

Unit 2

1. [Introducing Pandas](#)
2. [Reading files](#)
3. [Selecting data](#)
4. [Conditional selection](#)

1. Introducing Pandas



[Panda's documentation](#)

To begin we need to import pandas When you see pd, know it is referring to pandas

```
In [44]: import numpy as np
import pandas as pd
```

Pandas is a popular Python library used for working in tabular data (similar to the data stored in a spreadsheet).

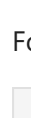
There are two main data structure used by pandas

- Series: equivalent to a vector or a list
- DataFrame: equivalent to a table.

Each column in a pandas DataFrame is a pandas Series data structure. We will mainly be looking at the DataFrame.

We can easily create a Pandas DataFrame by reading a .csv file

2. Reading files

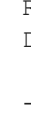


We will read the whole file at once using Pandas. Sometimes you might want to read the file line by line, and process each line. That's possible of course. See for example [here](#).

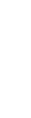
We will read data on [COVID-19 vaccinations](#)

In order to do that, I retrieved the raw data's url

Press on raw either here:



or here:



and retrieve the link:



```
In [2]: url = 'https://raw.githubusercontent.com/owid/covid-19-data/master/public/data/vaccinations/vaccinations.csv'
vacc_df = pd.read_csv(url)
```

read_csv has about 30 different options. See the [documentation](#)

For example, sep='t' is used for tab delimited files and 'usecol' reads only specific columns.

```
In [3]: type(vacc_df)
```

```
Out[3]: pandas.core.frame.DataFrame
```

view the shape of the dataframe:

```
In [47]: vacc_df.shape
```

```
Out[47]: (35163, 12)
```

view basic information:

```
In [5]: vacc_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 35163 entries, 0 to 35162
Data columns (total 12 columns):
 #   Column              Non-Null Count  Dtype  ---
 0   location            35163 non-null  object
 1   iso_code            35163 non-null  object
 2   date               35163 non-null  object
 3   total_vaccinations 20578 non-null  float64
 4   people_vaccinated   19765 non-null  float64
 5   people_fully_vaccinated 16719 non-null float64
 6   daily_vaccinations_raw 17397 non-null float64
 7   daily_vaccinations 34902 non-null float64
 8   total_vaccinations_per_hundred 20578 non-null float64
 9   people_vaccinated_per_hundred 19765 non-null float64
10   people_fully_vaccinated_per_hundred 16719 non-null float64
11   daily_vaccinations_per_million 34902 non-null float64
dtypes: float64(9), object(3)
memory usage: 3.2+ MB
```

```
In [6]: vacc_df.columns

Out[6]: Index(['location', 'iso_code', 'date', 'total_vaccinations', 'people_vaccinated', 'people_fully_vaccinated', 'daily_vaccinations_raw', 'daily_vaccinations', 'total_vaccinations_per_hundred'], dtype='object')
```

View the first few rows:

```
In [51]: vacc_df.head()
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations_per_hundred
0	Afghanistan	AFG	2021-02-22	0.0	0.0	NaN	NaN	NaN	NaN
1	Afghanistan	AFG	2021-02-23	NaN	NaN	NaN	NaN	NaN	1367.0
2	Afghanistan	AFG	2021-02-24	NaN	NaN	NaN	NaN	NaN	1367.0
3	Afghanistan	AFG	2021-02-25	NaN	NaN	NaN	NaN	NaN	1367.0
4	Afghanistan	AFG	2021-02-26	NaN	NaN	NaN	NaN	NaN	1367.0

What do you think that the 'tail' command does? Try it out!

What happens if we just type vacc_df, without a head or a tail?

```
In [53]: vacc_df.tail()
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations_per_hundred
35158	Zimbabwe	ZWE	2021-07-17	1771434.0	1132045.0	639389.0	39694.0	41958.0	20578.000000
35159	Zimbabwe	ZWE	2021-07-18	1785533.0	1144379.0	641154.0	14099.0	42019.0	20578.000000
35160	Zimbabwe	ZWE	2021-07-19	1827638.0	1184435.0	643203.0	42105.0	42253.0	20578.000000
35161	Zimbabwe	ZWE	2021-07-20	1897572.0	1247494.0	649843.0	69699.0	45971.0	20578.000000
35162	Zimbabwe	ZWE	2021-07-21	1949477.0	1292642.0	656830.0	52135.0	47976.0	20578.000000

```
In [9]: vacc_df.describe()
```

	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations_per_hundred
count	2.057800e+04	1.976500e+04	1.671900e+04	1.738700e+04	3.490200e+04	2.0578000000
mean	4.834277e+07	2.663580e+07	1.486042e+07	8.780545e+05	4.383342e+05	29.165589
std	2.276768e+08	1.267277e+08	6.475550e+07	3.552802e+06	2.474233e+06	35.025427
min	0.000000e+00	0.000000e+00	1.000000e+00	0.000000e+00	0.000000e+00	0.000000
25%	1.662121e+05	1.402700e+05	6.086300e+04	5.764000e+03	9.600000e+02	2.780000
50%	1.300552e+06	9.153020e+05	5.243060e+05	3.050300e+05	8.131000e+03	14.330000
75%	8.138526e+06	5.149939e+06	3.245194e+06	1.880345e+05	5.346750e+04	44.475000
max	3.787027e+09	2.092990e+09	1.048364e+09	4.794441e+07	4.329428e+07	232.650000

A summary of the functions so far:

- `pd.read_csv()` - Read data from a CSV file into a `Pandas DataFrame` object
- `df.info()` - View basic information about `Pandas` & data types
- `df.describe()` - View statistical information about numeric columns
- `df.columns` - Get the list of column names
- `df.shape` - Get the number of rows & columns as a tuple
- `df.head()`, `df.tail()` - View the beginning/end of the file

3. Selecting data

```
df = pd.read_csv('music.csv')
```

Pandas format is similar to a dictionary, not to a list

a list

```
In [18]: covid_data_list = [
{'date': '2020-08-30', 'new_cases': 1444, 'new_deaths': 1, 'new_tests': 53541},
{'date': '2020-09-01', 'new_cases': 1365, 'new_deaths': 4, 'new_tests': 42583},
{'date': '2020-09-01', 'new_cases': 996, 'new_deaths': 6, 'new_tests': 54395},
{'date': '2020-09-02', 'new_cases': 996, 'new_deaths': 8, 'new_tests': 54395},
{'date': '2020-09-03', 'new_cases': 1326, 'new_deaths': 6, 'new_tests': 54395}
]
covid_data_list
```

```
Out[18]: [{"date": "2020-08-30", "new_cases": 1444, "new_deaths": 1, "new_tests": 53541},
{"date": "2020-09-01", "new_cases": 1365, "new_deaths": 4, "new_tests": 42583},
{"date": "2020-08-31", "new_cases": 1365, "new_deaths": 4, "new_tests": 42583},
{"date": "2020-09-01", "new_cases": 996, "new_deaths": 6, "new_tests": 54395},
{"date": "2020-09-02", "new_cases": 996, "new_deaths": 8, "new_tests": 54395},
{"date": "2020-09-03", "new_cases": 1326, "new_deaths": 6, "new_tests": 54395}
]
```

a dictionary:

```
In [11]: covid_data_dict = {
'date': ['2020-08-30', '2020-09-01', '2020-09-01', '2020-09-02', '2020-09-03'],
'new_cases': [1444, 1365, 996, 975, 1326],
'new_deaths': [1, 4, 6, 8, 6],
'new_tests': [53541, 42583, 54395, 54395, 54395]
}
```

The index of a dataframe doesn't have to be numeric

```
In [54]: df = pd.DataFrame({'age': [30, 2, 12, 4, 32, 33, 69],
'color': ['blue', 'green', 'red', 'white', 'gray', 'black', 'red'],
'food': ['Steak', 'Lamb', 'Mango', 'Apple', 'Cheese', 'Melon', 'Beans'],
'height': [165, 70, 120, 80, 180, 172, 150],
'weight': [4.6, 8.3, 9.9, 3.3, 1.8, 9.5, 2.2],
'state': ['NY', 'TX', 'FL', 'AL', 'AK', 'TX', 'TX']
},
index=['Jane', 'Nick', 'Aaron', 'Penelope', 'Dean', 'Christina', 'Carmelia'])
df
```

```
Out[54]:   age  color  food  height  score  state
0     30  blue  Steak    165    4.6  NY
1     2   green  Lamb    70    8.3  TX
2     12  red  Mango    120    9.0  FL
3     4  white  Apple    80    3.3  AL
4     32  gray  Cheese   180    1.8  AK
5     33  black  Melon   172    9.5  TX
6     69  red  Bears   150    2.2  TX
```

In our file the index is numeric:

```
In [13]: vacc_df.head()
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations_per_hundred
0	Afghanistan	AFG	2021-02-22	0.0	0.0	NaN	NaN	NaN	NaN
1	Afghanistan	AFG	2021-02-23	NaN	NaN	NaN	NaN	NaN	1367.0
2	Afghanistan	AFG	2021-02-24	NaN	NaN	NaN	NaN	NaN	1367.0
3	Afghanistan	AFG	2021-02-25	NaN	NaN	NaN	NaN	NaN	1367.0
4	Afghanistan	AFG	2021-02-26	NaN	NaN	NaN	NaN	NaN	1367.0

return a single column as a Series:

note: using the `loc` notation is possible only for columns whose names do not contain spaces or special characters.

```
In [14]: vacc_df.location
vacc_df['location']
```

```
Out[14]: 0    Afghanistan
1    Afghanistan
2    Afghanistan
3    Afghanistan
4    Afghanistan
...
35158  Zimbabwe
35159  Zimbabwe
35160  Zimbabwe
35161  Zimbabwe
35162  Zimbabwe
Name: location, Length: 35163, dtype: object
```

return a single column as a dataframe:

```
In [15]: type(vacc_df.location)
```

```
Out[15]: pandas.core.series.Series
```

```
In [56]: vacc_df[['location']]
```

	location	date
0	Afghanistan	2021-02-22
1	Afghanistan	2021-02-23
2	Afghanistan	2021-02-24
3	Afghanistan	2021-02-25
4	Afghanistan	2021-02-26
...
35158	Zimbabwe	2021-07-17
35159	Zimbabwe	2021-07-18
35160	Zimbabwe	2021-07-19
35161	Zimbabwe	2021-07-20
35162	Zimbabwe	2021-07-21

35163 rows x 2 columns

```
In [58]: vacc_df.location[600]
```

```
Out[58]: 'Algeria'
```

retrieve two columns

```
In [18]: vacc_df[['location', 'date']]
```

	location	date
0	Afghanistan	2021-02-22
1	Afghanistan	2021-02-23
2	Afghanistan	2021-02-24
3	Afghanistan	2021-02-25
4	Afghanistan	2021-02-26
...
35158	Zimbabwe	2021-07-17
35159	Zimbabwe	2021-07-18
35160	Zimbabwe	2021-07-19
35161	Zimbabwe	2021-07-20
35162	Zimbabwe	2021-07-21

35163 rows x 2 columns

Selecting subsets of rows and columns

One way to do that is `iloc`.

`df.iloc` - selects subsets of rows and columns by integer location only

```
In [69]: #vacc_df.iloc[0] #first row
#vacc_df.iloc[1:3] #rows 1 and 2 as a series
vacc_df.iloc[3:5] #last row
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations_per_hundred
3	Afghanistan	AFG	2021-02-25	NaN	NaN	NaN	NaN	NaN	1367.0
4	Afghanistan	AFG	2021-02-26	NaN	NaN	NaN	NaN	NaN	1367.0

The operator

- when used alone it means "everything"
- also used to indicate a slice of values

```
In [28]: vacc_df.iloc[1:2, 2:22] # second and third row
vacc_df.iloc[1:2, 2:22] # a few specific rows

# Columns:
vacc_df.iloc[:, 0] # first column of data frame
vacc_df.iloc[:, 1] # second column of data frame
vacc_df.iloc[:, :] # all columns of data frame
```

```
Out[28]: #Rows and columns:
vacc_df.iloc[0:5, :] # first five rows of data frame
vacc_df.iloc[:, 0:2] # first two columns of data frame with all rows
vacc_df.iloc[0:3, 6:24], [0:5, 6:] # 1st, 4th, 7th, 25th row + 1st 6th 7th column.
```

	location	people_fully_vaccinated	daily_vaccinations_raw
0	Afghanistan	NaN	NaN
3	Afghanistan	NaN	NaN
6	Afghanistan	NaN	NaN
24	Afghanistan	NaN	NaN

```
In [21]: vacc_df.iloc[0:19]
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations_per_hundred
0	Afghanistan	AFG	2021-02-22	0.0	0.0	NaN	NaN	NaN	NaN
1	Afghanistan	AFG	2021-02-23	NaN	NaN	NaN	NaN	NaN	1367.0
2	Afghanistan	AFG	2021-02-24	NaN	NaN	NaN	NaN	NaN	1367.0
3	Afghanistan	AFG	2021-02-25	NaN	NaN	NaN	NaN	NaN	1367.0
4	Afghanistan	AFG	2021-02-26	NaN	NaN	NaN	NaN	NaN	1367.0
5	Afghanistan	AFG	2021-02-27	NaN	NaN	NaN	NaN	NaN	1367.0
6	Afghanistan	AFG	2021-02-28	8200.0	8200.0	NaN	NaN	NaN	1367.0
7	Afghanistan	AFG	2021-03-01	NaN	NaN	NaN	NaN	NaN	1580.0
8	Afghanistan	AFG	2021-03-02	NaN	NaN	NaN	NaN	NaN	1794.0
9	Afghanistan	AFG	2021-03-03	NaN	NaN	NaN	NaN	NaN	2008.0
10	Afghanistan	AFG	2021-03-04	NaN	NaN	NaN	NaN	NaN	2221.0
11	Afghanistan	AFG	2021-03-05	NaN	NaN	NaN	NaN	NaN	2435.0
12	Afghanistan	AFG	2021-03-06	NaN	NaN	NaN	NaN	NaN	2649.0
13	Afghanistan	AFG	2021-03-07	NaN	NaN	NaN	NaN	NaN	2862.0
14	Afghanistan	AFG	2021-03-08	NaN	NaN	NaN	NaN	NaN	2862.0
15	Afghanistan	AFG	2021-03-09	NaN	NaN	NaN	NaN	NaN	2862.0
16	Afghanistan	AFG	2021-03-10	NaN	NaN	NaN	NaN	NaN	2862.0
17	Afghanistan	AFG	2021-03-11	NaN	NaN	NaN	NaN	NaN	2862.0
18	Afghanistan	AFG	2021-03-12	NaN	NaN	NaN	NaN	NaN	2862.0

What if I want to select the 'daily_vaccinations' column, but I don't remember which column it is?

Use `df.loc`

`df.loc` - selects subsets of rows and columns by label only. Allowed inputs are:

- A single label, e.g. 5 or 'a', (note that 5 is interpreted as a label of the index, and never as an integer position along the index).
- A list or array of labels, e.g. ['a', 'b', 'c'].
- A slice object with labels, e.g. 'a':'f'.

```
In [75]: vacc_df.loc[2:3, ['daily_vaccinations', 'date']]
```

	daily_vaccinations	date
2	1367.0	2021-02-24
3	1367.0	2021-02-25

I'm missing the location. Let's add it:

```
In [23]: vacc_df.loc[0:3, ['location', 'daily_vaccinations']]

Out[23]:   location  daily_vaccinations
0  Afghanistan          1367.0
1  Afghanistan          1367.0
2  Afghanistan          1367.0
3  Afghanistan          1367.0
```

Semantics are similar to `loc`. But note:

- `df.iloc` excludes the last element. `df.iloc[0:10000]` will return entries 0..9999
- `df.loc` includes the last element. `df.loc[0:10000]` will return entries 0..10000

you try it! What is the difference between:

```
vacc_df.iloc[0:5]
vacc_df.loc[0:5]
```

```
In [24]: vacc_df.iloc[0:5]
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations_per_hundred
0	Afghanistan	AFG	2021-02-22	0.0	0.0	NaN	NaN	NaN	NaN
1	Afghanistan	AFG	2021-02-23	NaN	NaN	NaN	NaN	NaN	1367.0
2	Afghanistan	AFG	2021-02-24	NaN	NaN	NaN	NaN	NaN	1367.0
3	Afghanistan	AFG	2021-02-25	NaN	NaN	NaN	NaN	NaN	1367.0
4	Afghanistan	AFG	2021-02-26	NaN	NaN	NaN	NaN	NaN	1367.0

```
In [25]: vacc_df.loc[0:5]
```

	location	iso_code	date	total_vaccinations	people_vaccinated
--	----------	----------	------	--------------------	-------------------


```
Out[36]: location      433
iso_code      433
date          433
total_vaccinations      376
people_vaccinated      431
people_fully_vaccinated 399
daily_vaccinations_raw  371
daily_vaccinations      431
total_vaccinations_per_hundred 376
people_vaccinated_per_hundred 431
people_fully_vaccinated_per_hundred 399
daily_vaccinations_per_million 431
dtype: int64
```

At the end of the file we have some world data.

Use str.contains if you're not sure how this location is called

```
In [86]: vacc_df[vacc_df['location'].str.contains('World')]

Out[86]:
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations
34612	World	OWID_WRL	2020-12-02	0.000000e+00	0.000000e+00	NaN	NaN	NaN	NaN
34613	World	OWID_WRL	2020-12-03	0.000000e+00	0.000000e+00	NaN	0.0	0.0	0.0
34614	World	OWID_WRL	2020-12-04	1.000000e+00	1.000000e+00	NaN	1.0	0.0	0.0
34615	World	OWID_WRL	2020-12-05	1.000000e+00	1.000000e+00	NaN	0.0	0.0	0.0
34616	World	OWID_WRL	2020-12-06	1.000000e+00	1.000000e+00	NaN	0.0	0.0	0.0
...
34840	World	OWID_WRL	2021-07-19	3.662774e+09	2.055135e+09	1.015653e+09	19296764.0	30735269.0	30735269.0
34841	World	OWID_WRL	2021-07-20	3.695078e+09	2.067338e+09	1.025288e+09	32304332.0	31101056.0	31101056.0
34842	World	OWID_WRL	2021-07-21	3.727059e+09	2.079488e+09	1.033198e+09	31980155.0	30874061.0	30874061.0
34843	World	OWID_WRL	2021-07-22	3.760632e+09	2.087375e+09	1.043255e+09	33573399.0	30413169.0	30413169.0
34844	World	OWID_WRL	2021-07-23	3.787027e+09	2.092904e+09	1.048364e+09	26394993.0	29076881.0	29076881.0

233 rows × 10 columns

Remove the world data:

```
In [38]: vacc_df_noWorld = vacc_df.loc[vacc_df.location != 'World']
vacc_df_noWorld.tail()

Out[38]:
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations
35158	Zimbabwe	ZWE	2021-07-17	1771434.0	1132045.0	639389.0	39694.0	41958.0	41958.0
35159	Zimbabwe	ZWE	2021-07-18	1785533.0	1144379.0	641154.0	14099.0	42019.0	42019.0
35160	Zimbabwe	ZWE	2021-07-19	1827638.0	1184435.0	643203.0	42105.0	42253.0	42253.0
35161	Zimbabwe	ZWE	2021-07-20	1897337.0	1247494.0	649843.0	69699.0	45971.0	45971.0
35162	Zimbabwe	ZWE	2021-07-21	1949472.0	1292642.0	656830.0	52135.0	47976.0	47976.0

Find the country with the maximum vaccinations:

```
In [87]: max_vacc = vacc_df_noWorld['total_vaccinations'].max()
max_vacc

Out[87]: 2383659444.0
```

```
In [40]: vacc_df_noWorld.loc[vacc_df_noWorld.total_vaccinations == max_vacc]

Out[40]:
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations
15671	Asia	OWID_AS	2021-07-22	2.383659e+09	1.246089e+09	461410827.0	19637719.0	19995909.0	19995909.0

What do you think this function does?

```
In [41]: vacc_df_noWorld.total_vaccinations.mean()

Out[41]: 36918988.861587614
```

Your turn:

Select the number of daily vaccinations in Israel on date 2021-02-06 (hint: use &)

```
In [88]: vacc_df.loc[(vacc_df.location == 'Israel') & (vacc_df.date == '2021-02-20')]

Out[88]:
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations
15671	Israel	ISR	2021-07-20	11008624.0	5750067.0	5258557.0	14231.0	10060.0	10060.0

Find all the countries with more than 3000000 vaccinations

```
In [43]: vacc_df_noWorld.loc[(vacc_df_noWorld.daily_vaccinations > 300000)]

Out[43]:
```

	location	iso_code	date	total_vaccinations	people_vaccinated	people_fully_vaccinated	daily_vaccinations_raw	daily_vaccinations	total_vaccinations
224	Africa	OWID_AFR	2021-03-27	1.030312e+07	6962741.0	3606163.0	325837.0	301268.0	301268.0
225	Africa	OWID_AFR	2021-03-28	1.056197e+07	7218927.0	3608826.0	258849.0	335169.0	335169.0
226	Africa	OWID_AFR	2021-03-29	1.077488e+07	7305693.0	3743035.0	212915.0	348122.0	348122.0
227	Africa	OWID_AFR	2021-03-30	1.097167e+07	7387045.0	3870510.0	196785.0	355345.0	355345.0
267	Africa	OWID_AFR	2021-05-09	2.119979e+07	16063514.0	5518752.0	729471.0	351787.0	351787.0
...
33674	Upper middle income	OWID_UMC	2021-07-18	1.928636e+09	943605954.0	379618270.0	11998743.0	16733054.0	16733054.0
33675	Upper middle income	OWID_UMC	2021-07-19	1.943552e+09	945911507.0	381584744.0	14915465.0	16626595.0	16626595.0
33676	Upper middle income	OWID_UMC	2021-07-20	1.961319e+09	950175010.0	384027832.0	17767193.0	16728652.0	16728652.0
33677	Upper middle income	OWID_UMC	2021-07-21	1.980014e+09	953111895.0	386240484.0	18694961.0	16468485.0	16468485.0
33678	Upper middle income	OWID_UMC	2021-07-22	2.000667e+09	955873275.0	388296478.0	20622943.0	16617879.0	16617879.0

3614 rows × 10 columns

Summary of the functions in this unit:

- `.index.values` - the row indexes of this part of the dataframes
- `.str.contains` - selects rows and columns that contain a string
- `.max` - maximum value
- `.mean` - average value
- `.count` - the number of rows that contain a value
- `.len()` - dataframe length