



Data Science Lab

Pandas

Andrea Pasini Flavio Giobergia Elena Baralis

DataBase and Data Mining Group



Introduction to Pandas





Pandas

- Provides useful data structures (Series and DataFrames) and data analysis tools
- Based on Numpy arrays

Tools:

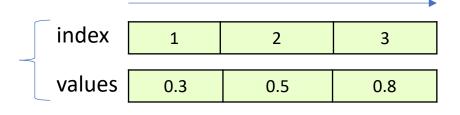
- Managing tables and series
 - data selection
 - grouping, pivoting
- Managing missing data
- Statistics on data







- Series: 1-Dimensional sequence of homogeneous elements
- Elements are associated to an explicit index
 - index elements can be either strings or integers
- Examples:



index	'3-July'	'4-July'	'5-July'
values	0.3	0.5	0.8







Creation from list



When not specified, index is set automatically with a progressive number

```
In [1]: import pandas as pd
    s1 = pd.Series([2.0, 3.1, 4.5])
    print(s1)
```

```
Out[1]: 0 2.0
1 3.1
2 4.5
```







Creation from list, specifying index









Creation from dictionary





```
In [1]: pd.Series({'c':2.0, 'b':3.1, 'a':4.5})
```

```
Out[1]: 'c' 2.0 'b' 3.1 'a' 4.5
```







Obtaining values and index from a Series



```
In [1]: s1 = pd.Series([2.0, 3.1, 4.5], index=['mon', 'tue', 'wed'])
    print(s1.values) # Numpy array
    print(s1.index)

Out[1]: [2.0, 3.1, 4.5]
    Index(['mon', 'tue', 'wed'], dtype='object')
```

Index is a custom Python object defined in Pandas





- Accessing Series elements
- Access by Index
 - Explicit: the one specified while creating a Series
 - Use the Series.loc attribute
 - Implicit: number associated to the element order (similarly to Numpy arrays)
 - Use the Series.iloc attribute



In [1]:

Pandas Series





Accessing Series elements



```
print(s1.loc['a']) # With explicit index
         print(s1.iloc[0]) # With implicit index
         s1.loc['b'] = 10  # Allows editing values
         print(f"Series:\n{s1}")
Out[1]:
         2.0
         2.0
         Series:
         'a'
               2.0
         'b'
               10
         'c'
               4.5
```

s1 = pd.Series([2.0, 3.1, 4.5], index=['a', 'b', 'c'])







Accessing Series elements: slicing



```
In [1]:
    s1 = pd.Series([2.0, 3.1, 4.5], index=['a', 'b', 'c'])
    print(s1.loc['b':'c']) # explicit index (stop element included)
    print(s1.iloc[1:3]) # implicit index (stop element excluded)
```

```
Out[1]: b 3.1 c 4.5 b 3.1 c 4.5
```







Accessing Series elements: masking



```
In [1]: s1 = pd.Series([2.0, 3.1, 4.5], index=['a', 'b', 'c'])
    print(s1[(s1>2) & (s1<10)])</pre>
```

```
Out[1]: b 3.1 c 4.5
```







Accessing Series elements: fancy indexing



```
In [1]:
    s1 = pd.Series([2.0, 3.1, 4.5], index=['a', 'b', 'c'])
    print(s1.loc[['a', 'c']])
    print(s1.iloc[[0, 2]])
```

```
Out[1]: a 2.0 c 4.5 a 2.0 c 4.5 c 4.5
```







- DataFrame: 2-Dimensional array
 - Can be thought as a table where columns are
 Series objects that share the same index
 - Each column has a name

Index	'Price'	'Quantity'	'Liters'
'Water'	1.0	5	1.5
'Beer'	1.4	10	0.3
'Wine'	5.0	8	1







Creation from Series



Use a dictionary to set column names

```
Out[1]:
                       Ouantity
             Price
                                   Liters
               1.0
                              5
                                      1.5
          а
          h
               1.4
                             10
                                      0.3
               5.0
                              8
                                      1.0
          C
```







Creation from dictionary of key-list pairs



- Each value (list) is associated to a column
 - Column name given by the key
- Index is automatically set to a progressive number
 - Unless explicitly passed as parameter (index=...)

```
In [1]: dct = { "c1": [0, 1, 2], "c2": [0, 2, 4] }
    df = pd.DataFrame(dct)
    print(df)
```











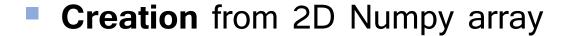
- Each dictionary is associated to a row
- Index is automatically set to a progressive number
 - Unless explicitly passed as parameter (index=...)

```
In [1]: dic_list = [{'c1':i, 'c2':2*i} for i in range(3)]
    df = pd.DataFrame(dic_list)
    print(df)
```











```
Out[1]: c1 c2
a 0 1
b 2 3
c 4 5
```

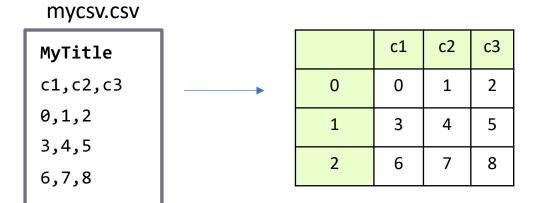






- Load DataFrame from csv file
 - Allows specifying the column delimiter (sep)
 - Automatically read header from first line of the file (after skipping the specified number of rows)
 - Column data types are inferred

```
df = pd.read_csv('./mycsv.csv', sep=',', skiprows=1)
```









- Load DataFrame from csv file
 - If it contains **null** values, you can specify how to recognize them
 - Empty columns are converted to "NaN" (Not a Number)
 - Using np.nan (NumPy's representation of NaN)
 - The string 'NaN' is automatically recognized

c1,c2,c3
0,no info,
3,4,5
6,x,NaN

	c1	c2	c3 <u>/</u>
0	0	NaN	NaN
1	3	4.0	5.0
2	6	NaN	NaN

type(np.nan) > float,
hence cz and c3 are floats







Save DataFrame to csv

	c1	c2	c3
0	0	NaN	2
1	3	4	5
2	6	NaN	NaN

savedcsv.csv

Use index=False to avoid writing the index



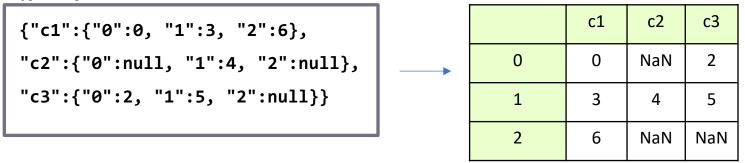




Load DataFrame from json file

```
df = pd.read_json('./myjson.json')
```

myjson.json



Use pd.to_json(path) to save a DataFrame in json format







- Many other data types are supported
 - Excel, HTML, HDF5, SAS, ...
- Check the pandas documentation
 - https://pandas.pydata.org/pandasdocs/stable/user guide/io.html







Obtaining column names and index from a DataFrame

Index	Price	Quantity	Liters
а	1.0	5	1.5
b	1.4	10	0.3
С	5.0	8	1







Accessing DataFrame data

Get a 2D Numpy array

Index	Price	Quantity	Liters
а	1.0	5	1.5
b	1.4	10	0.3
С	5.0	8	1







- Accessing DataFrames
 - Access a DataFrame column
 - Access rows and columns with indexing
 - df.loc
 - Explicit index
 - Slicing, masking, fancy indexing
 - df.iloc
 - Implicit index
- Whether a copy or view will be returned it depends on the context
 - Usually it is difficult to make assumptions
 - https://pandas-docs.github.io/pandas-docstravis/user guide/indexing.html







Accessing DataFrame columns





Index	Price	Quantity	Liters
а	1.0	5	1.5
b	1.4	10	0.3
С	5.0	8	1

In [1]: df['Quantity']

Out[1]: a 5
 b 10
 c 8







Accessing single DataFrame row by index



- loc (explicit), iloc (implicit)
- Return a Series with an element for each column

```
In [1]:
           print(df.loc['a'])
                                        # Get the first row (explicit)
           print(df.iloc[0])
                                        # Get the first row
Out[1]:
           Price
                    1.0
           Quantity 5.0
           Liters
                    1.5
           Price
                    1.0
           Quantity 5.0
           Liters
                    1.5
```







Accessing DataFrames with slicing



Allows selecting rows and columns







Accessing DataFrames with masking



Select rows based on a condition

Index	Price	Quantity	Liters
а	1.0	5	1.5
b	1.4	10	0.3
С	5.0	8	1

```
In [1]: mask = (df['Quantity']<10) & (df['Liters']>1)
    df.loc[mask, 'Quantity':] # Use masking and slicing
```

Out[1]: Quantity Liters
a 5 1.5







Accessing DataFrames with fancy indexing



To select columns...

Index	Price	Quantity	Liters
а	1.0	5	1.5
b	1.4	10	0.3
С	5.0	8	1

```
In [1]: mask = (df['Quantity']<10) & (df['Liters']>1)
    df.loc[mask, ['Price', 'Liters']] # Use masking and fancy
```

Out[1]: Price Liters
a 1.0 1.5







Accessing DataFrames with fancy indexing



To select **rows** and **columns**...

Index	Price	Quantity	Liters
а	1.0	5	1.5
b	1.4	10	0.3
С	5.0	8	1

```
In [1]:
           df.loc[['a', 'c'], ['Price','Liters']]
```

1.5

Out[1]: Price Liters 1.0 a

1.0 5.0 C







Assign value to selected items

Index	Price	Quantity	Liters
а	0.0	5	0.0
b	1.4	10	0.3
С	0.0	8	0.0







Add new column to DataFrame

DataFrame is modified inplace

Index	Price	Quantity	Liters
а	0.0	5	0.0
b	1.4	10	0.3
С	0.0	8	0.0

Index	Price	Quantity	Liters	Available
а	1.0	5	1.5	True
b	1.4	10	0.3	False
С	5.0	8	1	True

If the DataFrame already has a column with the specified name, then this is replaced







Add new column to DataFrame

It is also possible to assign directly a list

Index	Price	Quantity	Liters
а	0.0	5	0.0
b	1.4	10	0.3
С	0.0	8	0.0

Index	Price	Quantity	Liters	Available
а	1.0	5	1.5	True
b	1.4	10	0.3	False
С	5.0	8	1	True







Drop column(s)

- Returns a copy of the updated DataFrame
 - Unless inplace=True, in which case the original DataFrame is modified
 - This applies to many pandas methods -- always check the documentation!

Index	Price	Quantity	Liters/	Available
а	1.0	5	1.5	True
b	1.4	10	Ø.3	False
С	5.0	8	1	True







Rename column(s)

- Use a dictionary which maps old names with new names
- Returns a copy of the updated DataFrame

Index	Price	Quantity	Liters	Available
а	1.0	5	1.5	True
b	1.4	10	0.3	False
С	5.0	8	1	True

Index	Price	nItems	[L]	Available
а	1.0	5	1.5	True
b	1.4	10	0.3	False
С	5.0	8	1	True







- Unary operations on Series and DataFrames
 - exponentiation, logarithms, ...
- Operations between Series and DataFrames
 - Operations are performed element-wise, being aware of their indices/columns
- Aggregations (min, max, std, ...)







- Unary operations on Series and DataFrames
 - They work with any Numpy ufunc
 - The operation is applied to each element of the Series/DataFrame

Examples:

```
res = my series/4 + 1
```

- res = np.abs(my_series)
- res = np.exp(my dataframe)
- res = np.sin(my_series/4)
- • •







- Operations between Series (+,-,*,/)
 - Applied element-wise after aligning indices
 - Index elements which do not match are set to NaN (Not a Number)
 After index alignment
 - Example:

res = my series1 + my series2

Index	
b	3
а	1
С	10

mv	seri	es1
'''y_		c_{3}

Index	
а	1
b	3
d	30

my series2

Arter muex alignment
index in the result is sorted

Index	
а	2
b	6
С	NaN
d	NaN







- Operations between DataFrames
 - Applied element-wise after aligning indices and columns
 - Example (align index):
 - res = my_dataframe1 + my_dataframe2

Index in the result is **sorted**

Index	Total	Quantity
b	3	4
а	1	2
С	10	20

my_c	datafra	me1
------	---------	-----

Index	Total	Quantity
а	1	2
b	3	4
d	30	40

my_dataframe2

Index	Total	Quantity
а	2	4
b	6	8
С	NaN NaN	
d	NaN	NaN







- Operations between DataFrames
 - Example (align columns)
 - res = my_dataframe1 + my_dataframe2

Columns in the result are **sorted**

Index	Total	Quantity
а	1 2	
b	3 4	
С	5	6

Index	Total	Price
a	1 2	
b	3	4
С	5	6

Index	Price	Quantity	Total
а	NaN	NaN	2
b	NaN	NaN	6
С	NaN	NaN	10

my_dataframe1

my_dataframe2







- Operations between DataFrames and Series
 - The operation is applied between the Series and each row of the DataFrame
 - Follows broadcasting rules
 - Example:
 - res = my dataframe1 + my series1

Index	Total	Quantity
а	1	2
b	3	4
С	5	6

Index	
Total	1
Quantity	2

Index	Total	Quantity
а	2	4
b	4	6
С	6	8

my_dataframe1

my_series1







- Pandas Series and DataFrames allow performing aggregations
 - mean, std, min, max, sum
- Examples

```
In [1]: my_series.mean() # Return the mean of Series elements
```

 For DataFrames, aggregate functions are applied column-wise and return a Series

```
In [1]: my_df.mean() # Return a Series
```







Example of aggregations with DataFrames: z-score normalization

Index	Total	Quantity
а	1	2
b	3	4
С	5	6

Index	
Total	3.0
Quantity	4.0

Index	
Total	2.0
Quantity	2.0

my_dataframe1

mean_series

std_series







- Represented with sentinel value
 - None: Python null value
 - np.nan: Numpy Not A Number
- None is a Python object:
 - np.array([4, None, 5]) has dtype=Object
- np.NaN is a Floating point number
 - np.array([4, np.nan, 5]) has dtype=Float
- Using nan achieves better performances when performing numerical computations







- Pandas supports both None and NaN, and automatically converts between them when appropriate
- Example:







- Operating on missing values (for Series and DataFrames)
 - isnull()
 - Return a boolean mask indicating null values
 - notnull()
 - Return a boolean mask indicating not null values
 - dropna()
 - Return filtered data containing null values
 - fillna()
 - Return new data with filled or input missing values





- Operating on missing values: isnull, notnull
 - Return a new Series/DataFrame with the same shape as the input



dtype=float64





- Operating on missing values: dropna
 - For Series it removes null elements

```
In [1]: s1 = pd.Series([4, None, 5, np.nan])
    s1.dropna()

Out[1]: 0     4.0
     2     5.0
```







- Operating on missing values: dropna
 - For DataFrames it removes **rows** that contain at least a missing value (default behaviour)
 - Passing how=all removes rows if they contain all NaN's

Index	Total	Quantity
а	1	2
b	3	NaN
С	5	6

Index	Total	Quantity
а	1	2
С	5	6

Alternatively, it is possible to remove columns







- Operating on missing values: fillna
 - Fill null fields with a specified value (for both Series and DataFrames)

```
In [1]: s1 = pd.Series([4, None, 5, np.nan])
s1.fillna(0)

Out[1]: 0    4.0
    1    0.0
    2    5.0
    3    0.0
    dtype=float64
```







- Operating on missing values: fillna
 - The parameter **method** allows specifying different filling techniques
 - ffill: propagate last valid observation forward
 - bfill: use next valid observation to fill gap

```
In [1]: s1 = pd.Series([4, None, 5, np.nan])
s1.fillna(method='ffill')

Out[1]: 0    4.0
    1    4.0
    2    5.0
    3    5.0
```



Notebook Examples

- 3-PandasExamples.ipynb
 - 1. AccessingDataFrames and Series









- Pandas provides 2 methods for combining Series and DataFrames
 - concat()
 - Concatenate a sequence of Series/DataFrames
 - append()
 - Append a Series/DataFrame to the specified object







- Concatenating 2 Series
 - Index is preserved, even if duplicated
 - There is nothing that prevents duplicate indices in pandas!

```
In [1]:
    s1 = pd.Series(['a', 'b'], index=[1,2])
    s2 = pd.Series(['c', 'd'], index=[1,2])
    pd.concat((s1, s2))
```

```
Out[1]: 1 a 2 b 1 c 2 d d dtype=object
```







- Concatenating 2 Series
 - To avoid duplicates use ignore_index

```
In [1]:
    s1 = pd.Series(['a', 'b'], index=[1,2])
    s2 = pd.Series(['c', 'd'], index=[1,2])
    pd.concat((s1, s2), ignore_index=True)
```

```
Out[1]: 0 a

1 b

2 c

3 d

dtype=object
```







- Concatenating 2 DataFrames
 - Concatenate vertically by default

Index	Total	Quantity
а	1	2
b	3	4

Index	Total	Quantity
С	5	6
d	7	8

Index	Total	Quantity
а	1	2
b	3	4
С	5	6
d	7	8







- Concatenating 2 DataFrames
 - Missing columns are filled with NaN

Index	Total	Quantity
а	1	2
b	3	4

Index	Total	Quantity	Liters
С	5	6	1
d	7	8	2

Index	Total	Quantity	Liters
а	1	2	NaN
b	3	4	NaN
С	5	6	1.0
d	7	8	2.0







- The append() method is a shortcut for concatenating DataFrames
 - Returns the result of the concatenation

```
In [1]: df_concat = df1.append(df2)
```

is equivalent to:

```
In [1]: df_concat = pd.concat((df1, df2))
```







- Joining DataFrames with relational algebra: merge()
 - Merge on:
 - The column(s) with same name in the two DFs, by default
 - Specific columns, by specifying on=columns
 - left_on and right_on may also be used
 - The indices, if left_index/right_index are True
 - This preserves the indices (discarded otherwise)
 - Depending on the DataFrames, a one-to-one, many-to-one or many-to-many join can be performed
 - validate='1:1'|'1:m'|'m:1'|'m:m' to enforce the specific merge

```
In [1]: joined_df = pd.merge(df1, df2)
```





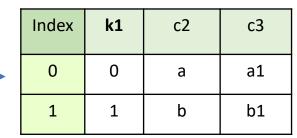


Examples (1)

pd.merge(df1, df2) → merge on columns in common, ["k1"]

Index	k1	c2
i1	0	a
i2	1	b

Index	k1	c3
i1	1	b1
i2	0	a1



pd.merge(df1, df2, right_index=True, left_index=True) → merge on index

Index	k1	c2
i1	0	а
i2	1	b
i3	0	С
i4	1	d

Index	k1	с3
i1	1	b1
i2	0	a1

Index	k1_x	c2	k1_y	c3
i1	0	а	1	b1
i2	1	b	0	a1







Examples (2)

pd.merge(df1, df2) → performs a one-to-one merge

Index	k1	c2
i1	0	a
i2	1	b

Index	k1	с3
i1	1	b1
i2	0	a1

Index	k1	c2	сЗ
0	0	а	a1
1	1	b	b1

pd.merge(df1, df2) → performs a many-to-one merge

Index	k1	c2
i1	0	а
i2	1	b
i3	0	С
i4	1	d

Index	k1	c3
i1	1	b1
i2	0	a1

Index	k1	c2	c3
0	0	а	a1
1	0	С	a1
2	1	b	b1
3	1	d	b1







- Pandas provides the equivalent of the SQL group by statement
- It allows the following operations:
 - Iterating on groups
 - Aggregating the values of each group (mean, min, max, ...)
 - Filtering groups according to a condition







Applying group by

- Specify the column(s) where you want to group (key)
- Obtain a DataFrameGroupBy object

Index	k	c1	c2
0	а	2	4
1	b	10	20
2	а	3	5
3	b	15	30

Index	k	c1	c2
0	а	2	4
2	а	3	5
1	b	10	20
3	b	15	30







Iterating on groups

Each group is a subset of the original DataFrame

Out[1]:

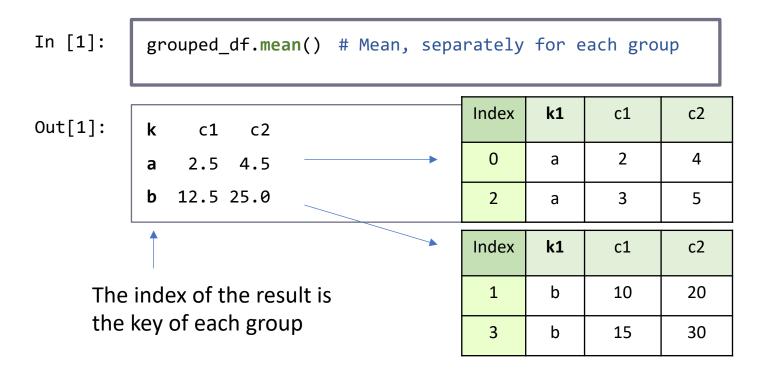
а					Index	k1	c1	c2
	k1	c1	c2	-	0	a	2	4
0	а	2	4		2	а	3	5
2	а	3	5					
b					Index	k1	c1	c2
	k1	c1	c2	-	1	b	10	20
1	b	10	20		3	b	15	30
3	b	15	30		3	D D	13	30







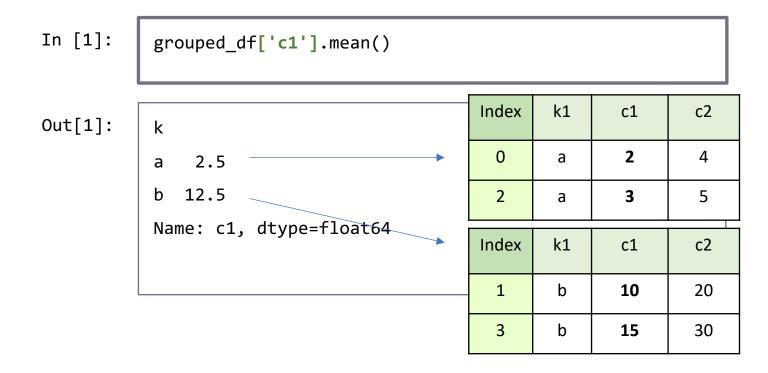
- Aggregating by group (min, max, mean, std)
 - The output is a DataFrame with the result of the aggregation for each group







- Aggregating a single column by group
 - The output is a Series with the result of the aggregation for each group







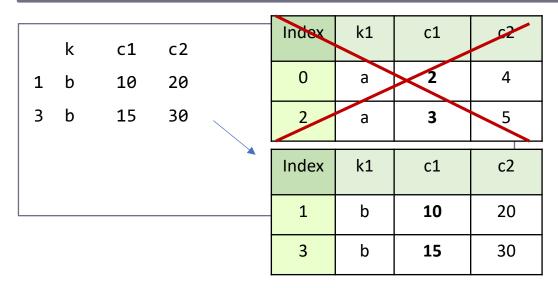


Filtering data by group

 The filter is expressed with a lambda function working with each group DataFrame (x)

```
In [1]: # Keep groups for which column c1 has a mean > 5
grouped_df.filter(lambda x: x['c1'].mean()>5)
```

Out[1]:



mean = 2.5
x: filtered
out

mean = 12.5
x: kept in
the result



Pivoting





- Pivoting allows inspecting relationships within a dataset
- Suppose to have the following dataset:

that shows failures for sensors of a given type and class during some test

Index	type	class	fail
0	а	3	1
1	b	2	1
2	b	3	1
3	а	3	0
4	b	2	1
5	а	1	0
6	b	1	0
7	а	2	0

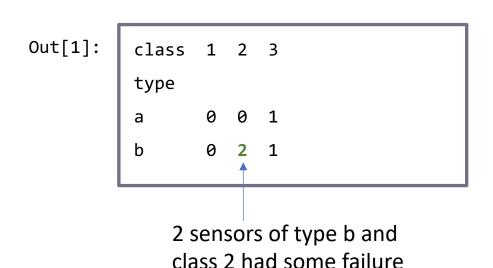


Pivoting





Shows the number of failures for all the combinations of type and class



Index	type	class	fail
0	а	3	1
1	b	2	1
2	b	3	1
3	а	3	0
4	b	2	1
5	а	1	0
6	b	1	0
7	а	2	0



Pivoting

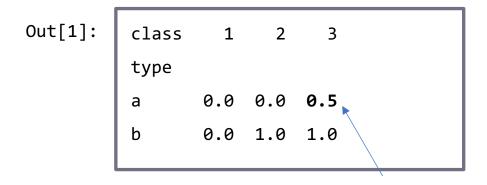




Shows the percentage of failures for all the combinations of type and class

50% of sensors of type a

and class 3 had some



failure

Index	type	class	fail
0	а	3	1
1	b	2	1
2	b	3	1
3	а	3	0
4	b	2	1
5	а	1	0
6	b	1	0
7	а	2	0

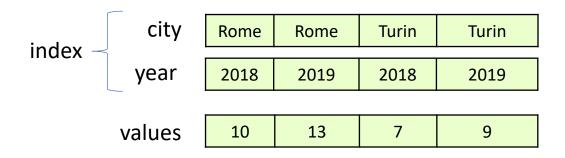


Multi-Index





- Multi-Index allows specifying an index hierarchy for
 - Series
 - DataFrames
- Example: index a Series by city and year









Building a multi-indexed Series

```
Rome 2018 10
2019 13
Turin 2018 7
2019 9
```







Naming index levels



```
In [1]: s1.index.names=['city', 'year']
    print(s1)
```

```
      city
      year

      Rome
      2018
      10

      2019
      13

      Turin
      2018
      7

      2019
      9
```







Accessing index levels



- Slicing and simple indexing are allowed
- Slicing on index levels follows Numpy rules

```
In [1]:
           print(s1.loc['Rome']) # Outer index level
           print(s1.loc[:,'2018']) # All cities, only 2018
Out[1]:
           year
                                                          Turin
                                                                   Turin
                                          Rome
                                                 Rome
           2018
                   10
                                          2018
                                                  2019
                                                          2018
                                                                   2019
           2019
                   13
                                           10
                                                   13
                                                           7
                                                                     9
           city
           Rome
                     10
           Turin
```







Accessing index levels (Examples)

```
In [1]: print(s1.loc['Turin', '2018':'2019'])
    print(s1[s1>10]) # Masking
```

						_
city	year		Rome	Rome	Turin	Turin
Turin	2018	7	2018	2019	2018	2019
	2019	9	2010	2013	2010	2013
			10	13	7	9
city	year					
Rome	2019	13				







Multi-indexed DataFrame

- Specify a multi-index for rows
- Columns can be multi-indexed as well

		Humidity		Temperat	ure
		max	min	max	min
Turkin	2018	33	48	6	33
Turin	2019	35	45	5	35
Domo	2018	40	59	2	33
Rome	2019	41	57	3	34







Multi-indexed DataFrame: creation

```
Out[1]:
```

```
c1 c2

a b a b

Rome 2018 0 1 2 3

2019 4 5 6 7

Turin 2018 8 9 10 11

2019 12 13 14 15
```







Multi-indexed DataFrame: access with outer index level

```
In [1]: print(df.loc[:, 'c1'])  # Access by column (all rows)
    print(df.loc['Rome', 'c1']) # Access rows and cols
```

```
Out[1]:
```

```
a b

Rome 2018 0 1
2019 4 5

Turin 2018 8 9
2019 12 13
```

0 1

2018

2019 4 5

		С	1	C	:2
		а	b	а	b
Rome	2018	0	1	2	3
	2019	4	5	6	7
Turin	2018	8	9	10	11
	2019	12	13	14	15







Multi-indexed DataFrame: access with outer and inner column levels using tuples

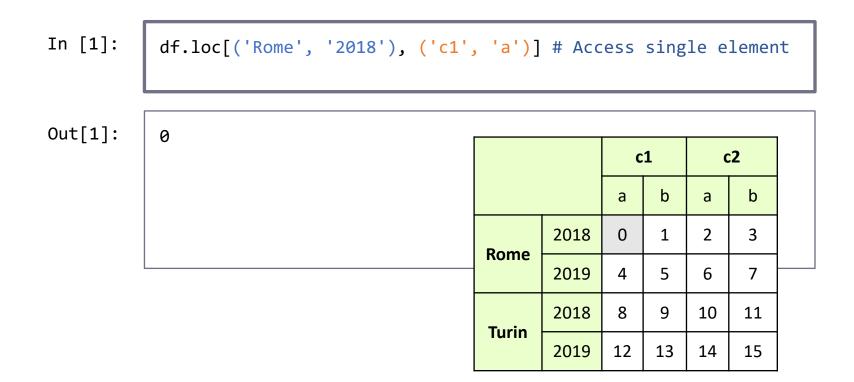
```
In [1]:
           df.loc[:, ('c1', 'a')] # Access by column
Out[1]:
           Rome
                  2018
                                                                       c2
                                                              c1
                  2019
                                                                b
                                                                          b
                                                                     a
                                                            a
           Turin 2018
                                                     2018
                                                            0
                                                                     2
                                                                          3
                  2019
                        12
                                                                1
                                              Rome
                                                     2019
                                                            4
                                                                5
                                                                          7
                                                     2018
                                                            8
                                                                9
                                                                    10
                                                                         11
                                              Turin
                                                     2019
                                                           12
                                                                13
                                                                         15
                                                                    14
```







• Multi-indexed DataFrame: access with outer and inner column and index levels using tuples









Multi-indexed DataFrame: slicing

pd.IndexSlice: Pandas object to make indexing easier

		c1	c2
		a	a
Rome	2018	0	2
Turin	2018	8	10

		С	1	C	:2
		а	b	а	b
Rome	2018	0	1	2	3
	2019	4	5	6	7
Turin	2018	8	9	10	11
	2019	12	13	14	15







Reset Index: transform index to DataFrame columns and create new (single level) index

```
In [1]:
               df.index.names = ['city', 'year']
               df_reset = df.reset_index()
               print(df reset)
     Out[1]:
                    city
                           year
                                 c1
                                         c2
                                         a b
                                     b
                                 0 1 2 3
                    Rome
                           2018
                           2019 4 5
                                         6 7
                    Rome
                   Turin
                           2018
                                 8
                                     9
                                        10
                                            11
                   Turin
                           2019
                                 12
                                    13
                                        14
                                            15
New index
```







- Set Index: transform columns to Multi-Index
 - Inverse function of reset_index()

	city	voar	c1	L	C	2
	City	year	а	Ь	а	b
0	Rome	2018	0	1	2	3
1	Rome	2019	4	5	6	7
2	Turin	2018	8	9	10	11
3	Turin	2019	12	13	14	15

oitu	VOOR .	c1		c2	
city	year -	а	b	а	b
Rome	2018	0	1	2	3
	2019	4	5	6	7
Turin	2018	8	9	10	11
	2019	12	13	14	15

New index







Unstack: transform multi-indexed Series to a Dataframe

myseries.unstack()

city	year	
Rome	2018	0
	2019	4
Turin	2018	8
	2019	12

	2018	2019
Rome	0	4
Turin	8	12







- Stack: inverse function of unstack()
 - From DataFrame to multi-indexed Series

mydataframe.stack()

	2018	2019
Rome	0	4
Turin	8	12

Rome	2018	0
Nome	2019	4
Turin	2018	8
Turin	2019	12







Aggregates on multi-indices

- Allowed by passing the level parameter
- Level specifies the row granularity at which the result is computed

my_dataframe.max(level='city')

city	voor	c1		c2	
city	year	а	b	а	b
Rome	2018	0	1	2	3
	2019	4	5	6	7
Turin	2018	8	9	10	11
	2019	12	13	14	15

city	c1	L	c2		
city	а	b	а	b	
Rome	4	5	6	7	
Turin	12	13	14	15	







Aggregates on multi-indices

my_dataframe.max(level='year')

city	year	c1		c2	
		а	b	а	b
Rome	2018	0	1	2	3
	2019	4	5	6	7
Turin	2018	8	9	10	11
	2019	12	13	14	15

voar	c1	<u>L</u>	c2		
year	а	b	а	b	
2018	8	9	10	11	
2019	12	13	14	15	







Aggregates on multi-indices

- Can also aggregate columns
 - Specify axis=1

my_dataframe.max(axis=1, level=0)

city	year	c1		c2	
		а	b	а	b
Rome	2018	0	1	2	3
	2019	4	5	6	7
Turin	2018	8	9	10	11
	2019	12	13	14	15

city	year	c1	c2
Rome	2018	1	3
Rome	2019	5	7
Turin	2018	9	11
Turin	2019	13	15