

Curs 4: Pandas - elemente avansate

Lucrul cu valori lipsa in Pandas

Reprezentarea valorilor lipsa in Pandas

Pandas foloseste doua variante pentru reprezentarea de valori lipsa: None si NaN. NaN este utilizat pentru tipuri numerice in virgula mobila. None este convertit la NaN daca seria este numerica; daca seria este ne-numerica, se considera de tip `object` :

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

```
In [2]: print(f'pandas version: {pd.__version__}')
print(f'numpy version: {np.__version__}')
```

```
# pandas version: 0.24.1
# numpy version: 1.16.2
```

```
pandas version: 1.0.1
numpy version: 1.18.1
```

NaN si None sunt echivalente in context numeric, in Pandas:

```
In [3]: pd.Series([1, np.nan, 2, None])
```

```
Out[3]: 0    1.0
1    NaN
2    2.0
3    NaN
dtype: float64
```

```
In [4]: pd.Series(['John', 'Danny', None])
```

```
Out[4]: 0    John
1    Danny
2    None
dtype: object
```

Intrucat doar tipurile numerice floating point suporta valoare de NaN, conform standardului IEEE 754, se va face transformarea unei serii de tip intreg intr-una de tip floating point daca se insereaza sau adauga un NaN:

```
In [5]: # creare de serie cu valori intregi
x = pd.Series([10, 20], dtype=int)
x
```

```
Out[5]: 0    10
        1    20
        dtype: int32
```

```
In [6]: x[1] = np.nan
x
```

```
Out[6]: 0    10.0
        1     NaN
        dtype: float64
```

```
In [7]: # adaugare cu append
x = pd.Series([10, 20], dtype=int)
print(f'Serie de intregi:\n{x}')
x = x.append(pd.Series([100, np.nan]))
print(f'Dupa adaugare:\n{x}')
```

```
Serie de intregi:
0    10
1    20
dtype: int32
Dupa adaugare:
0    10.0
1    20.0
0    100.0
1     NaN
dtype: float64
```

Operatii cu valori lipsa in Pandas

Metodele ce se pot folosi pentru operarea cu valori lipsa sunt:

- `isnull()` - genereaza o matrice de valori logice, ce specifica daca pe pozitiile corespunzatoare sunt valori lipsa
- `notnull()` - complementara lui `isnull()`
- `dropna()` - returneaza o versiune filtrata a datelor, doar acele linii care nu au null
- `fillna()` - returneaza o copie a obiectului initial, in care valorile lipsa sunt umplute cu ceva specificat

`isnull()` si `notnull()`

```
In [8]: data = pd.Series([1, np.nan, 'hello', None])
data
```

```
Out[8]: 0      1
1      NaN
2     hello
3      None
dtype: object
```

```
In [9]: data.isnull()
```

```
Out[9]: 0     False
1      True
2     False
3      True
dtype: bool
```

Selectarea doar acelor valori din obiectul Series care sunt ne-nule se face cu:

```
In [10]: # filtrare
data[data.notnull()]
```

```
Out[10]: 0      1
2     hello
dtype: object
```

Funcțiile `isnull()` și `notnull()` funcționează la fel și pentru obiecte `DataFrame`:

```
In [11]: df = pd.DataFrame({'Name': ['Will', 'Mary', 'Joan'], 'Age': [20, 25, 30]})
df
```

```
Out[11]:
```

	Name	Age
0	Will	20
1	Mary	25
2	Joan	30

```
In [12]: df.loc[2, 'Age'] = np.NaN
df
```

```
Out[12]:
```

	Name	Age
0	Will	20.0
1	Mary	25.0
2	Joan	NaN

```
In [13]: df.isnull()
```

```
Out[13]:
```

	Name	Age
0	False	False
1	False	False
2	False	True

```
In [14]: df.notnull()
```

```
Out[14]:
```

	Name	Age
0	True	True
1	True	True
2	True	False

In cazul obiectelor DataFrame, aplicarea lui `notnull()` nu lasa afara elemente din dataframe:

```
In [15]: df[df.notnull()]
```

```
Out[15]:
```

	Name	Age
0	Will	20.0
1	Mary	25.0
2	Joan	NaN

Stergerea de elemente cu `dropna()`

Pentru un obiect Series, metoda `dropna()` produce un alt obiect in care liniile cu valori de null sunt sterse:

```
In [16]: data
```

```
Out[16]:
```

```
0      1
1     NaN
2   hello
3     None
dtype: object
```

```
In [17]: data2 = data.dropna()
data2
```

```
Out[17]:
```

```
0      1
2   hello
dtype: object
```

Pentru un obiect DataFrame se pot sterge doar linii sau coloane intregi - obiectul care ramane trebuie sa fie tot un DataFrame:

```
In [18]: df = pd.DataFrame([[1, np.nan, 2],[2, 3, 5],[np.nan, 4, 6]])  
df
```

Out[18]:

	0	1	2
0	1.0	NaN	2
1	2.0	3.0	5
2	NaN	4.0	6

```
In [19]: # Implicit: eliminare de linii care contin null  
df2 = df.dropna()  
df2
```

Out[19]:

	0	1	2
1	2.0	3.0	5

Mai sus s-a ales implicit stergerea de linii, datorita faptului ca parametrul `axis` are implicit valoarea 0:

In [20]: `help(df.dropna)`

Help on method dropna in module pandas.core.frame:

dropna(axis=0, how='any', thresh=None, subset=None, inplace=False) method of pandas.core.frame.DataFrame instance

Remove missing values.

See the :ref:`User Guide <missing_data>` for more on which values are considered missing, and how to work with missing data.

Parameters

axis : {0 or 'index', 1 or 'columns'}, default 0

Determine if rows or columns which contain missing values are removed.

* 0, or 'index' : Drop rows which contain missing values.

* 1, or 'columns' : Drop columns which contain missing value.

.. versionchanged:: 1.0.0

Pass tuple or list to drop on multiple axes.
Only a single axis is allowed.

how : {'any', 'all'}, default 'any'

Determine if row or column is removed from DataFrame, when we have at least one NA or all NA.

* 'any' : If any NA values are present, drop that row or column.

* 'all' : If all values are NA, drop that row or column.

thresh : int, optional

Require that many non-NA values.

subset : array-like, optional

Labels along other axis to consider, e.g. if you are dropping rows these would be a list of columns to include.

inplace : bool, default False

If True, do operation inplace and return None.

Returns

DataFrame

DataFrame with NA entries dropped from it.

See Also

DataFrame.isna: Indicate missing values.

DataFrame.notna : Indicate existing (non-missing) values.

DataFrame.fillna : Replace missing values.

Series.dropna : Drop missing values.

Index.dropna : Drop missing indices.

Examples

```
>>> df = pd.DataFrame({"name": ['Alfred', 'Batman', 'Catwoman'],
...                     "toy": [np.nan, 'Batmobile', 'Bullwhip'],
...                     "born": [pd.NaT, pd.Timestamp("1940-04-25"),
...                               pd.NaT]})
```

```
>>> df
      name      toy      born
0  Alfred      NaN      NaT
1  Batman  Batmobile 1940-04-25
2  Catwoman  Bullwhip      NaT
```

Drop the rows where at least one element is missing.

```
>>> df.dropna()
      name      toy      born
1  Batman  Batmobile 1940-04-25
```

Drop the columns where at least one element is missing.

```
>>> df.dropna(axis='columns')
      name
0  Alfred
1  Batman
2  Catwoman
```

Drop the rows where all elements are missing.

```
>>> df.dropna(how='all')
      name      toy      born
0  Alfred      NaN      NaT
1  Batman  Batmobile 1940-04-25
2  Catwoman  Bullwhip      NaT
```

Keep only the rows with at least 2 non-NA values.

```
>>> df.dropna(thresh=2)
      name      toy      born
1  Batman  Batmobile 1940-04-25
2  Catwoman  Bullwhip      NaT
```

Define in which columns to look for missing values.

```
>>> df.dropna(subset=['name', 'born'])
      name      toy      born
1  Batman  Batmobile 1940-04-25
```

Keep the DataFrame with valid entries in the same variable.

```
>>> df.dropna(inplace=True)
>>> df
      name      toy      born
1  Batman  Batmobile 1940-04-25
```

Se poate opta pentru stergerea de coloane care contin null:

In [21]: df

Out[21]:

	0	1	2
0	1.0	NaN	2
1	2.0	3.0	5
2	NaN	4.0	6

```
In [22]: # stergere de coloane cu null
# df3 = df.dropna(axis=1) # functioneaza
df3 = df.dropna(axis='columns')
df3
```

Out[22]:

	2
0	2
1	5
2	6

Operatiile de mai sus sterg o linie sau o coloana daca ea contine cel putin o valoare de null. Se poate cere stergerea doar in cazul in care intreaga linie sau coloana e plina cu null, folosind parametrul `how` :

In [23]: df

Out[23]:

	0	1	2
0	1.0	NaN	2
1	2.0	3.0	5
2	NaN	4.0	6

```
In [24]: df2 = df.dropna(how='all')
df2
```

Out[24]:

	0	1	2
0	1.0	NaN	2
1	2.0	3.0	5
2	NaN	4.0	6

De remarcat ca `dropna()` nu modifica obiectul original, decat daca se specifica parametrul `inplace=True` .

Umplerea de valori nule cu `fillna()`

```
In [25]: data = pd.Series([1, np.nan, 2, None, 3], index=list('abcde'))
```

```
In [26]: # umplere cu valoare constanta
data2 = data.fillna(0)
data2
```

```
Out[26]: a    1.0
         b    0.0
         c    2.0
         d    0.0
         e    3.0
         dtype: float64
```

```
In [27]: # Umplere cu copierea ultimei valori cunoscute:
data2 = data.fillna(method='ffill')
data2
```

```
Out[27]: a    1.0
         b    1.0
         c    2.0
         d    2.0
         e    3.0
         dtype: float64
```

```
In [28]: # Umplere 'inapoi':
data2 = data.fillna(method='bfill')
data2
```

```
Out[28]: a    1.0
         b    2.0
         c    2.0
         d    3.0
         e    3.0
         dtype: float64
```

```
In [29]: # umplerea cu valoare calculata:
print(f'Media valorilor non-nan este: {data.mean()}')
data2 = data.fillna(data.mean())
data2
```

Media valorilor non-nan este: 2.0

```
Out[29]: a    1.0
         b    2.0
         c    2.0
         d    2.0
         e    3.0
         dtype: float64
```

Agregare si grupare

Agregari simple

```
In [30]: np.random.seed(100)
ser = pd.Series(np.random.rand(10))
ser
```

```
Out[30]: 0    0.543405
1    0.278369
2    0.424518
3    0.844776
4    0.004719
5    0.121569
6    0.670749
7    0.825853
8    0.136707
9    0.575093
dtype: float64
```

```
In [31]: ser.sum(), ser.max(), ser.min()
```

```
Out[31]: (4.425757785871915, 0.8447761323199037, 0.004718856190972565)
```

Pentru obiecte DataFrame, operatiile de agregare opereaza pe coloane:

```
In [32]: df = pd.DataFrame({'A': np.random.rand(10), 'B': -np.random.rand(10) }, index=
['line ' + str(i) for i in range(1, 11)])
df
```

```
Out[32]:
```

	A	B
line 1	0.891322	-0.431704
line 2	0.209202	-0.940030
line 3	0.185328	-0.817649
line 4	0.108377	-0.336112
line 5	0.219697	-0.175410
line 6	0.978624	-0.372832
line 7	0.811683	-0.005689
line 8	0.171941	-0.252426
line 9	0.816225	-0.795663
line 10	0.274074	-0.015255

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

```
Out[33]: A    0.466647
B   -0.414277
dtype: float64
```

.. si daca se doreste calculul pe linii, se poate indica via parametrul `axis` :

```
In [34]: # df.mean(axis=1)
df.mean(axis='columns')
```

```
Out[34]: line 1      0.229809
line 2     -0.365414
line 3     -0.316161
line 4     -0.113868
line 5      0.022144
line 6      0.302896
line 7      0.402997
line 8     -0.040243
line 9      0.010281
line 10     0.129409
dtype: float64
```

Exista o metoda utila, care pentru un obiect DataFrame calculeaza statisticile:

```
In [35]: df.describe()
```

```
Out[35]:
```

	A	B
count	10.000000	10.000000
mean	0.466647	-0.414277
std	0.356280	0.333688
min	0.108377	-0.940030
25%	0.191297	-0.704673
50%	0.246886	-0.354472
75%	0.815089	-0.194664
max	0.978624	-0.005689

Operatiile nu iau in considerare valorile lipsa:

```
In [36]: df.iloc[0, 0] = df.iloc[0,1] = np.nan
df.iloc[5, 0] = df.iloc[7, 1] = df.iloc[9, 1] = np.nan
df
```

Out[36]:

	A	B
line 1	NaN	NaN
line 2	0.209202	-0.940030
line 3	0.185328	-0.817649
line 4	0.108377	-0.336112
line 5	0.219697	-0.175410
line 6	NaN	-0.372832
line 7	0.811683	-0.005689
line 8	0.171941	NaN
line 9	0.816225	-0.795663
line 10	0.274074	NaN

```
In [37]: df.count()
```

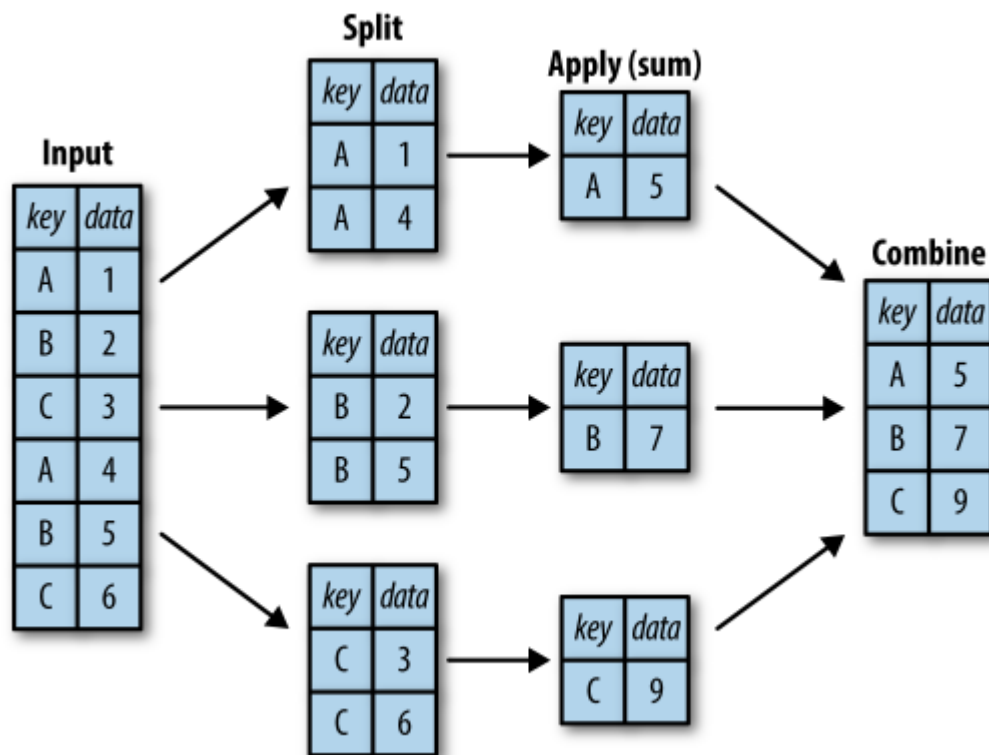
Out[37]: A 8
B 7
dtype: int64

Metoda de agregare	Descriere
count()	Numarul total de elemente
first(), last()	primul si ultimul element
mean(), median()	Media si mediana
min(), max()	Minimul si maximul
std(), var()	Deviatia standard si varianta
mad()	Deviatia absoluta medie
prod(), sum()	Produsul si suma elementelor

Gruparea datelor: split(), apply(), combine()

Pasii care se fac pentru agregarea datelor urmeaza secventa: imparte, aplica operatie, combina:

1. imparte - via metoda `split()` : separa datele initiale in grupuri, pe baza unei chei
2. aplica, via metoda `apply()` : calculeaza o functie pentru fiecare grup: agregare, transformare, filtrare
3. combina, via metoda `combine()` : concateneaza rezultatele si produ raspunsul final



```
In [38]: df = pd.DataFrame({'key': ['A', 'B', 'C', 'A', 'B', 'C'], 'data': range(6)}, columns=['key', 'data'])
df
```

Out[38]:

	key	data
0	A	0
1	B	1
2	C	2
3	A	3
4	B	4
5	C	5

```
In [39]: groups = df.groupby('key')
type(groups)
```

Out[39]: pandas.core.groupby.generic.DataFrameGroupBy

In [40]: `print(groups)`

<pandas.core.groupby.generic.DataFrameGroupBy object at 0x00000111D740BC08>

In [41]: `groups.sum()`

Out[41]:

data	
key	
A	3
B	5
C	7

Ca functie de agregare se poate folosi orice functie Pandas sau NumPy.

In [42]: `import seaborn as sns`
`planets = sns.load_dataset('planets')`

In [43]: `planets.head()`

Out[43]:

	method	number	orbital_period	mass	distance	year
0	Radial Velocity	1	269.300	7.10	77.40	2006
1	Radial Velocity	1	874.774	2.21	56.95	2008
2	Radial Velocity	1	763.000	2.60	19.84	2011
3	Radial Velocity	1	326.030	19.40	110.62	2007
4	Radial Velocity	1	516.220	10.50	119.47	2009

In [44]: `planets.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1035 entries, 0 to 1034
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   method          1035 non-null   object
1   number          1035 non-null   int64
2   orbital_period   992 non-null    float64
3   mass            513 non-null    float64
4   distance        808 non-null    float64
5   year            1035 non-null   int64
dtypes: float64(3), int64(2), object(1)
memory usage: 48.6+ KB
```

In [45]: `planets.describe()`

Out[45]:

	number	orbital_period	mass	distance	year
count	1035.000000	992.000000	513.000000	808.000000	1035.000000
mean	1.785507	2002.917596	2.638161	264.069282	2009.070531
std	1.240976	26014.728304	3.818617	733.116493	3.972567
min	1.000000	0.090706	0.003600	1.350000	1989.000000
25%	1.000000	5.442540	0.229000	32.560000	2007.000000
50%	1.000000	39.979500	1.260000	55.250000	2010.000000
75%	2.000000	526.005000	3.040000	178.500000	2012.000000
max	7.000000	730000.000000	25.000000	8500.000000	2014.000000

In [46]: `planets.method.unique()`

Out[46]: `array(['Radial Velocity', 'Imaging', 'Eclipse Timing Variations', 'Transit', 'Astrometry', 'Transit Timing Variations', 'Orbital Brightness Modulation', 'Microlensing', 'Pulsar Timing', 'Pulsation Timing Variations'], dtype=object)`

Pentru grupurile rezultate se poate alege o coloana, pentru care sa se calculeze valori agregate:

In [47]: `planets.groupby('method')['orbital_period'].median()`

Out[47]:

method	
Astrometry	631.180000
Eclipse Timing Variations	4343.500000
Imaging	27500.000000
Microlensing	3300.000000
Orbital Brightness Modulation	0.342887
Pulsar Timing	66.541900
Pulsation Timing Variations	1170.000000
Radial Velocity	360.200000
Transit	5.714932
Transit Timing Variations	57.011000

Name: orbital_period, dtype: float64

Grupurile pot fi iterate, returnand pentru fiecare grup un obiect de tip Series sau DataFrame:


```
In [48]: print(f'Number of columns: {len(planets.columns)}')

for (method, group) in planets.groupby('method'):
    print("{0:30s} shape={1}".format(method, group.shape))
```

```
Number of columns: 6
Astrometry           shape=(2, 6)
Eclipse Timing Variations shape=(9, 6)
Imaging              shape=(38, 6)
Microlensing         shape=(23, 6)
Orbital Brightness Modulation shape=(3, 6)
Pulsar Timing        shape=(5, 6)
Pulsation Timing Variations shape=(1, 6)
Radial Velocity      shape=(553, 6)
Transit              shape=(397, 6)
Transit Timing Variations shape=(4, 6)
```

Fiecare grup rezultat, fiind vazut ca un Series sau DataFrame, suporta apel de metode aferente acestor obiecte:

```
In [49]: planets.groupby('method')['year'].describe()
```

Out[49]:

	count	mean	std	min	25%	50%	75%	max
method								
Astrometry	2.0	2011.500000	2.121320	2010.0	2010.75	2011.5	2012.25	2013.0
Eclipse Timing Variations	9.0	2010.000000	1.414214	2008.0	2009.00	2010.0	2011.00	2012.0
Imaging	38.0	2009.131579	2.781901	2004.0	2008.00	2009.0	2011.00	2013.0
Microlensing	23.0	2009.782609	2.859697	2004.0	2008.00	2010.0	2012.00	2013.0
Orbital Brightness Modulation	3.0	2011.666667	1.154701	2011.0	2011.00	2011.0	2012.00	2013.0
Pulsar Timing	5.0	1998.400000	8.384510	1992.0	1992.00	1994.0	2003.00	2011.0
Pulsation Timing Variations	1.0	2007.000000	NaN	2007.0	2007.00	2007.0	2007.00	2007.0
Radial Velocity	553.0	2007.518987	4.249052	1989.0	2005.00	2009.0	2011.00	2014.0
Transit	397.0	2011.236776	2.077867	2002.0	2010.00	2012.0	2013.00	2014.0
Transit Timing Variations	4.0	2012.500000	1.290994	2011.0	2011.75	2012.5	2013.25	2014.0

Metodele aggregate(), filter(), transform(), apply()

Inainte de pasul de combinare a datelor se pot folosi metode care implementeaza operatii pe grupuri inainte de a face in final gruparea rezultatelor din grupuri.

```
In [50]: df = pd.DataFrame({'key': ['A', 'B', 'C', 'A', 'B', 'C'],
'data1': range(6),
'data2': np.random.randint(0, 10, 6)},
columns = ['key', 'data1', 'data2'])
df
```

Out[50]:

	key	data1	data2
0	A	0	5
1	B	1	8
2	C	2	1
3	A	3	0
4	B	4	7
5	C	5	6

Metoda `aggregate()` permite specificare de functii prin numele lor (string sau referinta la functie):

```
In [51]: df.groupby('key').aggregate(['min', np.median, max])
```

Out[51]:

	data1			data2		
	min	median	max	min	median	max
key						
A	0	1.5	3	0	2.5	5
B	1	2.5	4	7	7.5	8
C	2	3.5	5	1	3.5	6

Filtrarea cu `filter()` permite selectarea doar acelor grupuri care satisfac o anumita conditie:

```
In [52]: def filter_func(x): # x este o linie, corespunzand fiecarui grup
return x['data2'].std() > 4
```

```
In [53]: df.groupby('key').std()
```

Out[53]:

	data1	data2
key		
A	2.12132	3.535534
B	2.12132	0.707107
C	2.12132	3.535534

```
In [54]: df.groupby('key').filter(filter_func)
```

```
Out[54]:
```

	key	data1	data2
--	-----	-------	-------

Acelasi efect se obtine cu lambda functii:

```
In [55]: df.groupby('key').filter(lambda row: row['data2'].std() > 4)
```

```
Out[55]:
```

	key	data1	data2
--	-----	-------	-------

Transformarea cu `transform()` produce un dataframe cu acelasi numar de linii ca si cel initial, dar cu valorile calculate prin aplicarea unei operatii la nivelul fiecarui grup:

```
In [56]: df
```

```
Out[56]:
```

	key	data1	data2
0	A	0	5
1	B	1	8
2	C	2	1
3	A	3	0
4	B	4	7
5	C	5	6

Media pe fieare grup este:

```
In [57]: df.groupby('key').mean()
```

```
Out[57]:
```

	data1	data2
key		
A	1.5	2.5
B	2.5	7.5
C	3.5	3.5

Centrarea valorilor pentru fiecare grup - adica: in fiecare grup sa fie media 0 - se face cu:

```
In [58]: df.groupby('key').transform(lambda x: x - x.mean())
```

Out[58]:

	data1	data2
0	-1.5	2.5
1	-1.5	0.5
2	-1.5	-2.5
3	1.5	-2.5
4	1.5	-0.5
5	1.5	2.5

```
In [59]: df.groupby('key').transform(lambda x: x - x.mean()).mean()
```

Out[59]: data1 0.0
data2 0.0
dtype: float64

Functia `apply()` permite calculul unei functii peste fiecare grup. Exemplul de mai jos calculeaza prima coloana impartita la suma elementelor din coloana `data2`, in cadrul fiecarui grup:

```
In [60]: def norm_by_data2(x):
# x is a DataFrame of group values
x['data1'] /= x['data2'].sum()
return x

df.groupby('key').apply(norm_by_data2)
```

Out[60]:

	key	data1	data2
0	A	0.000000	5
1	B	0.066667	8
2	C	0.285714	1
3	A	0.600000	0
4	B	0.266667	7
5	C	0.714286	6

Functia `apply()` se poate folosi si in afara lui `groupby`, permitand calcul vectorizat de mare viteza:

```
In [61]: data_len = 10000

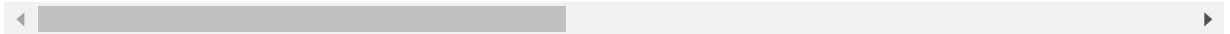
df_big = pd.DataFrame({'Noise_' + str(i) : np.random.rand(data_len) for i in range(1, 50)})

df_big.head()
```

Out[61]:

	Noise_1	Noise_2	Noise_3	Noise_4	Noise_5	Noise_6	Noise_7	Noise_8	Noise_9	No
0	0.030123	0.203968	0.706581	0.298033	0.534726	0.515900	0.258939	0.413919	0.026733	0.
1	0.776005	0.688731	0.204790	0.082986	0.053910	0.295277	0.478298	0.878959	0.426999	0.
2	0.550958	0.953967	0.185411	0.603051	0.411614	0.204954	0.782968	0.377960	0.100514	0.
3	0.381073	0.756840	0.121745	0.999780	0.766192	0.881829	0.667565	0.271940	0.286227	0.
4	0.529266	0.347373	0.184114	0.983282	0.353940	0.246467	0.866640	0.575963	0.430655	0.

5 rows × 49 columns



```
In [62]: all_noise_columns = [column for column in df_big.columns if column.startswith(
'Noise_')]

row = df_big.iloc[0]
row[all_noise_columns]
```

```
Out[62]: Noise_1      0.030123
Noise_2      0.203968
Noise_3      0.706581
Noise_4      0.298033
Noise_5      0.534726
Noise_6      0.515900
Noise_7      0.258939
Noise_8      0.413919
Noise_9      0.026733
Noise_10     0.547176
Noise_11     0.834616
Noise_12     0.631497
Noise_13     0.923611
Noise_14     0.551549
Noise_15     0.785927
Noise_16     0.280730
Noise_17     0.959686
Noise_18     0.287398
Noise_19     0.819674
Noise_20     0.756904
Noise_21     0.229681
Noise_22     0.050490
Noise_23     0.832008
Noise_24     0.982115
Noise_25     0.410147
Noise_26     0.856429
Noise_27     0.528605
Noise_28     0.577306
Noise_29     0.590815
Noise_30     0.147199
Noise_31     0.009771
Noise_32     0.625495
Noise_33     0.043671
Noise_34     0.914573
Noise_35     0.822432
Noise_36     0.405514
Noise_37     0.393812
Noise_38     0.769161
Noise_39     0.858692
Noise_40     0.461877
Noise_41     0.076768
Noise_42     0.700336
Noise_43     0.301304
Noise_44     0.381791
Noise_45     0.114720
Noise_46     0.870638
Noise_47     0.363271
Noise_48     0.828637
Noise_49     0.758510
Name: 0, dtype: float64
```

```
In [63]: # %%timeit

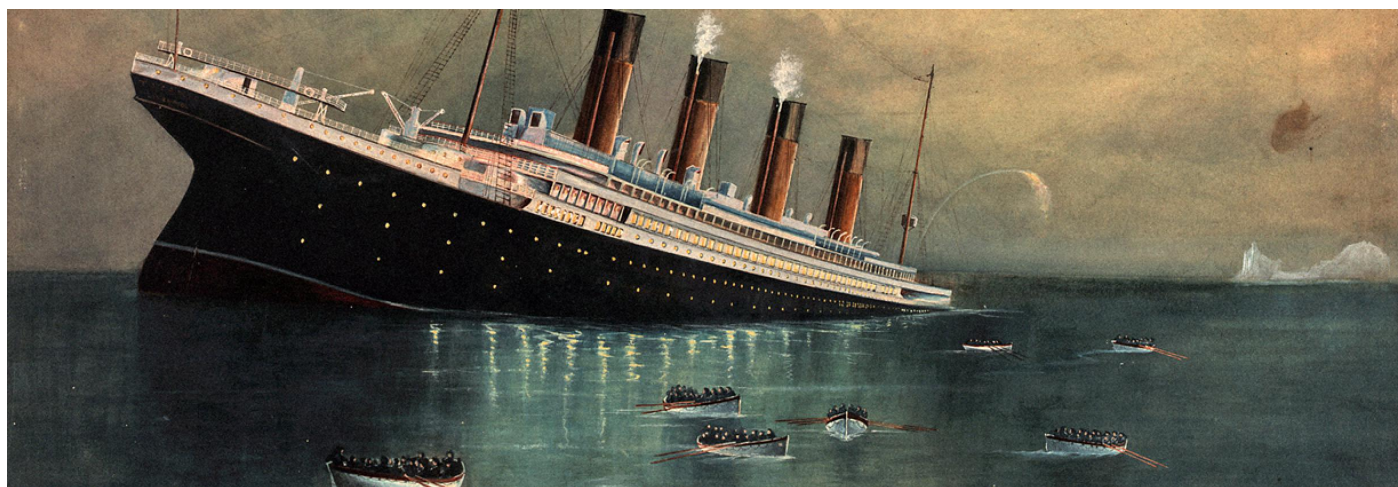
df_big['All_noises'] = df_big.apply(lambda row: np.mean(row[all_noise_columns]
]) > 0.1, axis=1)

# 7.65 s ± 135 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

```
In [64]: # %%timeit

for index in df_big.index:
    df_big.loc[index, 'All_noises'] = np.mean(df_big.loc[index, all_noise_columns]
]) > 0.1
# 16.1 s ± 2.43 s per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

Tabele pivot



```
In [65]: # Incarcarea datelor:

titanic = sns.load_dataset('titanic')
titanic.head()
```

Out[65]:

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True
1	1	1	female	38.0	1	0	71.2833	C	First	woman	False
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	False
3	1	1	female	35.0	1	0	53.1000	S	First	woman	False
4	0	3	male	35.0	0	0	8.0500	S	Third	man	True

Pornim de la urmatoarea problema: care este procentul de femei si barbati supravietuitori? Diferentierea de gen se face dupa coloana 'sex', iar supravietuirea este in coloana 'survived':

```
In [66]: titanic.groupby('sex')['survived'].mean()
```

```
Out[66]: sex
female    0.742038
male      0.188908
Name: survived, dtype: float64
```

Mai departe, se cere determinarea distributiei pe gen si clasa imbarcare, folosind `groupby()` :

```
In [67]: titanic.groupby(['sex', 'class'])['survived'].aggregate('mean').unstack()
```

```
Out[67]:
```

	class	First	Second	Third
sex				
female		0.968085	0.921053	0.500000
male		0.368852	0.157407	0.135447

Acest tip de operatii (grupare dupa doua attribute, calcul de valori agregate) este des intalnit si se numeste pivotare. Pandas introduce suport nativ pentru pivotare, simplificand codul:

```
In [68]: titanic.pivot_table('survived', index='sex', columns='class' )
```

```
Out[68]:
```

	class	First	Second	Third
sex				
female		0.968085	0.921053	0.500000
male		0.368852	0.157407	0.135447

Se poate face pivotare pe mai mult de doua niveluri (mai sus: sex si class). De exemplu, varsta poate fi adaugata pentru analiza, persoane sub 18 ani (copii) si cei peste 18 (adulti). In primul pas se poate face impartirea persoanelor pe cele doua subintervale de varsta (≤ 18 , > 18) folosind `cut` :


```
In [69]: age = pd.cut(titanic['age'], [0, 18, 80], labels=['child', 'adult'])
age.head(15)
```

```
Out[69]: 0      adult
1      adult
2      adult
3      adult
4      adult
5      NaN
6      adult
7      child
8      adult
9      child
10     child
11     adult
12     adult
13     adult
14     child
Name: age, dtype: category
Categories (2, object): [child < adult]
```

```
In [70]: titanic.pivot_table('survived', ['sex', age], 'class')
```

```
Out[70]:
```

	class	First	Second	Third
sex	age			
female	child	0.909091	1.000000	0.511628
	adult	0.972973	0.900000	0.423729
male	child	0.800000	0.600000	0.215686
	adult	0.375000	0.071429	0.133663

```
In [71]: fare_split = pd.cut(titanic.fare, 2, labels=['cheap fare', 'expensive fare'])
```

```
In [72]: fare_split
```

```
Out[72]: 0      cheap fare
1      cheap fare
2      cheap fare
3      cheap fare
4      cheap fare
...
886    cheap fare
887    cheap fare
888    cheap fare
889    cheap fare
890    cheap fare
Name: fare, Length: 891, dtype: category
Categories (2, object): [cheap fare < expensive fare]
```

```
In [73]: titanic.pivot_table('survived', ['sex', age, fare_split], 'class')
```

Out[73]:

		class	First	Second	Third
sex	age	fare			
female	child	cheap fare	0.900000	1.000000	0.511628
		expensive fare	1.000000	NaN	NaN
	adult	cheap fare	0.971429	0.900000	0.423729
		expensive fare	1.000000	NaN	NaN
male	child	cheap fare	0.800000	0.600000	0.215686
	adult	cheap fare	0.369565	0.071429	0.133663
		expensive fare	0.500000	NaN	NaN

In []: