

* Cleaning Data using Pandas:-

→ import pandas as pd
import numpy as np

```
people = {  
    "first": ["Essam", "Jannet", "John", np.nan, None, "NA"],  
    "last": ["Shenhav", "Doe", "Doe", np.nan, np.nan, "Missing"],  
    "age": [20, 21, 22, None, None, "Missing"]  
}
```

df = pd.DataFrame(people)

df →

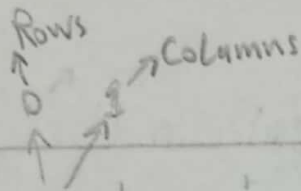
	first	Last	age
0	Essam	Shenhav	20
1	Jannet	Doe	21
2	John	Doe	22
3	NaN	NaN	None
4	None	NaN	None
5	NA	Missing	Missing

* Removing Rows/columns with missing values:-

→ df.dropna() → حذف الصفوف التي بها NaN أو الأعمدة التي بها NaN

→

	first	Last	age
0	Essam	Shenhav	20
1	Jannet	Doe	21
2	John	Doe	22
5	NA	Missing	Missing



→ `df.dropna(axis=0, how='any', inplace=True)`
 ↳ any: drop Rows if any of the values is missing
 ↳ all " " "all " " " " " " " " " " " "

* Dropping Rows with missing values in a certain Column

→ `df.dropna(axis=0, how='any', subset=['last'])`

→

	First	Last	age
0	Essam	Shenhav	20
1	Jannet	Doe	21
2	John	Doe	22
5	NA	Missing	Missing

→ `df.dropna(axis=0, how='all', subset=['last', 'age'])`

* Customizing ^{standard} missing values to Remove missing values Like (NA, missing)

→ `df.replace('NA', np.nan, inplace=True)`

`df.replace('Missing', np.nan, inplace=True)`

← كل قيمة ب NA أو Missing ستصبح NAN ليسهل حذفها

* identifying missing values using `isna()`

→ `df.isna()` → True و أي شيء آخر ب False يستبدل ب NAN

* Filling missing values

- `df.fillna('MISSING')` → `MISSING` بـ `NaN` يستبدل
وكن مفيد أكثر لاستبدال الـ `numerical data` بـ `mean` أو `حضر`
- `df['age'].mean()` → `21` أو `Error`

* Converting missing values to appropriate datatypes

- `df['age'] = df['age'].astype(float)`
- `df['Years Code'].unique()` → gives us all the unique of that column

* Removing duplicates

- `df.drop_duplicates(inplace=True)`

* Strip

- `df['Last name'] = df['Last name'].str.strip('...')`
الحروف التي سيتم حذفها ← `Left`

* Standardizing Phone numbers

- `df['Phone Number'] = df['Phone Number'].str.replace('[^a-zA-Z0-9]', '')`
- `df['Phone Number'].apply(lambda x: x[0:3] + '-' + x[3:6] + '-' + x[6:10])`
- `df['Phone Number'].apply(lambda x: str(x))`

* Splitting Columns

→ `df[['Street-Address', 'State', 'zip-code']] = df['Address'].str.
split(',', 2, expand=True)`

* Filtering Rows of Data

→ `for x in df.index:`
 `if df.loc[x, 'PhoneNumber'] == '':`
 `df.drop(x, inplace=True)`

df → ^{فأخذه} حذف كل rows التي بها بيانات

→ `df = df.dropna(subset="PhoneNumber", inplace=True)`

* Resetting index

→ `df = df.reset_index(drop=True)`

* Check For Null values

→ `df.isnull().any()`

→ `df[df['Phone Number'].isna() == 1]`

* Filling missing value equal to the value after it

→ `df.ffill`

↳ in front of it

→ `df.bfill`

↳ behind of it