

Московский Государственный Технический Университет
им. Н. Э. Баумана

Лабораторная работа №1
по курсу: «Технологии машинного обучения»

**Разведочный анализ данных. Исследование и
визуализация данных.**

Выполнила:
Студентка группы ИУ5-63
Нурлыева Д.Д.

Москва
2019

Задание:

- Выбрать набор данных (датасет).
- Создать ноутбук, который содержит следующие разделы:
 1. Текстовое описание выбранного Вами набора данных.
 2. Основные характеристики датасета.
 3. Визуальное исследование датасета.
 4. Информация о корреляции признаков.
- Сформировать отчет и разместить его в своем репозитории на github.

Текстовое описание набора данных

В качестве набора данных мы будем использовать набор данных Heart Disease UCI - <https://www.kaggle.com/ronitf/heart-disease-uci> В датасете отражено наличие сердечного заболевания у пациента в зависимости от разных признаков.

Датасет содержит следующие колонки:

age - возраст в годах

sex - (1 = мужчина; 0 = женщина)

cp - тип боли в груди

trestbps - артериальное давление в состоянии покоя (в мм рт. ст. при поступлении в стационар)

chol - холестерин в мг/дл

fbs - уровень сахара в крови натощак > 120 мг / дл) (1 = да; 0 = нет)

restecg - электрокардиографические результаты покоя

thalach - максимальная ЧСС

exang - стенокардия, вызванная физическими упражнениями (1 = да; 0 = Нет)

oldpeak - Депрессия, вызванная физическими упражнениями относительно покоя

slope - наклон пика упражнения сегмента

ca - количество крупных сосудов (0-3)

thal - 3 = нормальный; 6 = фиксированный дефект; 7 = реверзибельный дефект

target - заболевание 1-есть или 0-нет

Текст программы:

```
In [35]:
```

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(style="ticks")
```

```
In [36]:
```

```
data = pd.read_csv('/Users/user/Desktop/data.csv')
```

```
In [37]:
```

```
data.head()
```

```
Out[37]:
```

	age	sex	cp	trestbps	cho l	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

```
In [38]:
```

```
data.shape
```

```
Out[38]:
```

```
(303, 14)
```

```
In [39]:
```

```
total_count = data.shape[0]
```

```
print('Кол-во строк:', total_count)
```

Кол-во строк: 303

```
In [40]:
```

```
data.columns
```

```
Out[40]:
```

```
Index(['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs',  
      'restecg', 'thalach',  
      'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target'],  
      dtype='object')
```

```
In [41]:
```

```
data.dtypes
```

```
Out[41]:
```

```
age          int64  
sex          int64  
cp           int64  
trestbps     int64  
chol         int64  
fbs          int64  
restecg      int64  
thalach      int64  
exang        int64  
oldpeak      float64  
slope        int64  
ca           int64  
thal         int64
```

```
target          int64
dtype: object
In [42]:
for col in data.columns:
    print('{} - 
{}'.format(col,data[data[col].isnull()].shape[0]))
```

```
age - 0
sex - 0
cp - 0
trestbps - 0
chol - 0
fbs - 0
restecg - 0
thalach - 0
exang - 0
oldpeak - 0
slope - 0
ca - 0
thal - 0
target - 0
```

```
In [43]:
data.describe()
```

Out[43]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
count	303.00000	303.00000	303.00000	303.00000	303.00000	303.00000	303.00000	303.00000	303.00000	303.00000	303.00000	303.00000	303.00000	303.00000
mean	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.528053	149.646865	0.326733	1.039604	1.399340	0.729373	2.313531	0.544554
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860	22.905161	0.469794	1.161075	0.616226	1.022606	0.612277	0.498835
min	29.00000	0.00000	0.00000	94.00000	126.00000	0.00000	0.00000	71.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
25 %	47.50000	0.00000	0.00000	120.00000	211.00000	0.00000	0.00000	133.50000	0.00000	0.00000	1.00000	0.00000	2.00000	0.00000
50 %	55.00000	1.00000	1.00000	130.00000	240.00000	0.00000	1.00000	153.00000	0.00000	0.80000	1.00000	0.00000	2.00000	1.00000
75 %	61.00000	1.00000	2.00000	140.00000	274.50000	0.00000	1.00000	166.00000	1.00000	1.60000	2.00000	1.00000	3.00000	1.00000
max	77.00000	1.00000	3.00000	200.00000	564.00000	1.00000	2.00000	202.00000	1.00000	6.20000	2.00000	4.00000	3.00000	1.00000

```
In [50]:
# Уникальные значения для целевого признака
data['target'].unique()
Out[50]:
array([1, 0])
```

```
In [49]:
```

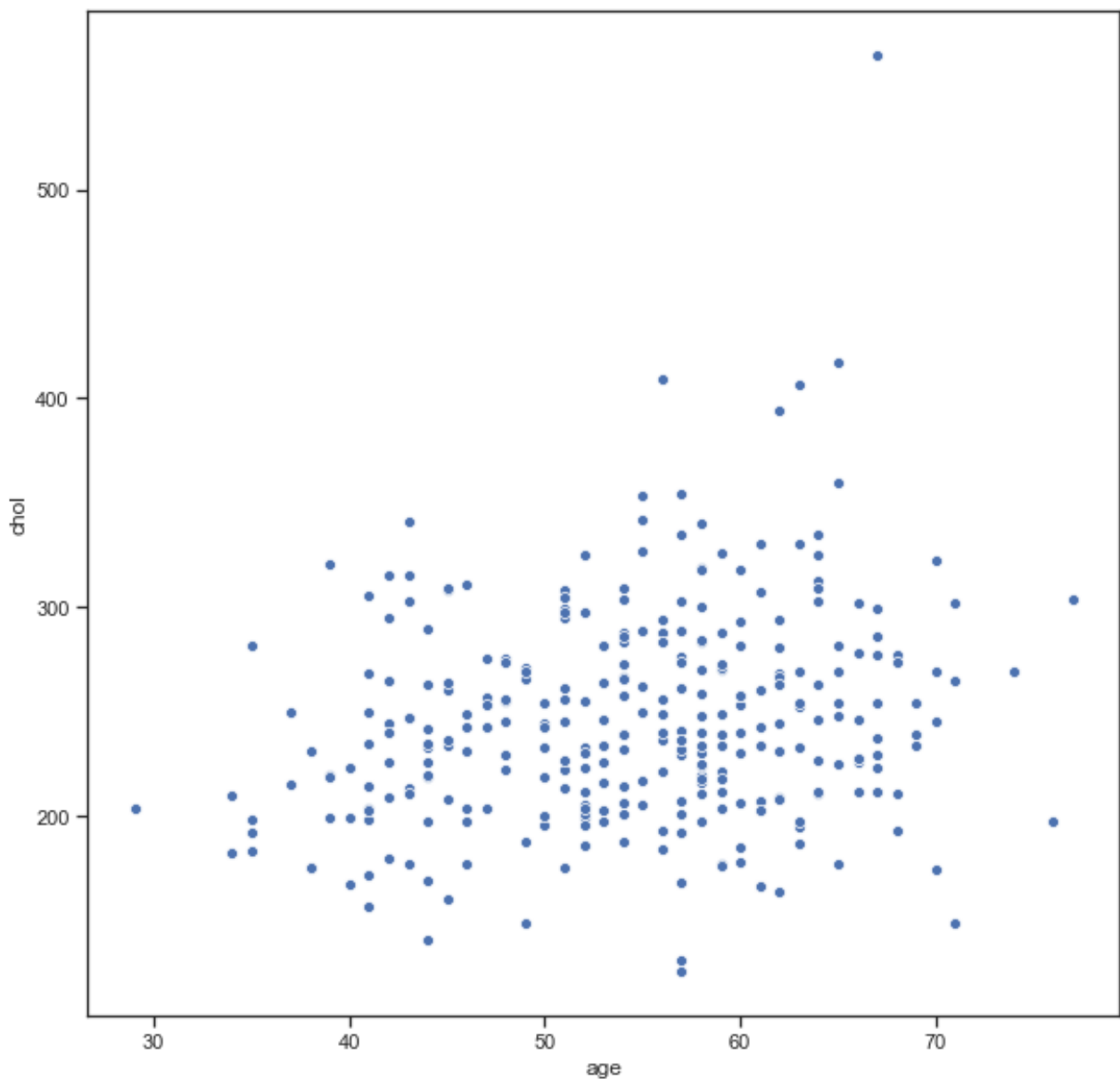
```
# Зависимость между возрастом пациента и содержанием  
холестерина в крови
```

```
fig, ax = plt.subplots(figsize=(10,10))
```

```
sns.scatterplot(ax=ax, x='age', y='chol', data=data)
```

```
Out[49]:
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a1e5e7470>
```



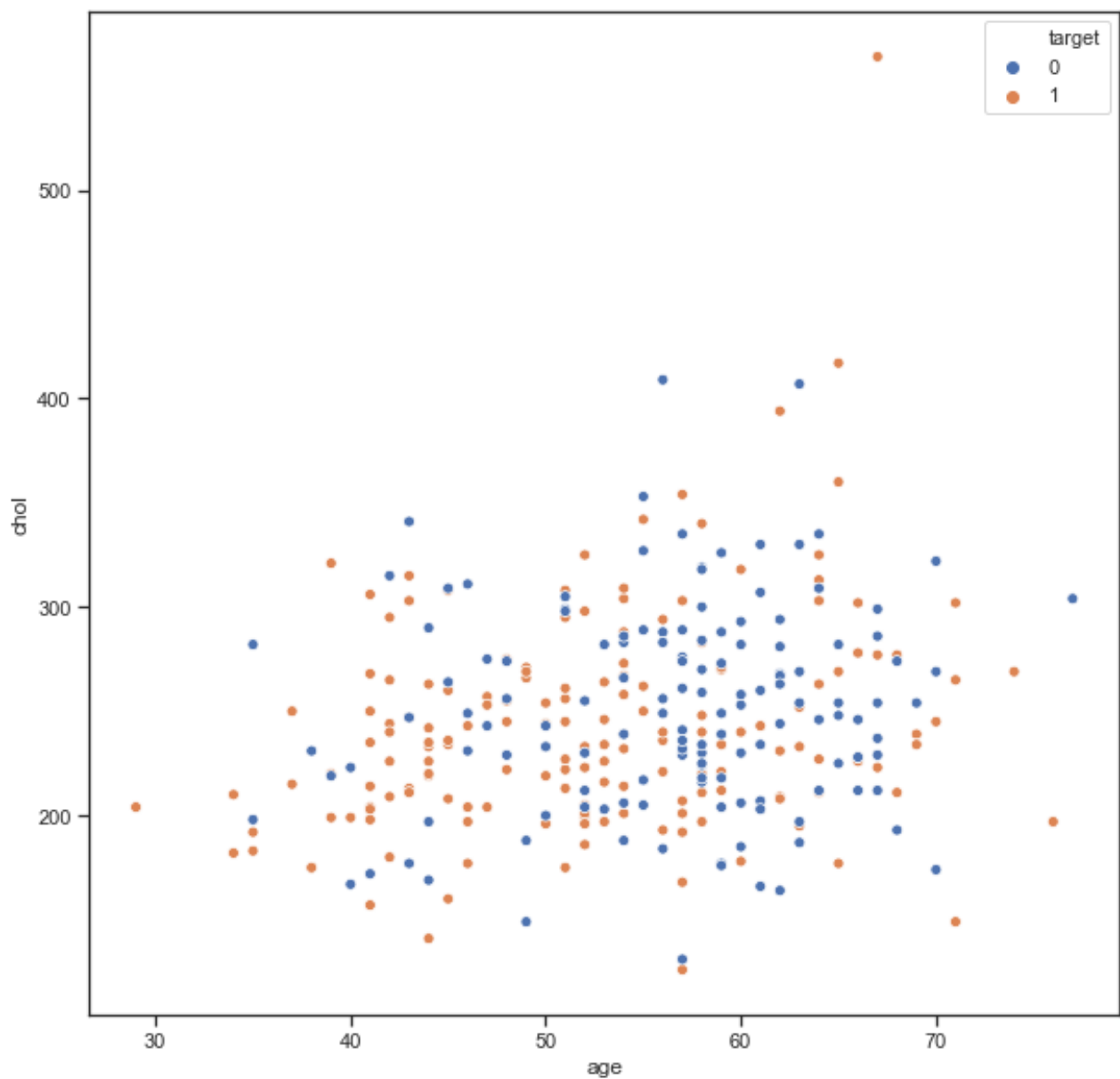
```
In [54]:
```

```
fig, ax = plt.subplots(figsize=(10,10))
```

```
sns.scatterplot(ax=ax, x='age', y='chol', data=data,  
hue='target')
```

```
Out[54]:
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a1e9e2518>
```



In [57]:

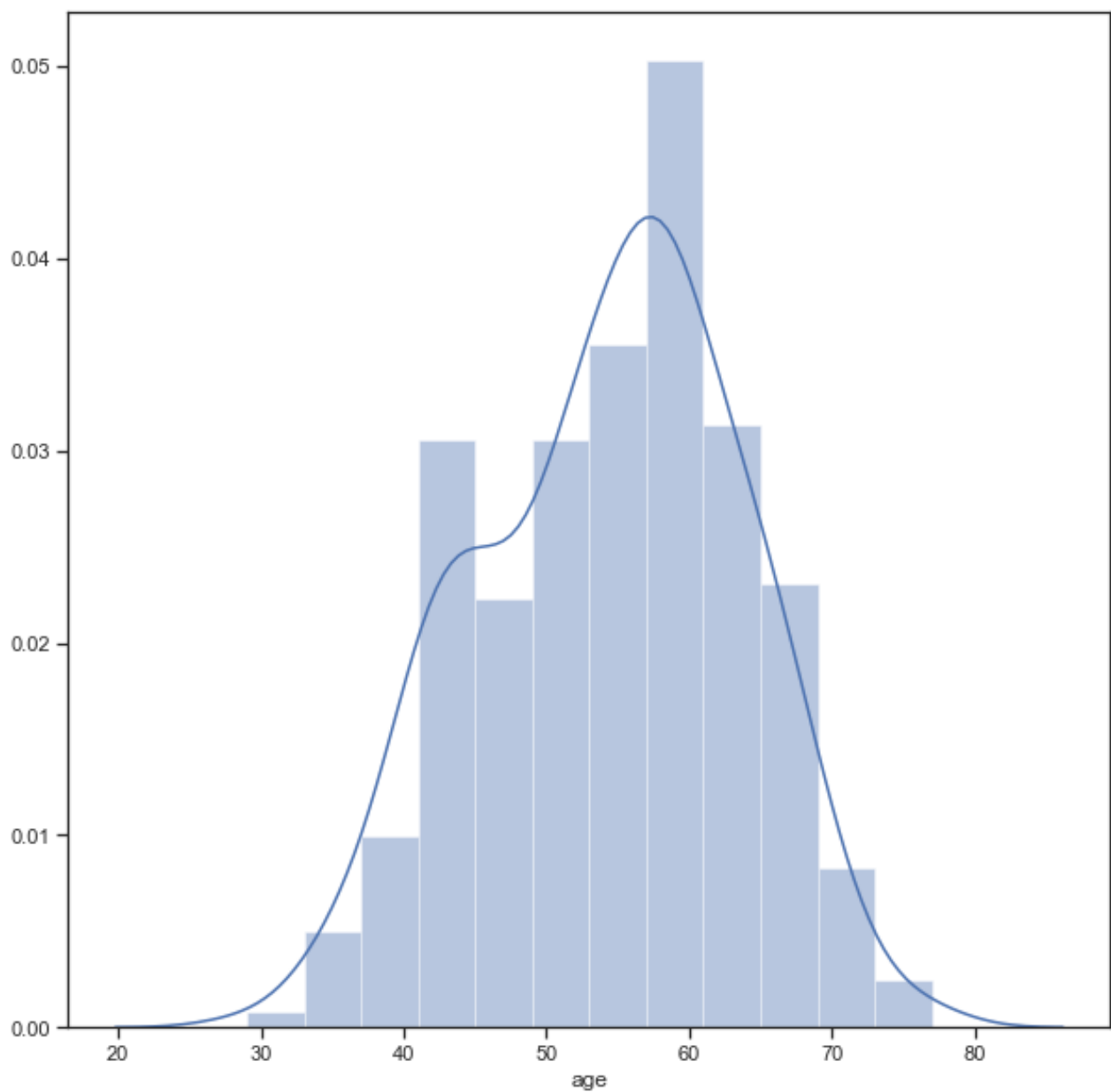
```
#Построение гистограммы
```

```
fig, ax = plt.subplots(figsize=(10,10))
```

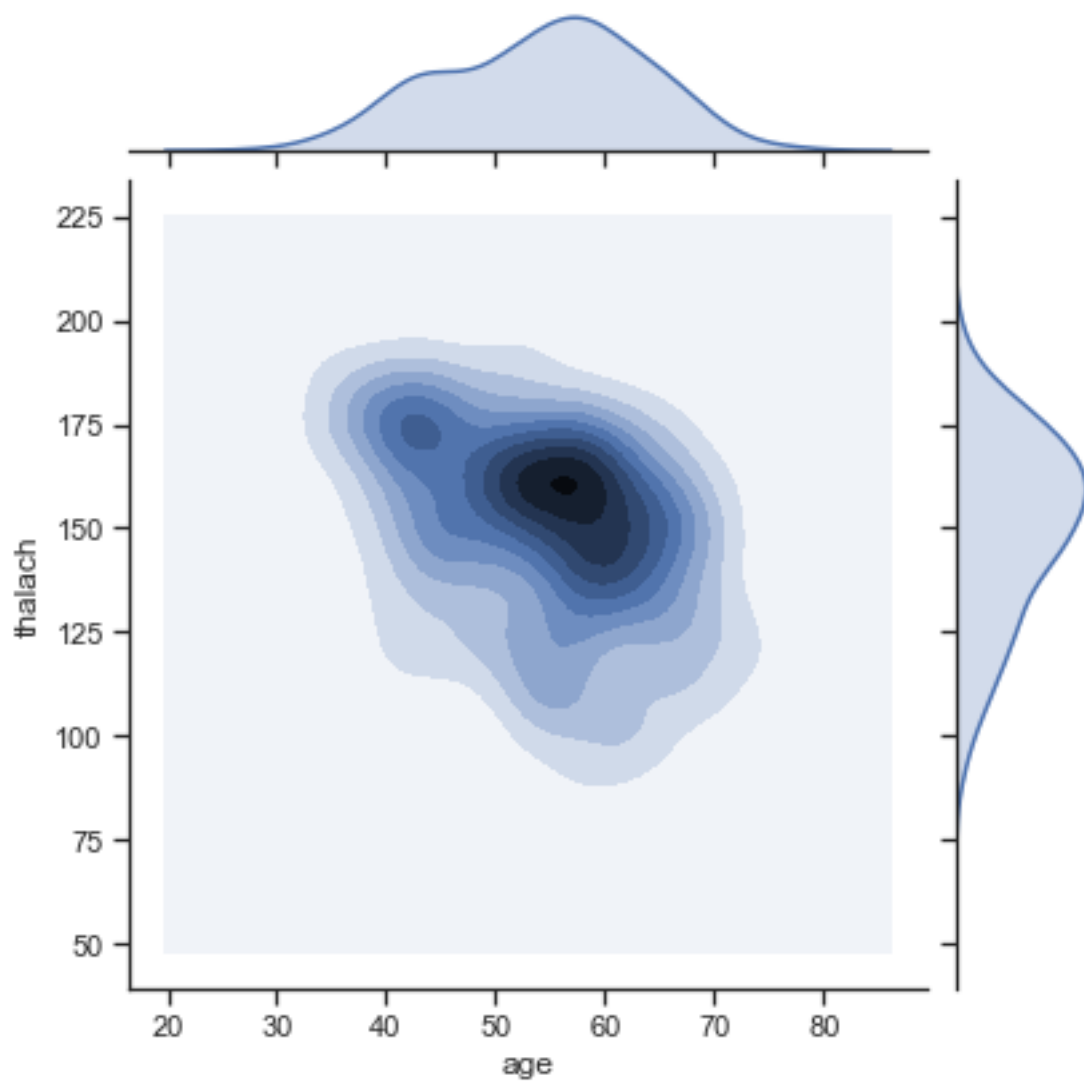
```
sns.distplot(data['age'])
```

Out[57]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a1f0210b8>
```



```
In [58]:  
#jointplot  
sns.jointplot(x='age', y='thalach', data=data, kind="kde")  
Out[58]:  
<seaborn.axisgrid.JointGrid at 0x1a1f3a9780>
```



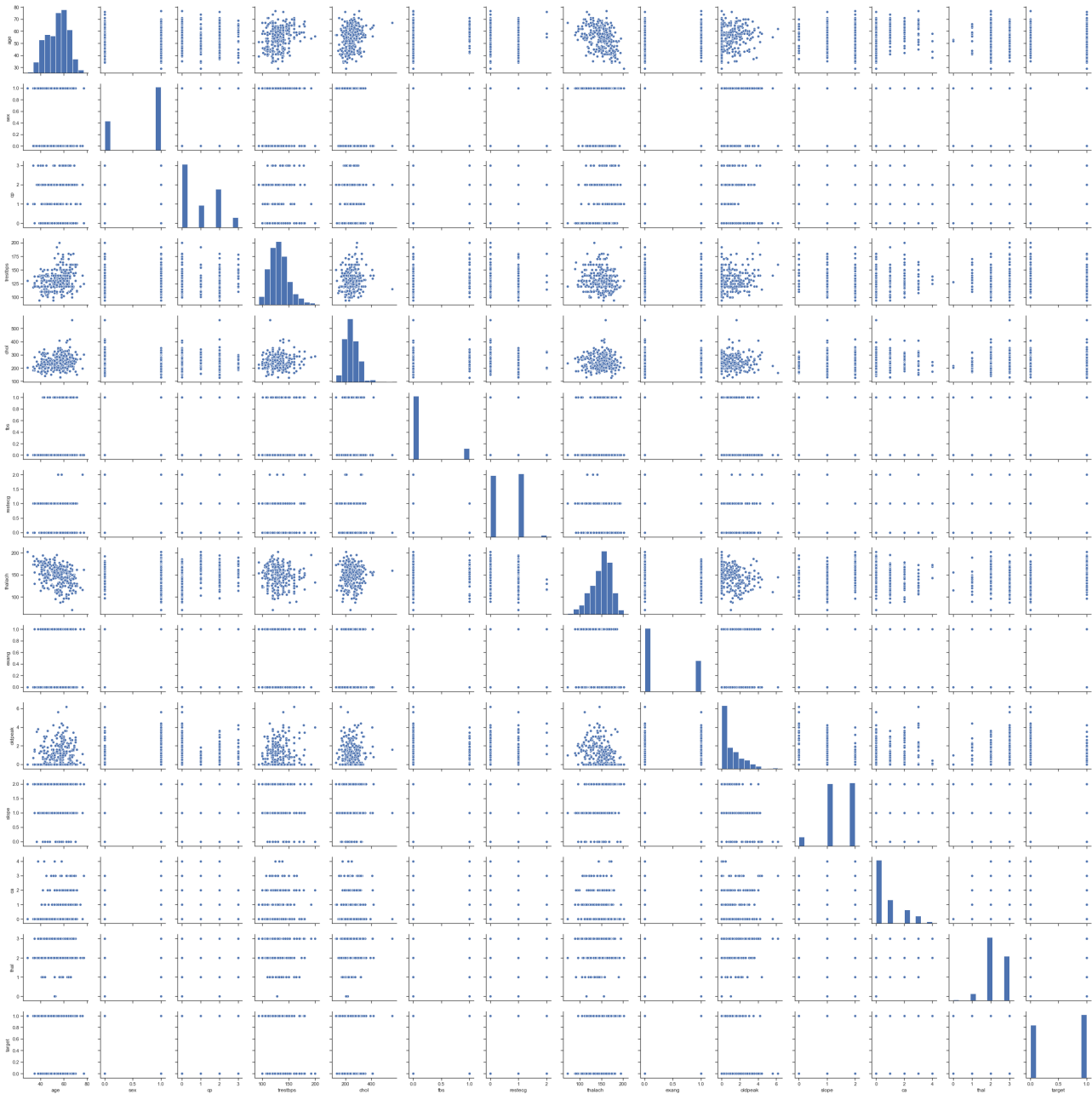
In [59]:

```
#парные диаграммы
```

```
sns.pairplot(data)
```

Out[59]:

```
<seaborn.axisgrid.PairGrid at 0x1a1d52e208>
```

In [60]:

```
sns.pairplot(data, hue="target")
```

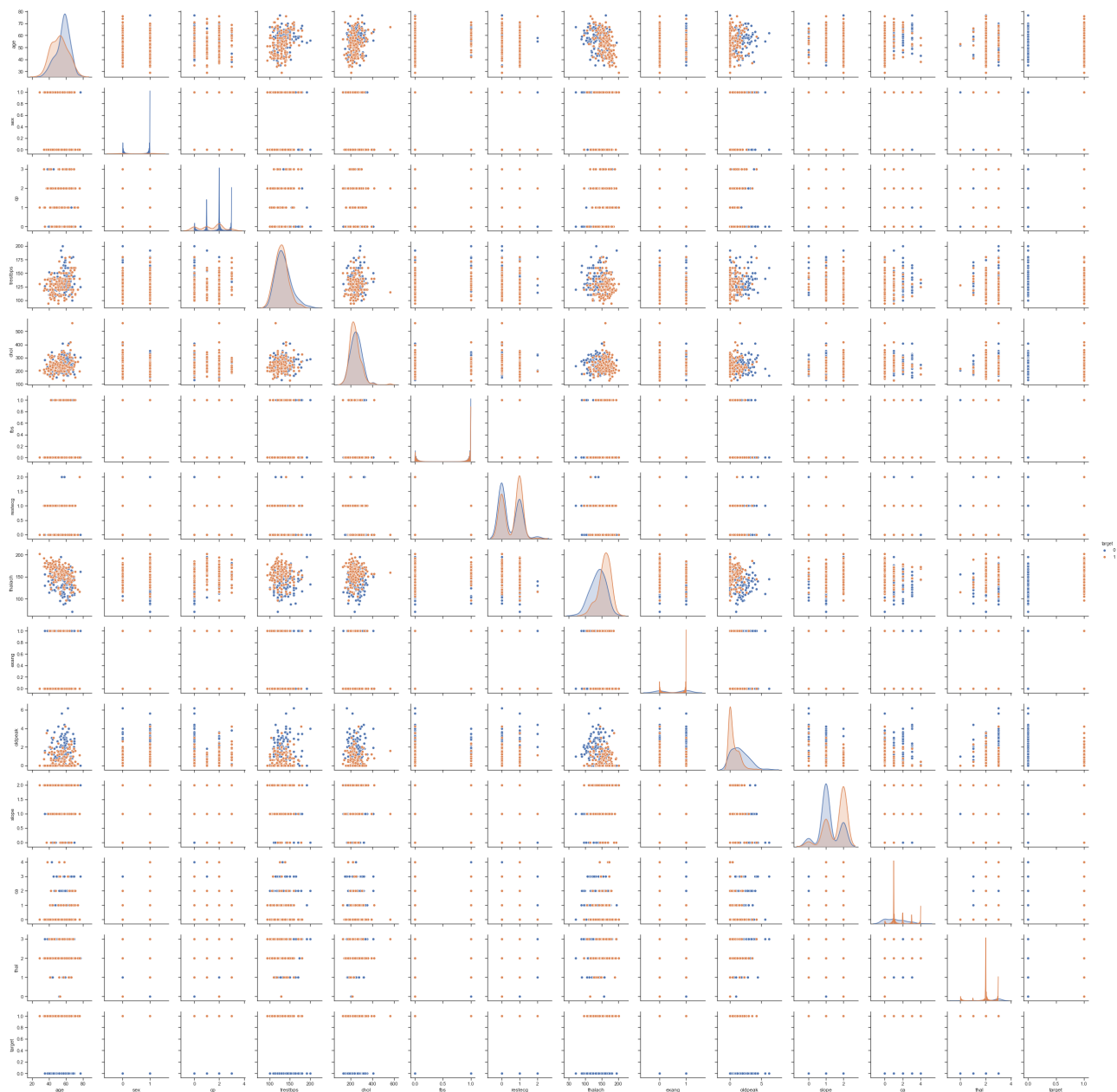
```
/anaconda3/lib/python3.6/site-packages/statsmodels/
nonparametric/kde.py:494: RuntimeWarning: invalid value
encountered in true_divide
```

```
    binned = fast_linbin(X,a,b,gridsize)/(delta*nobs)
/anaconda3/lib/python3.6/site-packages/statsmodels/
nonparametric/kdetools.py:34: RuntimeWarning: invalid value
encountered in double_scalars
```

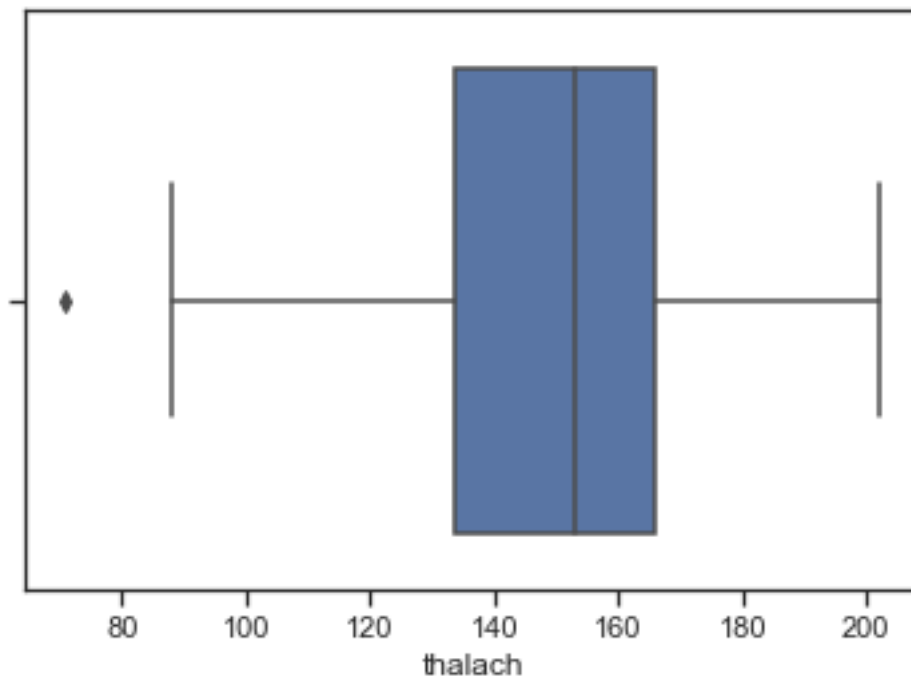
```
    FAC1 = 2*(np.pi*bw/RANGE)**2
/anaconda3/lib/python3.6/site-packages/numpy/core/
_methods.py:26: RuntimeWarning: invalid value encountered in
reduce
```

```
    return umr_maximum(a, axis, None, out, keepdims)
```

```
Out[60]:  
<seaborn.axisgrid.PairGrid at 0x1a23934710>
```



```
In [61]:  
# ящик с усами  
sns.boxplot(x=data['thalach'])  
Out[61]:  
<matplotlib.axes._subplots.AxesSubplot at 0x1a29136f28>
```



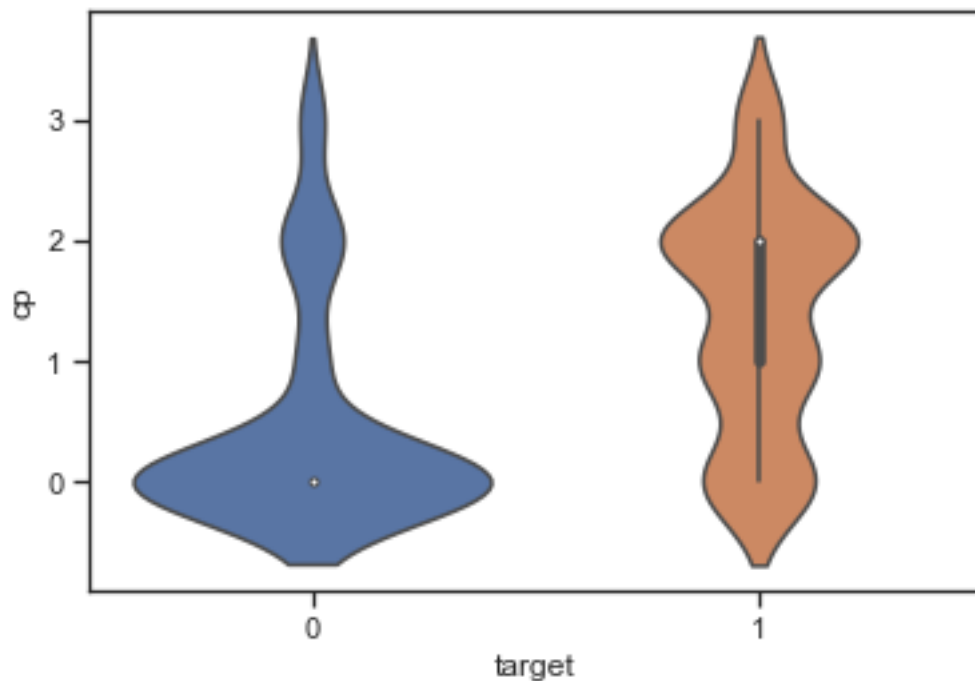
In [63]:

```
# Распределение параметра cp сгруппированные по target.
```

```
sns.violinplot(x='target', y='cp', data=data)
```

Out[63]:

<matplotlib.axes._subplots.AxesSubplot at 0x1a2c568860>



In [64]:

```
data.corr()
```

Out[64]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
--	-----	-----	----	----------	------	-----	---------	---------	-------	---------	-------	----	------	--------

age	1.000 000	-0.098 447	-0.068 653	0.279 351	0.213 678	0.121 308	-0.116 211	-0.398 522	0.096 801	0.210 013	-0.168 814	0.276 326	0.068 001	-0.225 439
sex	-0.098 447	1.000 000	-0.049 353	-0.056 769	-0.197 912	0.045 032	-0.058 196	-0.044 020	0.141 664	0.096 093	-0.030 711	0.118 261	0.210 041	-0.280 937
cp	-0.068 653	-0.049 353	1.000 000	0.047 608	-0.076 904	0.094 444	0.044 421	0.295 762	-0.394 280	-0.149 230	0.119 717	-0.181 053	-0.161 736	0.433 798
trest bps	0.279 351	-0.056 769	0.047 608	1.000 000	0.123 174	0.177 531	-0.114 103	-0.046 698	0.067 616	0.193 216	-0.121 475	0.101 389	0.062 210	-0.144 931
chol	0.213 678	-0.197 912	-0.076 904	0.123 174	1.000 000	0.013 294	-0.151 040	-0.009 940	0.067 023	0.053 952	-0.004 038	0.070 511	0.098 803	-0.085 239
fbs	0.121 308	0.045 032	0.094 444	0.177 531	0.013 294	1.000 000	-0.084 189	-0.008 567	0.025 665	0.005 747	-0.059 894	0.137 979	-0.032 019	-0.028 046
reste cg	-0.116 211	-0.058 196	0.044 421	-0.114 103	-0.151 040	-0.084 189	1.000 000	0.044 123	-0.070 733	-0.058 770	0.093 045	-0.072 042	-0.011 981	0.137 230
thala ch	-0.398 522	-0.044 020	0.295 762	-0.046 698	-0.009 940	-0.008 567	0.044 123	1.000 000	-0.378 812	-0.344 187	0.386 784	-0.213 177	-0.096 439	0.421 741
exan g	0.096 801	0.141 664	-0.394 280	0.067 616	0.067 023	0.025 665	-0.070 733	-0.378 812	1.000 000	0.288 223	-0.257 748	0.115 739	0.206 754	-0.436 757
oldp eak	0.210 013	0.096 093	-0.149 230	0.193 216	0.053 952	0.005 747	-0.058 770	-0.344 187	0.288 223	1.000 000	-0.577 537	0.222 682	0.210 244	-0.430 696
slope	-0.168 814	-0.030 711	0.119 717	-0.121 475	-0.004 038	-0.059 894	0.093 045	0.386 784	-0.257 748	-0.577 537	1.000 000	-0.080 155	-0.104 764	0.345 877
ca	0.276 326	0.118 261	-0.181 053	0.101 389	0.070 511	0.137 979	-0.072 042	-0.213 177	0.115 739	0.222 682	-0.080 155	1.000 000	0.151 832	-0.391 724
thal	0.068 001	0.210 041	-0.161 736	0.062 210	0.098 803	-0.032 019	-0.011 981	-0.096 439	0.206 754	0.210 244	-0.104 764	0.151 832	1.000 000	-0.344 029
targe t	-0.225 439	-0.280 937	0.433 798	-0.144 931	-0.085 239	-0.028 046	0.137 230	0.421 741	-0.436 757	-0.430 696	0.345 877	-0.391 724	-0.344 029	1.000 000

In [73]:

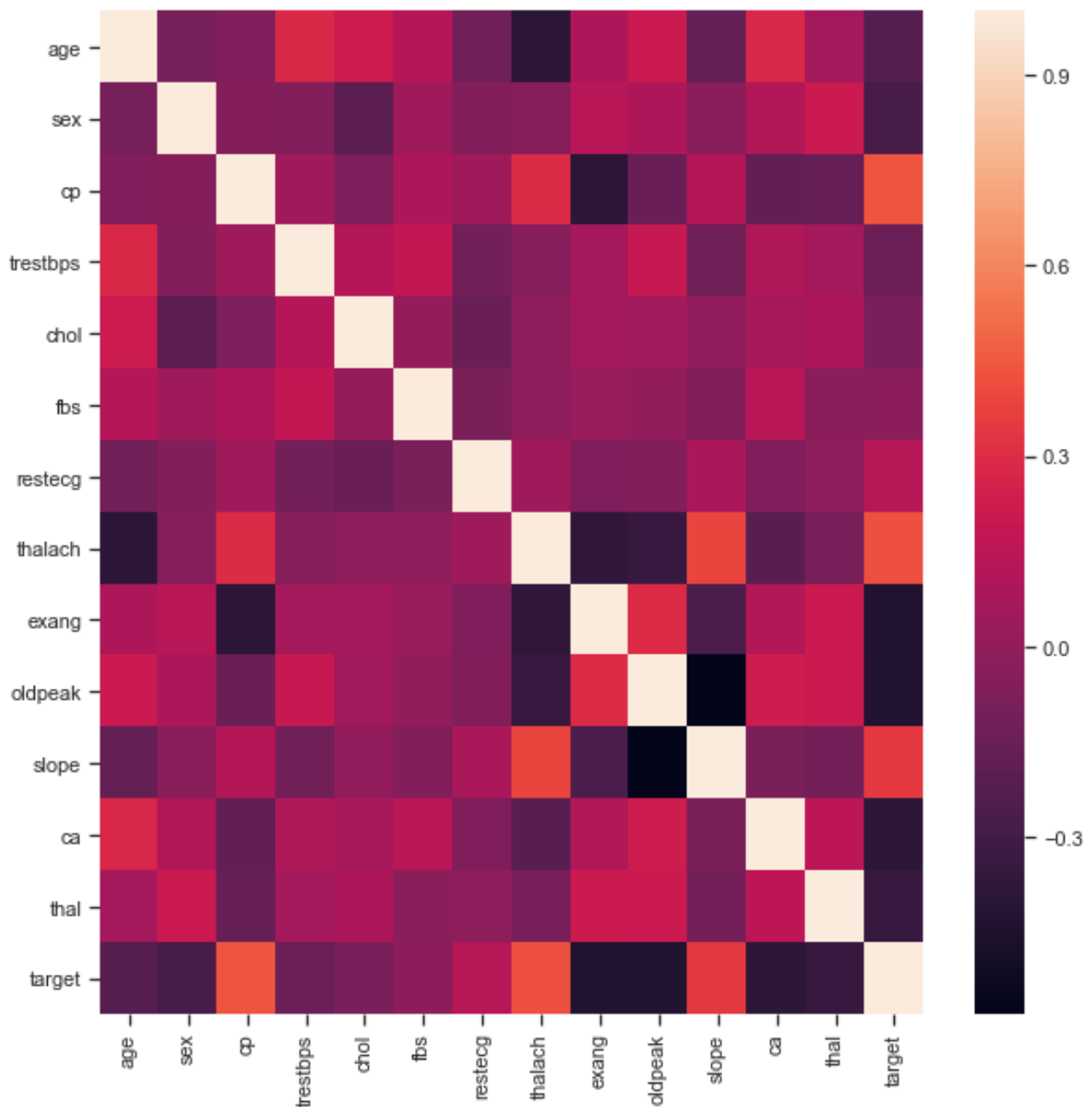
```
# тепловая карта со значениями в ячейках
```

```
plt.figure(figsize=(10, 10))
```

```
sns.heatmap(data.corr())
```

Out[73]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a31492278>
```



In [78]:

```
fig, ax = plt.subplots(1, 3, sharex='col', sharey='row',
figsize=(30,10))
sns.heatmap(data.corr(method='pearson'), ax=ax[0],
annot=True, fmt='.1f')
sns.heatmap(data.corr(method='kendall'), ax=ax[1],
annot=True, fmt='.1f')
sns.heatmap(data.corr(method='spearman'), ax=ax[2],
annot=True, fmt='.1f')
fig.suptitle('Корреляционные матрицы, построенные различными
методами')
ax[0].title.set_text('Pearson')
ax[1].title.set_text('Kendall')
ax[2].title.set_text('Spearman')
```

Корреляционные матрицы, построенные различными методами

