

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import matplotlib inline
import seaborn as sns

D:\Jupyter Notebook\lib\site-packages\scipy\_init_.py:146: UserWarning: A NumPy version >=1.16.5 and <1.23.0 is required for this version of SciPy (detected version 1.23.1)
warnings.warn(f"A NumPy version >=({np.minversion}) and <({np.maxversion})")

In [2]: data=pd.read_csv(r"D:\PwDS-Pooja Sharma\AI1 G-drive Files\Day15\IBM Attrition Data.csv")

In [3]: data.head()
```

	Age	Attrition	Department	DistanceFromHome	Education	EducationField	EnvironmentSatisfaction	JobSatisfaction	MaritalStatus	MonthlyIncome	NumCompaniesWorked	WorkLifeBalance	YearsAtCompany
0	41	Yes	Sales	1	2	Life Sciences	2	4	Single	5993	8	1	6
1	49	No	Research & Development	8	1	Life Sciences	3	2	Married	5130	1	3	10
2	37	Yes	Research & Development	2	2	Other	4	3	Single	2090	6	3	0
3	33	No	Research & Development	3	4	Life Sciences	4	3	Married	2909	1	3	8
4	27	No	Research & Development	2	1	Medical	1	2	Married	3468	9	3	2

```
In [4]: data.shape

Out[4]: (1470, 13)
```

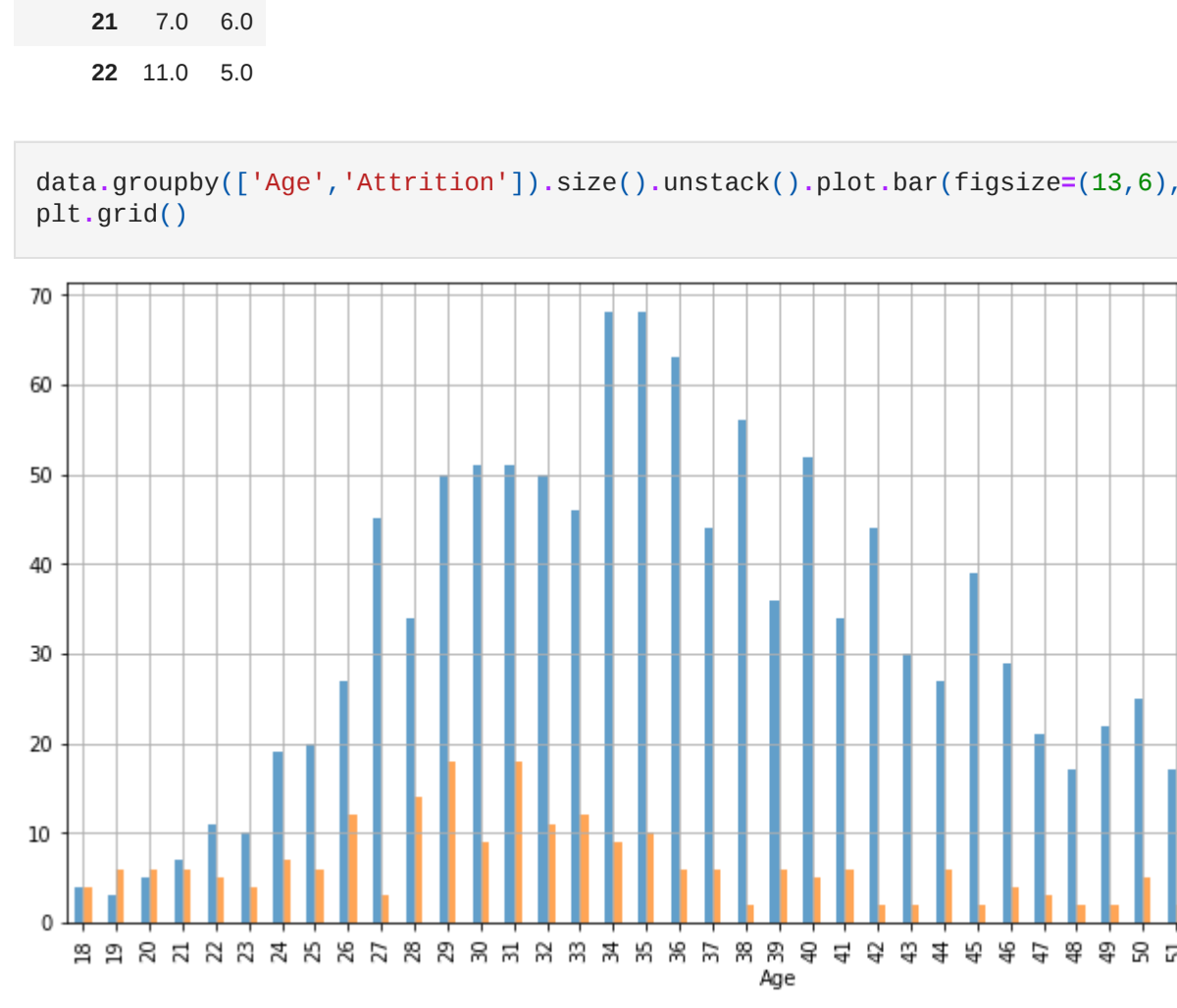
```
In [5]: data.info() #check for null values

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 13 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Age                   1470 non-null  int64
 1   Attrition              1470 non-null  object
 2   Department             1470 non-null  object
 3   DistanceFromHome       1470 non-null  int64
 4   Education              1470 non-null  int64
 5   EducationField         1470 non-null  object
 6   EnvironmentSatisfaction 1470 non-null  int64
 7   JobSatisfaction        1470 non-null  int64
 8   MaritalStatus         1470 non-null  object
 9   MonthlyIncome         1470 non-null  int64
10   NumCompaniesWorked    1470 non-null  int64
11   WorkLifeBalance       1470 non-null  int64
12   YearsAtCompany        1470 non-null  int64
dtypes: int64(9), object(4)
memory usage: 149.4+ KB
```

```
In [6]: data.isnull().sum() #sum of all null values must be zero

Attrition      0
Age             0
Department      0
DistanceFromHome 0
EducationField  0
EnvironmentSatisfaction 0
JobSatisfaction 0
MaritalStatus  0
MonthlyIncome  0
NumCompaniesWorked 0
WorkLifeBalance 0
YearsAtCompany 0
dtype: int64
```

```
In [7]: data['Age'].value_counts().plot.bar(figsize=(10,6),xlabel='Age',ylabel='Count',title='Age distribution',alpha=0.5)
plt.grid()
plt.show()
```



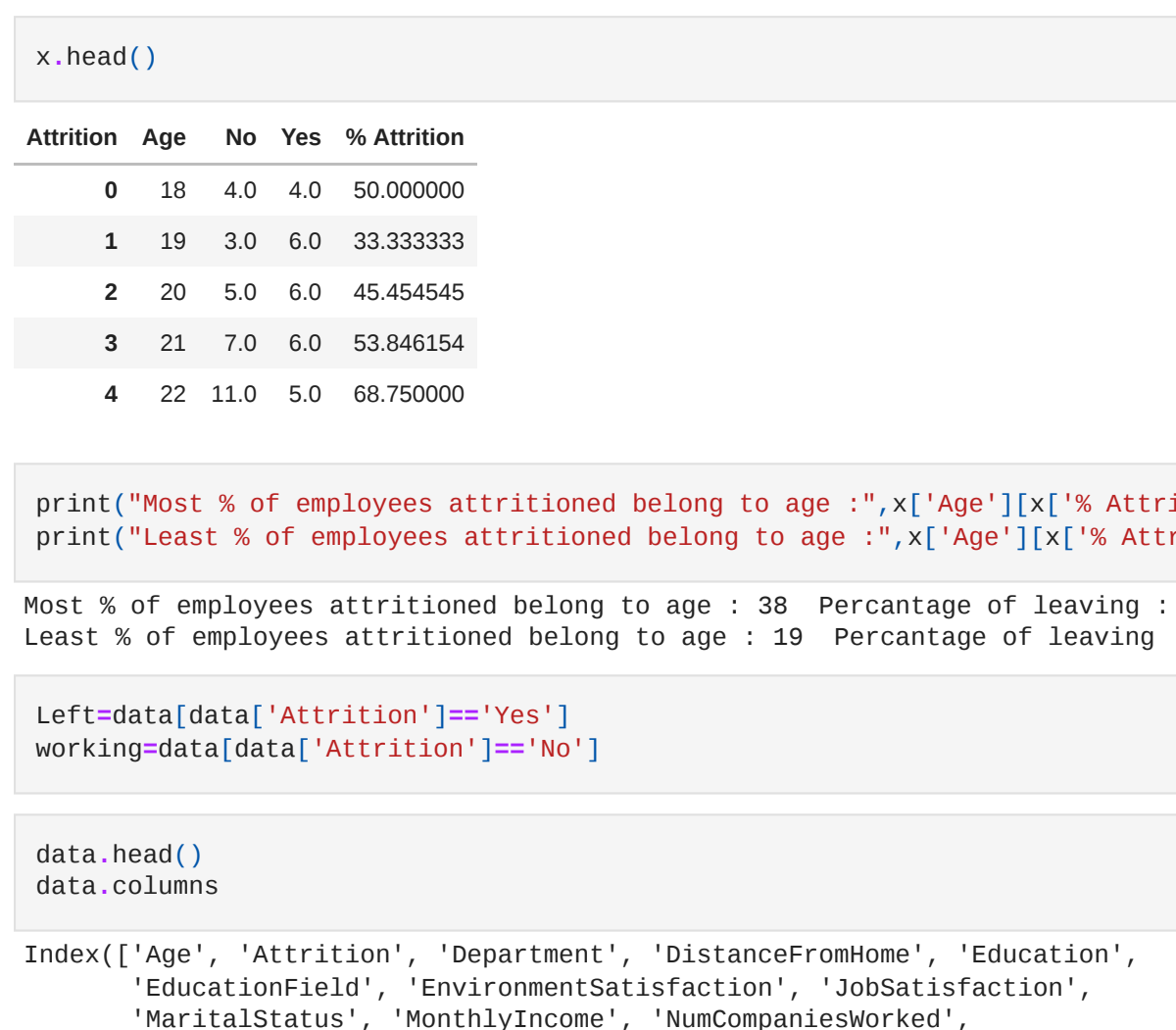
```
In [8]: data['Attrition'].value_counts()

Out[8]:
No    1225
Yes    237
Name: Attrition, dtype: int64
```

```
In [9]: data.groupby(['Age','Attrition']).size().unstack().head()
```

Age	No	Yes
18	4.0	4.0
19	3.0	6.0
20	5.0	6.0
21	7.0	6.0
22	11.0	5.0

```
In [10]: data.groupby(['Age','Attrition']).size().unstack().plot.bar(figsize=(10,6),alpha=0.7)
plt.grid()
```



```
In [11]: x=data.groupby(['Age','Attrition']).size().unstack().reset_index()
for i in x['Age']:
    x['Attrition']=[(x['No']*100)/(x['No']+x['Yes'])]
```

```
In [12]: x.head()
```

Attrition	Age	No	Yes	% Attrition
0	18	4.0	4.0	50.000000
1	19	3.0	6.0	33.333333
2	20	5.0	6.0	45.454545
3	21	7.0	6.0	53.846154
4	22	11.0	5.0	68.750000

```
In [13]: print("Most % of employees attritioned belong to age :",x['Age'][x['% Attrition'].idxmax()],"% Percentage of leaving :",x['% Attrition'].max())
print("Least % of employees attritioned belong to age :",x['Age'][x['% Attrition'].idxmin()],"% Percentage of leaving :",x['% Attrition'].min())

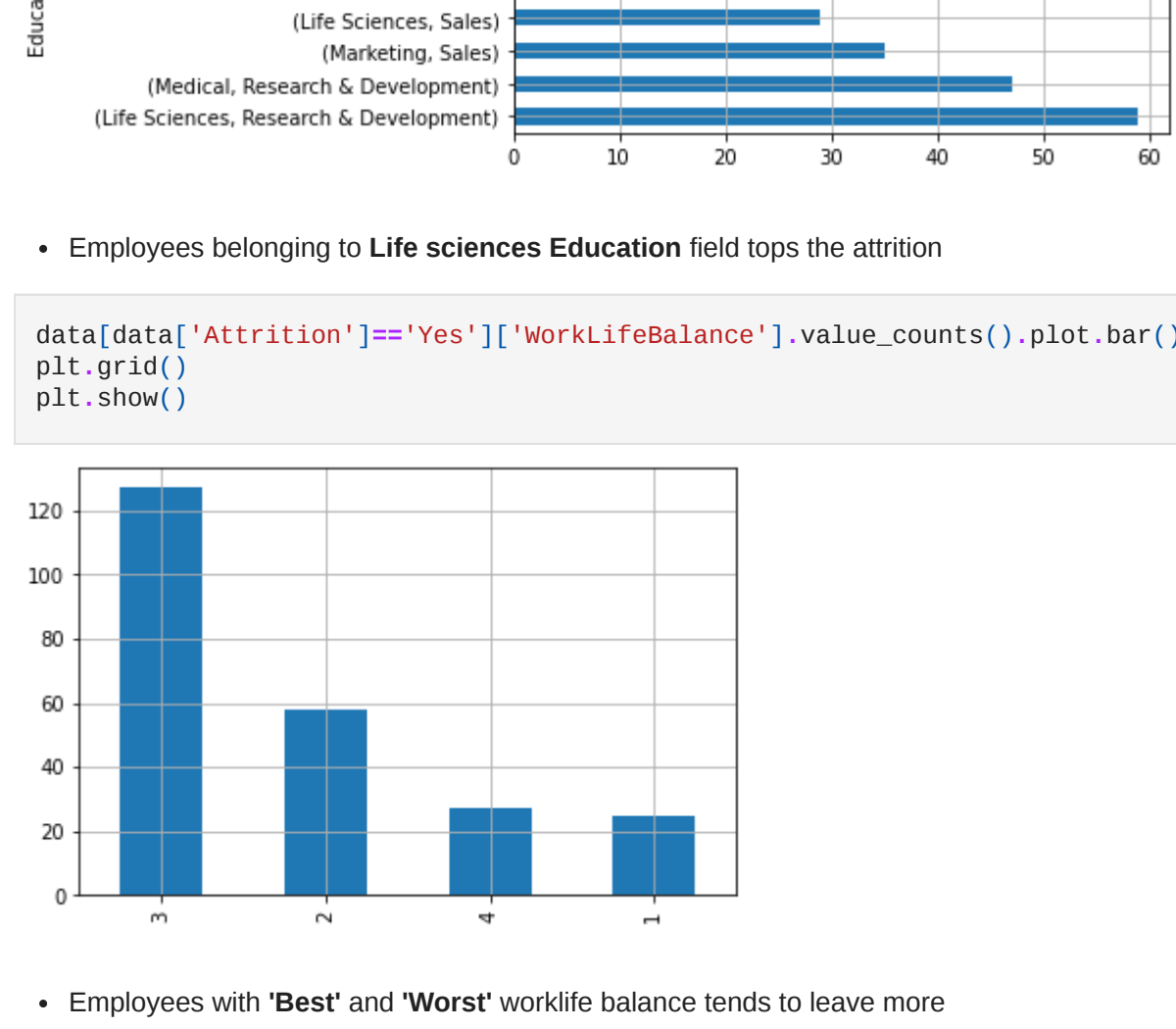
Most % of employees attritioned belong to age : 38 Percentage of leaving : 96.55172413793103
Least % of employees attritioned belong to age : 19 Percentage of leaving : 33.33333333333336
```

```
In [14]: Left=data[data['Attrition']=='Yes']
working=data[data['Attrition']=='No']

In [15]: data.head()

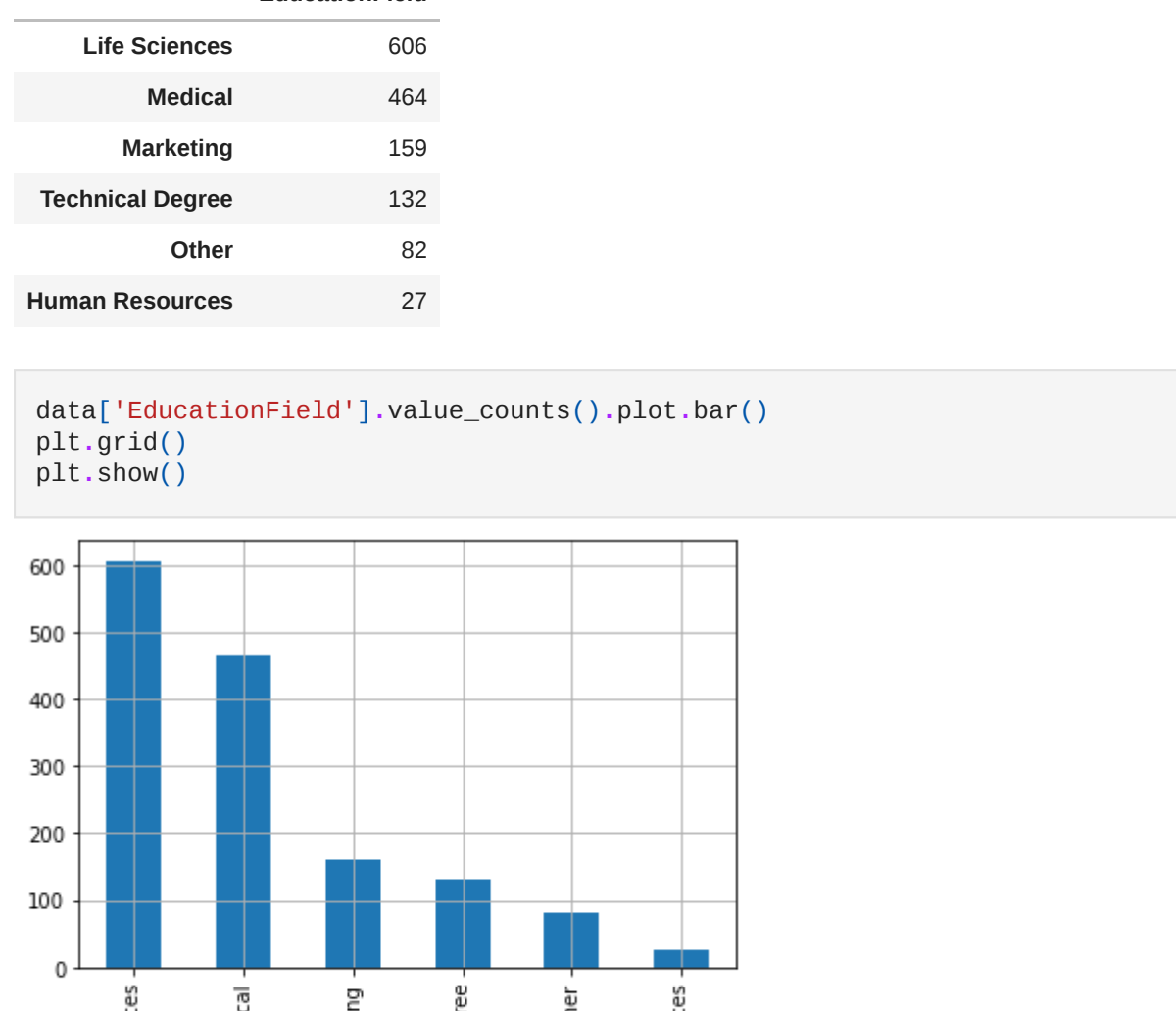
Out[16]: Index(['Age', 'Attrition', 'Department', 'DistanceFromHome', 'Education',
      'EducationField', 'EnvironmentSatisfaction', 'JobSatisfaction',
      'MaritalStatus', 'MonthlyIncome', 'NumCompaniesWorked',
      'WorkLifeBalance', 'YearsAtCompany'],
      dtype='object')
```

```
In [16]: data[data['Attrition']=='Yes'][['MaritalStatus']].value_counts().plot.bar(xlabel='Marital Status',ylabel='Count',title='Attrition vs Marital Status')
plt.grid()
plt.show()
```




- Single tops the number of employees leaving the company

```
In [17]: data[data['Attrition']=='Yes'][['EducationField','Department']].value_counts().plot.barh()
plt.grid()
plt.show()
```



- Employees belonging to Life sciences Education field tops the attrition

```
In [18]: data[data['Attrition']=='Yes'][['WorkLifeBalance']].value_counts().plot.bar()
plt.grid()
plt.show()
```



- Employees with 'Best' and 'Worst' worklife balance tends to leave more

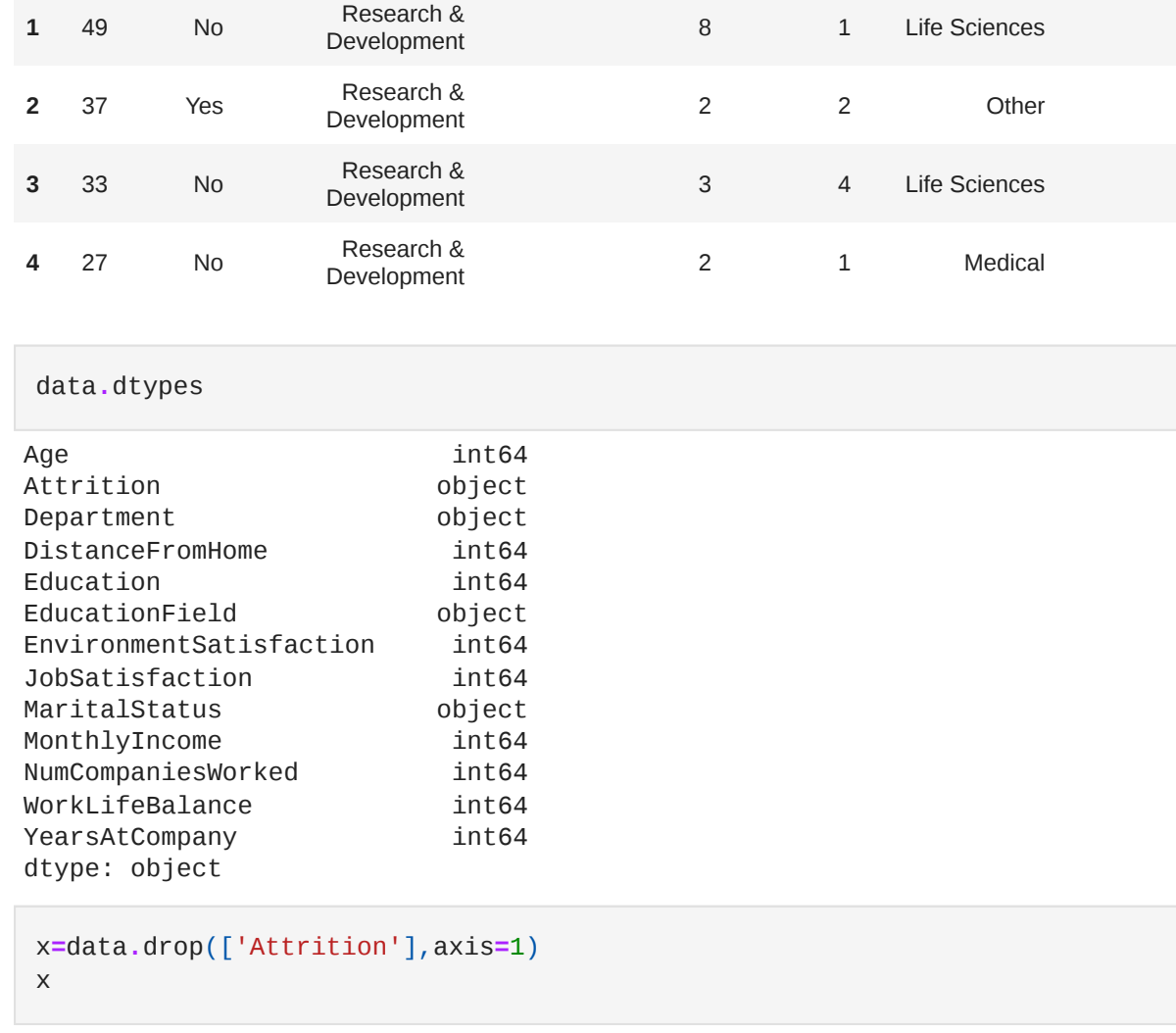
### Distribution of employees by Education Field

```
In [19]: pd.DataFrame(data['EducationField'].value_counts())

Out[19]:
```

EducationField	Count
Life Sciences	606
Medical	464
Marketing	159
Technical Degree	132
Other	82
Human Resources	27

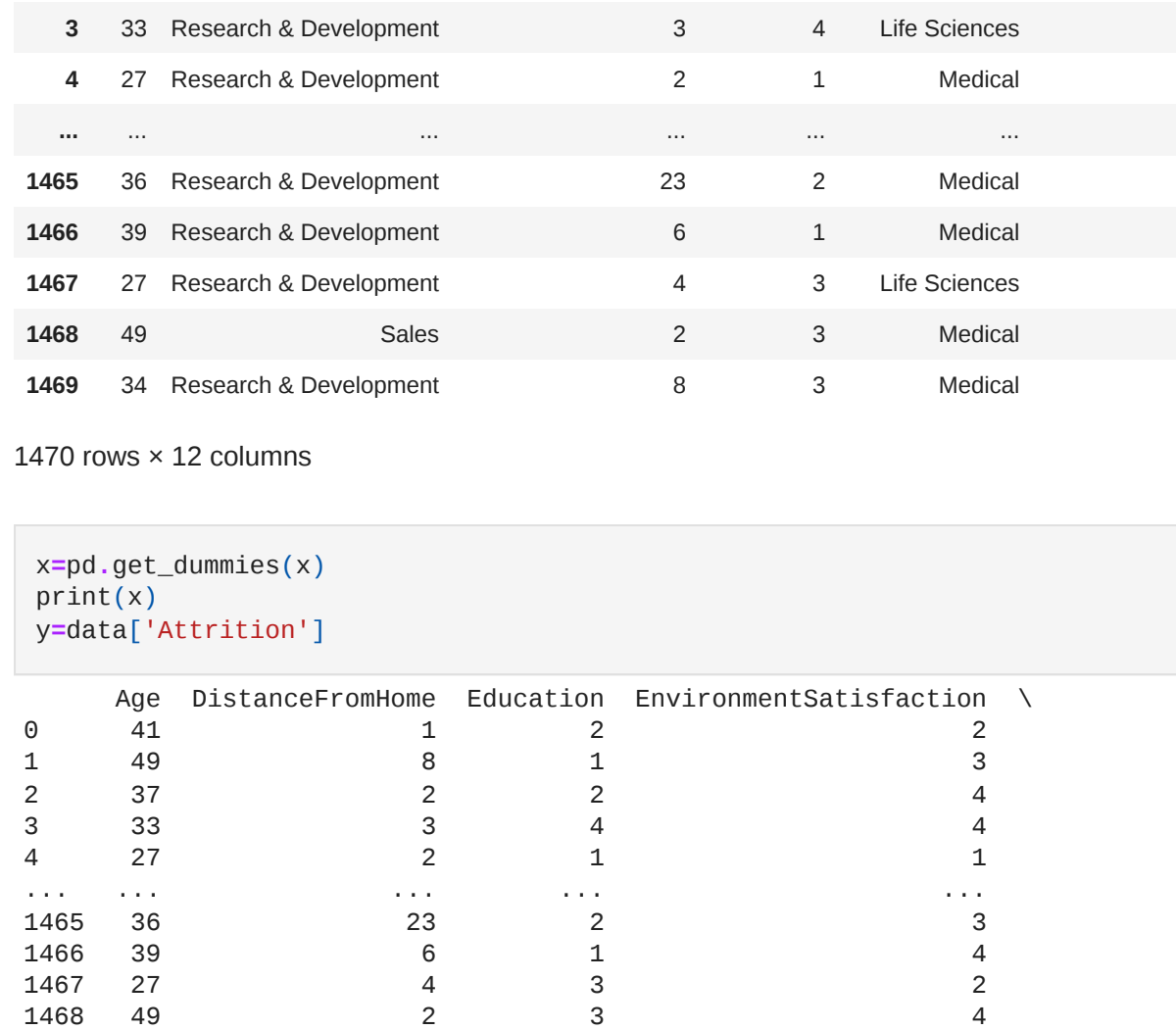
```
In [20]: data['EducationField'].value_counts().plot.bar()
plt.grid()
plt.show()
```



```
In [21]: print("Max Number of employees belong to :",data['EducationField'].value_counts().idxmax(),":",data['EducationField'].value_counts().max())

Max Number of employees belong to : Life Sciences : 606
```

```
In [22]: data['MaritalStatus'].value_counts().plot.bar()
plt.grid()
plt.show()
```



- Most of the employees are : Married

### Q. Build up a logistic regression model to predict which employees are likely to attrite.

```
In [23]: data.head()
```

	Age	Attrition	Department	DistanceFromHome	Education	EducationField	EnvironmentSatisfaction	JobSatisfaction	MaritalStatus	MonthlyIncome	NumCompaniesWorked	WorkLifeBalance	YearsAtCompany
0	41	Yes	Sales	1	2	Life Sciences	2	4	Single	5993	8	1	6
1	49	No	Research & Development	8	1	Life Sciences	3	2	Married	5130	1	3	10
2	37	Yes	Research & Development	2	2	Other	4	3	Single	2090	6	3	0
3	33	No	Research & Development	3	4	Life Sciences	4	3	Married	2909	1	3	8
4	27	No	Research & Development	2	1	Medical	1	2	Married	3468	9	3	2

```
In [24]: data.dtypes

Out[24]: Age                int64
Attrition              object
Department             object
DistanceFromHome       int64
Education              int64
EducationField         object
EnvironmentSatisfaction int64
JobSatisfaction         int64
MaritalStatus         object
MonthlyIncome         int64
NumCompaniesWorked    int64
WorkLifeBalance       int64
YearsAtCompany        int64
dtypes: object
```

```
In [25]: x=data.drop(['Attrition'],axis=1)
x
```

```
Out[25]:
```

	Age	Department	DistanceFromHome	Education</
--	-----	------------	------------------	-------------