



# Machine Learning for Networking

Pandas - Extra slides

Andrea Pasini Flavio Giobergia Elena Baralis **Gabriele Ciravegna** 

DataBase and Data Mining Group







- 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







- 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
dtype=float64
```







- 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

dropped\_df = df.dropna(axis='columns')







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







- 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
```



# **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

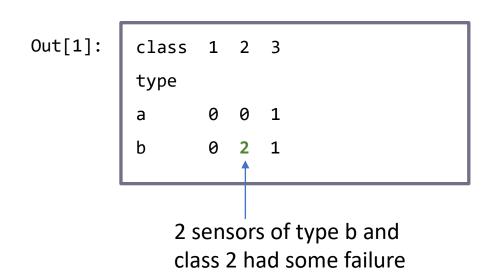


# **Pivoting**





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

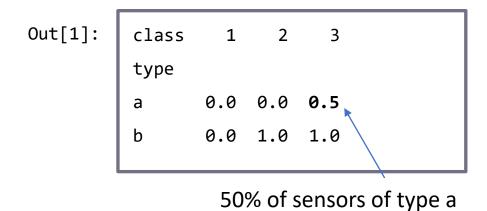


# **Pivoting**





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



failure

and class 3 had some

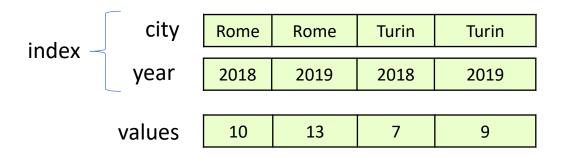
Index	type	class	fail
0	а	3	1
1	b	2	1
2	b	3	1
3	а	3	0
4	b	2	1
5	а	1	0
6	b	1	0
7	а	2	0







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









# Building a multi-indexed Series

Out[1]:

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







# Naming index levels



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

Out[1]:

```
      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
                     7
```







# Accessing index levels (Examples)

Turin 2018 7 2019 9

Rome	ome Rome Turin		Turin
2018	2019	2018	2019
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
<b>Turin</b> 2019	2019	35	45	5	35
Dome	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

Out[1]:

		a	b
Rome	2018	0	1
	2019	4	5
Turin	2018	8	9
	2019	12	13
	a h		

2018 **0** 1 2019 4 5

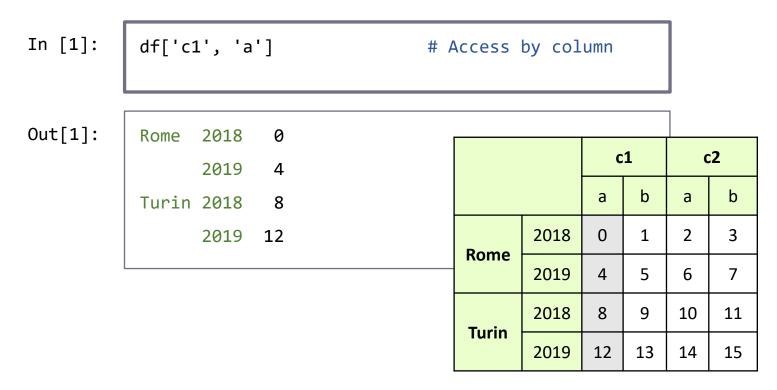
		C	1	C	:2
		а	b	а	b
Domo	2018	0	1	2	3
Rome	2019	4	5	6	7
Turin	2018	8	9	10	11
	2019	12	13	14	15







# Multi-indexed DataFrame: access with outer and inner index levels









# Multi-indexed DataFrame: access with outer and inner index levels

Out[1]:

c1	a	0
c2	а	2

		<b>c1</b>		<b>c2</b>	
		а	b	а	b
Domo	2018	0	1	2	3
Rome	2019	4	5	6	7
Turka	2018	8	9	10	11
Turin	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
                                        c2
                           year
                                 c1
                                     b
                                        a b
                                 0 1 2 3
                   Rome
                           2018
                           2019 4 5 6 7
                   Rome
                   Turin
                           2018 8
               2
                                       10 11
                   Turin
                           2019
                                12
                                    13
                                        14
                                           15
New index
```







- Set Index: transform columns to Multi-Index
  - Inverse function of reset\_index()

	city	oity year		c1		c2	
		year	а	b	а	р	
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	

city	year -	С	c1		c2	
		a	b	а	b	
Pomo	2018	0	1	2	3	
Rome	2019	4	5	6	7	
Turin	2018	8	9	10	11	
Turin	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

Domo	2018	0
Rome	2019	4
Turka	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	<b>11001</b>	c1		c2	
city	year	а	b	а	b
Domo	2018	0	1	2	3
Rome	2019	4	5	6	7
Turring	2018	8	9	10	11
Turin	2019	12	13	14	15

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

voor	c1	L	c2		
year	a 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	V00#	c1		c2	
city	year	а	b	а	b
Domo	2018	0	1	2	3
Rome	2019	4	5	6	7
Turin	2018	8	9	10	11
Turin	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