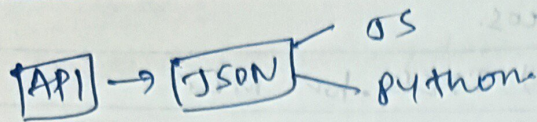


## working with JSON



code:

```
import pandas as pd
pd.read_json('_.json')
pd.read_json('url.com/json')
```

extract data from an API

## working with SQL

```
* !pip install mysql.connector
* import mysql.connector
* con = mysql.connector.connect(host='localhost', user='root',
                                password='', database='world')
* pd.read_sql_query('SELECT * FROM city', con)
```

citytable

## API Application Interface

code:

```
import pandas as pd
import requests
response = requests.get('_.url')
pd.DataFrame(response.json()['results'])
```

Rapid API

- when the site doesn't provide api's go for web scraping.

\* (use BeautifulSoup for web scraping)

## Understanding Your data

01. How big the data is?

- df.shape

02. How does data look like?

- df.head() | df.sample(5)

03. What is the data type or cols?

- df.info()



04. Are there any missing values?

- `df.isnull().sum()`

05. How does data look mathematically?

- `df.describe()`

06. Are there any duplicate values?

- `df.duplicated().sum()`

07. How is the correlation b/w cols?

- `df.corr()`

## EDA (Exploratory Data Analysis) (Titanic dataset)

univariate analysis: Analysis of single column

01. categorical data (plotting data)

`sns.countplot(df['survived'])`

`df['survived'].value_counts()`

} count how many have survived.

`df['survived'].value_counts().plot(kind='plot')`

02. Numerical data

Histogram

`plt.hist(df['Age'])`

- the distribution of data in that numerical column.

Distplot: Improvement of histogram.

`sns.distplot(df['Age'])`

Box plot:

outliers

min

Q1

median

Q3

max

$Q3 + 1.5 \times IQR$

$(Q1 - 1.5 \times IQR)$

$$IQR = Q3 - Q1$$

`sns.boxplot(df['fare'])`