

* Grouping and Aggregating:- Compining multiple pieces of data into a single result

* Median Function

→ `df['Salary'].median()` → 57287.0

→ `df.median()` → تحسب ال median لكل column بحول قيم عددية

Salary	57287
Age	29
work Week Hrs	40.0

* describe function

→ `df.describe()` → for a quick overview for the data
(count, mean, std, min, 25%, 50%, 75%, max) تحسب

* Count function

→ `df['Salary'].count()` → عدد العناصر الموجودة في ال column
بمعنى عدد الناس التي جاوبت على السؤال

* value_counts()

→ `df['Hobbyist'].value_counts()` → Yes 71257

No 17626

مثال عدد الناس التي جاوبوا بـ Yes أو No
أو FB أو YT أو insta

→ `df['Hobbyist'].value_count(normalize=True)`

→ To break it down by percentage

→ Yes 0.8017 No 0.1983

To each group independently

Dataframes into groups

* Grouping:- 1-split 2- Apply function 3- Combine Results

* groupby() function

→ country_grp = df.groupby(['country'])

country_grp.get_group('India') → تميل البيانات الأشخاص إلى

ال Country بتاعهم India زي القتر على اد Country

↓
filt = df['country'] == 'India'

df.loc[filt]

→ country_grp['Social Media'].value_counts() → تحسب ال counts
الخامس بال Social media ال Country

→ country_grp['Social Media'].value_counts().loc['india']
← نفس الحوار يكون مع India

→ country_grp['Social Media'].value_counts(normalize=True).loc['india']

← نفس الحوار لكن لميلام الأرقام كنسبة

→ country_grp['Salary'].median().loc['Germany']
← ال median الخامس بال Salary في Germany

* agg function for multiple function

→ `country_grp['Salary'].agg(['median', 'mean']).loc['India']`
India: المتوسط median والخاص mean

* using filtering to get python users by country

→ `filt = df['country'] == 'India'`
`df.loc[filt]['Language'].str.contains('Python').sum()`
→ 3105 → India: عدد مستخدمي Python في India

* لا يمكن استبدال `df.loc[filt]` بـ `country_grp` فقط بل نستعمل `apply`

→ `country_grp['Language'].apply(lambda x: x.str.contains('Python').sum())`

* % of developers using python in each country

→ `country_respondents = df['country'].value_counts()`
`country_uses_python = country_grp['Language'].apply(lambda x: x.str.contains('Python').sum())`
`python_df = pd.concat([country_respondents, country_uses_python], axis='columns', sort=False)`
`python_df['Pct Knows Python'] = (python_df['country_uses_python'] / python_df['country_respondents']) * 100`
`python_df.sort_values(by='Pct Knows Python', ascending=False, inplace=True)`

مثال آخر
مختصر → `country_grp['Language'].apply(lambda x: x.str.contains('Python').sum() / len(x) * 100)`

* When to use "groupby" in pandas?

→ Anytime you wanna Analyse some pandas series by some category.

→ when using the word 'For each...', what is something
continent ^{✓ex.} mean, max, ... ^{✓ex.}

* To display a groupby in visual Form

→ %matplotlib inline

drinks.groupby('continent').mean().plot(kind='bar')

bar plot لعرض البيانات على هيئة أعمدة

* To groupby multiple columns

→ df.groupby(['Salary', 'Language'])

* Working with Dates and Time Series Data:

* convert to datetime using to_datetime:

→ `df['Date'] = pd.to_datetime(df['Date'])`

Format) ممكن تدي error بـ Unknown Format لكن بتستخدم على

2020-03-13 08-PM 11:00 ←

→ `df['Date'] = pd.to_datetime(df['Date'], format='%Y-%m-%d %I-%p')`

→ `df.loc[0, 'Date'].day_name()` → `'Friday'`
 الجواب يومه في التاريخ

* Parse dates while loading data from csv

→ `d_parser = lambda x: pd.datetime.strptime(x, '%Y-%m-%d %I-%p')`
→ `df = pd.read_csv('DFname.csv', parse_dates=['Date'], date_parser=d_parser)`

* using day_name on whole series using dt class

→ `df['Date'].dt.day_name()` → `اليوم في Date`

* create column of dayname

→ `df['DayOfWeek'] = df['Date'].dt.day_name()`

* min & max methods on datetime series

→ `df['Date'].min()` → Timestamp('2017-07-01 11:00:00')

→ `df['Date'].max()` → Timestamp('2020-03-13 20:00:00')

time
delta → `df['Date'].max() - df['Date'].min()`

→ Timedelta('986 days 09:00:00')

* Filtering by dates as string

→ `Filt(df['Date'] >='2020')`

`df.loc[Filt]` → rows from 2020 onwards

→ `Filt(df['Date'] >='2019') & (df['Date'] < '2020')`

→ rows from 2019 to 2020

* Filter by to_datetime

→ `Filt = (df['Date'] >= pd.to_datetime('2019-01-01'))`

`& (df['Date'] < pd.to_datetime('2020-01-01'))`

* Set date as index

→ `df.set_index('Date', inplace=True)`

* Filtering by just passing the date in brackets

→ `df['2019']` → الطريقة دي أسهل من الفلتر

* Using Slices to get specific dates data

→ `df['2020-01':'2020-02']`

* Calculating average of a slice (timeLine)

→ `df['2020-1':'2020-02']['Close'].mean()` → 195.16569

* Getting max value of a column on a given day

→ `df['2020-01-01']['High'].max()`

* Resampling (breaking down by days) a whole column into a new variable

→ `df['High'].resample['D']` → Daily basis

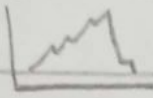
2 days ← 2D 3D W → Weekly
موجود في الـ documentation

→ `highs = df['High'].resample['D'].max()` → يعرف أعلى سعر لكل يوم

→ `highs['2020-01-01']` → يعرف أعلى سعر في يوم معين

* Plotting with matplotlib in pandas

→ %matplotlib inline

high.plot() → 

* resampling df with multiple columns

→ df.resample('W').mean() → columns ١١ mean ١١

* using agg to apply different function on different columns while resampling

→ df.resamp('w').agg({'Close': 'mean', 'High': 'max',
'Low': 'min', 'Volume': 'sum'})

→ ({'column ١١': 'Method ١١'})