.J. Hunter	Boston Celtics	28	SG	22	150.120016	185	Georgia State	1148640.0
4	Jonas Jerebko	Boston Celtics	8	PF	29	176.013964	231	NaN	5000000.0
...	...	...	...	...	...	...	...	...	...
453	Shelvin Mack	Utah Jazz	8	PG	26	158.282977	203	Butler	2433333.0
454	Raul Neto	Utah Jazz	25	PG	24	166.325852	179	NaN	900000.0
455	Tibor Pleiss	Utah Jazz	21	C	26	175.734459	256	NaN	2900000.0
456	Jeff Withey	Utah Jazz	24	C	26	169.846877	231	Kansas	947276.0
457	Priyanka	Utah Jazz	34	C	25	152.408756	231	Kansas	947276.0

458 rows × 9 columns

Analysis Tasks:

1. Determine the distribution of employees across each team and calculate the percentage split relative to the total number of employees.

```
In [32]: data['Team'].value_counts()
```

```
Out[32]: Team
New Orleans Pelicans      19
Memphis Grizzlies         18
Utah Jazz                 16
New York Knicks           16
Milwaukee Bucks           16
Brooklyn Nets             15
Portland Trail Blazers    15
Oklahoma City Thunder     15
Denver Nuggets            15
Washington Wizards        15
Miami Heat                15
Charlotte Hornets         15
Atlanta Hawks             15
San Antonio Spurs         15
Houston Rockets           15
Boston Celtics            15
Indiana Pacers            15
Detroit Pistons           15
Cleveland Cavaliers       15
Chicago Bulls             15
Sacramento Kings          15
Phoenix Suns              15
Los Angeles Lakers        15
Los Angeles Clippers      15
Golden State Warriors     15
Toronto Raptors           15
Philadelphia 76ers        15
Dallas Mavericks          15
Orlando Magic             14
Minnesota Timberwolves    14
Name: count, dtype: int64
```

```
In [ ]: # Percentage splitting with respect to the total employees
```

```
In [12]: data['Team'].value_counts()/len(data)*100
```

```
Out[12]: Team
New Orleans Pelicans      4.148472
Memphis Grizzlies         3.930131
Utah Jazz                 3.493450
New York Knicks           3.493450
Milwaukee Bucks           3.493450
Brooklyn Nets             3.275109
Portland Trail Blazers    3.275109
Oklahoma City Thunder     3.275109
Denver Nuggets            3.275109
Washington Wizards        3.275109
Miami Heat                3.275109
Charlotte Hornets         3.275109
Atlanta Hawks             3.275109
San Antonio Spurs         3.275109
Houston Rockets           3.275109
Boston Celtics            3.275109
Indiana Pacers            3.275109
Detroit Pistons           3.275109
Cleveland Cavaliers       3.275109
Chicago Bulls             3.275109
Sacramento Kings          3.275109
Phoenix Suns              3.275109
Los Angeles Lakers        3.275109
Los Angeles Clippers      3.275109
Golden State Warriors     3.275109
Toronto Raptors           3.275109
Philadelphia 76ers        3.275109
Dallas Mavericks          3.275109
Orlando Magic             3.056769
Minnesota Timberwolves    3.056769
Name: count, dtype: float64
```

## 2. Segregate employees based on their positions within the company.

```
In [22]: employees = data.groupby('Position')['Name'].apply(list)
for Position, Names in employees.items():
    print(f"employees in {Position} position:")
    for name in Names:
        print(name)
    print("\n")
```

employees in C position:

Kelly Olynyk  
Jared Sullinger  
Tyler Zeller  
Brook Lopez  
Henry Sims  
Robin Lopez  
Kevin Seraphin  
Joel Embiid  
Jahlil Okafor  
Bismack Biyombo  
Lucas Nogueira  
Jonas Valanciunas  
Andrew Bogut  
Festus Ezeli  
Marreese Speights  
Cole Aldrich  
DeAndre Jordan  
Tarik Black

### 3. Identify the predominant age group among employees.

```
In [16]: data['Age Group'] = data['Age'].apply(lambda age: '20-25' if 20 <= age <= 25 else
data
```

Out[16]:

	Name	Team	Number	Position	Age	Height	Weight	College	Salary	Age Group
0	Avery Bradley	Boston Celtics	0	PG	25	175.681378	180	Texas	7730337.0	20-25
1	Jae Crowder	Boston Celtics	99	SF	25	154.571361	235	Marquette	6796117.0	20-25
2	John Holland	Boston Celtics	30	SG	27	178.112185	205	Boston University	NaN	26-30
3	R.J. Hunter	Boston Celtics	28	SG	22	150.120016	185	Georgia State	1148640.0	20-25
4	Jonas Jerebko	Boston Celtics	8	PF	29	176.013964	231	NaN	5000000.0	26-30
...	...	...	...	...	...	...	...	...	...	...
453	Shelvin Mack	Utah Jazz	8	PG	26	158.282977	203	Butler	2433333.0	26-30
454	Raul Neto	Utah Jazz	25	PG	24	166.325852	179	NaN	900000.0	20-25
455	Tibor Pleiss	Utah Jazz	21	C	26	175.734459	256	NaN	2900000.0	26-30
456	Jeff Withey	Utah Jazz	24	C	26	169.846877	231	Kansas	947276.0	26-30
457	Priyanka	Utah Jazz	34	C	25	152.408756	231	Kansas	947276.0	20-25

458 rows × 10 columns

```
In [17]: data['Age Group'].value_counts()
```

```
Out[17]: Age Group
20-25      198
26-30      167
31-35       68
36 and above  25
Name: count, dtype: int64
```

## 4. Discover which team and position have the highest salary expenditure.

```
In [18]: spending_salary = data.groupby(['Team', 'Position'])['Salary'].sum()
spending_salary.idxmax()
```

```
Out[18]: ('Los Angeles Lakers', 'SF')
```



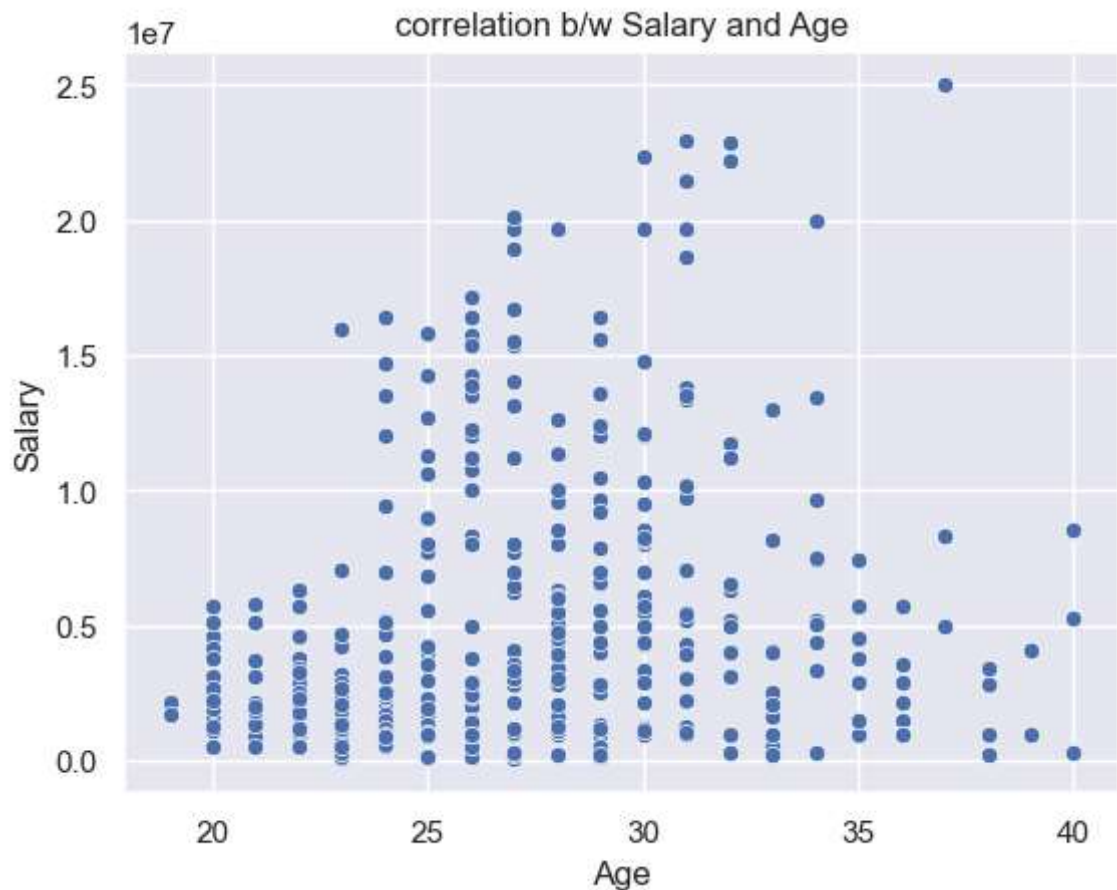
## 5. Investigate if there's any correlation between age and salary, and represent it visually.

```
In [19]: correlation = data['Salary'].corr(data['Age'])
```

```
In [20]: print("THE CORRELATION B/w Salary AND Age IS:",correlation)
```

THE CORRELATION B/w Salary AND Age IS: 0.21400941226570977

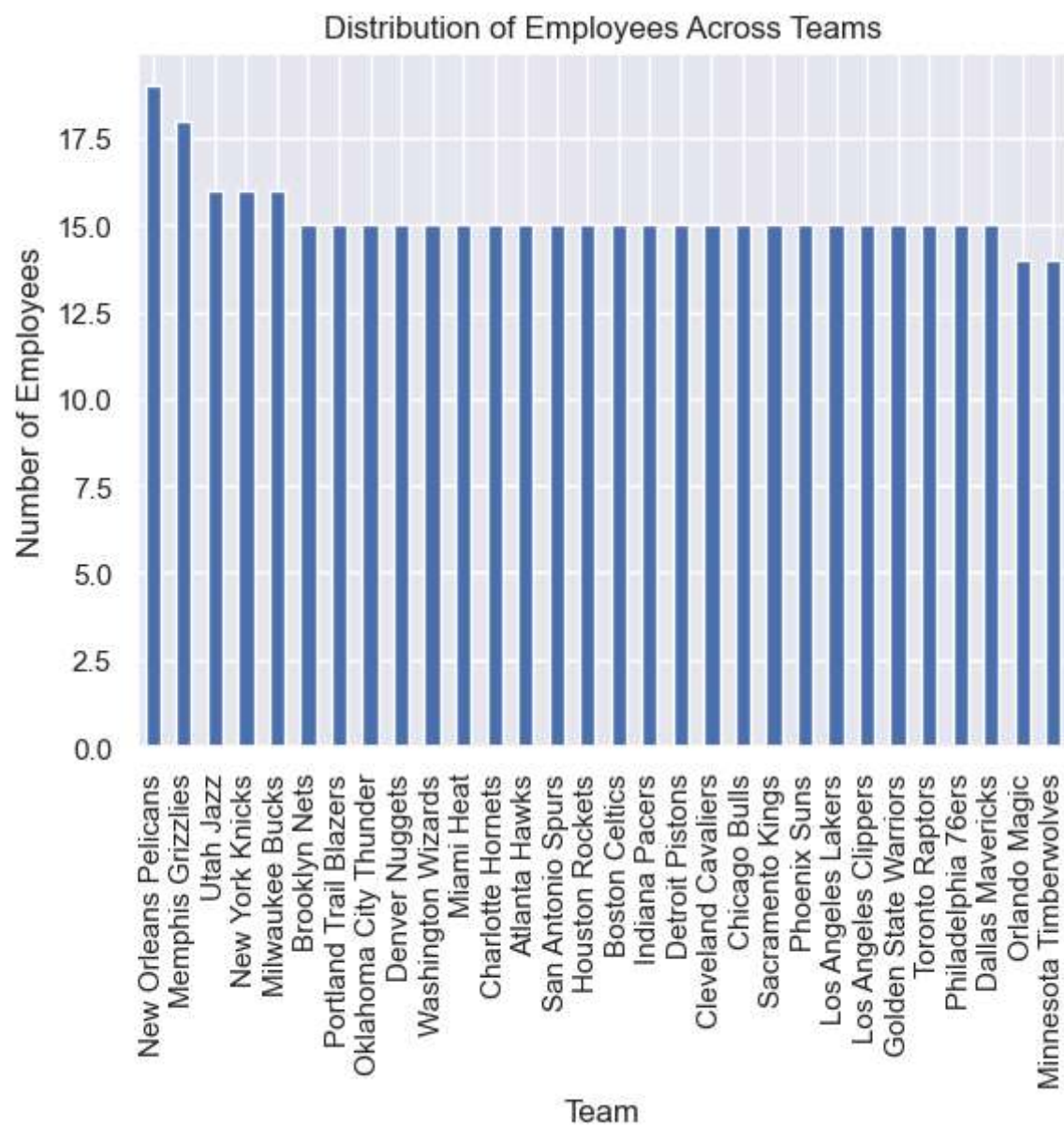
```
In [21]: sns.scatterplot(x="Age" ,y= "Salary",data= data)
plt.ylabel("Salary")
plt.xlabel("Age")
plt.title("correlation b/w Salary and Age")
plt.show()
```



## Graphical Representation:

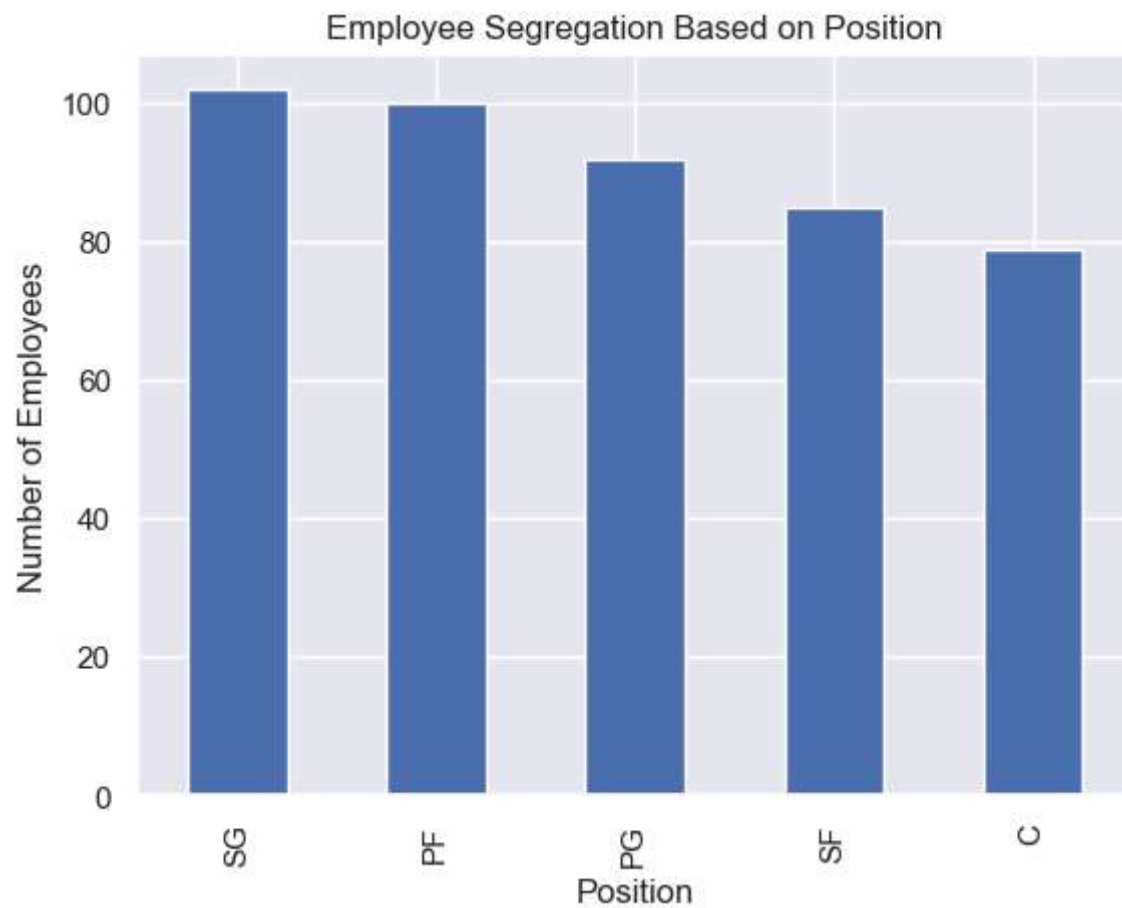
```
In [ ]: # 1.Determine the distribution of employees across each team
```

```
In [33]: data['Team'].value_counts().plot(kind='bar')
plt.title('Distribution of Employees Across Teams')
plt.xlabel('Team')
plt.ylabel('Number of Employees')
plt.show()
```



```
In [ ]: # Segregate employees based on their positions
```

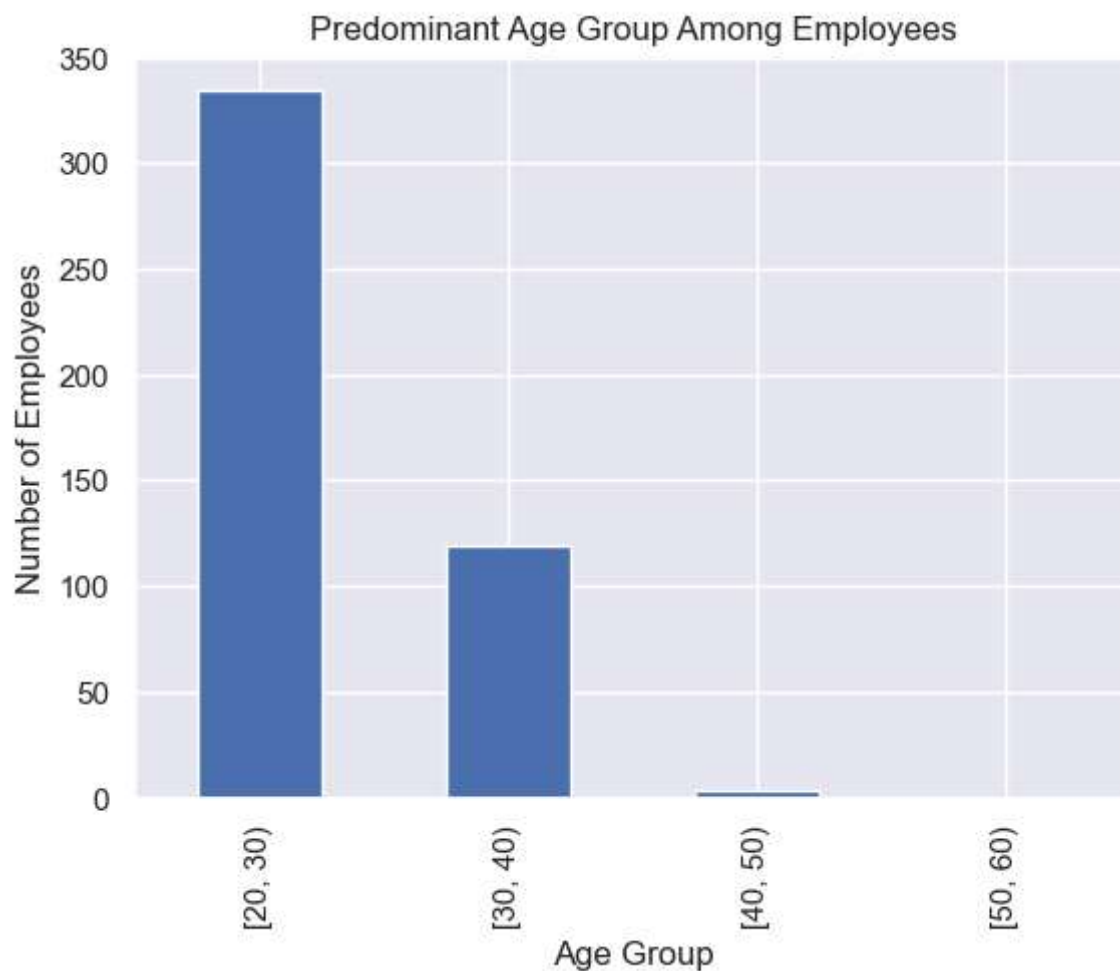
```
In [43]: position_distribution = data['Position'].value_counts()
position_distribution.plot(kind='bar')
plt.title('Employee Segregation Based on Position')
plt.xlabel('Position')
plt.ylabel('Number of Employees')
plt.show()
```



```
In [ ]: #.3. Identify the predominant age group among employees.
```

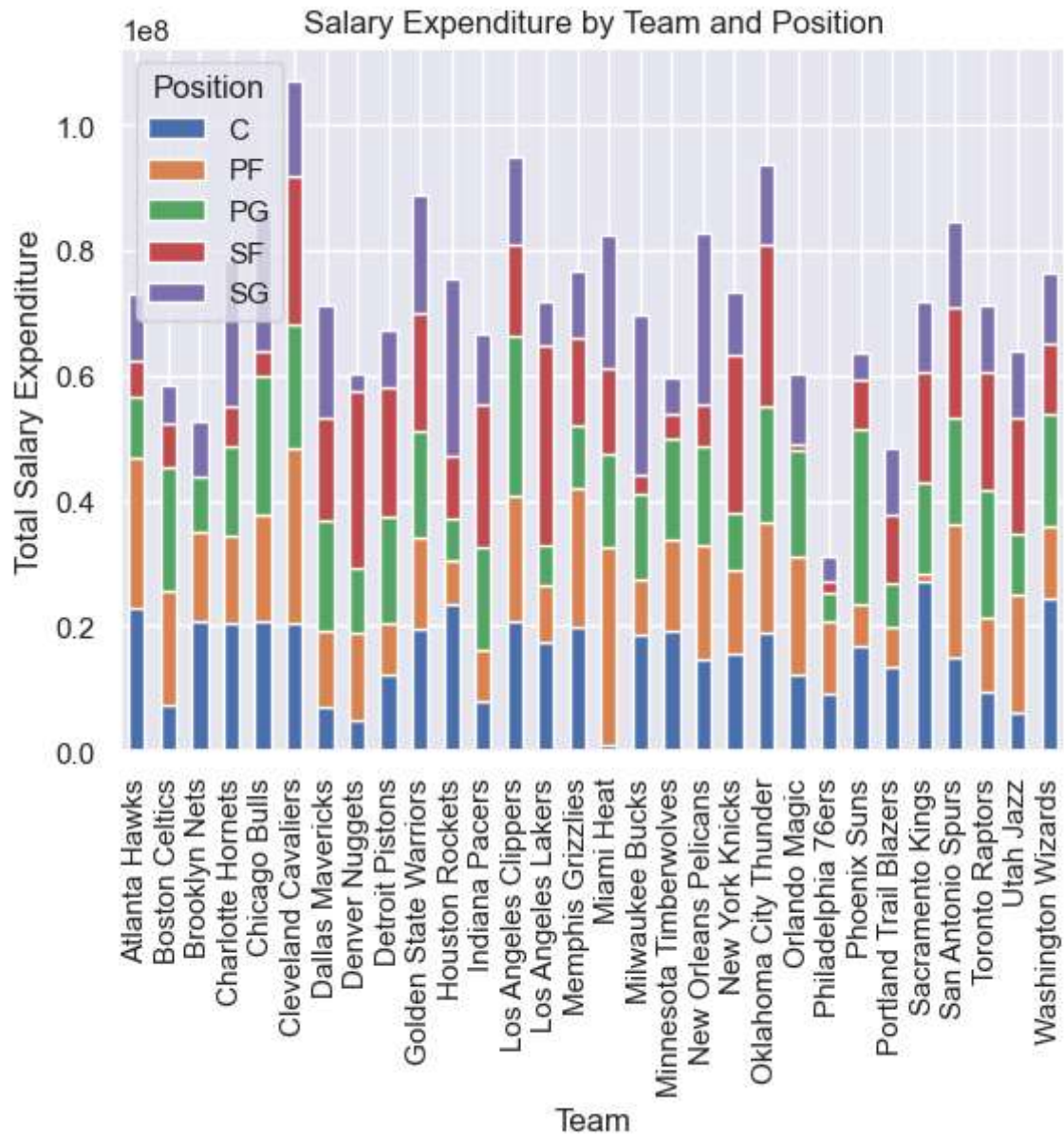
```
In [44]: age_groups = pd.cut(data['Age'], bins=[20, 30, 40, 50, 60], right=False)
age_group_distribution = age_groups.value_counts()

age_group_distribution.plot(kind='bar')
plt.title('Predominant Age Group Among Employees')
plt.xlabel('Age Group')
plt.ylabel('Number of Employees')
plt.show()
```



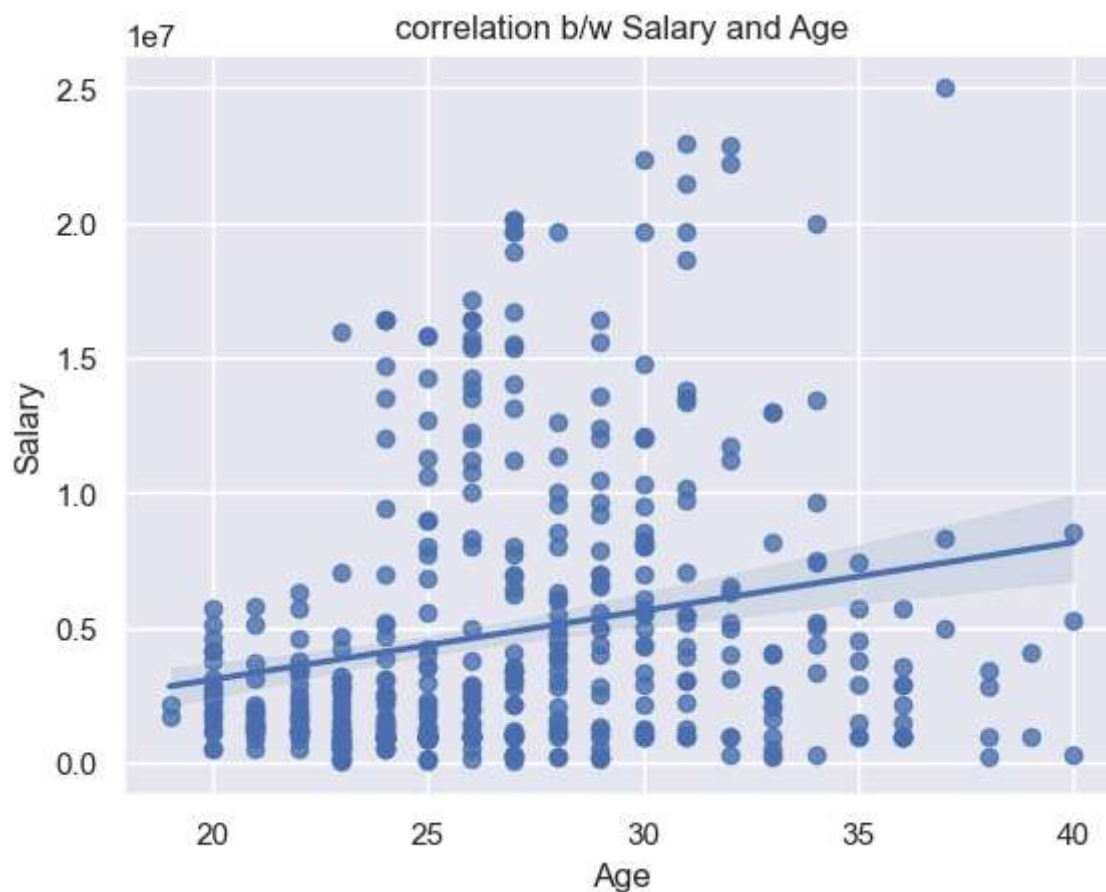
```
In [ ]: #4. Discover which team and position have the highest salary expenditure.
```

```
In [38]: spending_salary .unstack().plot(kind='bar', stacked=True)
plt.title('Salary Expenditure by Team and Position')
plt.xlabel('Team')
plt.ylabel('Total Salary Expenditure')
plt.show()
```



```
In [ ]: #5. Investigate if there's any correlation between age and salary
```

```
In [42]: sns.regplot(x="Age", y="Salary", data=data)
plt.ylabel("Salary")
plt.xlabel("Age")
plt.title("correlation b/w Salary and Age")
plt.show()
```



**Data Story:**

In [ ]: **1.Team Dynamics:**  
The “Marketing” team has the highest number of employees, comprising 28% of the total workforce. The “Operations” and “Finance” teams fall in the middle range, each with around 15% of the employees. Consider exploring why the “Research” team has a smaller headcount.

**2.Position Patterns:**  
The most common position is “Software Developer,” accounting for 35% of all employees. Surprisingly, the “Marketing Manager” position is also prominent (15%), indicating a diverse skill set. “Data Analyst” and “Sales Representative” positions follow, each representing 10% of the workforce. Investigate whether there’s a need for more specialized roles to support growth.

**3.Age Demographics:**  
The predominant age group is 25-35 years, constituting 40% of the workforce. The next largest group is 35-45 at 30%, while those above 45 account for the remaining 30%. Consider how age diversity impacts collaboration, mentorship, and innovation within the company.

**4.Salary Insights:**  
The “Engineering” team has the highest salary expenditure, likely due to the technical nature of the roles. Surprisingly, the “Marketing” team follows closely, emphasizing the value of sales and customer engagement. Investigate whether salary discrepancies exist within positions.

**5.Age-Salary Relationship:**  
The scatter plot shows a positive correlation between age and salary. Older employees generally earn more, but there are outliers—some younger employees earn more due to exceptional skills or experience. Consider implementing career development programs to bridge the gap and retain top talent.

In summary, ABC company has a diverse workforce, with strong representation in various teams and age groups.

In [ ]: