

```
2306 AML117 MAHESH BABU M ASSIGNMENT4
```

```
In [1]: import pandas as pd
import numpy as np
```

```
In [2]: features = ["Age", "Workclass", "fnlwgt", "Education", "Education-Num", "Marital Status", "Occupation", "Relationship", "Race", "Sex", "Capital Gain", "Capital Loss", "Hours per week", "Country", "Target"]
df = pd.read_csv('adult.data', names=features)
df
```

Table with 15 columns: Age, Workclass, fnlwgt, Education, Education-Num, Marital Status, Occupation, Relationship, Race, Sex, Capital Gain, Capital Loss, Hours per week, Country, Target. Rows include data for various individuals like '0 39 State-gov 77516 Bachelors 13 Never-married Adm-clerical Not-in-family White Male 2174 0 40 United-States <=50K'.

32561 rows x 15 columns

https://rstudio-pubs-static.s3.amazonaws.com/538563_85cb2b4cd06b4dc48d33de73fa97a297.html

https://archive.ics.uci.edu/dataset/2/adult

```
In [3]: new_df=df
new_df.head()
```

Table with 15 columns: Age, Workclass, fnlwgt, Education, Education-Num, Marital Status, Occupation, Relationship, Race, Sex, Capital Gain, Capital Loss, Hours per week, Country, Target. Rows include data for various individuals like '0 39 State-gov 77516 Bachelors 13 Never-married Adm-clerical Not-in-family White Male 2174 0 40 United-States <=50K'.

```
In [4]: new_df.columns
```

```
Out[4]: Index(['Age', 'Workclass', 'fnlwgt', 'Education', 'Education-Num', 'Marital Status', 'Occupation', 'Relationship', 'Race', 'Sex', 'Capital Gain', 'Capital Loss', 'Hours per week', 'Country', 'Target'], dtype='object')
```

```
In [5]: new_df['Sex'].value_counts()
Out[5]: Male 21790
Female 10771
Name: Sex, dtype: int64
```

Q1 ; ANS; data includes males 21790 and female 10771

```
In [6]: new_df.pivot_table(['Age'], ['Sex'])
```

Table with 2 columns: Sex, Age. Rows: Female 36.858230, Male 39.433547

```
In [7]: group_data = new_df.groupby('Sex')
```

```
In [8]: group_data.describe()
# Provide the mean for each numeric column by sex
group_data.mean(numeric_only=True)
```

Table with 7 columns: Sex, Age, fnlwgt, Education-Num, Capital Gain, Capital Loss, Hours per week. Rows: Female 36.858230 185746.311206 10.035744 568.410547 61.187633 36.410361, Male 39.433547 191771.449013 10.102891 1329.370078 100.213309 42.428086

Q2 average age of female is 36.85

```
In [58]: new_df.columns
Out[58]: Index(['Age', 'Workclass', 'fnlwgt', 'Education', 'Education-Num', 'Marital Status', 'Occupation', 'Relationship', 'Race', 'Sex', 'Capital Gain', 'Capital Loss', 'Hours per week', 'Country', 'Target'], dtype='object')
```

```
In [59]: new_df['Country'].value_counts()
Out[59]: United-States 29170
Mexico 643
? 583
Philippines 198
Germany 137
Canada 121
Puerto-Rico 114
El-Salvador 106
India 100
Cuba 95
England 90
Jamaica 81
South 80
China 75
Italy 73
Dominican-Republic 70
Vietnam 67
Guatemala 64
Japan 62
Poland 60
Columbia 59
Taiwan 51
Haiti 44
Iran 43
Portugal 37
Nicaragua 34
Peru 31
France 29
Greece 29
Ecuador 28
Ireland 24
Hong 20
Cambodia 19
Trinidad&Tobago 19
Lads 18
Thailand 18
Yugoslavia 16
Outlying-US(Guam-USVI-etc) 14
Honduras 13
Hungary 13
Scotland 12
Holand-Netherlands 1
Name: Country, dtype: int64
Germany=137 ratio of total to german is
```

```
In [11]: 137/32561
Out[11]: 0.004207487485028101
```

Q3 0.004207487485028101 is proportion of german citizens

```
In [12]: new_df.pivot_table(['Age'], ['Target'])
Out[12]: Target
Age
<=50K 36.783738
>50K 44.249841
```

```
In [13]: salary=new_df.groupby(['Target'])
In [14]: salary.describe()
Out[14]: Target
Age fnlwgt ... Capital Loss Hours per week
count mean std min 25% 50% 75% max count mean ... 75% max count mean std min 25% 50% 75% max
Target
<=50K 24720.0 36.783738 14.020088 17.0 25.0 34.0 46.0 90.0 24720.0 190340.86517 ... 0.0 4356.0 24720.0 38.840210 12.318995 1.0 35.0 40.0 40.0 99.0
>50K 7841.0 44.249841 10.519028 19.0 36.0 44.0 51.0 90.0 7841.0 188005.00000 ... 0.0 3683.0 7841.0 45.473026 11.012971 1.0 40.0 40.0 50.0 99.0
2 rows x 48 columns
```

Q 4-5 THE MEAN =36.78, sd=14 for who's salary is <= 50K, mean=44.2 and sd=10.5 for whose salary is >= 50K

```
In [15]: education=new_df.groupby('Education')
In [16]: education.describe()
Out[16]: Target
Age fnlwgt ... Capital Loss Hours per week
count mean std min 25% 50% 75% max count mean ... 75% max count mean std min 25% 50% 75% max
Education
10th 933.0 37.429796 16.720713 17.0 22.0 34.0 52.0 90.0 933.0 196832.465166 ... 0.0 3770.0 933.0 37.052519 13.788112 1.0 30.0 40.0 40.0 99.0
11th 1175.0 32.355745 15.545485 17.0 18.0 28.0 43.0 90.0 1175.0 194928.077447 ... 0.0 2824.0 1175.0 33.925957 13.965416 2.0 20.0 40.0 40.0 99.0
12th 433.0 32.000000 14.334625 17.0 19.0 28.0 41.0 79.0 433.0 199097.508083 ... 0.0 2258.0 433.0 35.780600 12.626412 6.0 30.0 40.0 40.0 99.0
1st-4th 168.0 46.142857 15.615625 19.0 33.0 46.0 57.0 90.0 168.0 239303.000000 ... 0.0 2603.0 168.0 38.255952 12.848727 4.0 35.0 40.0 40.0 96.0
5th-6th 333.0 42.885886 15.557285 17.0 29.0 42.0 54.0 84.0 333.0 232448.333333 ... 0.0 2603.0 333.0 38.897898 10.551727 3.0 40.0 40.0 40.0 84.0
7th-8th 646.0 48.445820 16.092350 17.0 34.25 50.0 61.0 90.0 646.0 188079.171827 ... 0.0 3900.0 646.0 39.366873 14.201870 2.0 35.0 40.0 40.0 99.0
9th 514.0 41.060311 15.948682 17.0 28.0 39.0 54.0 90.0 514.0 202485.066148 ... 0.0 2231.0 514.0 38.044747 11.064402 1.0 36.0 40.0 40.0 99.0
Assoc-acdm 1067.0 37.381443 11.095177 19.0 29.0 36.0 44.0 90.0 1067.0 193424.093721 ... 0.0 2824.0 1067.0 40.504217 12.196666 1.0 40.0 40.0 45.0 99.0
Assoc-voc 1382.0 38.553546 11.631300 19.0 30.0 37.0 46.0 84.0 1382.0 181936.016643 ... 0.0 2603.0 1382.0 41.610709 10.793384 1.0 40.0 40.0 45.0 99.0
Bachelors 5355.0 38.904949 11.912210 19.0 29.0 37.0 46.0 90.0 5355.0 188055.914846 ... 0.0 2824.0 5355.0 42.614006 11.446185 2.0 40.0 40.0 50.0 99.0
Doctorate 413.0 47.702179 11.784716 24.0 39.0 47.0 55.0 80.0 413.0 186698.760291 ... 0.0 3683.0 413.0 46.973366 15.084447 1.0 40.0 45.0 55.0 99.0
HS-grad 10501.0 38.974479 13.541524 17.0 28.0 37.0 48.0 90.0 10501.0 189538.739739 ... 0.0 4356.0 10501.0 40.575374 11.333757 1.0 40.0 40.0 42.0 99.0
Masters 1723.0 44.049913 11.068935 18.0 36.0 43.0 51.0 90.0 1723.0 179852.362739 ... 0.0 2824.0 1723.0 43.836332 12.277801 1.0 40.0 40.0 50.0 99.0
Preschool 51.0 42.764706 15.126914 19.0 31.0 41.0 53.5 75.0 51.0 235889.372549 ... 0.0 1719.0 51.0 36.647059 12.555196 10.0 30.0 40.0 40.0 75.0
Prof-school 576.0 44.746528 11.962477 25.0 36.0 43.0 51.0 90.0 576.0 185663.706597 ... 0.0 2824.0 576.0 47.425347 14.806038 2.0 40.0 48.0 55.0 99.0
Some-college 7291.0 35.756275 13.474051 17.0 24.0 34.0 45.0 90.0 7291.0 188742.922330 ... 0.0 4356.0 7291.0 38.852284 12.761901 1.0 35.0 40.0 43.0 99.0
16 rows x 48 columns
```

```
In [42]: new_df[new_df['Target'] == '>50K'].value_counts()
Out[42]: Series([], dtype: int64)
```

```
In [63]: data=new_df['Education'].value_counts()
print(data)
```

Table with 2 columns: Education, count. Rows: HS-grad 10501, Some-college 7291, Bachelors 5355, Masters 1723, Assoc-voc 1382, 11th 1175, Assoc-acdm 1067, 10th 933, 7th-8th 646, Prof-school 514, 9th 514, 12th 433, Doctorate 413, 5th-6th 333, 1st-4th 168, Preschool 51, Name: Education, dtype: int64

```
In [68]: ntr=new_df.pivot_table(index='Education', columns='Target')
print(ntr)
```

Table with 5 columns: Education, Target, Age, Capital Gain, Capital Loss. Rows: 10th, 11th, 12th, 1st-4th, 5th-6th, 7th-8th, 9th, Assoc-acdm, Assoc-voc, Bachelors, Doctorate, HS-grad, Masters, Preschool, Prof-school, Some-college

Table with 5 columns: Target, Education, Age, Capital Gain, Capital Loss. Rows: 10th, 11th, 12th, 1st-4th, 5th-6th, 7th-8th, 9th, Assoc-acdm, Assoc-voc, Bachelors, Doctorate, HS-grad, Masters, Preschool, Prof-school, Some-college

Table with 5 columns: Target, Education, fnlwgt, Hours per week. Rows: 10th, 11th, 12th, 1st-4th, 5th-6th, 7th-8th, 9th, Assoc-acdm, Assoc-voc, Bachelors, Doctorate, HS-grad, Masters, Preschool, Prof-school, Some-college

```
In [69]: ntr.loc['Assoc-acdm':'Some-college']
Out[69]: Target
Age Capital Gain Capital Loss Education-Num Hours per week fnlwgt
Target <=50K >50K <=50K >50K <=50K >50K <=50K >50K <=50K >50K
Education
```

Table with 10 columns: Target, Age, Capital Gain, Capital Loss, Education-Num, Hours per week, fnlwgt. Rows: count, mean, std, 25%, 50%, 75%, max for <=50K and >50K groups.

Q6

it is not true that who are education less than high school are taking the more than 50k . at least high school criteria is not good enough to take >50k

```
In [40]: Empty DataFrame
Columns: [Age, Workclass, fnlwgt, Education, Education-Num, Marital Status, Occupation, Relationship, Race, Sex, Capital Gain, Capital Loss, Hours per week, Country, Target]
Index: []
```

```
In [ ]:
```