

Assignment 01

Questions 01 and Answers

1. List all professions and products in the dataset.

```
1 # Import CSV file check the column and data inside of it
```

```
In [1]: 1 #import all the needed packages
        2 import pandas as pd
        3 import numpy as np
        4 import matplotlib.pyplot as plt
        5 import re
        6 import csv
```

```
In [2]: 1 #import the CSV file into the script to analyze
        2 insurance_data = pd.read_csv("Insurance_Data_and_Description_CW1 (1).csv")
        3
```

```
In [3]: 1 # check CSV file with its's initial file colum and data
        2 insurance_data
```

Out[3]:

	ID	Product	Region	Profession	Age	First Rate Paid	Contract sum	Closing Date
0	1	Life	150	Employee	26.0	Y	1,000	Jul-90
1	2	Life	150	Farmer	31.0	N	2,000	Jul-90
2	3	Life	150	Farmer	33.0	Y	2,000	Jul-90
3	4	Life	150	Farmer	33.0	Y	2,000	Jul-90
4	5	Life	120	Self-employed	34.0	Y	2,000	Jul-90

```
1 # Answering Assignment questions
```

```
1 # A) List all professions and products
```

```
In [4]: 1 # a) List all professions and products in the dataset
        2 insurance_product_and_profession = insurance_data[["Product", "Profession"]]
        3 insurance_product_and_profession
```

Out[4]:

	Product	Profession
0	Life	Employee
1	Life	Farmer
2	Life	Farmer
3	Life	Farmer
4	Life	Self-employed
...
14840	Household	Civil servant
14841	Life	Civil servant
14842	Health	Civil servant
14843	Health	Civil servant
14844	Household	Civil servant

Import all the libraries on the top and develop the script to get output for product and profession.

2. Show any problem with the data etc.

The main problem in this data set is Wrong data types and N/A, nan values – which mean Null values. The solutions are provided, then transformed into to correct data type also removed Null values from the table.

```
1 # B) Show problems & Handling problems

In [5]: 1 # Show any problem with the data etc.
2 # 1. check whether the csv file have any null values
3 missing_value = ["N/a", "na", np.nan]
4 insurance_data = pd.read_csv("Insurance_Data_and_Description_CW1 (1).csv", na_values = missing_value)
5 insurance_data.isnull().sum()
6 insurance_data = insurance_data.dropna()
7 insurance_data.isnull().sum()

Out[5]: ID          0
Product          0
Region           0
Profession       0
Age              0
First Rate Paid  0
Contract sum     0
Closing Date     0
dtype: int64
```

3. Clean and transform the data into the correct data types.

The abovementioned problem was sorted here.

```
1 # C) Clean and transform

In [8]: 1

In [6]: 1 #As this result there are so many Null values in both of the columns
2 #In order to clean these problem removing all the null values from dataset
3
4 insurance_data['Contract sum'] = insurance_data['Contract sum'].str.replace(',', '0')
5 insurance_data['Contract sum'] = insurance_data['Contract sum'].str.replace('K', '0')
6 insurance_data['Contract sum'] = insurance_data['Contract sum'].str.replace('k', '0')
7
8
9 insurance_data['Contract sum'] = insurance_data['Contract sum'].replace(',', '0')
10 insurance_data['Contract sum'] = insurance_data['Contract sum'].astype(int)
11
12 cleaned_insurance_data = insurance_data.dropna()
13
14 cleaned_insurance_data
```

4. Add a new column “Age Group” and fill up its data according to below criteria.

here I am using between and comparison operators to check the age and created the column for Age group.

D) Add a new column “Age Group”

```
In [8]: 1 cleaned_insurance_data.loc[cleaned_insurance_data['Age']<=16, 'Age Group'] = 'AG-1'
2 cleaned_insurance_data.loc[cleaned_insurance_data['Age'].between(17,24), 'Age Group'] = 'AG-2'
3 cleaned_insurance_data.loc[cleaned_insurance_data['Age'].between(25,29), 'Age Group'] = 'AG-3'
4 cleaned_insurance_data.loc[cleaned_insurance_data['Age'].between(30,39), 'Age Group'] = 'AG-4'
5 cleaned_insurance_data.loc[cleaned_insurance_data['Age'].between(40,49), 'Age Group'] = 'AG-5'
6 cleaned_insurance_data.loc[cleaned_insurance_data['Age'].between(50,59), 'Age Group'] = 'AG-6'
7 cleaned_insurance_data.loc[cleaned_insurance_data['Age']>=60, 'Age Group'] = 'AG-7'
8
9 # check the edited file
10 cleaned_insurance_data[["Age", "Age Group"]]
```

Out[8]:

	Age	Age Group
0	26.0	AG-4
1	31.0	AG-4
2	33.0	AG-4
3	33.0	AG-4
4	34.0	AG-4
...
14840	78.0	AG-7
14841	78.0	AG-7
14842	82.0	AG-7

5. Rank region by total insurance product.

using the count function to for region and group by the to identify the total insurance product for each Region.

1 # E) Rank region by total insurance product

```
In [9]: 1 Rank = cleaned_insurance_data.groupby('Region', sort=False).count()
2 Rank['Product']
```

Out[9]: Region
150 12281
120 16
160 2265
999 195
130 36
Name: Product, dtype: int64

- Find the mean, median, standard deviation of Age and Contract Sum.

Finding the summary of Age group and Contract sum in the dataset, which is showing up mean, median etc.

```
1 # F) Find the mean, median, standard deviation of Age and Contract Sum

In [10]: 1 # 1. For Age
        2 cleaned_insurance_data["Age"].describe()

Out[10]: count    14793.00000
        mean      44.93578
        std       14.66991
        min        1.00000
        25%       34.00000
        50%       45.00000
        75%       57.00000
        max       95.00000
        Name: Age, dtype: float64

In [11]: 1 # 2. For Contract Sum
        2 insurance_data['Contract sum'].astype(str).astype(int)
        3 insurance_data["Contract sum"].describe()

Out[11]: count    1.479300e+04
        mean     2.892438e+05
        std      2.584105e+05
        min      4.000000e+01
        25%      1.500000e+05
        50%      2.000000e+05
        75%      3.300000e+05
        max      4.250000e+06
        Name: Contract sum, dtype: float64
```

- Filter the rows of Age < 25 and Contract Sum between 50k and 100k and insurance type of Health or Life insurance.

Filtering the data with above condition, I used single line comparison to done the filtering.

```
1 # G) filtering age, contract sum and insurance type

In [13]: 1
        2 # selecting rows based on condition
        3 filterd_data = cleaned_insurance_data[(cleaned_insurance_data['Age'] < 25) & ((cleaned_insurance_data['Product'] == 'Life')
        4 filterd_data

Out[13]:
```

	ID	Product	Region	Profession	Age	First Rate Paid	Contract sum	Closing Date	Age Group	
	29	30	Health	150	Student	21.0	Y	50000	Aug-90	AG-4
	196	197	Health	150	Student	16.0	N	50000	Sep-90	AG-1
	1100	1101	Life	150	Executive employee	7.0	Y	100000	Nov-90	AG-1
	1164	1165	Health	150	Student	23.0	N	100000	Dec-90	AG-4
	1165	1166	Life	160	Student	24.0	Y	100000	Dec-90	AG-4

	14506	14507	Health	150	Student	23.0	N	1000000	Dec-96	AG-4
	14507	14508	Health	150	Student	23.0	N	1000000	Dec-96	AG-4
	14509	14510	Health	160	Farmer	24.0	N	1000000	Dec-96	AG-4
	14510	14511	Health	150	Student	24.0	N	1000000	Dec-96	AG-4
	14511	14512	Health	150	Student	24.0	N	1000000	Dec-96	AG-4

8. Find total insurance product by profession.

Finding product over the profession, by grouping profession and calculate product in the professions.

Find total insurance product by profession

```
In [13]: 1 total_insurance_product = cleaned_insurance_data.groupby('Profession').count()
         2 total_insurance_product['Product']
```

```
Out[13]: Profession
Civil servant    2044
Employee        5602
Executive employee  575
Farmer          3686
Pupil           190
Self-employed    20
Student         2552
Worker          124
Name: Product, dtype: int64
```

9. Find correlations between product, profession, and contract sum. What are your conclusions?

The correlation is can only calculate with numerical values, so other variable could not do for the correlation.

Find correlations between product, profession and contract sum

```
In [14]: 1 dataset = cleaned_insurance_data[['Product', 'Profession', 'Contract sum']]
         2
         3 print( dataset.corr())
         4
```

```
Contract sum
Contract sum    1.0
```

10. Calculate and Visualize products on age group. Describe your findings

Calculate and Visualise products on age group

```
In [24]: 1 x = []
2 y = []
3
4 visualize_data = total_insurance_product = cleaned_insurance_data.groupby('Profession').count()
5 visualize_data['Product'] = visualize_data['Product'].astype(int)
6
7
8 new_data = visualize_data[['Product', 'Age Group']]
9 new_data = new_data.to_csv('newfile.csv')
10
11 with open('newfile.csv', 'r') as csvfile:
12     plots = csv.reader(csvfile, delimiter = ',')
13
14     for row in plots:
15         x.append(int(row[0]))
16         y.append(row[1])
17
18
19 plt.bar(x, y, color = 'g', width = 0.72, label = "Age")
20 plt.xlabel('Products')
21 plt.ylabel('Age group')
22 plt.title('Ages of different product')
23 plt.legend()
24 plt.show()
25
26 visualize_data
```