	* Python Pandas: -a Library used for working with	1
- Marie Care	data sets	
	It has functions for analyzing, cleaning, exploring and manipulating data	
	Pandas allows us to analyze big data and make conclusions based on statistical theories	
	LONG LAST-113 DAJEM DI, DEGLISBICAL DIRECTIES	
-	Pandas can clean messy data sets, and make them readable and relevant	
<u>*</u>	Getting started:-	
<del>-</del> >	import pandas as pd	The second second
	× Loading data	
	df = od. read -csy ('csy File Location') -	
	df = pd. read -csv ('csv File Location') Zex.  ('data/servey-result-public.	(evs)
	df > dota  frame sicilizio upe upe	
	df.shape -> (89184,84)	
	df.info () -> dete types ) Wise general jesepen in 15	
	pd. set_option ('display. max_rows', 86)dF	
	سيعرض جميع ال والماعة عند الطياعة ·	
	display.max_columns' = coloums )1 to Utill 9	
-	لملاعة أول وسمع فقط أوعد معين إذا كتب حر ( ) df. head	
	dF. tail() -> ( ماماعة آخود عنوا أوسر مين ( ) الماماعة آخود الماماعة الماماعة آخود الماماعة آخود الماماعة آخود الماماعة آخود الماماعة الماماعة آخود الماماعة آخود الماماعة آخود الماماعة آخود الماماعة	

* Pandas Data Frame & Series
people = f
"First": ["Essam", "Jannet", "Will"],
"Last": ("Shenhab", "Ope", "Smith"),
"99e": [20,21,22]
7
people['age'] -> [20,21,22]
* Creating a Dataframe:-
import pardes as pol
dF = pd. Data Frame (people) * Load data into a Data Frame
dF-
First Last age
0 Essam Shenhab 20
1 Jannet Doe 21
2 Will Smith 22
The state of the s
df['age'] > 0 20mm
1 dF.age 1 21
2 22 14 6
method مساخی نام طلا می می می می می است ال ا کا را کا کا بی می است می
The state of the s

LE REAL CONTRACTOR OF THE PROPERTY OF THE PROP
* Loc: Accessing a row by Label Location
= Cot recessing a row of sales
df. Loc [[0,1], 'Last'] , o Shenhab
1 Doe
When used in a survey
df[''].value counts() =x Yes 7128
1896
df-loc[0] -> prints out one person's entire survey
dF. Loc[0, 'Hobbyist'] -> 'Yes' -> Prints out gresponse
For aspecific Row
X SLicing (Not-inclusive)
dF. Loc [0:2, 'Hobbyist': 'Employment']
Hobbyist Age Employment
o Yes 18-24 years old Not Employed
1 No 25-34 11 WOYK From home
2 Yes 45-54,,, Employed full-time

	in the same		- Charles
	* Loc: Accessing	a row by Label Loca	tion
		The state of the second	
	df. Loc [[0,1],	'Last') o Shen	hab
	20 / 3 /	1 Ooe	
	* When used		
		A	7128
			1896
	df-100507 > 1	vints out one person's e	ntive survey
-		yist'] -> 'Yes' -> Prints	
	L	for as	pecific Row
	X SLicing (Not:	-inclusive)	
			1-31
	dF. Loc [0:2. 1	Hobbyist': 'Employment']	
	Hobbyis	,	Emologment
	o Yes		Not Employed
	1 No	25-34 // ///	Work From home
	2 Yes	45-54,,,,	Employed full-time
		F 68 22	Ten III
	100 100		
Telephone I	The second second second		The second section is a second section of

Indexes:- Joint deall column 13 49 [1]  Indexes:- List Age  df > 0 Essam Shenhab 20  1 Jannet Que 21  2 Will smith 22  X set something as the index For the data Frame  df. set_index(!last') -> last first Age  Shenhab Essam 20  Que Janet 21  Smith Will 22  culp 19 line and olicids of griefling df clip index 11
df > 0 Essam shenhab 20  1 Jannet Doe 21  2 Will smith 22  Set something as the index Forthedata Frame  df. set_index('last') -> last first Age  Shenhab Essam 20  Doe Janet 21  Smith Will 22
Jannet Doe 21  2) Will smith 22  Set something as the index For the data Frame  df. set_index('last') -> last first Age  Shenhab Essam 20  Doe Janet 21  Smith Will 22
1 Jannet Doe 21 2 Will smith 22  x set something as the index Forthedge Frame  df. set_index ('last') -> last first Age  Shenhab Essam 20  Doe Janet 21  Smith Will 22
set something as the index Forthedge Frame  df. set_index('last') -> last first Age  Shenhab Essam 20  Doe: Janet 21  Smith Will 22
df. set_index('last') -> last first Age Shenhab Essam 20 Doe: Janet 21 Smith Will 22
df. set_index('Last') -> Last first Age Shenhab Essam 20 Doe Janet 21 Smith Will 22
Shenhab Essam 20 Doe Janet 21 Smith Will 22
Smith Will 22
Smith Will 22
Leus 19 line and olice of a greet index
df. set_index('last' simplace = True) list signification of
Last 1132 Age
Shenhob Essam 20
Doe Jannet 21
Smith Will 22
Maria Company of the
df.index -> Index (['Shenhab', 'Doe', 'smith'], dtype
= 'object', name='last')
df.loc['Shenhab'] > First Essam
Last Shenhab
Age 20
df. Loc ['Shenhab', 'First'] -> 'Essam

dF. Loc [o] -> Error	The second second
dF. i(oc [o] -> First	Essam
Lasti	Shenhab
Age	20
The state of the s	The second secon
* Re-setting the inc	lex
df. reset sindex (inpla	
0.7	First Age
O shenhab	Essam 20
1 00e	Jannet 21
2 Smith	Will 22
· CSV calo	عن قرادة الم
	July Western
dF=pd. read_csv('da	ta/survey_results_public.cov',
	index cal = 'Response ID')
df. head (3) →	Country Hobbyist
	2) - 1 - 1
1	Egypt Yes
2	USA No
113	UK Yes
schema_df.sort_inder	
ا من المعدد	Age 1st code Better Life
	THE RESERVE OF THE PARTY OF THE

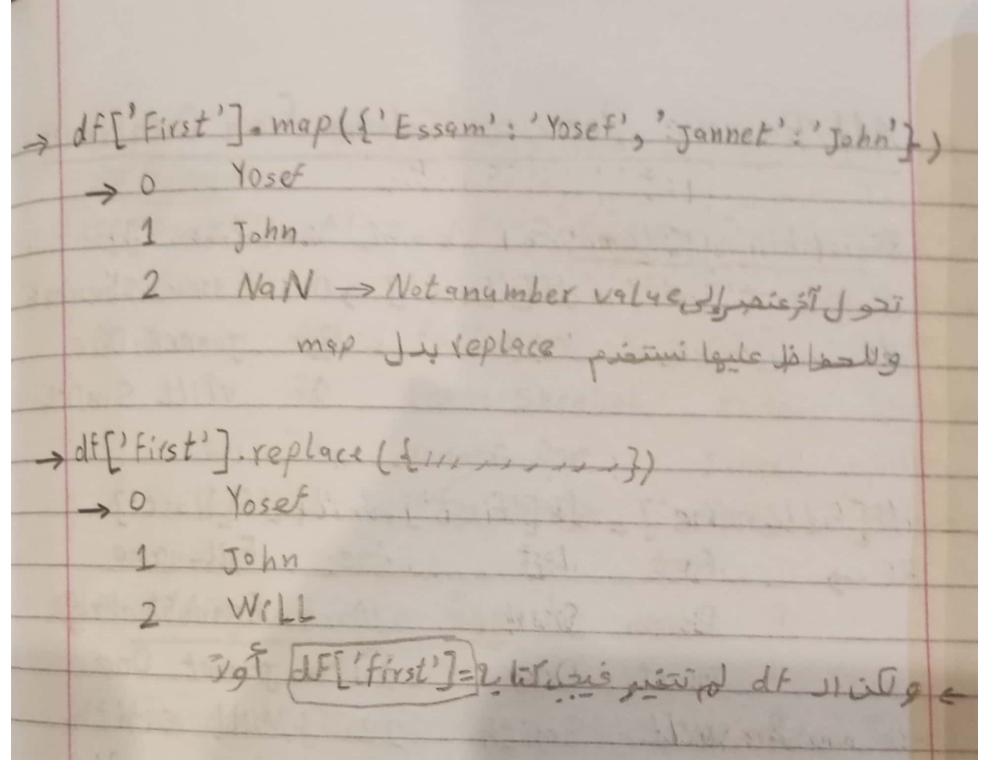
* Filtering:-
* O False
filt = (df['Last'] - 1 Poe') + 1 True
2 False
df[Filt] 1 Jannet Doe 21
df[filt] Jannet Doe 21  df.la[filt] 2 John Doe 21  df.la[filt] 2 John Doe 21
df. Loc [filt, 'Age'] > 1 21
3 24
* '8' > and '1' > or
Filt=(df['Last']=='Doe') & (df['first']=='John')
df. Loc [Filt, 'Age'] > 3 24
df. loc[~filt, 'Age'] > 0 20
False siedlinger 1 21
2 122
* Examples:-
countries = ['USA'; India', 'UK']
Filt = df['Country'], isin (countries)
df.loc[Filt, 'Country'] > 2 USA
6 TIK
12 India
27 India
107 USA
The property of the party of th

Filt = df ['Language Worked With'] . str. contains ('Python' ging = False) df. Loc [ Filt, 'Language Worked With'] HTML/CSS; Java; Python; C C++;HTML/css; python Bash/Shell/ Power Shell; Python; Java

\* Updating Rows and Columns: df. columns -> Index (['first', 'last', 'Age'], dtype = 'object') Ight colymns I pout busil x df.columns=['First\_name', 'Last\_name', 'Age'] df => Index ([' First - name', 'Last - name', 'Age'], dtype = 'object') methods observe columns I de Jucil \* -> df. columns = [x method() For x in df. columns]

upper -> Index(['FIRST\_NAME','L']) -> df-columns= df. columns-str. replace (' ', '-') 2000 colymn Just & , df. rename (columns = { First name : 'First', 'Last name': 'last'); inplace=True) so column clib duci x ~ df. Loc[2] = { "WILL", 'Ferrell', "45'} > df. Loc[2, ['Last', 'Age']]=['Ferrel', '45'] >df. Loc[2, 'Last'] = "Ferrel' > df. qt[[,,,,]] = 'Ferrel' & When updating a single object

* Common Mistakes	
Filt=(df['Last'] == 'smith')	
JETEN TON Setting	With
df[Filt]['last'] = 'Ferrel' -> Error-Setting	warning
ENI df-Loc[Filt, 'Last'] = 'Ferrel'	
avaloc [ Tive, Last ] = 10100	
columne Il dis citilul de ser	4
df['last']=df['last'].str.upper()	
df['Last'] > 0 SHENHAB	
1 DOE	
2 FERREL	
× Some Methods	
def update email (email):	
veturn email upper ()	
df['email'].apply (update_email)	
df['Last'] = df['email'].apply(Lambda x: x.lower	(1)
	31
df.apply (len) 0 3	
1 3	
2 3	
dF.applymap (Len), first Last email	L
0 5 7 23	35-
*Age (int) has no Len() 1 6 3 17	40
2 4 5 17	-



	Add/Remove Rows and Columns From Dataframes:
*	Add/ Vernore Lous du cochin
	X Combining Columns
	df ['First'] + ' + df [ last ] > C Essam Shenhab
	2 Jannet Doe 2 Will Smith
-	df['full_name'] = df['first']+ '+df['last']
	dF > first last Age Full-name
	Essam Shenhab 20 Essam Shenhab
	Jannet Doe 21 Jannet Doe
	2 will Smith 22 Will Smith
	* Removing columns
	df.drop (columns = ['First', 'Last'], inplace = True)
	df = Age fall-name
	0 20 Essam Shenhab
	1 21 Jannet Doe
-	2 22 Will Smith
	X Splitting Strings
	df['full_name].str.split(' ')
	> 0 [Essam, Shenhab]
	1 [Januet, Doe]
	2 [Will, Smith]

* Expand
df[['First', 'Last']] = df['Full_name'].str split!
expand = Tru
df > Age full-name First Last
0 20 Essam Shenhab Essam Shenha
1 21 Jannet Doe Jannet Doe
2 22 Will smith Will Smith
Maria Maria
* Adding a Single Row
df.append ({'Ficst': 'Tony'}, ignoreindex = True)
Age Full-name First Last
0 20 EssamShenhab Essam Shenhab
1 21 Jannet Doe Jannet Doe
2 22 WILL Smith WILL Smith
3 NaN NaN Tony NaN
*Append Data Frames
people = {
'First': ['Tong', 'STeve'],
Last': ['STark', Rogers'],
df2 = Pd. DataFrame (people)

15 115	2 innove index=True) > Hatala
>df-df-appena (a)	2, ignore_index=True) > distibling
-> df=pd.concql(	([df,df2], ignore_index=True, sort=true, sort=true)
JE 990	Elect Last
df - 990 0 20	
1 21	
2 22	C. C.
3 NaN	
4 NaN	
1 / 4 / 4	, or
* Remove Row	MS .
> df. drop (index	
	0 20 Essem Shenhab
	1 21 Jannet Doe
	2 22 Will Smith
	4 NON STEVE Rogers
af. drop(index = c	df[df['Last'] == 'Doe'].index)
	First Last age
	Essam Shenhab 20
	Will Smith 22
0	
	STeve Rogers NaN

	Sorting Data
*	500
	df. sort = values (by = 'Last') > Dataframe is sorted
	alphabitically depending on the (last
	and if it was numbers it would be sorted
Will.	by the highest number
	> First Last age
	1 Jannet Doe 21
	0 Essam Shenhab 20
	2 Will Smith 22
	of an indicated Para (COP) hard of a familiar to be
-	* Sorting in descending order
>	df. sort_values (by='lest', ascending = False)
	*If there duplicates ->
-	df. sort_values(by=['last','Fixst'])
	of the first 'First' 7 ascending folk. True
7	df.sort_values (by=['last', 'First'], ascending[false, True  First last Age inplace=True
	O Covey Barlog 45
	3 Adam Doe 16 1 Jane Doe 20
1	Other
	df.sort_index -> خاتمی الترتیب تانی در
1000	

	df['last'] sort values() > 3	Doe
1	1	Doe
	2	Doe
	0	Barlog
	df['Salary USD']. n Largest(3) -> 28	dex salarguso 307 20000 16 20000 219 20000
→ →	dfonlargest (3, 'Salary USD') -> 9051101	المعرض جميع البيانا الانامادة