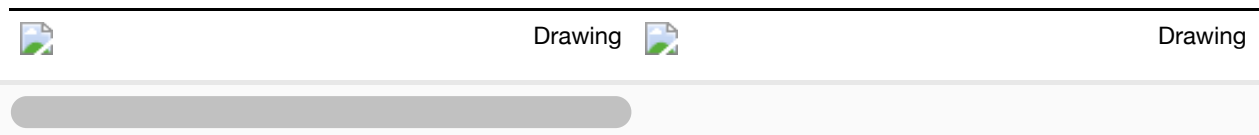


# Pandas

`pandas` is the Python library that structures and simplifies data manipulation and analysis. The name, `pandas` is derived by **panel** and **data**. The library indeed focusses on representing data as tables ( `DataFrames` ), indices ( `Index` ) and data series ( `Series` ).



`pandas` uses `NumPy` arrays to represent data. It provides specific functions to manipulated data as `code` objects . This means that to use `pandas` you need to import also `NumPy`. `pandas` requires and builds on top of `NumPy` .

There are three main data object types in `pandas`:

- `Series` - A mutable, one-dimensional array of indexed data.
- `DataFrames` - A two-dimensional, size-mutable, potentially heterogeneous tabulated data structure.
- `Index` - An one-dimensional, immutable array or ordered set (technically a multi-set, as `Index` objects may contain repeated values).

Below we will learn a little bit more about these three data objects.

## ***A short history of Pandas.***

Wes McKinney started developing what then became `pandas` while working at the capital management firm Applied Quantitative Research (AQR). `Pandas` was developed initially as a closed-source project and was made open source in 2009. `Pandas` is sponsored by [NumFOCUS, Inc.](https://numfocus.org/) (<https://numfocus.org/>), that promotes support and sponsorships of python based open source code.

```
In [1]: import numpy as np
import pandas as pd
```

## **The pandas Series object**

`Series` are dictionary-like objects. `Pandas Series` is a one-dimensional array that comes with labels assigned. They similar to `NumPy` one-dimensional arrays but they are labeled or indexed. So they are built on top of `NumPy` arrays but come with additional functionality as well as constrains. They can be thought of as a specification of `NumPy` arrays. For example, let's evaluate the code below.

First we define a `pandas Series` objects and assign a set of string values to the elements of the

object:

```
In [2]: data = pd.Series(['a', 'b', 'c', 'd'])
```

After defining the object let's explore it.

```
In [3]: print (data)
```

```
0    a
1    b
2    c
3    d
dtype: object
```

The data type is defined as `object`, the values we assigned are stored in the second column. The first column is the index automatically assigned by pandas to each value. Now, each value comes with an index. This is a fundamental data organization aspect of pandas. Values always have indices, because pandas deals with panel data, i.e., tables. So even a one-dimensional array as a Series is assigned indices just like a table is.

We will discuss pandas index objects later.

A Pandas Series can be created directly by assigning values into an array (using `[]`), and that array to a series (`pd.Series()`). Yet, the final product of that series definition creates something different than an Array. It creates a set of pairs of values where a label is associated to a corresponding value.

Series are capable of holding data of any type (integer, string, float, python objects, etc.). For example:

```
In [4]: data = pd.Series(['string', 10, np.nan, 0.01])
print(data)
```

```
0    string
1         10
2         NaN
3         0.01
dtype: object
```

Once the Series object is created, the Index is created and it becomes part of the methods of the Series object. This means that the Index to the entries can be retrieved and used to address the corresponding entry. For example, we can call the `Index` as a method and it will return the range, with the min value, the max value and the steps in between them:

```
In [5]: data.index
```

```
Out[5]: RangeIndex(start=0, stop=4, step=1)
```

The index can be used to address the corresponding value in the Series object:

```
In [6]: data[0:3]
```

```
Out[6]: 0    string
        1      10
        2     NaN
        dtype: object
```

Because the Index in pandas is explicitly defined and it becomes a method (this is in contrast with NumPy arrays where indices are implicitly defined and hence not addressable or callable), the pandas Series object becomes much more useful and structured.

For example, we can create a Series object by explicitly assigning values and indices:

```
In [7]: data1 = pd.Series(['string',10, np.nan,0.01],index=[0,3,2,4])
        print(data1)
```

```
0    string
3      10
2     NaN
4    0.01
dtype: object
```

Even though `data` and `data1` might look the same, they are not. The indices are different:

```
In [8]: print(data.index)
        print(data1.index)
```

```
RangeIndex(start=0, stop=4, step=1)
Int64Index([0, 3, 2, 4], dtype='int64')
```

So, that, if we address the second entry in either object we will get different values:

```
In [9]: print(data[3])
        print(data1[3])
```

```
0.01
10
```

One way to think about pandas Series objects is that they are dictionaries where a label is paired to a value, say label 1 is assigned to value 1, or label 2 to value 3 etc. Indeed, a

pandas Series object can be constructed by hand building a dictionary a set of pairing of labels and values.

For example, let's build a [Python dictionary](https://docs.python.org/3/tutorial/datastructures.html#dictionaries).

(<https://docs.python.org/3/tutorial/datastructures.html#dictionaries>) of cognitive health. The dictionary pairs a label to a value:

```
In [10]: cognitive_health = {'happyness':10,  
                             'language': 2,  
                             'energy': 5,  
                             'memory': 3}
```

Note that python dictionaries have attributes (you can type `cognitive_health.` and press tab twice to get a list of attributes), yet the python dictionary does not have `index` as an attribute.

The following line will return an error.

```
In [12]: cognitive_health.index
```

```
-----  
--  
AttributeError                                Traceback (most recent call las  
t)  
Input In [12], in <cell line: 1>()  
----> 1 cognitive_health.index  
  
AttributeError: 'dict' object has no attribute 'index'
```

Python dictionaries can be used to set pandas Series directly:

```
In [13]: cognitive_health_Series = pd.Series(cognitive_health)
```

The dictionary has been made into a panda Series and it is now ordered and labelled. So, the following operation will not return an error, but the labels of the Series (the indices):

```
In [14]: cognitive_health_Series.index
```

```
Out[14]: Index(['happyness', 'language', 'energy', 'memory'], dtype='object')
```

Another important property that the pandas Series have and python dictionaries do not is the ability to allow slicing. Whereas, a dictionary would return an error if called as follows:

```
In [15]: cognitive_health['happyness':'energy']
```

```
-----
--
TypeError                                Traceback (most recent call last)
Input In [15], in <cell line: 1>()
----> 1 cognitive_health['happyness':'energy']

TypeError: unhashable type: 'slice'
```

A pandas Series do allow slicing:

```
In [16]: cognitive_health_Series['happyness':'energy']
```

```
Out[16]: happyness      10
         language       2
         energy         5
         dtype: int64
```

To summarize what we have learned about pandas Series:

- They are ordered and labelled one-dimensional arrays.
- They can be populated by assigning
  - values only (indices are automatically assigned): `series = pd.Series(['a','b','c'])`
  - values and indices explicitly: `series = pd.Series(['a','b','c'],index = [1,3,2])`
  - python dictionaries directly: `dic = {1:'a',2:'b',3:'c'}, series = pd.Series(dic)`
- They always come with the property `index`

## Pandas index object

Let's briefly discuss the `index` object in pandas.

An `Index` is the pandas object that hosts information regarding the ordering and labels of the arrays inside other objects such as `Series` and `DataFrames`.

Index objects also have many of the attributes familiar from NumPy arrays (e.g., `shape`, `dimensions`, `size`, etc):

```
In [17]: ind = pd.Index([1, 2, 3, 4, 5])

print(ind.size, ind.shape, ind.ndim, ind.dtype)

print(ind)
```

```
5 (5,) 1 int64
Int64Index([1, 2, 3, 4, 5], dtype='int64')
```

Pandas Index objects are immutable.

An index is similar to a NumPy array, but it is immutable. This means that once an Index is defined the values inside the index cannot be changed.

Where arrays can be modified after definition, a pandas index cannot. Let's try this. Above we defined `ind` as a pandas index and set in the third position the value `3`.

Let's try to change that value, evaluate the following operation in which we attempt to set the value `10` in the third position of `ind`:

```
In [18]: ind[2] = 10
```

```
-----
--
TypeError                                Traceback (most recent call last)
Input In [18], in <cell line: 1>()
----> 1 ind[2] = 10

File ~/opt/anaconda3/lib/python3.8/site-packages/pandas/core/indexes/base.py:5021, in Index.__setitem__(self, key, value)
    5019 @final
    5020 def __setitem__(self, key, value):
-> 5021     raise TypeError("Index does not support mutable operations")

TypeError: Index does not support mutable operations
```

The error should return the following line at the end:

```
TypeError: Index does not support mutable operations
```

A pandas Index does not allow changes. This is helpful. One way to think about the index is that it is a specialized NumPy array. Specialized means that it has a more narrow scope than the more general goal of a NumPy array. The scope is that to define the (ahem) index of a data frame. Because of this scope (store an index) changes to the values of the array are not allowed, in other words the Index is immutable, or cannot be changed after definition by assigning a different value to any of its elements.

If we think about it, this makes sense. Changing a value inside an index of a data frame would change the definition of the data frame and really invalidating the purpose of the index and of the data frame.

Pandas `Index` objects are designed to facilitate operations on the array and serve the task of keeping track of positions of data entries in the objects.

They support and facilitate operations such as joining datasets. Because of this the pandas index object follows many operations of the built in python datatype `set` (<https://docs.python.org/3/library/stdtypes.html#set-types-set-frozenset>).

Because of this the `Index` object allows many of operations also served by Python set data structure, such as unions, intersections, differences, and other combinations can be computed in a familiar way.

For example, two pandas index object can be united. This means that the unique indices are combined:

```
In [19]: index1 = pd.Index([1, 2, 3, 4, 5])
index2 = pd.Index([1, 3, 4, 6, 20])

index1.union(index2)
```

```
Out[19]: Int64Index([1, 2, 3, 4, 5, 6, 20], dtype='int64')
```

Other logical operations can be performed using pandas index objects, for example intersection (find the common elements):

```
In [20]: index1.intersection(index2)
```

```
Out[20]: Int64Index([1, 3, 4], dtype='int64')
```

In sum, pandas index objects allow performing slicing, indexing operations and facilitate keeping track of, ahem, the indices.

## Pandas DataFrame objects

Whereas pandas `Index` and `Series` objects are the backbone of the pandas library, the `DataFrame` object is the workhorse of the library. DataFrames have effectively made pandas the library for data science.

The `DataFrame` is an extension of the `Series` object. It is a 2-dimensional, mutable array that it can be conceptualized as either a more general NumPy array, or a specialized Python dictionary.

A `DataFrame` can be defined most simply by setting some values to it. The indexing and labelling is automatically assigned by pandas. For example, we can create a one-dimensional list and assign the values of the list to a `DataFrame`:

```
In [21]: list = ['a', 'b', 'c', 'd']  
data_frame1 = pd.DataFrame(list)  
print(data_frame1)
```

```
0  
0  a  
1  b  
2  c  
3  d
```

Even though the result might seem very similar to that obtained about with the Series object, in reality, the DataFrame object has assigned not one but two labels, one label for the rows dimension ( `[ 0 : 3 ]` ) and one for the column dimension ( `0` ).

Indeed, we can notice from the `print()` output that two dimensions were automatically labelled (indexed) with numbers. Whereas the Series is created as a one-dimensional object, the DataFrame is created as a two-dimensional object.

Another way to think about a DataFrame is that it is a sequence of *aligned* Series objects. What do we mean by that?

Each column in the DataFrame can be thought of as a Series. Each series is labelled by the column index. Importantly, the series are aligned, this means that the set of Series (the columns of the DataFrame) are indexed by a common `Index` object.

Let's take a look at all this.

Let's construct a new Series object similar to `cognitive_health_Series`, the object we created above (make sure that object is still in memory, if needed re-run that section of the cells). Let's assume that that original series represented the data collected on a subject.

```
In [22]: cognitive_health_Series
```

```
Out[22]: happiness    10  
language             2  
energy               5  
memory               3  
dtype: int64
```

The new Series will represent data on a second subject. To build the series we will use steps similar to the ones used above. We will first make a python dictionary and then make a pandas Series object out of it.



```
In [23]: cognitive_health_subject2_dict = {'happyness': 15,
                                           'language': 4,
                                           'energy': 9,
                                           'memory': 6}
cognitive_health_subject2_series = pd.Series(cognitive_health_subject2_dict)
cognitive_health_subject2_series
```

```
Out[23]: happyness    15
language      4
energy        9
memory        6
dtype: int64
```

Now that we have two pandas Series, we can construct a pandas DataFrame by combining the two Series objects.

To do so, we will first create a single dictionary containing the two series, labelled as `subject 1` and `subject 2` and then use that dictionary to make a DataFrame.

The dictionary for the DataFrame is formed by assigning each subject to a label:

```
In [24]: dict = {'subject 1' : cognitive_health_Series, 'subject 2' : cognitive_health_Series}
print(dict)
```

```
{'subject 1': happyness    10
language      2
energy        5
memory        3
dtype: int64, 'subject 2': happyness    15
language      4
energy        9
memory        6
dtype: int64}
```

The dictionary can now be used to build a DataFrame:

```
In [25]: sample = pd.DataFrame(dict)
sample
```

```
Out[25]:
```

	subject 1	subject 2
<b>happyness</b>	10	15
<b>language</b>	2	4
<b>energy</b>	5	9
<b>memory</b>	3	6

Excellent. Pandas did its Kung Fu. The dictionary comprising two pandas Series, was organized

into a pandas DataFrame.

Next try adding two more subjects to the same DataFrame ( `sample` ). Let's practice this with a subject that has all values higher than subject 2 by a value of 3 (just add 3) and another subject that has all values lower than subject 2 by a value of 2 (just subtract 3).

```
In [26]: cognitive_health_subject4_dict = {'happyness': 15-2,  
                                           'language': 4-2,  
                                           'energy': 9-2,  
                                           'memory': 6-2}  
cognitive_health_subject4_series = pd.Series(cognitive_health_subject4_dict)  
cognitive_health_subject4_series
```

```
Out[26]: happyness    13  
         language     2  
         energy       7  
         memory       4  
         dtype: int64
```

```
In [27]: cognitive_health_subject3_dict = {'happyness': 15+3,  
                                           'language': 4+3,  
                                           'energy': 9+3,  
                                           'memory': 6+3}  
cognitive_health_subject3_series = pd.Series(cognitive_health_subject3_dict)  
cognitive_health_subject3_series
```

```
Out[27]: happyness    18  
         language     7  
         energy      12  
         memory       9  
         dtype: int64
```

```
In [28]: dict = {'subject 1' : cognitive_health_Series,
                'subject 2' : cognitive_health_subject2_series,
                'subject 3' : cognitive_health_subject3_series,
                'subject 4' : cognitive_health_subject4_series,
                }
print(dict)
```

```
{'subject 1': happiness      10
 language      2
 energy        5
 memory        3
 dtype: int64, 'subject 2': happiness      15
 language      4
 energy        9
 memory        6
 dtype: int64, 'subject 3': happiness      18
 language      7
 energy       12
 memory        9
 dtype: int64, 'subject 4': happiness      13
 language      2
 energy        7
 memory        4
 dtype: int64}
```

```
In [29]: sample = pd.DataFrame(dict)
sample
```

Out[29]:

	subject 1	subject 2	subject 3	subject 4
<b>happiness</b>	10	15	18	13
<b>language</b>	2	4	7	2
<b>energy</b>	5	9	12	7
<b>memory</b>	3	6	9	4

The newly created DataFrame has various attributes that we can explore. We can address and extract the columns for example:

```
In [30]: cols = sample.columns
print(cols)
```

```
Index(['subject 1', 'subject 2', 'subject 3', 'subject 4'], dtype='object')
```

We can address the rows, which in technical terms are called the labels or the index:

```
In [31]: rows = sample.index  
print(rows)
```

```
Index(['happyness', 'language', 'energy', 'memory'], dtype='object')
```

## Summary

We have learned about the library `pandas` and the three fundamental objects of the library:

- Index
- Series
- DataFrames

These last one are the most used, versatile data representation objects and are most commonly used for data science projects.

Pandas DataFrames are 2-dimensional objects composed by collections of one-dimensional objects called Series. The one-dimensional objects in a DataFrame are `aligned` meaning that all entries of a series are matched, they are related, each row is indexed by the same index no matter what series.

## Exercise

To practice with pandas Series and DataFrames, build a new dataset that has 10 series (subject 1-10) representing measurements from the following measures of brain health:

```
['neuronal activity', 'blood oxygenation', 'blood pulsation rate',  
'cortical thickness']
```

```
In [32]: idx = ['neuronal activity', 'blood oxygenation', 'blood pulsation rate', 'cortical thickness']
measurements = [.5, 1, 1.5, 2]
nseries = 10
data_dict = {}
label = "subject"
for i in range(1, nseries + 1):
    name = label + str(i)
    data_dict[name] = pd.Series([val * i for val in measurements], index=idx)
df = pd.DataFrame(data_dict)
print(df)
```

	subject1	subject2	subject3	subject4	subject5	\
neuronal activity	0.5	1.0	1.5	2.0	2.5	
blood oxygenation	1.0	2.0	3.0	4.0	5.0	
blood pulsation rate	1.5	3.0	4.5	6.0	7.5	
cortical thickness	2.0	4.0	6.0	8.0	10.0	

	subject6	subject7	subject8	subject9	subject10
neuronal activity	3.0	3.5	4.0	4.5	5.0
blood oxygenation	6.0	7.0	8.0	9.0	10.0
blood pulsation rate	9.0	10.5	12.0	13.5	15.0
cortical thickness	12.0	14.0	16.0	18.0	20.0

```
In [33]: df.transpose()
```

Out[33]:

	neuronal activity	blood oxygenation	blood pulsation rate	cortical thickness
subject1	0.5	1.0	1.5	2.0
subject2	1.0	2.0	3.0	4.0
subject3	1.5	3.0	4.5	6.0
subject4	2.0	4.0	6.0	8.0
subject5	2.5	5.0	7.5	10.0
subject6	3.0	6.0	9.0	12.0
subject7	3.5	7.0	10.5	14.0
subject8	4.0	8.0	12.0	16.0
subject9	4.5	9.0	13.5	18.0
subject10	5.0	10.0	15.0	20.0

```
In [ ]:
```