



Data Preprocessing For Machine learning



About



Ayham Ramadan

Machine Learning Engineer



+201229492980

Mail: ay_ramadan@hotmail.com

Work Experience

DevisionX , Cairo **Oct 2021 - Present**
Computer Vision Engineer

Vodafone UK VIS , Cairo **Aug 2017 - Sep 2018**
Crisis Management Team Leader

NCCI - for Insurance , Riyadh **Jun 2015 - Sep 2015**
Oracle Developer

Information Technology Institute - MCIT , Alexandria
Branch Apr 2021 – Jan 2021
AI – PRO EPITA Certified

International School of Scandinavia, Copenhagen
Dec 2019 – Present
Project Management and Supply Chain optimization

Alexandria University of Engineering
Jan 2013 - Jan 2017
BSc in Computer and Communication

Agenda For Today

01 Pandas

What is Pandas and Why is it used ?

02 Dataframe

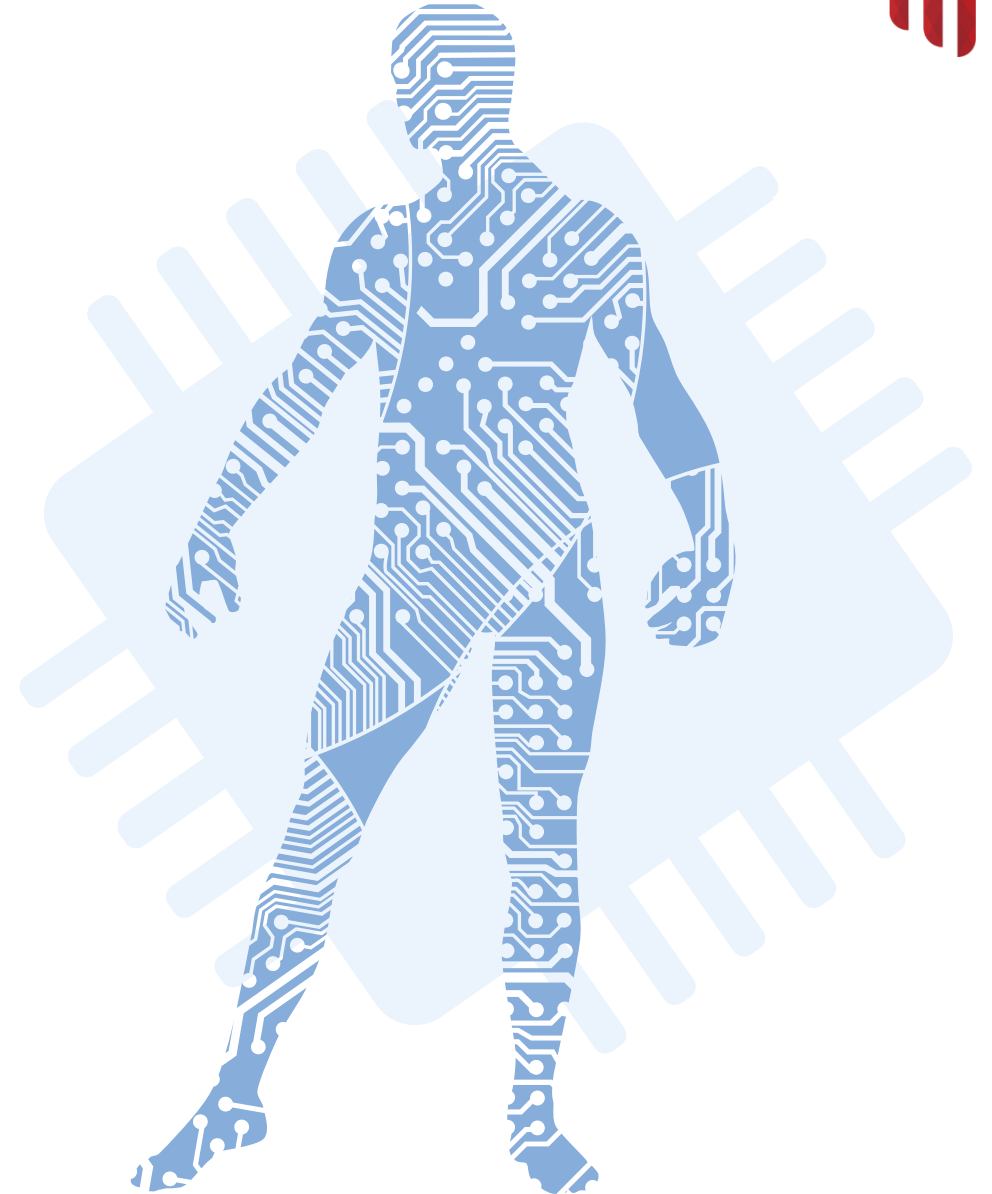
How do we apply changes and Preprocessing on Dataframes?

03 Loading Datasets (Practical)

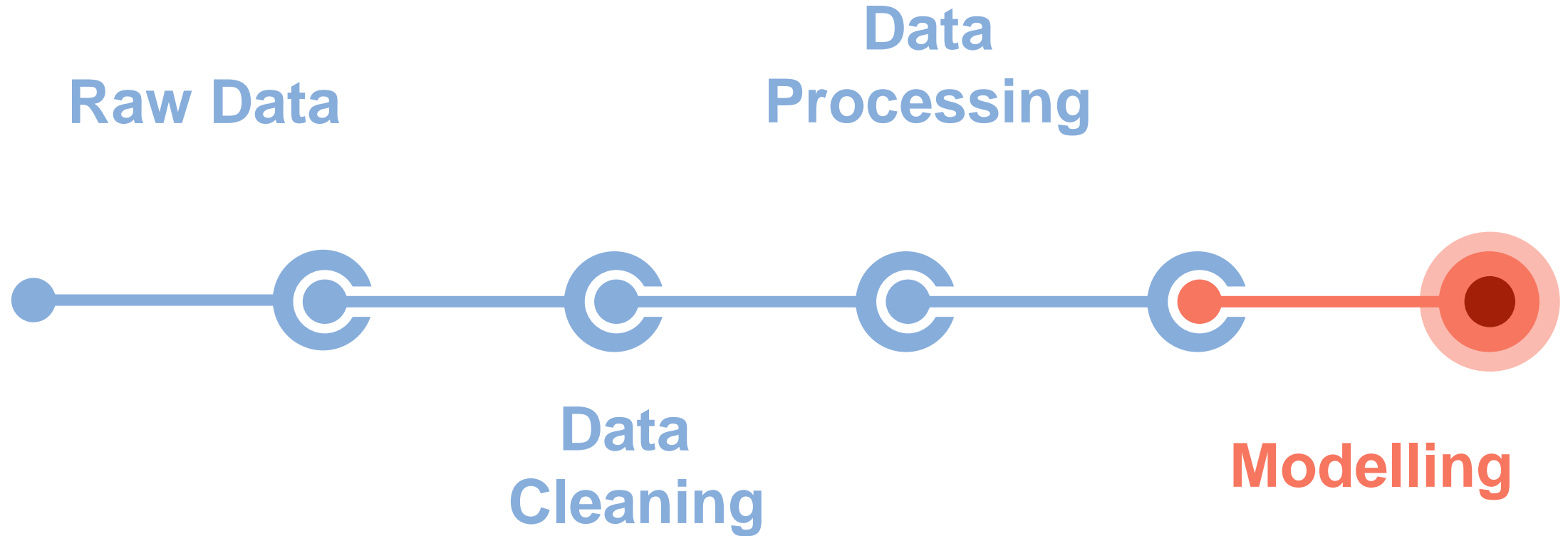
How to load different Datasets and How to use the documentation ?

04 Basic Preprocessing

Applying what was learned today?

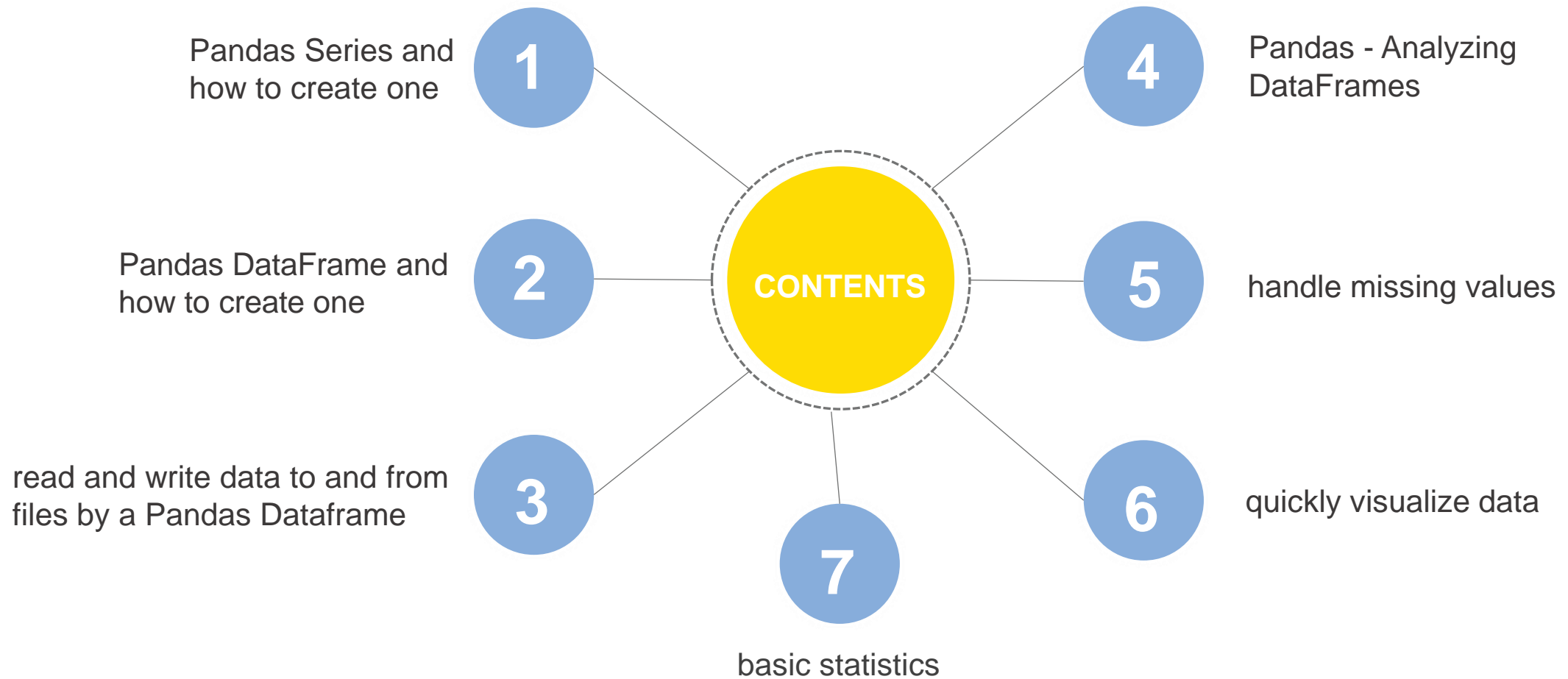


Data Pipeline





Pandas





What is Pandas?

- Pandas is a Python library used for working with data sets.
- It has functions for analyzing, cleaning, exploring, and manipulating data.
- The name "Pandas" has a reference to both "Panel Data", and "Python Data Analysis".
- Why Use Pandas?
Pandas allows us to analyze big data and make conclusions based on statistical theories.
- Pandas can clean messy datasets and make them readable and relevant.
- Relevant data is very important in data science.

Install it using this command:

```
1 !pip install pandas
```

Checking Pandas Version

The version string is stored under **version** attribute.

```
: 1 import pandas as pd
  2
  3 print(pd.__version__)
```

1.1.5

❑ Pandas Series



A Pandas Series is like a column in a table.
It is a one-dimensional array holding data of any type.

1. Creating a series from array

```
1 #Create a simple Pandas Series from a list
2 import pandas as pd
3
4 Name = ["Avery Bradley", "John Holland", "Jonas Jerebko", "Jordan Mickey", "Terry Rozier",
5         "Jared Sullinger", "Evan Turner"]
6
7 myvar = pd.Series(Name)
8
9 print(myvar)
```

```
0    Avery Bradley
1     John Holland
2    Jonas Jerebko
3    Jordan Mickey
4     Terry Rozier
5   Jared Sullinger
6      Evan Turner
dtype: object
```





1. Accessing element of Series

There are two ways through which we can access element of series, they are :

- Accessing Element from Series with Position
- Accessing Element Using Label (index)

2.1 Accessing Element from Series with Position : In order to access the series element refers to the index number. Use the index operator [] to access an element in a series. The index must be an integer. In order to access multiple elements from a series, we use Slice operation.

Accessing first 5 elements of Series

```
: 1 # import pandas and numpy
  2 import pandas as pd
  3 import numpy as np
  4
  5 # creating simple array
  6 data = np.array(['g','e','e','k','s','f','o','r','g','e','e','k','s'])
  7 ser = pd.Series(data)
  8 #retrieve the first 5 elements
  9 print("\n",ser[:5])
10 print("\n",ser[[0,1,2,4]])
11 #retrieve element
12 print("\n- 5th ele. is ",ser[5])
13
```

```
0    g
1    e
2    e
3    k
4    s
dtype: object

0    g
1    e
2    e
4    s
dtype: object

- 5th ele. is  f
```

2.2 Accessing Element Using Label (index) :

In order to access an element from series, we have to set values by index label. A Series is like a fixed-size dictionary in that you can get and [set values by index label](#).

```
1 # import pandas and numpy
2 import pandas as pd
3 import numpy as np
4
5 # creating simple array
6 data = np.array(['g','e','e','k','s','f','o','r','g','e','e','k','s'])
7 ser = pd.Series(data,index=["10","11","12","13","14","15","16","17","18","19","20","21","22"])
8
9 # accessing a element using label "index"
10 print(ser["16"])
11 print(ser[["16","15"]])
```

```
o
16    o
15    f
dtype: object
```

What will be the result if you pass dict into pd.Series() ?

```

1 # Create a simple Pandas Series from a dictionary:
2
3 import pandas as pd
4
5 calories = {"day1": [420,50], "day2": [380,34], "day3": [390,43]}
6
7 myvar = pd.Series(calories)
8
9 print(myvar)
10 print("\nDay1 data =",myvar["day1"])

```

```

day1    [420, 50]
day2    [380, 34]
day3    [390, 43]
dtype: object

```

Day1 data = [420, 50]

```

: 1 import pandas as pd
2
3 calories = {"day1": 420, "day2": 380, "day3": 390}
4
5 myvar = pd.Series(calories, index = ["day1", "day2"])
6
7 print(myvar)

```

```

day1    420
day2    380
dtype: int64

```

Binary operation methods on series:

FUNCTION	DESCRIPTION
<u>add()</u>	Method is used to add series or list like objects with same length to the caller series
<u>sub()</u>	Method is used to subtract series or list like objects with same length from the caller series
<u>mul()</u>	Method is used to multiply series or list like objects with same length with the caller series
<u>div()</u>	Method is used to divide series or list like objects with same length by the caller series
<u>sum()</u>	Returns the sum of the values for the requested axis
<u>prod()</u>	Returns the product of the values for the requested axis
<u>mean()</u>	Returns the mean of the values for the requested axis
<u>pow()</u>	Method is used to put each element of passed series as exponential power of caller series and returned the results
<u>abs()</u>	Method is used to get the absolute numeric value of each element in Series/DataFrame
<u>cov()</u>	Method is used to find covariance of two series

Pandas series method



FUNCTION	DESCRIPTION
Series()	A pandas Series can be created with the Series() constructor method. This constructor method accepts a variety of inputs
<u>combine_first()</u>	Method is used to combine two series into one
count()	Returns number of non-NA/null observations in the Series
size()	Returns the number of elements in the underlying data
name()	Method allows to give a name to a Series object, i.e. to the column
is_unique()	Method returns boolean if values in the object are unique
idxmax()	Method to extract the index positions of the highest values in a Series
idxmin()	Method to extract the index positions of the lowest values in a Series
sort_values()	Method is called on a Series to sort the values in ascending or descending order
sort_index()	Method is called on a pandas Series to sort it by the index instead of its values
head()	Method is used to return a specified number of rows from the beginning of a Series. The method returns a brand new Series
tail()	Method is used to return a specified number of rows from the end of a Series. The method returns a brand new Series
<u>le()</u>	Used to compare every element of Caller series with passed series.It returns True for every element which is Less than or Equal to the element in passed series
<u>ne()</u>	Used to compare every element of Caller series with passed series. It returns True for every element which is Not Equal to the element in passed series
<u>ge()</u>	Used to compare every element of Caller series with passed series. It returns True for every element which is Greater than or Equal to the element in passed series
<u>eq()</u>	Used to compare every element of Caller series with passed series. It returns True for every element which is Equal to the element in passed series

Pandas series method (Continue)



<u>lt()</u>	Used to compare two series and return Boolean value for every respective element
<u>clip()</u>	Used to clip value below and above to passed Least and Max value
<u>clip_lower()</u>	Used to clip values below a passed least value
<u>clip_upper()</u>	Used to clip values above a passed maximum value
<u>astype()</u>	Method is used to change data type of a series
<u>tolist()</u>	Method is used to convert a series to list
<u>get()</u>	Method is called on a Series to extract values from a Series. This is alternative syntax to the traditional bracket syntax
<u>unique()</u>	Pandas unique() is used to see the unique values in a particular column
<u>nunique()</u>	Pandas nunique() is used to get a count of unique values
<u>value_counts()</u>	Method to count the number of the times each unique value occurs in a Series
<u>factorize()</u>	Method helps to get the numeric representation of an array by identifying distinct values
<u>map()</u>	Method to tie together the values from one object to another
<u>between()</u>	Pandas between() method is used on series to check which values lie between first and second argument
<u>apply()</u>	Method is called and feeded a Python function as an argument to use the function on every Series value. This method is helpful for executing custom operations that are not included in pandas or numpy



❑ Pandas DataFrame

Pandas DataFrame is a 2-dimensional labeled data structure with columns of potentially different types. It is generally the most commonly used pandas object.

Pandas DataFrame can be created in multiple ways. Let's discuss different ways to create a DataFrame one by one.

Method #1: Creating Pandas DataFrame from lists of lists.

```
1 # Import pandas library
2 import pandas as pd
3
4 # initialize list of lists
5 data = [['tom', 10], ['nick', 15], ['juli', 14]]
6 # Create the pandas DataFrame
7 df = pd.DataFrame(data, columns = ['Name', 'Age'])
8 # print dataframe.
9 df
10
```

	Name	Age
0	tom	10
1	nick	15
2	juli	14



Method #2: Creating DataFrame from dict of narray/lists

To create DataFrame from dict of narray/list, all the narray must be of same length. If index is passed then the length index should be equal to the length of arrays. If no index is passed, then by default, index will be range(n) where n is the array length.

```
1 # Python code demonstrate creating
2 # DataFrame from dict narray / lists
3 # By default addresses.
4
5 import pandas as pd
6
7 # initialise data of lists.
8 data = {'Name':['Tom', 'nick', 'krish', 'jack'],
9         'Age':[20, 21, 19, 18]}
10 # Create DataFrame
11 df = pd.DataFrame(data)
12 # Print the output.
13 df
14
```

	Name	Age
0	Tom	20
1	nick	21
2	krish	19
3	jack	18

Method #3: Creates a indexes DataFrame using arrays.

```
1 # Python code demonstrate creating
2 # pandas DataFrame with indexed by
3
4 # DataFrame using arrays.
5 import pandas as pd
6 # initialise data of lists.
7 data = {'Name':['Tom', 'Jack', 'nick', 'juli'],'marks':[99, 98, 95, 90]}
8 # Creates pandas DataFrame.
9 df = pd.DataFrame(data, index=['rank1',
10 'rank2','rank3','rank4'])
11 # print the data
12 df
13
```

	Name	marks
rank1	Tom	99
rank2	Jack	98
rank3	nick	95
rank4	juli	90



Method #4: Creating Dataframe from list of dicts

Pandas DataFrame can be created by passing lists of dictionaries as a input data. By default dictionary keys taken as columns.

```
1 # Python code demonstrate how to create
2 # Pandas DataFrame by lists of dicts.
3 import pandas as pd
4 # Initialise data to lists.
5 data = [{'a': 1, 'b': 2, 'c':3},
6         {'a':10, 'b': 20, 'c': 30}]
7 # Creates DataFrame.
8 df = pd.DataFrame(data)
9 # Print the data
10 df
11
```

	a	b	c
0	1	2	3
1	10	20	30



Method #5: Creating DataFrame using zip() function.

Two lists can be merged by using list(zip()) function. Now, create the pandas DataFrame by calling pd.DataFrame() function.

```
1 # Python program to demonstrate creating
2 # pandas Datadaframe from lists using zip.
3
4 import pandas as pd
5 # List1
6 Name = ['tom', 'krish', 'nick', 'juli']
7 # List2
8 Age = [25, 30, 26, 22]
9 # get the list of tuples from two lists.
10 # and merge them by using zip().
11 list_of_tuples = list(zip(Name, Age))
12 # Assign data to tuples.
13 list_of_tuples
14 # Converting lists of tuples into
15 # pandas Dataframe.
16 df = pd.DataFrame(list_of_tuples,
17 columns = ['Name', 'Age'])
18 # Print data.
19 df
20
```

	Name	Age
0	tom	25
1	krish	30
2	nick	26
3	juli	22

```
[('tom', 25), ('krish', 30), ('nick', 26), ('juli', 22)]
```

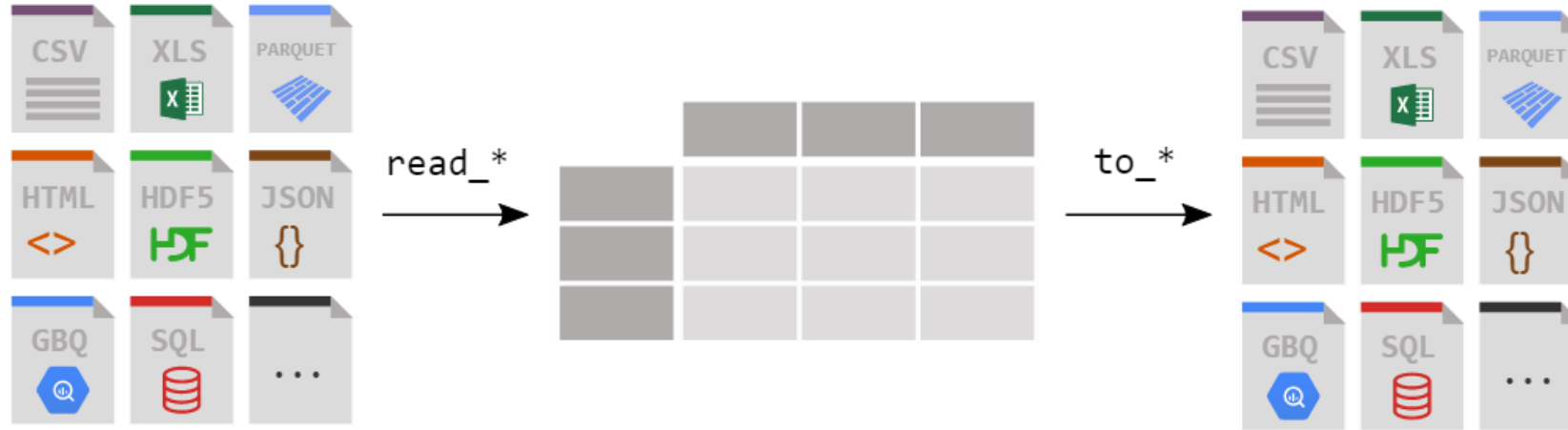
Method #6: Creating DataFrame from Dicts of series.

To create DataFrame from Dicts of series, dictionary can be passed to form a DataFrame. The resultant index is the union of all the series of passed indexed.

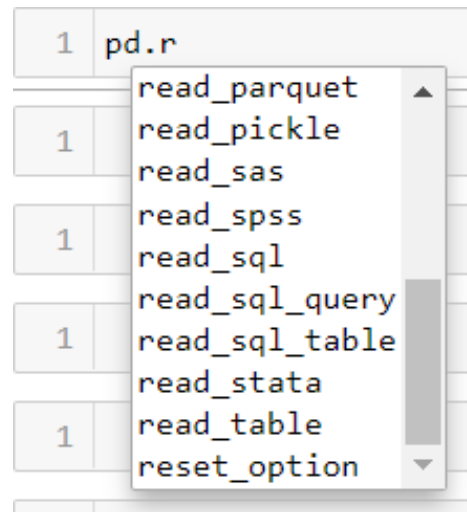
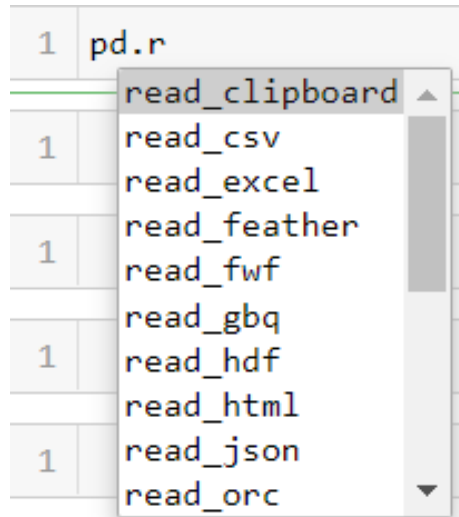
```
1 # Python code demonstrate creating
2 # Pandas Dataframe from Dicts of series.
3
4 import pandas as pd
5
6 # Initialise data to Dicts of series.
7 d = {'one' : pd.Series([10, 20, 30, 40],
8 index =['a', 'b', 'c', 'd']), 'two' : pd.Series([10, 20, 30, 40], index =['a', 'b', 'c', 'd'])}
9
10 # creates Dataframe.
11 df = pd.DataFrame(d)
12
13 # print the data.
14 df
15
```

	one	two
a	10	10
b	20	20
c	30	30
d	40	40

❑ Read and write data to and from files by a Pandas Dataframe



- pandas provides the [read_csv\(\)](#) function to read data stored as a csv file into a pandas DataFrame. pandas supports many different file formats or data sources out of the box (csv, excel, sql, json, parquet, ...), each of them with the prefix read_*.



<https://pandas.pydata.org/docs/reference/io.html>



1. Pandas Read CSV

A simple way to store big data sets is to use CSV files (comma separated files).

CSV files contains plain text and is a well know format that can be read by everyone including Pandas.

https://pandas.pydata.org/docs/reference/api/pandas.read_csv.html

```
1 #Load the CSV into a DataFrame:  
2 import pandas as pd  
3 df = pd.read_csv('dataset/data.csv')  
4 print(df)
```

	Duration	Pulse	Maxpulse	Calories
0	60	110	130	409.1
1	60	117	145	479.0
2	60	103	135	340.0
3	45	109	175	282.4
4	45	117	148	406.0
..
164	60	105	140	290.8
165	60	110	145	300.0
166	60	115	145	310.2
167	75	120	150	320.4
168	75	125	150	330.4

[169 rows x 4 columns]



Make sure to always have a check on the data after reading in the data. When displaying a DataFrame, the first and last 5 rows will be shown by default

- `df.head()`: by default, display first 5 rows from data.
- `df.tail()` : by default, display last 5 rows from data.

```
1 df.head()
```

	Duration	Pulse	Maxpulse	Calories
0	60	110	130	409.1
1	60	117	145	479.0
2	60	103	135	340.0
3	45	109	175	282.4
4	45	117	148	406.0

```
1 df.tail()
```

	Duration	Pulse	Maxpulse	Calories
164	60	105	140	290.8
165	60	110	145	300.0
166	60	115	145	310.2
167	75	120	150	320.4
168	75	125	150	330.4

```
1 df.head(10)
```

	Duration	Pulse	Maxpulse	Calories
0	60	110	130	409.1
1	60	117	145	479.0
2	60	103	135	340.0
3	45	109	175	282.4
4	45	117	148	406.0
5	60	102	127	300.0
6	60	110	136	374.0
7	45	104	134	253.3
8	30	109	133	195.1
9	60	98	124	269.0

```
1 df.tail(10)
```

	Duration	Pulse	Maxpulse	Calories
159	30	80	120	240.9
160	30	85	120	250.4
161	45	90	130	260.4
162	45	95	130	270.0
163	45	100	140	280.9
164	60	105	140	290.8
165	60	110	145	300.0
166	60	115	145	310.2
167	75	120	150	320.4
168	75	125	150	330.4

Reading in files with encoding problems

Most files you'll encounter will probably be encoded with UTF-8. This is what Python expects by default, so most of the time you won't run into problems. However, sometimes you'll get an error like this:

```
: 1 # modules we'll use
2 import numpy as np
3 # helpful character encoding module
4 import chardet

: 1 police_killings = pd.read_csv("PoliceKillingsUS.csv")
2039
21 self._first_chunk

pandas\_libs\parsers.pyx in pandas._libs.parsers.TextReader.read()
pandas\_libs\parsers.pyx in pandas._libs.parsers.TextReader._read_low_memory()
pandas\_libs\parsers.pyx in pandas._libs.parsers.TextReader._read_rows()
pandas\_libs\parsers.pyx in pandas._libs.parsers.TextReader._convert_column_data()
pandas\_libs\parsers.pyx in pandas._libs.parsers.TextReader._convert_tokens()
pandas\_libs\parsers.pyx in pandas._libs.parsers.TextReader._convert_with_dtype()
pandas\_libs\parsers.pyx in pandas._libs.parsers.TextReader._string_convert()
pandas\_libs\parsers.pyx in pandas._libs.parsers._string_box_utf8()

UnicodeDecodeError: 'utf-8' codec can't decode byte 0x96 in position 2: invalid start byte
```



chardet.detect()

The **detect** function takes one argument, a non-Unicode string. It returns a dictionary containing the auto-detected character encoding and a confidence level from 0 to 1

```
1 with open("PoliceKillingsUS.csv", 'rb') as rawdata:
2     result = chardet.detect(rawdata.read(100000))
3
4 # check what the character encoding might be
5 print(result)
```

```
{'encoding': 'Windows-1252', 'confidence': 0.73, 'language': ''}
```

encoding attribute : <https://docs.python.org/3/library/codecs.html#standard-encodings>

```
1 police_killings = pd.read_csv("PoliceKillingsUS.csv", encoding='Windows-1252')
2 police_killings
```

	id	name	date	manner_of_death	armed	age	gender	race	city	state	signs_of_mental_illness	threat_level	flee	body_camera
0	3	Tim Elliot	02/01/15	shot	gun	53.0	M	A	Shelton	WA	True	attack	Not fleeing	False
1	4	Lewis Lee Lembke	02/01/15	shot	gun	47.0	M	W	Aloha	OR	False	attack	Not fleeing	False
2	5	John Paul Quintero	03/01/15	shot and Tasered	unarmed	23.0	M	H	Wichita	KS	False	other	Not fleeing	False



Save the DataFrame to csv file

`DataFrame.to_csv(path_or_buf=None, sep=',', na_rep='', float_format=None, columns=None, header=True, index=True, index_label=None, mode='w', encoding=None, compression='infer', quoting=None, quotechar='"', line_terminator=None, chunksize=None, date_format=None, doublequote=True, escapechar=None, decimal='.', errors='strict', storage_options=None)`

https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.to_csv.html#pandas.DataFrame.to_csv

```
: 1 police_killings.to_csv("encoded_PK.csv")
```

```
: 1 data=pd.read_csv("encoded_PK.csv")
  2 data
```

```
: 
```

	id	name	date	manner_of_death	armed	age	gender	race	city	state	signs_of_mental_illness	threat_level	flee	body_camera
0	3	Tim Elliot	02/01/15	shot	gun	53.0	M	A	Shelton	WA	True	attack	Not fleeing	False
1	4	Lewis Lee Lembke	02/01/15	shot	gun	47.0	M	W	Aloha	OR	False	attack	Not fleeing	False
2	5	John Paul Quintero	03/01/15	shot and Tasered	unarmed	23.0	M	H	Wichita	KS	False	other	Not fleeing	False
3	8	Matthew Hoffman	04/01/15	shot	toy weapon	32.0	M	W	San Francisco	CA	True	attack	Not fleeing	False
4	9	Michael Rodriguez	04/01/15	shot	nail gun	39.0	M	H	Evans	CO	False	attack	Not fleeing	False
...

❑ Pandas - Analyzing DataFrames



access, modify, add, sort, filter, and delete data

1. Viewing and getting information about data

1. DataFrame.head(**N=5**): view the first 5 rows by default
2. DataFrame.tail(**N=5**): view the last 5 rows by default
3. DataFrame.info() :The DataFrames object has a method called info(), that gives you more information about the data set
4. DataFrame.index : to view the DataFrame indexes
5. DataFrame.columns : to view the dataframe columns name

2. Access, modify

data Access

1. DataFrame["column name"]: access certain column from the main dataframe
2. DataFrame[["column names list "]]: access multiples columns from the main dataframe
3. DataFrame. Loc[row's name, column's name] : access data point by row and column names
4. DataFrame. iLoc[row's index, column's index] :access data point by row and column indexes

modify

1. DataFrame["column name"]= New value
2. DataFrame[["column names list "]] =New value
3. DataFrame. Loc[row's name, column's name] = New value
4. DataFrame. iLoc[row's index, column's index] =New value
5. DataFrame.index = new index array
- 6.DataFrame.columns = new column names array

Add new column/s

DataFrame["new column name"]= value/array of values

3. Delete row/column

DataFrame.drop(labels=None, axis=0, index=None, columns=None, level=None, inplace=False, errors='raise')



```
1 file_dir="default of credit card clients.csv"
2 Raw_data = pd.read_csv(file_dir,index_col=1,header=1)
3 Raw_data.head()
```

	0	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	...	BILL_AMT4	BILL_AMT5	BILL_AMT6	PAY_AMT1	PAY_AMT2	PAY_AMT3
ID																	
1	1	20000	female	university	married	24	2	2	-1	-1	...	0	0	0	0	689	
2	2	120000	female	university	single	26	-1	2	0	0	...	3272	3455	3261	0	1000	
3	3	90000	female	university	single	34	0	0	0	0	...	14331	14948	15549	1518	1500	
4	4	50000	female	university	married	37	0	0	0	0	...	28314	28959	29547	2000	2019	
5	5	50000	male	university	married	57	-1	0	-1	0	...	20940	19146	19131	2000	36681	

5 rows × 25 columns

```
1 Raw_data.drop(labels="0",axis=1)
```

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	...	BILL_AMT4	BILL_AMT5	BILL_AMT6	PAY_AMT1	PAY_AMT2	PAY_AMT3
ID																	
1	20000	female	university	married	24	2	2	-1	-1	-2	...	0	0	0	0	689	
2	120000	female	university	single	26	-1	2	0	0	0	...	3272	3455	3261	0	1000	
3	90000	female	university	single	34	0	0	0	0	0	...	14331	14948	15549	1518	1500	

```
1 Raw_data.drop(columns=["0"])
```

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	...	BILL_AMT4	BILL_AMT5	BILL_AMT6	PAY_AMT1	PAY_AMT2	PAY_AMT3
ID																	
1	20000	female	university	married	24	2	2	-1	-1	-2	...	0	0	0	0	689	
2	120000	female	university	single	26	-1	2	0	0	0	...	3272	3455	3261	0	1000	
3	90000	female	university	single	34	0	0	0	0	0	...	14331	14948	15549	1518	1500	
4	50000	female	university	married	37	0	0	0	0	0	...	28314	28959	29547	2000	2019	

1. Data cleaning :

1. Duplicated Data

I. check the existence of duplicated data using `DataFrame.duplicated()`

`DataFrame.duplicated(subset=None, keep='first')`

II. then, remove it using `DataFrame.drop_duplicates(inplace=True)`

`DataFrame.drop_duplicates(subset=None, keep='first', inplace=False, ignore_index=False)`

https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.drop_duplicates.html

<https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.duplicated.html>

```
: 1 Raw_data.duplicated()
: 1      False
: 2      False
: 3      False
: 4      False
: 5      False
: ...
: 29996     False
: 29997     False
: 29998     False
: 29999     False
: 30000     False
: Length: 30000, dtype: bool
```

```
: 1 Raw_data.duplicated().sum()
```

```
: 35
```

```
1 Raw_data.drop_duplicates(inplace=True)
```

```
1 Raw_data.duplicated().sum()
```

```
0
```

1. Filtration : `DataFrame.filter()` it is used to select subset of data

<https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.filter.html>

2. missing values



3. Check the features datatype : DataFrame.dtypes()

```
1 data_A_filteration.dtypes
2 ### this line represent that all dataset features saved as object
```

```
ID
LIMIT_BAL          object
SEX                object
EDUCATION           object
MARRIAGE            object
AGE                object
PAY_0               object
PAY_2               object
PAY_3               object
PAY_4               object
PAY_5               object
PAY_6               object
BILL_AMT1           object
BILL_AMT2           object
BILL_AMT3           object
BILL_AMT4           object
BILL_AMT5           object
BILL_AMT6           object
PAY_AMT1            object
PAY_AMT2            object
PAY_AMT3            object
PAY_AMT4            object
PAY_AMT5            object
PAY_AMT6            object
default payment next month  object
dtype: object
```

4. Convert data from one type to another:

DataFrame.astype(dtype)

DataFrame[column names array].astype(dtype)

DataFrame.astype({"col1_name": "int32","column2_name":"int64"})

```
: 1 data_A_filteration[numerical_columns]=data_A_filteration[numerical_columns].astype("int64")
```

```
: 1 data_A_filteration.dtypes
```

```
: ID
LIMIT_BAL          int64
SEX                object
EDUCATION           object
MARRIAGE            object
AGE                int64
PAY_0               int64
PAY_2               int64
PAY_3               int64
PAY_4               int64
PAY_5               int64
PAY_6               int64
BILL_AMT1           int64
BILL_AMT2           int64
BILL_AMT3           int64
BILL_AMT4           int64
BILL_AMT5           int64
BILL_AMT6           int64
PAY_AMT1            int64
```

How to analyze , visualize and deal with categorical data ?

ID	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	...	BILL_AMT4	BILL_AMT5	BILL_AMT6	PAY_AMT1	PAY_AMT2
1	20000	female	university	married	24	2	2	-1	-1	-2	...	0	0	0	0	689
2	120000	female	university	single	26	-1	2	0	0	0	...	3272	3455	3261	0	1000
3	90000	female	university	single	34	0	0	0	0	0	...	14331	14948	15549	1518	1500
4	50000	female	university	married	37	0	0	0	0	0	...	28314	28959	29547	2000	2019
5	50000	male	university	married	57	-1	0	-1	0	0	...	20940	19146	19131	2000	36681

5 rows × 24 columns

```

: 1 data_A_filteration["EDUCATION"].unique()
: array(['university', 'graduate school', 'others', 'high school', 0],
      dtype=object)

: 1 data_A_filteration["EDUCATION"].value_counts()
: university      13857
  graduate school  10513
  high school      4811
  others           121
  0                 14
  Name: EDUCATION, dtype: int64

: 1 print("University =",(data_A_filteration["EDUCATION"]== 'university' ).sum())
  2 print("graduate school =",(data_A_filteration["EDUCATION"]== 'graduate school' ).sum())
  3 print("others =",(data_A_filteration["EDUCATION"]== 'others' ).sum())
  4 print("high school =",(data_A_filteration["EDUCATION"]== 'high school' ).sum())
  5 print("0 =",(data_A_filteration["EDUCATION"]== 0 ).sum())

University = 13857
graduate school = 10513
others = 121
high school = 4811
0 = 14

```

DataFrame.unique()
 DataFrame.value_counts()
 DataFrame.values



```
1 # Let's see the data types and non-null values for each column
2 data_After_pro.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
Index: 29316 entries, 1 to 30000
```

```
Data columns (total 24 columns):
```

#	Column	Non-Null Count	Dtype
0	LIMIT_BAL	29316 non-null	int64
1	SEX_male	29316 non-null	int64
2	EDUCATION_encoded	29316 non-null	int32
3	Marital_state	29316 non-null	int64
4	AGE	29316 non-null	int64
5	PAY_0	29316 non-null	int64
6	PAY_2	29316 non-null	int64
7	PAY_3	29316 non-null	int64
8	PAY_4	29316 non-null	int64
9	PAY_5	29316 non-null	int64
10	PAY_6	29316 non-null	int64
11	BILL_AMT1	29316 non-null	int64
12	BILL_AMT2	29316 non-null	int64
13	BILL_AMT3	29316 non-null	int64
14	BILL_AMT4	29316 non-null	int64
15	BILL_AMT5	29316 non-null	int64
16	BILL_AMT6	29316 non-null	int64
17	PAY_AMT1	29316 non-null	int64
18	PAY_AMT2	29316 non-null	int64
19	PAY_AMT3	29316 non-null	int64
20	PAY_AMT4	29316 non-null	int64
21	PAY_AMT5	29316 non-null	int64
22	PAY_AMT6	29316 non-null	int64
23	default payment next month	29316 non-null	int64

```
dtypes: int32(1), int64(23)
```

```
memory usage: 6.7+ MB
```



□ Handling missing values

The first thing to do when you get a new dataset is take a look at some of it. This lets you see that it all read in correctly and gives an idea of what's going on with the data. In this case, let's see if there are any missing values, which will be represented with NaN or None.

	Date	GameID	Drive	qtr	down	time	TimeUnder	TimeSecs	PlayTimeDiff	SideofField	...	yacEPA
0	2009-09-10	2009091000	1	1	NaN	15:00	15	3600.0	0.0	TEN	...	NaN
1	2009-09-10	2009091000	1	1	1.0	14:53	15	3593.0	7.0	PIT	...	1.146076
2	2009-09-10	2009091000	1	1	2.0	14:16	15	3556.0	37.0	PIT	...	NaN
3	2009-09-10	2009091000	1	1	3.0	13:35	14	3515.0	41.0	PIT	...	-5.03142
4	2009-09-10	2009091000	1	1	4.0	13:27	14	3507.0	8.0	PIT	...	NaN

How many missing data points do we have?

Let's see how many we have in each column using `dataframe.isnull().sum()` , `dataframe.isna()`



```
1 data.isna().sum()
```

ID	0
Name	0
Age	0
Photo	0
Nationality	0
...	
GKHandling	48
GKKicking	48
GKPositioning	48
GKReflexes	48
Release Clause	1564
Length: 88, dtype: int64	

```
1 data.isnull().sum()
```

ID	0
Name	0
Age	0
Photo	0
Nationality	0
...	
GKHandling	48
GKKicking	48
GKPositioning	48
GKReflexes	48
Release Clause	1564
Length: 88, dtype: int64	

```
1 data.isna()
```

	ID	Name	Age	Photo	Nationality	Flag	Overall	Potential	Club	Club Logo	...	Composure	Marking	StandingTackle	SlidingTackle	GKDividing	GKHar
0	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
...
18202	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
18203	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
18204	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
18205	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
18206	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False

18207 rows x 88 columns

```
1 data.isnull()
```

	ID	Name	Age	Photo	Nationality	Flag	Overall	Potential	Club	Club Logo	...	Composure	Marking	StandingTackle	SlidingTackle	GKDividing	GKHar
0	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
...
18202	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
18203	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
18204	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
18205	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False
18206	False	False	False	False	False	False	False	False	False	False	...	False	False	False	False	False	False

18207 rows x 88 columns



It's very important to figure out why the data is missing

Is this value missing because it wasn't recorded or because it doesn't exist?

If a value is missing because it doesn't exist (like the height of the oldest child of someone who doesn't have any children) then it doesn't make sense to try and guess what it might be. These values you probably do want to keep as NaN or change it to zero . On the other hand, if a value is missing because it wasn't recorded, then you can try to guess what it might have been based on the other values in that column and row. This is called imputation, and we'll learn how to do it next! :)

How to deal with missing values ?

1. drop columns with missing values
2. drop rows with missing values
3. replace them with another value
4. imputation :
 1. replace them with mean , median values if you deal with numerical data or most frequent value in case of categorical data
 2. KNN supervised ML model
 3. k-mean Un-supervised ML model



1. Drop missing values :

using `DataFrame.dropna()`

`DataFrame.dropna(axis=0, how='any', thresh=None, subset=None, inplace=False)`

- By default, drop rows with missing value. But, to drop columns set `axis=1`
- This function return new data after removing missing values. But, to override the old data use `inplace = True`
- Subset: Define in which columns to look for missing values in case of dropping rows and visa verse in case of dropping columns
- Thresh: Keep only the rows/columns with at least `N` non-NA values
- how: “any ” >> drop if any missing value exist , “all” >> drop if all the row / column values is nan

```
In [133]: 1 df = pd.DataFrame({"name": [np.nan, 'Batman', 'Catwoman'],
2                               "toy": [np.nan, 'Batmobile', 'Bullwhip'],
3                               "born": [pd.NaT, pd.Timestamp("1940-04-25"),pd.NaT],
4                               "null":[np.nan,np.nan,np.nan]})
5 df
```

Out[133]:

	name	toy	born	null
0	NaN	NaN	NaT	NaN
1	Batman	Batmobile	1940-04-25	NaN
2	Catwoman	Bullwhip	NaT	NaN

Dropping rows

```
In [140]: 1 df.dropna(how="all")
```

Out[140]:

	name	toy	born	null
1	Batman	Batmobile	1940-04-25	NaN
2	Catwoman	Bullwhip	NaT	NaN

```
In [141]: 1 df.dropna(how="any")
```

Out[141]:

	name	toy	born	null
--	------	-----	------	------

Dropping columns

```
1 df.dropna(axis=1,how="all")
```

	name	toy	born
0	NaN	NaN	NaT
1	Batman	Batmobile	1940-04-25
2	Catwoman	Bullwhip	NaT

```
1 df.dropna(axis=1,how="any")
```

0
1
2

thresh:

```
|:
```

	name	toy	born	null
0	NaN	NaN	NaT	NaN
1	Batman	Batmobile	1940-04-25	NaN
2	Catwoman	Bullwhip	NaT	NaN

```
|: 1 #Keep only the rows with at least 2 non-NA values.
   2 df.dropna(thresh=2)
```

```
|:
```

	name	toy	born	null
1	Batman	Batmobile	1940-04-25	NaN
2	Catwoman	Bullwhip	NaT	NaN

```
|: 1 #Keep only the rows with at least 3 non-NA values.
   2 df.dropna(thresh=3)
```

```
|:
```

	name	toy	born	null
1	Batman	Batmobile	1940-04-25	NaN

subset:

```
1 #Define in which columns to look for missing values.
2 df.dropna(subset=['name', 'toy'])
```

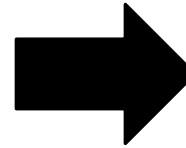
	name	toy	born	null
1	Batman	Batmobile	1940-04-25	NaN
2	Catwoman	Bullwhip	NaT	NaN

1. Filling in missing values:

We can use the Panda's `fillna()` function to fill in missing values in a dataframe for us

```
1 df
```

	name	toy	born	null
0	NaN	NaN	NaT	NaN
1	Batman	Batmobile	1940-04-25	NaN
2	Catwoman	Bullwhip	NaT	NaN



```
1 df.fillna(0)
```

	name	toy	born	null
0	0	0	0	0.0
1	Batman	Batmobile	1940-04-25 00:00:00	0.0
2	Catwoman	Bullwhip	0	0.0

DataFrame.fillna(value=None, method=None, axis=None, inplace=False, limit=None, downcast=None)

<https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.fillna.html>



❑ Pandas - Data Correlations

Finding Relationships

A great aspect of the Pandas module is the `corr()` method.

The `corr()` method calculates the relationship between each column in your data set.

• **Dataframe.corr()**

Types of correlations:

1. Perfect correlation
2. Positive correlation
3. Negative correlation
4. Bad correlation

How to visualize the correlation between features?

1. `sns.heatmap(correlation_result_data, annot=True)`
2. `sns.pairplot(dataframe)`
3. `sns.scatterplot(x=feature1, y=feature2)`



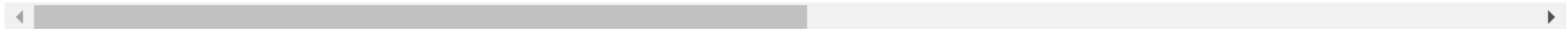
calculate summary statistics using **DataFrame.describe()**

```
|: 1 data_A_filteration[numerical_columns].describe()
```

|:

	ID	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	BILL_AMT1	BILL_AMT2	BILL_AMT3	...
count	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	2.931600e+04	...
mean	35.426695	-0.017465	-0.131259	-0.164074	-0.219232	-0.264224	-0.288580	51042.246316	49045.626995	4.691128e+04
std	9.497365	1.125777	1.199962	1.199591	1.171496	1.136187	1.151949	73480.513879	71051.572267	6.925505e+04
min	21.000000	-2.000000	-2.000000	-2.000000	-2.000000	-2.000000	-2.000000	-165580.000000	-69777.000000	-1.572640e+05
25%	28.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	3519.500000	2975.750000	2.646750e+03
50%	34.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	22282.000000	21095.500000	2.006850e+04
75%	41.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	66807.250000	63736.250000	5.995375e+04
max	267.000000	8.000000	8.000000	8.000000	8.000000	8.000000	8.000000	964511.000000	983931.000000	1.664089e+06	...	9

8 rows × 21 columns



```
: 1 data_A_filteration[categorical_columns].describe()
```

:

	ID	SEX	EDUCATION	MARRIAGE
count	29316	29316	29316	29316
unique	2	5	3	3
top	female	university	single	single
freq	17692	13857	15797	15797

Numerical data

Categorical data

Aggregating statistics



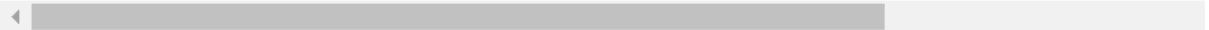
I. Numerical data

```
|: 1 data_A_filteration[numerical_columns].describe()
```

```
|:
```

	ID	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	
count	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	29316.000000	29316
mean	35.426695	-0.017465	-0.131259	-0.164074	-0.219232	-0.264224	-0.288580	-0.288580	5
std	9.497365	1.125777	1.199962	1.199591	1.171496	1.136187	1.151949	1.151949	7
min	21.000000	-2.000000	-2.000000	-2.000000	-2.000000	-2.000000	-2.000000	-2.000000	-16
25%	28.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	5
50%	34.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	2
75%	41.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	6
max	267.000000	8.000000	8.000000	8.000000	8.000000	8.000000	8.000000	8.000000	96

8 rows × 21 columns



```
|: 1 data_A_filteration["AGE"].mean()
```

```
|: 35.426695319961794
```

```
|: 1 data_A_filteration["AGE"].count()
```

```
|: 29316
```

```
|: 1 data_A_filteration["AGE"].std()
```

```
|: 9.497365305617228
```

```
|: 1 data_A_filteration["AGE"].max()
```

```
|: 267
```

II. Categorical data

```
1 data_A_filteration[categorical_columns].describe()
```

	ID	SEX	EDUCATION	MARRIAGE
count	29316	29316	29316	29316
unique	2	5	3	3
top	female	university	single	single
freq	17692	13857	15797	15797

```
1 data_A_filteration["SEX"].count()
```

29316

```
1 data_A_filteration["SEX"].value_counts().count()
```

2

```
1 data_A_filteration["SEX"].value_counts().idxmax()
```

'female'

```
1 data_A_filteration["SEX"].value_counts().max()
```

17692

```
1 data_A_filteration["SEX"].nunique()
```

2