



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:

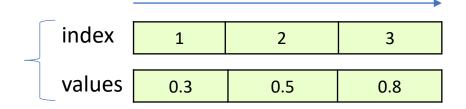
- Managing tables and series
 - data selection
 - grouping, pivoting
- Managing missing data
- Statistics on data







- 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'])



In [1]:

Pandas Series

c 4.5





Accessing Series elements: slicing



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







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>
```







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]:
                      Quantity
             Price
                                  Liters
               1.0
                              5
                                      1.5
          a
          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)
```







Creation from list of dictionaries



- 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)
```







Creation from 2D Numpy array



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







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



Returns a Series with column data

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







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)
                                     # Get the first row
          print(df.iloc[0])
Out[1]:
                   1.0
          Price
          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']]
```

Out[1]: Price Liters

a 1.0 1.5

c 5.0 1.0







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	6 .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



Notebook Examples

3.1 Pandas Series and DataFrame.ipynb

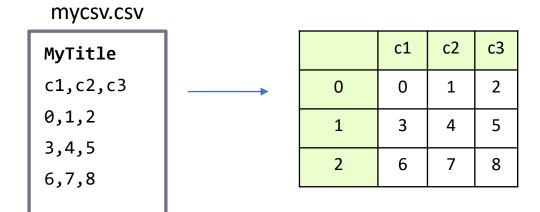








- Load DataFrame from csv file
 - Allows specifying the column delimiter (sep)
 - Automatically read header from first line of the file (after skipping the specified number of rows)
 - Column data types are inferred









- 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	
CSV,	0,no info,	
mycsv.csv	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 c2 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



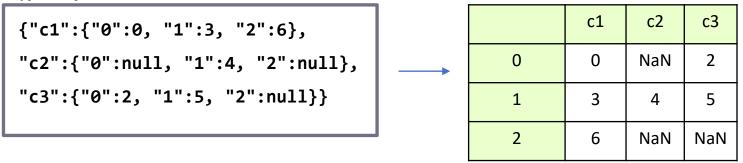




Load DataFrame from json file

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

myjson.json



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



DataFrames and I/O





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







- 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)
 - Example:
 - res = my series1 + my series2

Index	
b	3
а	1
С	10

mv	series	- 1
IIIY_	_361163	ΣТ

Index	
а	1
b	3
d	30

my series2

After index alignment index in the result is **sorted**

Index	
a	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

Ind	lex	in	the	resu	t
is s	ort	ed	l		

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

my_dataframe1	-
---------------	---

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

my dataframe2

Index	Total Quantit	
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	
а	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
a	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



Notebook Examples

3.2 PandasOperations.ipynb









- 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

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

Index	Index Total Qua	
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

Index	k1	c2	c3
0	0	а	a1
1	1	b	b1

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





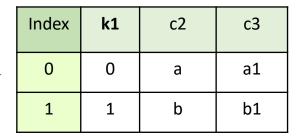


Examples (2)

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



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	с3
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]:

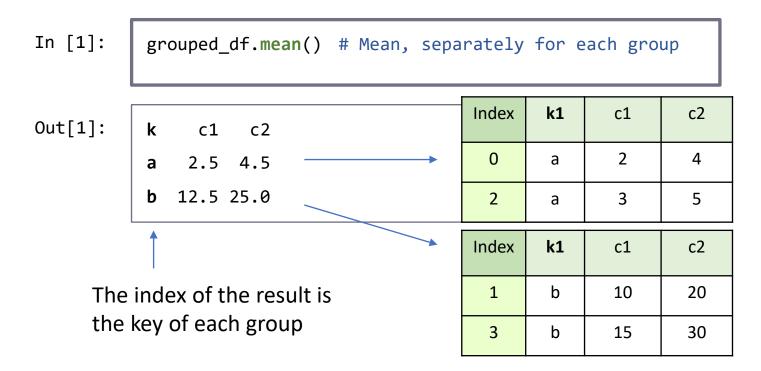
ā	a					Index	k1	c1	c2
		k1	c1	c2	-	0	а	2	4
(9	a	2	4		2	а	3	5
2	2	a	3	5					
ł)					Index	k1	c1	c2
		k1	c1	c2	-	1	b	10	20
1	L	b	10	20		3	b	15	30
3	3	b	15	30		3	D D	13	30







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









- Many Aggregations by group with .agg([..., ...])
 - The output is a DataFrame with the results of the aggregations for each retained column

Out[1]: c1 c2

k max min max min

a 3 2 5 4

b 15 10 30 20

The results now have 4 columns: 'max' and 'min' for each of the previous columns

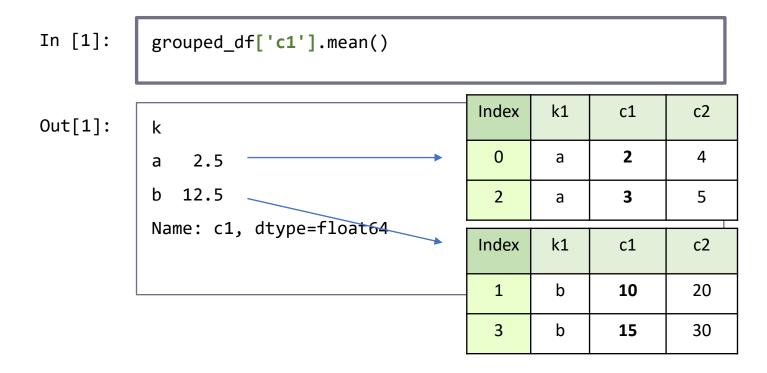
Index	k1	c1	c2
0	а	2	4
2	а	3	5
Index	k1	c 1	c2

Index	k1	c1	c2
1	b	10	20
3	b	15	30





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







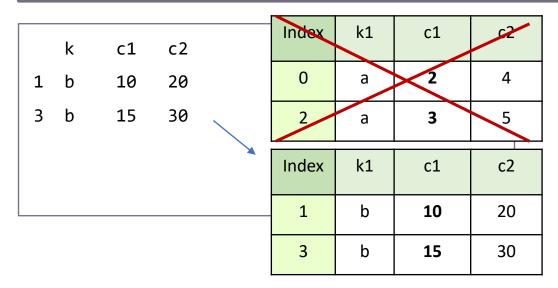


Filtering data by group

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



Notebook Examples

3.3 PandasGrouping.ipynb

