

## PK1

Мартынова П.В. ИУ5-61Б вариант 13 задача 2, датасет 5

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
import seaborn as sns
```

```
df = pd.read_csv('data/states_all_extended.csv')
df.head(10)
```

	PRIMARY_KEY	STATE	YEAR	ENROLL	\
0	1992_ALABAMA	ALABAMA	1992	NaN	
1	1992_ALASKA	ALASKA	1992	NaN	
2	1992_ARIZONA	ARIZONA	1992	NaN	
3	1992_ARKANSAS	ARKANSAS	1992	NaN	
4	1992_CALIFORNIA	CALIFORNIA	1992	NaN	
5	1992_COLORADO	COLORADO	1992	NaN	
6	1992_CONNECTICUT	CONNECTICUT	1992	NaN	
7	1992_DELAWARE	DELAWARE	1992	NaN	
8	1992_DISTRICT_OF_COLUMBIA	DISTRICT_OF_COLUMBIA	1992	NaN	
9	1992_FLORIDA	FLORIDA	1992	NaN	

	TOTAL_REVENUE	FEDERAL_REVENUE	STATE_REVENUE	LOCAL_REVENUE	\
0	2678885.0	304177.0	1659028.0	715680.0	
1	1049591.0	106780.0	720711.0	222100.0	
2	3258079.0	297888.0	1369815.0	1590376.0	
3	1711959.0	178571.0	958785.0	574603.0	
4	26260025.0	2072470.0	16546514.0	7641041.0	
5	3185173.0	163253.0	1307986.0	1713934.0	
6	3834302.0	143542.0	1342539.0	2348221.0	
7	645233.0	45945.0	420942.0	178346.0	
8	709480.0	64749.0	0.0	644731.0	
9	11506299.0	788420.0	5683949.0	5033930.0	

	TOTAL_EXPENDITURE	INSTRUCTION_EXPENDITURE	...	
0	2653798.0	1481703.0	...	NaN
1	972488.0	498362.0	...	NaN
2	3401580.0	1435908.0	...	NaN
3	1743022.0	964323.0	...	NaN
4	27138832.0	14358922.0	...	NaN

5	3264826.0	1642466.0	...	NaN
6	3721338.0	2148041.0	...	NaN
7	638784.0	372722.0	...	NaN
8	742893.0	329160.0	...	NaN
9	11305642.0	5166374.0	...	NaN

	G08_HI_A_MATHEMATICS	G08_AS_A_READING	G08_AS_A_MATHEMATICS	\
0	NaN	NaN	NaN	
1	NaN	NaN	NaN	
2	NaN	NaN	NaN	
3	NaN	NaN	NaN	
4	NaN	NaN	NaN	
5	NaN	NaN	NaN	
6	NaN	NaN	NaN	
7	NaN	NaN	NaN	
8	NaN	NaN	NaN	
9	NaN	NaN	NaN	

	G08_AM_A_READING	G08_AM_A_MATHEMATICS	G08_HP_A_READING	\
0	NaN	NaN	NaN	
1	NaN	NaN	NaN	
2	NaN	NaN	NaN	
3	NaN	NaN	NaN	
4	NaN	NaN	NaN	
5	NaN	NaN	NaN	
6	NaN	NaN	NaN	
7	NaN	NaN	NaN	
8	NaN	NaN	NaN	
9	NaN	NaN	NaN	

	G08_HP_A_MATHEMATICS	G08_TR_A_READING	G08_TR_A_MATHEMATICS
0	NaN	NaN	NaN
1	NaN	NaN	NaN
2	NaN	NaN	NaN
3	NaN	NaN	NaN
4	NaN	NaN	NaN
5	NaN	NaN	NaN
6	NaN	NaN	NaN
7	NaN	NaN	NaN
8	NaN	NaN	NaN
9	NaN	NaN	NaN

[10 rows x 266 columns]

```
df.describe()
```

	YEAR	ENROLL	TOTAL_REVENUE	FEDERAL_REVENUE \
count	1715.000000	1.224000e+03	1.275000e+03	1.275000e+03
mean	2002.075219	9.175416e+05	9.102045e+06	7.677799e+05
std	9.568621	1.066514e+06	1.175962e+07	1.146992e+06
min	1986.000000	4.386600e+04	4.656500e+05	3.102000e+04
25%	1994.000000	2.645145e+05	2.189504e+06	1.899575e+05
50%	2002.000000	6.499335e+05	5.085826e+06	4.035480e+05
75%	2010.000000	1.010532e+06	1.084516e+07	8.279320e+05
max	2019.000000	6.307022e+06	8.921726e+07	9.990221e+06

	STATE_REVENUE	LOCAL_REVENUE	TOTAL_EXPENDITURE \
count	1.275000e+03	1.275000e+03	1.275000e+03
mean	4.223743e+06	4.110522e+06	9.206242e+06
std	5.549735e+06	5.489562e+06	1.199279e+07
min	0.000000e+00	2.209300e+04	4.816650e+05
25%	1.165776e+06	7.151210e+05	2.170404e+06
50%	2.537754e+06	2.058996e+06	5.242672e+06
75%	5.055548e+06	4.755293e+06	1.074420e+07
max	5.090457e+07	3.610526e+07	8.532013e+07

	INSTRUCTION_EXPENDITURE	SUPPORT_SERVICES_EXPENDITURE \
count	1.275000e+03	1.275000e+03
mean	4.768010e+06	2.682587e+06
std	6.300569e+06	3.357214e+06
min	2.655490e+05	1.399630e+05
25%	1.171336e+06	6.380760e+05
50%	2.658253e+06	1.525471e+06
75%	5.561959e+06	3.222924e+06
max	4.396452e+07	2.605802e+07

	OTHER_EXPENDITURE	...	G08_HI_A_READING	G08_HI_A_MATHEMATICS
\count	1.224000e+03	...	246.000000	248.000000
mean	4.299509e+05	...	254.845528	269.995968
std	5.347893e+05	...	5.617077	5.992909
min	1.154100e+04	...	239.000000	255.000000
25%	1.034492e+05	...	251.000000	266.000000
50%	2.717040e+05	...	255.000000	270.000000
75%	5.172222e+05	...	258.000000	274.000000
max	3.995951e+06	...	277.000000	287.000000

	G08_AS_A_READING	G08_AS_A_MATHEMATICS	G08_AM_A_READING \
count	153.000000	157.000000	61.000000
mean	279.803922	306.898089	246.688525
std	9.017570	11.034436	7.630074
min	257.000000	277.000000	229.000000
25%	274.000000	301.000000	242.000000
50%	281.000000	309.000000	247.000000
75%	286.000000	314.000000	252.000000
max	298.000000	332.000000	261.000000

	G08_AM_A_MATHEMATICS	G08_HP_A_READING	G08_HP_A_MATHEMATICS \
count	60.000000	14.000000	13.000000
mean	261.016667	248.285714	264.846154
std	6.652365	6.568322	7.861787
min	240.000000	234.000000	250.000000
25%	257.750000	244.250000	258.000000
50%	261.000000	248.500000	266.000000
75%	265.000000	253.500000	269.000000
max	275.000000	259.000000	276.000000

	G08_TR_A_READING	G08_TR_A_MATHEMATICS
count	141.000000	145.000000
mean	268.085106	284.082759
std	7.218517	7.684313
min	249.000000	266.000000
25%	263.000000	279.000000
50%	267.000000	284.000000
75%	272.000000	289.000000
max	291.000000	305.000000

[8 rows x 264 columns]

### Обработка пропусков в данных

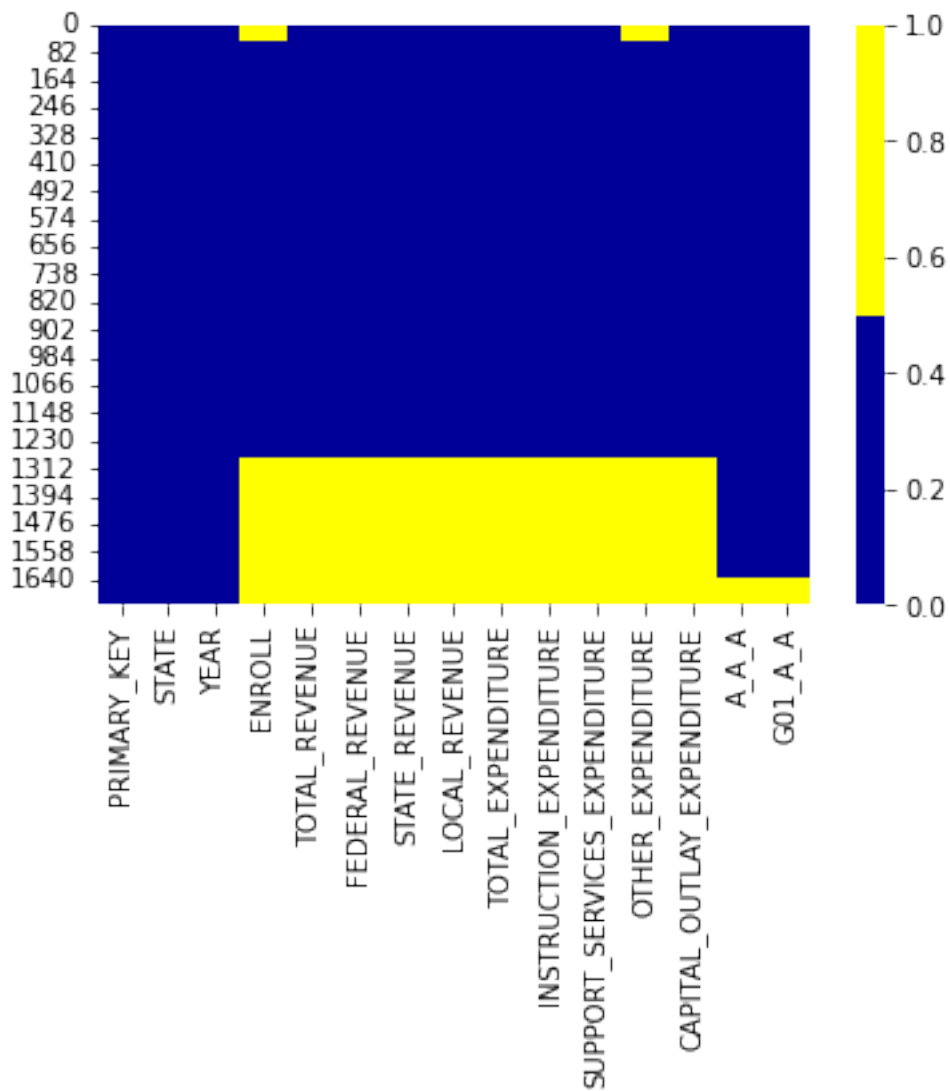
```
cols = df.columns[:15]
```

*# желтый - пропущенные данные, синий - не пропущенные*

```
colours = ['#000099', '#ffff00']
```

```
sns.heatmap(df[cols].isnull(), cmap=sns.color_palette(colours))
```

<AxesSubplot:>



*#Количество пустых ячеек в колонках:*

```
df.isnull().sum()
```

```
PRIMARY_KEY      0
STATE            0
YEAR            0
ENROLL          491
TOTAL_REVENUE    440
...
G08_AM_A_MATHEMATICS  1655
G08_HP_A_READING    1701
G08_HP_A_MATHEMATICS  1702
G08_TR_A_READING    1574
G08_TR_A_MATHEMATICS  1570
Length: 266, dtype: int64
```

*#Типы данных в колонках:*

df.dtypes

```
PRIMARY_KEY      object
STATE            object
YEAR             int64
ENROLL           float64
TOTAL_REVENUE     float64
...
G08_AM_A_MATHEMATICS float64
G08_HP_A_READING   float64
G08_HP_A_MATHEMATICS float64
G08_TR_A_READING   float64
G08_TR_A_MATHEMATICS float64
Length: 266, dtype: object
```

*#Количество пустых числовых значений*

num\_cols = []

total\_count = df.shape[0]

**for** col **in** df.columns:

*# Количество пустых значений*

temp\_null\_count = df[df[col].isnull()].shape[0]

dt = str(df[col].dtype)

**if** temp\_null\_count>0 **and** (dt=='float64' **or** dt=='int64'):

num\_cols.append(col)

temp\_perc = round((temp\_null\_count / total\_count) \* 100.0, 2)

print('Колонка {}. Тип данных {}. Количество пустых значений {}, {}%.'.format(col, dt, temp\_null\_count, temp\_perc))

Колонка ENROLL. Тип данных float64. Количество пустых значений 491, 28.63%.

Колонка TOTAL\_REVENUE. Тип данных float64. Количество пустых значений 440, 25.66%.

Колонка FEDERAL\_REVENUE. Тип данных float64. Количество пустых значений 440, 25.66%.

Колонка STATE\_REVENUE. Тип данных float64. Количество пустых значений 440, 25.66%.

Колонка LOCAL\_REVENUE. Тип данных float64. Количество пустых значений 440, 25.66%.

Колонка TOTAL\_EXPENDITURE. Тип данных float64. Количество пустых значений 440, 25.66%.

Колонка INSTRUCTION\_EXPENDITURE. Тип данных float64. Количество пустых значений 440, 25.66%.

Колонка SUPPORT\_SERVICES\_EXPENDITURE. Тип данных float64. Количество пустых значений 440, 25.66%.

Колонка OTHER\_EXPENDITURE. Тип данных float64. Количество пустых значений 491, 28.63%.

Колонка CAPITAL\_OUTLAY\_EXPENDITURE. Тип данных float64. Количество пустых значений 440, 25.66%.

Колонка A\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G01\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G02\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G03\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G04\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G05\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G07\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G08\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G09\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G10\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G11\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка G12\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка KG\_A\_A. Тип данных float64. Количество пустых значений 83, 4.84%.

Колонка PK\_A\_A. Тип данных float64. Количество пустых значений 173, 10.09%.

Колонка G01-G08\_A\_A. Тип данных float64. Количество пустых значений 695, 40.52%.

Колонка G09-G12\_A\_A. Тип данных float64. Количество пустых значений 644, 37.55%.

Колонка G01\_AM\_F. Тип данных float64. Количество пустых значений 1308, 76.27%.

Колонка G01\_AM\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.

Колонка G01\_AS\_F. Тип данных float64. Количество пустых значений 1307, 76.21%.

Колонка G01\_AS\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.

Колонка G01\_BL\_F. Тип данных float64. Количество пустых значений 1307, 76.21%.

Колонка G01\_BL\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.

Колонка G01\_HI\_F. Тип данных float64. Количество пустых значений 1308, 76.27%.

Колонка G01\_HI\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.

Колонка G01\_HP\_F. Тип данных float64. Количество пустых значений 1351, 78.78%.

Колонка G01\_HP\_M. Тип данных float64. Количество пустых значений 1352, 78.83%.

[illegible]



[illegible]

[illegible]

[illegible]

[illegible]



76.21%.  
Колонка G12\_HI\_F. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка G12\_HI\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка G12\_HP\_F. Тип данных float64. Количество пустых значений 1352, 78.83%.  
Колонка G12\_HP\_M. Тип данных float64. Количество пустых значений 1352, 78.83%.  
Колонка G12\_TR\_F. Тип данных float64. Количество пустых значений 1344, 78.37%.  
Колонка G12\_TR\_M. Тип данных float64. Количество пустых значений 1344, 78.37%.  
Колонка G12\_WH\_F. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка G12\_WH\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка KG\_AM\_F. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка KG\_AM\_M. Тип данных float64. Количество пустых значений 1308, 76.27%.  
Колонка KG\_AS\_F. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка KG\_AS\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка KG\_BL\_F. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка KG\_BL\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка KG\_HI\_F. Тип данных float64. Количество пустых значений 1308, 76.27%.  
Колонка KG\_HI\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка KG\_HP\_F. Тип данных float64. Количество пустых значений 1349, 78.66%.  
Колонка KG\_HP\_M. Тип данных float64. Количество пустых значений 1350, 78.72%.  
Колонка KG\_TR\_F. Тип данных float64. Количество пустых значений 1344, 78.37%.  
Колонка KG\_TR\_M. Тип данных float64. Количество пустых значений 1344, 78.37%.  
Колонка KG\_WH\_F. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка KG\_WH\_M. Тип данных float64. Количество пустых значений 1307, 76.21%.  
Колонка PK\_AM\_F. Тип данных float64. Количество пустых значений 1332, 77.67%.  
Колонка PK\_AM\_M. Тип данных float64. Количество пустых значений 1321, 77.03%.  
Колонка PK\_AS\_F. Тип данных float64. Количество пустых значений 1321,

77.03%.

Колонка PK\_AS\_M. Тип данных float64. Количество пустых значений 1323, 77.14%.

Колонка PK\_BL\_F. Тип данных float64. Количество пустых значений 1321, 77.03%.

Колонка PK\_BL\_M. Тип данных float64. Количество пустых значений 1321, 77.03%.

Колонка PK\_HI\_F. Тип данных float64. Количество пустых значений 1321, 77.03%.

Колонка PK\_HI\_M. Тип данных float64. Количество пустых значений 1321, 77.03%.

Колонка PK\_HP\_F. Тип данных float64. Количество пустых значений 1387, 80.87%.

Колонка PK\_HP\_M. Тип данных float64. Количество пустых значений 1384, 80.7%.

Колонка PK\_TR\_F. Тип данных float64. Количество пустых значений 1357, 79.13%.

Колонка PK\_TR\_M. Тип данных float64. Количество пустых значений 1357, 79.13%.

Колонка PK\_WH\_F. Тип данных float64. Количество пустых значений 1321, 77.03%.

Колонка PK\_WH\_M. Тип данных float64. Количество пустых значений 1321, 77.03%.

Колонка G04\_A\_A\_READING. Тип данных float64. Количество пустых значений 1065, 62.1%.

Колонка G04\_A\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1150, 67.06%.

Колонка G04\_A\_M\_READING. Тип данных float64. Количество пустых значений 1065, 62.1%.

Колонка G04\_A\_M\_MATHEMATICS. Тип данных float64. Количество пустых значений 1150, 67.06%.

Колонка G04\_A\_F\_READING. Тип данных float64. Количество пустых значений 1065, 62.1%.

Колонка G04\_A\_F\_MATHEMATICS. Тип данных float64. Количество пустых значений 1150, 67.06%.

Колонка G04\_WH\_A\_READING. Тип данных float64. Количество пустых значений 1450, 84.55%.

Колонка G04\_WH\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1450, 84.55%.

Колонка G04\_BL\_A\_READING. Тип данных float64. Количество пустых значений 1489, 86.82%.

Колонка G04\_BL\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1486, 86.65%.

Колонка G04\_HI\_A\_READING. Тип данных float64. Количество пустых значений 1465, 85.42%.

Колонка G04\_HI\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1465, 85.42%.

Колонка G04\_AS\_A\_READING. Тип данных float64. Количество пустых значений 1551, 90.44%.

Колонка G04\_AS\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений

значений 1547, 90.2%.  
Колонка G04\_AM\_A\_READING. Тип данных float64. Количество пустых значений 1651, 96.27%.  
Колонка G04\_AM\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1652, 96.33%.  
Колонка G04\_HP\_A\_READING. Тип данных float64. Количество пустых значений 1699, 99.07%.  
Колонка G04\_HP\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1700, 99.13%.  
Колонка G04\_TR\_A\_READING. Тип данных float64. Количество пустых значений 1532, 89.33%.  
Колонка G04\_TR\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1532, 89.33%.  
Колонка G08\_A\_A\_READING. Тип данных float64. Количество пустых значений 1153, 67.23%.  
Колонка G08\_A\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1113, 64.9%.  
Колонка G08\_A\_M\_READING. Тип данных float64. Количество пустых значений 1153, 67.23%.  
Колонка G08\_A\_M\_MATHEMATICS. Тип данных float64. Количество пустых значений 1113, 64.9%.  
Колонка G08\_A\_F\_READING. Тип данных float64. Количество пустых значений 1153, 67.23%.  
Колонка G08\_A\_F\_MATHEMATICS. Тип данных float64. Количество пустых значений 1113, 64.9%.  
Колонка G08\_WH\_A\_READING. Тип данных float64. Количество пустых значений 1450, 84.55%.  
Колонка G08\_WH\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1450, 84.55%.  
Колонка G08\_BL\_A\_READING. Тип данных float64. Количество пустых значений 1493, 87.06%.  
Колонка G08\_BL\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1494, 87.11%.  
Колонка G08\_HI\_A\_READING. Тип данных float64. Количество пустых значений 1469, 85.66%.  
Колонка G08\_HI\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1467, 85.54%.  
Колонка G08\_AS\_A\_READING. Тип данных float64. Количество пустых значений 1562, 91.08%.  
Колонка G08\_AS\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1558, 90.85%.  
Колонка G08\_AM\_A\_READING. Тип данных float64. Количество пустых значений 1654, 96.44%.  
Колонка G08\_AM\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1655, 96.5%.  
Колонка G08\_HP\_A\_READING. Тип данных float64. Количество пустых значений 1701, 99.18%.  
Колонка G08\_HP\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1702, 99.24%.  
Колонка G08\_TR\_A\_READING. Тип данных float64. Количество пустых



значений 1574, 91.78%.

Колонка G08\_TR\_A\_MATHEMATICS. Тип данных float64. Количество пустых значений 1570, 91.55%.

Возьмем в качестве количественного признака признак G06\_A\_A - общее количество учащихся шестого класса. Заменяем пропуски на медианное значение:

```
med = df['G06_A_A'].median()
print(med)
df['G06_A_A'] = df['G06_A_A'].fillna(med)
```

48672.0

```
for col in df.columns:
    pct_missing = np.mean(df[col].isnull())
    print('{} - {}'.format(col, round(pct_missing*100)))
```

```
for col in df.columns:
    missing = df[col].isnull()
    num_missing = np.sum(missing)
```

PRIMARY\_KEY - 0%

STATE - 0%

YEAR - 0%

ENROLL - 29%

TOTAL\_REVENUE - 26%

FEDERAL\_REVENUE - 26%

STATE\_REVENUE - 26%

LOCAL\_REVENUE - 26%

TOTAL\_EXPENDITURE - 26%

INSTRUCTION\_EXPENDITURE - 26%

SUPPORT\_SERVICES\_EXPENDITURE - 26%

OTHER\_EXPENDITURE - 29%

CAPITAL\_OUTLAY\_EXPENDITURE - 26%

A\_A\_A - 5%

G01\_A\_A - 5%

G02\_A\_A - 5%

G03\_A\_A - 5%

G04\_A\_A - 5%

G05\_A\_A - 5%

G06\_A\_A - 0%

G07\_A\_A - 5%

G08\_A\_A - 5%

G09\_A\_A - 5%

G10\_A\_A - 5%

G11\_A\_A - 5%

G12\_A\_A - 5%

KG\_A\_A - 5%

PK\_A\_A - 10%

G01-G08\_A\_A - 41%

G09-G12\_A\_A - 38%

G01\_AM\_F - 76%

G01\_AM\_M - 76%

G01\_AS\_F - 76%

G01\_AS\_M - 76%

G01\_BL\_F - 76%

G01\_BL\_M - 76%

G01\_HI\_F - 76%

G01\_HI\_M - 76%

G01\_HP\_F - 79%

G01\_HP\_M - 79%

G01\_TR\_F - 78%

G01\_TR\_M - 78%

G01\_WH\_F - 76%

G01\_WH\_M - 76%

G02\_AM\_F - 76%

G02\_AM\_M - 76%

G02\_AS\_F - 76%

G02\_AS\_M - 76%

G02\_BL\_F - 76%

G02\_BL\_M - 76%

G02\_HI\_F - 76%

G02\_HI\_M - 76%

G02\_HP\_F - 79%

G02\_HP\_M - 79%

G02\_TR\_F - 78%

G02\_TR\_M - 78%

G02\_WH\_F - 76%

G02\_WH\_M - 76%

G03\_AM\_F - 76%

G03\_AM\_M - 76%

G03\_AS\_F - 76%

G03\_AS\_M - 76%

G03\_BL\_F - 76%

G03\_BL\_M - 76%

G03\_HI\_F - 76%

G03\_HI\_M - 76%

G03\_HP\_F - 79%

G03\_HP\_M - 79%

G03\_TR\_F - 78%

G03\_TR\_M - 78%

G03\_WH\_F - 76%

G03\_WH\_M - 76%

G04\_AM\_F - 76%

G04\_AM\_M - 76%

G04\_AS\_F - 76%

G04\_AS\_M - 76%

G04\_BL\_F - 76%

G04\_BL\_M - 76%

G04\_HI\_F - 76%

G04\_HI\_M - 76%  
G04\_HP\_F - 79%  
G04\_HP\_M - 79%  
G04\_TR\_F - 78%  
G04\_TR\_M - 78%  
G04\_WH\_F - 76%  
G04\_WH\_M - 76%  
G05\_AM\_F - 76%  
G05\_AM\_M - 76%  
G05\_AS\_F - 76%  
G05\_AS\_M - 76%  
G05\_BL\_F - 76%  
G05\_BL\_M - 76%  
G05\_HI\_F - 76%  
G05\_HI\_M - 76%  
G05\_HP\_F - 79%  
G05\_HP\_M - 79%  
G05\_TR\_F - 78%  
G05\_TR\_M - 78%  
G05\_WH\_F - 76%  
G05\_WH\_M - 76%  
G06\_AM\_F - 76%  
G06\_AM\_M - 76%  
G06\_AS\_F - 76%  
G06\_AS\_M - 76%  
G06\_BL\_F - 76%  
G06\_BL\_M - 76%  
G06\_HI\_F - 76%  
G06\_HI\_M - 76%  
G06\_HP\_F - 79%  
G06\_HP\_M - 79%  
G06\_TR\_F - 78%  
G06\_TR\_M - 78%  
G06\_WH\_F - 76%  
G06\_WH\_M - 76%  
G07\_AM\_F - 76%  
G07\_AM\_M - 76%  
G07\_AS\_F - 76%  
G07\_AS\_M - 76%  
G07\_BL\_F - 76%  
G07\_BL\_M - 76%  
G07\_HI\_F - 76%  
G07\_HI\_M - 76%  
G07\_HP\_F - 79%  
G07\_HP\_M - 79%  
G07\_TR\_F - 78%  
G07\_TR\_M - 78%  
G07\_WH\_F - 76%  
G07\_WH\_M - 76%  
G08\_AM\_F - 76%

G08\_AM\_M - 76%  
G08\_AS\_F - 76%  
G08\_AS\_M - 76%  
G08\_BL\_F - 76%  
G08\_BL\_M - 76%  
G08\_HI\_F - 76%  
G08\_HI\_M - 76%  
G08\_HP\_F - 79%  
G08\_HP\_M - 79%  
G08\_TR\_F - 78%  
G08\_TR\_M - 78%  
G08\_WH\_F - 76%  
G08\_WH\_M - 76%  
G09\_AM\_F - 76%  
G09\_AM\_M - 76%  
G09\_AS\_F - 76%  
G09\_AS\_M - 76%  
G09\_BL\_F - 76%  
G09\_BL\_M - 76%  
G09\_HI\_F - 76%  
G09\_HI\_M - 76%  
G09\_HP\_F - 79%  
G09\_HP\_M - 79%  
G09\_TR\_F - 78%  
G09\_TR\_M - 78%  
G09\_WH\_F - 76%  
G09\_WH\_M - 76%  
G10\_AM\_F - 76%  
G10\_AM\_M - 76%  
G10\_AS\_F - 76%  
G10\_AS\_M - 76%  
G10\_BL\_F - 76%  
G10\_BL\_M - 76%  
G10\_HI\_F - 76%  
G10\_HI\_M - 76%  
G10\_HP\_F - 79%  
G10\_HP\_M - 79%  
G10\_TR\_F - 78%  
G10\_TR\_M - 78%  
G10\_WH\_F - 76%  
G10\_WH\_M - 76%  
G11\_AM\_F - 76%  
G11\_AM\_M - 76%  
G11\_AS\_F - 76%  
G11\_AS\_M - 76%  
G11\_BL\_F - 76%  
G11\_BL\_M - 76%  
G11\_HI\_F - 76%  
G11\_HI\_M - 76%  
G11\_HP\_F - 79%

G11\_HP\_M - 79%  
G11\_TR\_F - 78%  
G11\_TR\_M - 78%  
G11\_WH\_F - 76%  
G11\_WH\_M - 76%  
G12\_AM\_F - 76%  
G12\_AM\_M - 76%  
G12\_AS\_F - 76%  
G12\_AS\_M - 76%  
G12\_BL\_F - 76%  
G12\_BL\_M - 76%  
G12\_HI\_F - 76%  
G12\_HI\_M - 76%  
G12\_HP\_F - 79%  
G12\_HP\_M - 79%  
G12\_TR\_F - 78%  
G12\_TR\_M - 78%  
G12\_WH\_F - 76%  
G12\_WH\_M - 76%  
KG\_AM\_F - 76%  
KG\_AM\_M - 76%  
KG\_AS\_F - 76%  
KG\_AS\_M - 76%  
KG\_BL\_F - 76%  
KG\_BL\_M - 76%  
KG\_HI\_F - 76%  
KG\_HI\_M - 76%  
KG\_HP\_F - 79%  
KG\_HP\_M - 79%  
KG\_TR\_F - 78%  
KG\_TR\_M - 78%  
KG\_WH\_F - 76%  
KG\_WH\_M - 76%  
PK\_AM\_F - 78%  
PK\_AM\_M - 77%  
PK\_AS\_F - 77%  
PK\_AS\_M - 77%  
PK\_BL\_F - 77%  
PK\_BL\_M - 77%  
PK\_HI\_F - 77%  
PK\_HI\_M - 77%  
PK\_HP\_F - 81%  
PK\_HP\_M - 81%  
PK\_TR\_F - 79%  
PK\_TR\_M - 79%  
PK\_WH\_F - 77%  
PK\_WH\_M - 77%  
G04\_A\_A\_READING - 62%  
G04\_A\_A\_MATHEMATICS - 67%  
G04\_A\_M\_READING - 62%

G04\_A\_M\_MATHEMATICS - 67%  
 G04\_A\_F\_READING - 62%  
 G04\_A\_F\_MATHEMATICS - 67%  
 G04\_WH\_A\_READING - 85%  
 G04\_WH\_A\_MATHEMATICS - 85%  
 G04\_BL\_A\_READING - 87%  
 G04\_BL\_A\_MATHEMATICS - 87%  
 G04\_HI\_A\_READING - 85%  
 G04\_HI\_A\_MATHEMATICS - 85%  
 G04\_AS\_A\_READING - 90%  
 G04\_AS\_A\_MATHEMATICS - 90%  
 G04\_AM\_A\_READING - 96%  
 G04\_AM\_A\_MATHEMATICS - 96%  
 G04\_HP\_A\_READING - 99%  
 G04\_HP\_A\_MATHEMATICS - 99%  
 G04\_TR\_A\_READING - 89%  
 G04\_TR\_A\_MATHEMATICS - 89%  
 G08\_A\_A\_READING - 67%  
 G08\_A\_A\_MATHEMATICS - 65%  
 G08\_A\_M\_READING - 67%  
 G08\_A\_M\_MATHEMATICS - 65%  
 G08\_A\_F\_READING - 67%  
 G08\_A\_F\_MATHEMATICS - 65%  
 G08\_WH\_A\_READING - 85%  
 G08\_WH\_A\_MATHEMATICS - 85%  
 G08\_BL\_A\_READING - 87%  
 G08\_BL\_A\_MATHEMATICS - 87%  
 G08\_HI\_A\_READING - 86%  
 G08\_HI\_A\_MATHEMATICS - 86%  
 G08\_AS\_A\_READING - 91%  
 G08\_AS\_A\_MATHEMATICS - 91%  
 G08\_AM\_A\_READING - 96%  
 G08\_AM\_A\_MATHEMATICS - 97%  
 G08\_HP\_A\_READING - 99%  
 G08\_HP\_A\_MATHEMATICS - 99%  
 G08\_TR\_A\_READING - 92%  
 G08\_TR\_A\_MATHEMATICS - 92%

```
print(df['G06_A_A'])
```

0	59929.0
1	9542.0
2	53832.0
3	35017.0
4	399776.0
	...
1710	48672.0
1711	48672.0
1712	48672.0
1713	48672.0

```
1714      48672.0
Name: G06_A_A, Length: 1715, dtype: float64
```

В качестве категориального признака можно было бы взять, например, state, но так как в этом столбце нет пропущенных значений, обработаем пропуски еще для одного количественного признака: GRADES 9\_12 - Число учащихся с девятого по двенадцатый классы. Заполним пропуски средним значением.

```
mean = df['G09-G12_A_A'].mean()
print(mean)
df['G09-G12_A_A'] = df['G09-G12_A_A'].fillna(mean)
```

```
282069.08496732026
```

```
for col in df.columns:
    pct_missing = np.mean(df[col].isnull())
    print('{} - {}'.format(col, round(pct_missing*100)))
```

```
for col in df.columns:
    missing = df[col].isnull()
    num_missing = np.sum(missing)
```

```
PRIMARY_KEY - 0%
STATE - 0%
YEAR - 0%
ENROLL - 29%
TOTAL_REVENUE - 26%
FEDERAL_REVENUE - 26%
STATE_REVENUE - 26%
LOCAL_REVENUE - 26%
TOTAL_EXPENDITURE - 26%
INSTRUCTION_EXPENDITURE - 26%
SUPPORT_SERVICES_EXPENDITURE - 26%
OTHER_EXPENDITURE - 29%
CAPITAL_OUTLAY_EXPENDITURE - 26%
A_A_A - 5%
G01_A_A - 5%
G02_A_A - 5%
G03_A_A - 5%
G04_A_A - 5%
G05_A_A - 5%
G06_A_A - 0%
G07_A_A - 5%
G08_A_A - 5%
G09_A_A - 5%
G10_A_A - 5%
G11_A_A - 5%
G12_A_A - 5%
KG_A_A - 5%
PK_A_A - 10%
```

G01-G08\_A\_A - 41%

G09-G12\_A\_A - 0%

G01\_AM\_F - 76%

G01\_AM\_M - 76%

G01\_AS\_F - 76%

G01\_AS\_M - 76%

G01\_BL\_F - 76%

G01\_BL\_M - 76%

G01\_HI\_F - 76%

G01\_HI\_M - 76%

G01\_HP\_F - 79%

G01\_HP\_M - 79%

G01\_TR\_F - 78%

G01\_TR\_M - 78%

G01\_WH\_F - 76%

G01\_WH\_M - 76%

G02\_AM\_F - 76%

G02\_AM\_M - 76%

G02\_AS\_F - 76%

G02\_AS\_M - 76%

G02\_BL\_F - 76%

G02\_BL\_M - 76%

G02\_HI\_F - 76%

G02\_HI\_M - 76%

G02\_HP\_F - 79%

G02\_HP\_M - 79%

G02\_TR\_F - 78%

G02\_TR\_M - 78%

G02\_WH\_F - 76%

G02\_WH\_M - 76%

G03\_AM\_F - 76%

G03\_AM\_M - 76%

G03\_AS\_F - 76%

G03\_AS\_M - 76%

G03\_BL\_F - 76%

G03\_BL\_M - 76%

G03\_HI\_F - 76%

G03\_HI\_M - 76%

G03\_HP\_F - 79%

G03\_HP\_M - 79%

G03\_TR\_F - 78%

G03\_TR\_M - 78%

G03\_WH\_F - 76%

G03\_WH\_M - 76%

G04\_AM\_F - 76%

G04\_AM\_M - 76%

G04\_AS\_F - 76%

G04\_AS\_M - 76%

G04\_BL\_F - 76%

G04\_BL\_M - 76%



G04\_HI\_F - 76%  
G04\_HI\_M - 76%  
G04\_HP\_F - 79%  
G04\_HP\_M - 79%  
G04\_TR\_F - 78%  
G04\_TR\_M - 78%  
G04\_WH\_F - 76%  
G04\_WH\_M - 76%  
G05\_AM\_F - 76%  
G05\_AM\_M - 76%  
G05\_AS\_F - 76%  
G05\_AS\_M - 76%  
G05\_BL\_F - 76%  
G05\_BL\_M - 76%  
G05\_HI\_F - 76%  
G05\_HI\_M - 76%  
G05\_HP\_F - 79%  
G05\_HP\_M - 79%  
G05\_TR\_F - 78%  
G05\_TR\_M - 78%  
G05\_WH\_F - 76%  
G05\_WH\_M - 76%  
G06\_AM\_F - 76%  
G06\_AM\_M - 76%  
G06\_AS\_F - 76%  
G06\_AS\_M - 76%  
G06\_BL\_F - 76%  
G06\_BL\_M - 76%  
G06\_HI\_F - 76%  
G06\_HI\_M - 76%  
G06\_HP\_F - 79%  
G06\_HP\_M - 79%  
G06\_TR\_F - 78%  
G06\_TR\_M - 78%  
G06\_WH\_F - 76%  
G06\_WH\_M - 76%  
G07\_AM\_F - 76%  
G07\_AM\_M - 76%  
G07\_AS\_F - 76%  
G07\_AS\_M - 76%  
G07\_BL\_F - 76%  
G07\_BL\_M - 76%  
G07\_HI\_F - 76%  
G07\_HI\_M - 76%  
G07\_HP\_F - 79%  
G07\_HP\_M - 79%  
G07\_TR\_F - 78%  
G07\_TR\_M - 78%  
G07\_WH\_F - 76%  
G07\_WH\_M - 76%

G08\_AM\_F - 76%  
G08\_AM\_M - 76%  
G08\_AS\_F - 76%  
G08\_AS\_M - 76%  
G08\_BL\_F - 76%  
G08\_BL\_M - 76%  
G08\_HI\_F - 76%  
G08\_HI\_M - 76%  
G08\_HP\_F - 79%  
G08\_HP\_M - 79%  
G08\_TR\_F - 78%  
G08\_TR\_M - 78%  
G08\_WH\_F - 76%  
G08\_WH\_M - 76%  
G09\_AM\_F - 76%  
G09\_AM\_M - 76%  
G09\_AS\_F - 76%  
G09\_AS\_M - 76%  
G09\_BL\_F - 76%  
G09\_BL\_M - 76%  
G09\_HI\_F - 76%  
G09\_HI\_M - 76%  
G09\_HP\_F - 79%  
G09\_HP\_M - 79%  
G09\_TR\_F - 78%  
G09\_TR\_M - 78%  
G09\_WH\_F - 76%  
G09\_WH\_M - 76%  
G10\_AM\_F - 76%  
G10\_AM\_M - 76%  
G10\_AS\_F - 76%  
G10\_AS\_M - 76%  
G10\_BL\_F - 76%  
G10\_BL\_M - 76%  
G10\_HI\_F - 76%  
G10\_HI\_M - 76%  
G10\_HP\_F - 79%  
G10\_HP\_M - 79%  
G10\_TR\_F - 78%  
G10\_TR\_M - 78%  
G10\_WH\_F - 76%  
G10\_WH\_M - 76%  
G11\_AM\_F - 76%  
G11\_AM\_M - 76%  
G11\_AS\_F - 76%  
G11\_AS\_M - 76%  
G11\_BL\_F - 76%  
G11\_BL\_M - 76%  
G11\_HI\_F - 76%  
G11\_HI\_M - 76%

G11\_HP\_F - 79%  
G11\_HP\_M - 79%  
G11\_TR\_F - 78%  
G11\_TR\_M - 78%  
G11\_WH\_F - 76%  
G11\_WH\_M - 76%  
G12\_AM\_F - 76%  
G12\_AM\_M - 76%  
G12\_AS\_F - 76%  
G12\_AS\_M - 76%  
G12\_BL\_F - 76%  
G12\_BL\_M - 76%  
G12\_HI\_F - 76%  
G12\_HI\_M - 76%  
G12\_HP\_F - 79%  
G12\_HP\_M - 79%  
G12\_TR\_F - 78%  
G12\_TR\_M - 78%  
G12\_WH\_F - 76%  
G12\_WH\_M - 76%  
KG\_AM\_F - 76%  
KG\_AM\_M - 76%  
KG\_AS\_F - 76%  
KG\_AS\_M - 76%  
KG\_BL\_F - 76%  
KG\_BL\_M - 76%  
KG\_HI\_F - 76%  
KG\_HI\_M - 76%  
KG\_HP\_F - 79%  
KG\_HP\_M - 79%  
KG\_TR\_F - 78%  
KG\_TR\_M - 78%  
KG\_WH\_F - 76%  
KG\_WH\_M - 76%  
PK\_AM\_F - 78%  
PK\_AM\_M - 77%  
PK\_AS\_F - 77%  
PK\_AS\_M - 77%  
PK\_BL\_F - 77%  
PK\_BL\_M - 77%  
PK\_HI\_F - 77%  
PK\_HI\_M - 77%  
PK\_HP\_F - 81%  
PK\_HP\_M - 81%  
PK\_TR\_F - 79%  
PK\_TR\_M - 79%  
PK\_WH\_F - 77%  
PK\_WH\_M - 77%  
G04\_A\_A\_READING - 62%  
G04\_A\_A\_MATHEMATICS - 67%

G04\_A\_M\_READING - 62%  
 G04\_A\_M\_MATHEMATICS - 67%  
 G04\_A\_F\_READING - 62%  
 G04\_A\_F\_MATHEMATICS - 67%  
 G04\_WH\_A\_READING - 85%  
 G04\_WH\_A\_MATHEMATICS - 85%  
 G04\_BL\_A\_READING - 87%  
 G04\_BL\_A\_MATHEMATICS - 87%  
 G04\_HI\_A\_READING - 85%  
 G04\_HI\_A\_MATHEMATICS - 85%  
 G04\_AS\_A\_READING - 90%  
 G04\_AS\_A\_MATHEMATICS - 90%  
 G04\_AM\_A\_READING - 96%  
 G04\_AM\_A\_MATHEMATICS - 96%  
 G04\_HP\_A\_READING - 99%  
 G04\_HP\_A\_MATHEMATICS - 99%  
 G04\_TR\_A\_READING - 89%  
 G04\_TR\_A\_MATHEMATICS - 89%  
 G08\_A\_A\_READING - 67%  
 G08\_A\_A\_MATHEMATICS - 65%  
 G08\_A\_M\_READING - 67%  
 G08\_A\_M\_MATHEMATICS - 65%  
 G08\_A\_F\_READING - 67%  
 G08\_A\_F\_MATHEMATICS - 65%  
 G08\_WH\_A\_READING - 85%  
 G08\_WH\_A\_MATHEMATICS - 85%  
 G08\_BL\_A\_READING - 87%  
 G08\_BL\_A\_MATHEMATICS - 87%  
 G08\_HI\_A\_READING - 86%  
 G08\_HI\_A\_MATHEMATICS - 86%  
 G08\_AS\_A\_READING - 91%  
 G08\_AS\_A\_MATHEMATICS - 91%  
 G08\_AM\_A\_READING - 96%  
 G08\_AM\_A\_MATHEMATICS - 97%  
 G08\_HP\_A\_READING - 99%  
 G08\_HP\_A\_MATHEMATICS - 99%  
 G08\_TR\_A\_READING - 92%  
 G08\_TR\_A\_MATHEMATICS - 92%

При обработке пропусков в данных для категориального признака можно было бы использовать стратегии "most\_frequent" или "constant"

#### Диаграмма рассеивания

```

fig, ax = plt.subplots(figsize=(10,10))
plt.xticks(rotation=90)
sns.scatterplot(ax=ax, x='STATE', y='G09-G12_A_A', data=df)

<AxesSubplot:xlabel='STATE', ylabel='G09-G12_A_A'>
  
```

