### TASK 1

### • Installing Pandas: Python Data Analysis Library

```
Requirement already satisfied: pandas in c:\users\sanaj\appdata\local\programs\python\python312\lib\site-packages (2.2.3)
Requirement already satisfied: numpy>=1.26.0 in c:\users\sanaj\appdata\local\programs\python\python312\lib\site-packages (from pandas) (2.2.1)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\users\sanaj\appdata\local\programs\python\python312\lib\site-packages (from pandas) (2.9.0.po st0)
Requirement already satisfied: pytz>=2020.1 in c:\users\sanaj\appdata\local\programs\python\python312\lib\site-packages (from pandas) (2024.2)
Requirement already satisfied: tzdata>=2022.7 in c:\users\sanaj\appdata\local\programs\python\python312\lib\site-packages (from pandas) (2024.2)
Requirement already satisfied: six>=1.5 in c:\users\sanaj\appdata\local\programs\python\python312\lib\site-packages (from python-dateutil>=2.8.2->pandas)
(1.16.0)

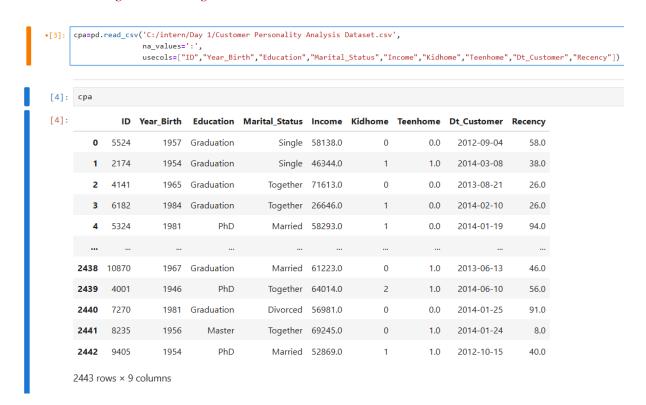
[notice] A new release of pip is available: 24.0 -> 25.0.1
[notice] To update, run: python.exe -m pip install --upgrade pip
```

#### Importing the pandas module

Pandas5 provides high-performance data structures and data analysis tools. The main feature of Pandas is a fast and efficient Data Frame object for data manipulation with integrated indexing. The structure of Data Frame can be seen as a spreadsheet which offers very flexible ways of working with it. Pandas offers features like adding or removing rows and columns and reshaping it the way you want. It also provides many high- performance functions such as aggregation , merging , and joining datasets. Pandas tools for exporting and importing data from different formats file like .CSV , text files , Microsoft Excel , SQL database etc.

[2]: import pandas as pd

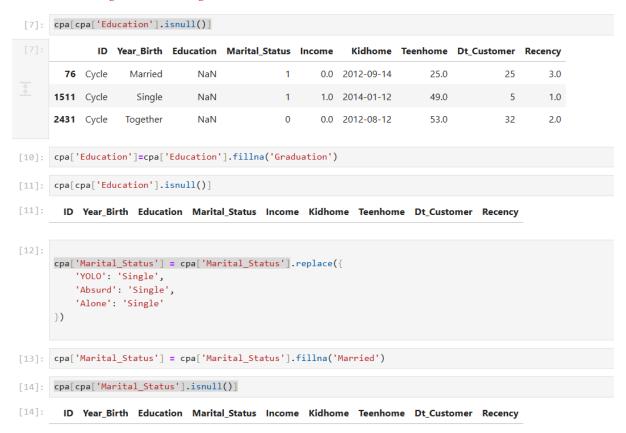
### • Loading and Reading the dataset



First of all, I have created a new notebook as I want. Here, I have created a notebook with the name 'intern' then I have downloaded the data file and stored it in the same directory to retrieve the same. The above code has to be written to read the .csv file in pandas is by using read\_csv method with the file path, na\_value parameter is used to tell pandas they consider ":" as NaN value print NaN at the place of ":" and the other method which is usecols are used to tell the pandas that which cols are we going to use.

The output of the execution shows that ad DataFrame size is 2443 rows x 9 columns.

### • Filtering and handling the null values



: [	сра									
]:		ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome	Dt_Customer	Recency
	0	5524	1957	Graduation	Single	58138.0	0	0.0	2012-09-04	58.0
	1	2174	1954	Graduation	Single	46344.0	1	1.0	2014-03-08	38.0
	2	4141	1965	Graduation	Together	71613.0	0	0.0	2013-08-21	26.0
	3	6182	1984	Graduation	Together	26646.0	1	0.0	2014-02-10	26.0
	4	5324	1981	PhD	Married	58293.0	1	0.0	2014-01-19	94.0
	2438	10870	1967	Graduation	Married	61223.0	0	1.0	2013-06-13	46.0
	2439	4001	1946	PhD	Together	64014.0	2	1.0	2014-06-10	56.0
	2440	7270	1981	Graduation	Divorced	56981.0	0	0.0	2014-01-25	91.0
	2441	8235	1956	Master	Together	69245.0	0	1.0	2014-01-24	8.0
	2442	9405	1954	PhD	Married	52869.0	1	1.0	2012-10-15	40.0
					a(cpa['Income'			SettingWith	o ConvWarning.	
: (	C:\Use A valu Try us See th cpa[	ers\sana e is tr ing .lo e cavea 'Income	nj\AppData\ rying to be oc[row_inde  ots in the '] =cpa['I	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil	a(cpa['Income' ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Incom	\1220748 from a l instead	429.py:1: DataFrame. ata.org/pa	-		uide/inde
: (	C:\Use A valu Try us See th cpa[ cpa[cr	ers\sana e is tr ing .lo e cavea 'Income	nj\AppData\ rying to be oc[row_inde  ots in the ''] =cpa['I  ome'].isnul	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil	ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Incom	\1220748 from a l instead ndas.pyd e'].mean	429.py:1: DataFrame. ata.org/pa	ndas-docs/s	stable/user_go	uide/inde
	C:\Use A valu Try us See th cpa[ cpa[cr	ers\sana e is tr ing .lo e cavea 'Income	nj\AppData\ rying to be oc[row_inde  ots in the ''] =cpa['I  ome'].isnul	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil	ipykernel_3380 opy of a slice exer] = value on: https://pa	\1220748 from a l instead ndas.pyd e'].mean	429.py:1: DataFrame. ata.org/pa	ndas-docs/s	stable/user_go	uide/inde
	C:\Use A valu Try us See th cpa[ cpa[cr	ers\sana e is tr ing .lo e cavea 'Income a['Inco	nj\AppData\ rying to be oc[row_inde  ots in the ''] =cpa['I  ome'].isnul	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil  1()]  n Marital_S	ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Incom	\1220748 from a l instead ndas.pyd e'].mean	429.py:1: DataFrame. ata.org/pa	ndas-docs/s	stable/user_go	uide/inde
	C:\Use A valu Try us See th cpa[cpa[cp	ers\sanaae is tr ing .lo ee cavea 'Income aa['Inco	ry/AppData\ rying to be oc[row_inde  ots in the ''] =cpa['I  ome'].isnul  th Educatio  whome'].isn	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil  1()] on Marital_S ull()]	ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Incom	\12207480 from a linstead ndas.pyda e'].mean	429.py:1: DataFrame. ata.org/pa ())  Teenhome	ndas-docs/s e Dt_Custor	stable/user_go	uide/inde
	C:\Use A valu Try us See th cpa[cp ID cpa[cp	ers\sanaae is tr ing .lo ee cavea 'Income aa['Inco Year_Birt	ry/AppData\ rying to be oc[row_inde  ots in the ''] =cpa['I  ome'].isnul  th Educatio  whome'].isn	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil  1()] on Marital_S ull()] on Marital_S	<pre>ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Income status Income</pre>	\12207480 from a linstead ndas.pyda e'].mean	429.py:1: DataFrame. ata.org/pa ())  Teenhome	ndas-docs/s e Dt_Custor	stable/user_go	uide/inde
	C:\Use A valu Try us See th cpa[cp ID  cpa[cp ID  cpa[cp	ers\sanaae is tr ing .lo ee cavea 'Income aa['Inco Year_Birt Year_Birt	nj\AppData\ rying to be oc[row_inde  ots in the ots in the i'] =cpa['I  ome'].isnul  th Educatio  the Educatio  ome'].isnu	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil  1()] on Marital_S  ull()] on Marital_S	<pre>ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Income status Income</pre>	\1220748- from a linstead  ndas.pyd. e'].mean  Kidhome	429.py:1: DataFrame.  ata.org/pa ())  Teenhomo	ndas-docs/s  e Dt_Custor	mer Recency	uide/inde
	C:\Use A valu Try us See th cpa[cp ID  cpa[cr ID  cpa[cr ID	ers\sanae is triing lose cavea 'Income cavea 'Income ca	nj\AppData\ rying to be oc[row_inde  ots in the o'] =cpa['I  ome'].isnul  th Educatio  the Educatio  the Educatio  the Educatio  the Educatio  the Educatio	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil  1()] on Marital_S  ull()] on Marital_S  11()]	ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Income	\1220748- from a linstead  ndas.pyd. e'].mean  Kidhome	429.py:1: DataFrame.  ata.org/pa ())  Teenhomo	ndas-docs/s  e Dt_Custor	mer Recency	uide/inde
	C:\Use A valu Try us See th cpa[cp ID  cpa[cp ID  cpa[cp ID  cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp	ers\sana e is tr ing .lo e cavea 'Income a['Inco Year_Birt a['Kidh Year_Birt a['Rece	nj\AppData\ rying to be oc[row_inde  ots in the o'] =cpa['I  ome'].isnul  th Educatio  the Educatio  the Educatio  come'].isnu  the Educatio  come'].isnu	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil  1()] on Marital_S  ull()] on Marital_S  11()]	ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Income status Income	\1220748- from a linstead ndas.pyd. e'].mean Kidhome Kidhome	429.py:1: DataFrame.  ata.org/pa ())  Teenhomo	ndas-docs/s  e Dt_Custon  e Dt_Custon  e Dt_Custon	mer Recency mer Recency	uide/inde
	C:\Use A valu Try us See th cpa[cp ID  cpa[cp ID  cpa[cp ID  cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp Cpa[cp	ers\sana e is tr ing .lo e cavea 'Income a['Inco Year_Birt a['Kidh Year_Birt a['Rece	nj\AppData\ rying to be oc[row_inde  ots in the o'] =cpa['I  ome'].isnul  th Educatio  the Educatio  the Educatio  come'].isnu  the Educatio  come'].isnu	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil  1()] on Marital_S  ull()] on Marital_S  11()]	ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Income	\1220748- from a linstead ndas.pyd. e'].mean Kidhome Kidhome	429.py:1: DataFrame.  ata.org/pa ())  Teenhomo	ndas-docs/s  e Dt_Custon  e Dt_Custon  e Dt_Custon	mer Recency mer Recency	uide/inde
	C:\Use A valu Try us See th cpa[cp ID  cpa[cp ID  cpa[cp ID  cpa[cp ID	ers\sanae is triing lo ee cavea 'Income aa['Inco Year_Birt aa['Kidh Year_Birt aa['Rece	nj\AppData\ rying to be oc[row_inde  ots in the o'] =cpa['I  ome'].isnul  th Educatio  the Educatio  the Educatio  come'].isnu  the Educatio  come'].isnu	Local\Temp\ set on a c xer,col_ind documentati ncome'].fil  1()] on Marital_S  ull()] on Marital_S  11()] on Marital_S  11()] on Marital_S	ipykernel_3380 opy of a slice exer] = value on: https://pa lna(cpa['Income status Income	\1220748- from a linstead ndas.pyd. e'].mean Kidhome Kidhome	429.py:1: DataFrame.  ata.org/pa ())  Teenhomo	ndas-docs/s  e Dt_Custon  e Dt_Custon  e Dt_Custon	mer Recency mer Recency	uide/inde

# • Dropping or removing the duplicate values

cpa=c	pa.drop	_duplicate	s()						
сра									
	ID	Year_Birth	Education	Marital_Status	Income	Kidhome	Teenhome	Dt_Customer	Recency
0	5524	1957	Graduation	Single	58138.0	0	0.0	2012-09-04	58.0
1	2174	1954	Graduation	Single	46344.0	1	1.0	2014-03-08	38.0
2	4141	1965	Graduation	Together	71613.0	0	0.0	2013-08-21	26.0
3	6182	1984	Graduation	Together	26646.0	1	0.0	2014-02-10	26.0
4	5324	1981	PhD	Married	58293.0	1	0.0	2014-01-19	94.0
				on Together 26646.0 1 0.0 2014-02-10 26.0 nD Married 58293.0 1 0.0 2014-01-19 94.0					
2438	10870	1967	Graduation	Married	61223.0	0	1.0	2013-06-13	46.0
2439	4001	1946	PhD	Together	64014.0	2	1.0	2014-06-10	56.0
2440	7270	1981	Graduation	Divorced	56981.0	0	0.0	2014-01-25	91.0
2441	8235	1956	Master	Together	69245.0	0	1.0	2014-01-24	8.0
2442	9405	1954	PhD	Married	52869.0	1	1.0	2012-10-15	40.0

2223 rows × 9 columns

# • Checking is there any null value remaining

### • Renaming column headers to be clean and uniform

]: cp	а=ср	oa.rename(co			_id","Year_Bir me","Kidhome":		_			
]: <b>c</b> p	a									
		customer_id	year_birth	education	marital_status	income	kidhome	teenhome	dt_customer	recency
	0	5524	1957	Graduation	Single	58138.0	0	0.0	2012-09-04	58.0
	1	2174	1954	Graduation	Single	46344.0	1	1.0	2014-03-08	38.0
	2	4141	1965	Graduation	Together	71613.0	0	0.0	2013-08-21	26.0
	3	6182	1984	Graduation	Together	26646.0	1	0.0	2014-02-10	26.0
	4	5324	1981	PhD	Married	58293.0	1	0.0	2014-01-19	94.0
24	138	10870	1967	Graduation	Married	61223.0	0	1.0	2013-06-13	46.0
24	139	4001	1946	PhD	Together	64014.0	2	1.0	2014-06-10	56.0
24	40	7270	1981	Graduation	Divorced	56981.0	0	0.0	2014-01-25	91.0
24	41	8235	1956	Master	Together	69245.0	0	1.0	2014-01-24	8.0
24	142	9405	1954	PhD	Married	52869.0	1	1.0	2012-10-15	40.0
222		0 1								

2223 rows × 9 columns

### • Converting data in date format and removing invalid data after conversion

[33]:	cpa['	dt_customer'	] = pd.to_0	datetime(cpa	a['dt_customer	'], erro	rs='coerce	e')						
[35]:	сра =	<pre>cpa = cpa[cpa['dt_customer'].notnull()]</pre>												
[36]:	сра	сра												
[36]:		customer_id	year_birth	education	marital_status	income	kidhome	teenhome	dt_customer	recency				
	0	5524	1957	Graduation	Single	58138.0	0	0.0	2012-09-04	58.0				
	1	2174	1954	Graduation	Single	46344.0	1	1.0	2014-03-08	38.0				
	2	4141	1965	Graduation	Together	71613.0	0	0.0	2013-08-21	26.0				
	3	6182	1984	Graduation	Together	26646.0	1	0.0	2014-02-10	26.0				
	4	5324	1981	PhD	Married	58293.0	1	0.0	2014-01-19	94.0				
	2438	10870	1967	Graduation	Married	61223.0	0	1.0	2013-06-13	46.0				
	2439	4001	1946	PhD	Together	64014.0	2	1.0	2014-06-10	56.0				
	2440	7270	1981	Graduation	Divorced	56981.0	0	0.0	2014-01-25	91.0				
	2441	8235	1956	Master	Together	69245.0	0	1.0	2014-01-24	8.0				
	2442	9405	1954	PhD	Married	52869.0	1	1.0	2012-10-15	40.0				

2037 rows × 9 columns

### • Ensuring correct datatypes

```
[38]: cpa['customer_id'] = cpa['customer_id'].astype(int)
    cpa['year_birth'] = cpa['year_birth'].astype(int)
    cpa['education'] = cpa['education'].astype(str)
    cpa['marital_status'] = cpa['marital_status'].astype(str)
    cpa['income'] = cpa['income'].astype(float)
    cpa['recency'] = cpa['recency'].astype(int)
    cpa['kidhome'] = cpa['kidhome'].astype(int)
    cpa['teenhome'] = cpa['teenhome'].astype(int)
```

[39]:	сра									
[39]:		customer_id	year_birth	education	marital_status	income	kidhome	teenhome	dt_customer	recency
	0	5524	1957	Graduation	Single	58138.0	0	0	2012-09-04	58
	1	2174	1954	Graduation	Single	46344.0	1	1	2014-03-08	38
	2	4141	1965	Graduation	Together	71613.0	0	0	2013-08-21	26
	3	6182	1984	Graduation	Together	26646.0	1	0	2014-02-10	26
	4	5324	1981	PhD	Married	58293.0	1	0	2014-01-19	94
	2438	10870	1967	Graduation	Married	61223.0	0	1	2013-06-13	46
	2439	4001	1946	PhD	Together	64014.0	2	1	2014-06-10	56
	2440	7270	1981	Graduation	Divorced	56981.0	0	0	2014-01-25	91
	2441	8235	1956	Master	Together	69245.0	0	1	2014-01-24	8
	2442	9405	1954	PhD	Married	52869.0	1	1	2012-10-15	40

2037 rows × 9 columns

## • Convert 'dt\_customer' to consistent date format (dd-mm-yyyy)

```
[40]: cpa['dt_customer'] = cpa['dt_customer'].dt.strftime('%d-%m-%Y')
|
```

[41]:	сра									
[41]:		customer_id	year_birth	education	marital_status	income	kidhome	teenhome	dt_customer	recency
	0	5524	1957	Graduation	Single	58138.0	0	0	04-09-2012	58
	1	2174	1954	Graduation	Single	46344.0	1	1	08-03-2014	38
	2	4141	1965	Graduation	Together	71613.0	0	0	21-08-2013	26
	3	6182	1984	Graduation	Together	26646.0	1	0	10-02-2014	26
	4	5324	1981	PhD	Married	58293.0	1	0	19-01-2014	94
	2438	10870	1967	Graduation	Married	61223.0	0	1	13-06-2013	46
	2439	4001	1946	PhD	Together	64014.0	2	1	10-06-2014	56
	2440	7270	1981	Graduation	Divorced	56981.0	0	0	25-01-2014	91
	2441	8235	1956	Master	Together	69245.0	0	1	24-01-2014	8
	2442	9405	1954	PhD	Married	52869.0	1	1	15-10-2012	40
	2037 rd	ows × 9 colum	ns							

### • Standardrize the text value

[43]:	сра									
[43]:		customer_id	year_birth	education	marital_status	income	kidhome	teenhome	dt_customer	recency
	0	5524	1957	graduation	single	58138.0	0	0	04-09-2012	58
	1	2174	1954	graduation	single	46344.0	1	1	08-03-2014	38
	2	4141	1965	graduation	together	71613.0	0	0	21-08-2013	26
	3	6182	1984	graduation	together	26646.0	1	0	10-02-2014	26
	4	5324	1981	phd	married	58293.0	1	0	19-01-2014	94
	2438	10870	1967	graduation	married	61223.0	0	1	13-06-2013	46
	2439	4001	1946	phd	together	64014.0	2	1	10-06-2014	56
	2440	7270	1981	graduation	divorced	56981.0	0	0	25-01-2014	91
	2441	8235	1956	master	together	69245.0	0	1	24-01-2014	8
	2442	9405	1954	phd	married	52869.0	1	1	15-10-2012	40

2037 rows × 9 columns

### • Saved clean data

[44]: cpa.to\_csv("cleaned\_customer\_personality\_analysis\_data.csv", index=False)