



**Министерство науки и высшего образования Российской Федерации  
Федеральное государственное бюджетное образовательное учреждение  
высшего образования  
«Московский государственный технический университет  
имени Н.Э. Баумана  
(национальный исследовательский университет)»  
(МГТУ им. Н.Э. Баумана)**

**Факультет «Информатика и системы управления»  
Кафедра ИУ5 «Системы обработки информации и управления»**

Лабораторная работа №1  
по дисциплине «Технология машинного обучения» на тему:

Разведочный анализ данных. Исследование и визуализация данных.

---

Выполнил:  
студент группы № ИУ5-62  
Морозенков О.Н.  
подпись, дата

Проверил:  
Ю.Е. Гапанюк  
подпись, дата

2020 г.

## Задание:

- Выбрать набор данных (датасет).

Для лабораторных работ не рекомендуется выбирать датасеты большого размера.

- 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.

In [18]:

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

In [19]:

```
data = pd.read_csv('heart.csv', sep=",")
data.head()
```

Out[19]:

	age	sex	cp	trestbps	chol	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 [20]:

```
data.tail()
```

Out[20]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
298	57	0	0	140	241	0	1	123	1	0.2	1	0	3	0
299	45	1	3	110	264	0	1	132	0	1.2	1	0	3	0
300	68	1	0	144	193	1	1	141	0	3.4	1	2	3	0
301	57	1	0	130	131	0	1	115	1	1.2	1	1	3	0
302	57	0	1	130	236	0	0	174	0	0.0	1	1	2	0

In [21]:

```
data.shape
stroki = data.shape[0]
stolbec = data.shape[1]
print('Всего строк: {}'.format(stroki), ', всего столбцов: {}'.format(stolbec))
```

Всего строк: 303 , всего столбцов: 14

In [22]:

```
data.columns
```

Out[22]:

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

In [23]:

```
data.dtypes
```

Out[23]:

age int64
sex 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 [24]:

```
for col in data.columns:
    # Количество пустых значений - все значения заполнены
    temp_null_count = data[data[col].isnull()].shape[0]
    print('{} - {}'.format(col, temp_null_count))
```

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 [25]:

```
data.describe()
```

Out[25]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca
count	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
mean	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.528053	149.646865	0.326733	1.039604	1.399340	0.729379
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860	22.905161	0.469794	1.161075	0.616226	1.022608
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000	71.000000	0.000000	0.000000	0.000000	0.000000
25%	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000	133.500000	0.000000	0.000000	1.000000	0.000000
50%	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000	153.000000	0.000000	0.800000	1.000000	0.000000
75%	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000	166.000000	1.000000	1.600000	2.000000	1.000000
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000	202.000000	1.000000	6.200000	2.000000	4.000000

In [26]:

```
data['target'].unique()
```

Out[26]:

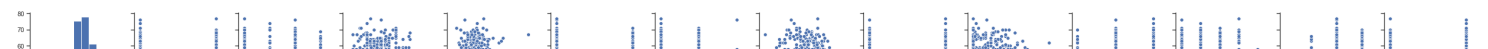
array([1, 0], dtype=int64)

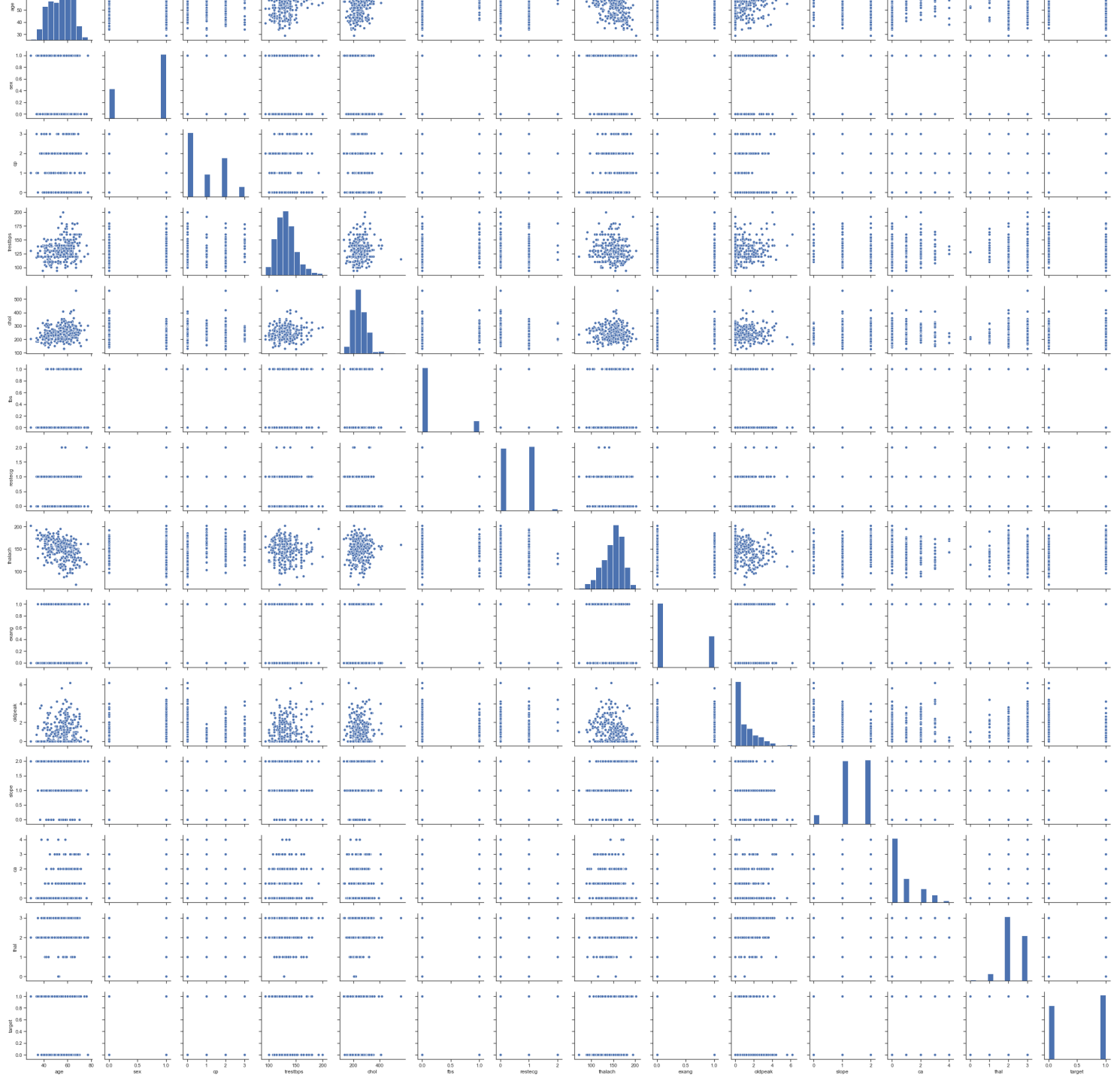
In [27]:

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

Out[27]:

<seaborn.axisgrid.PairGrid at 0x250ea414278>





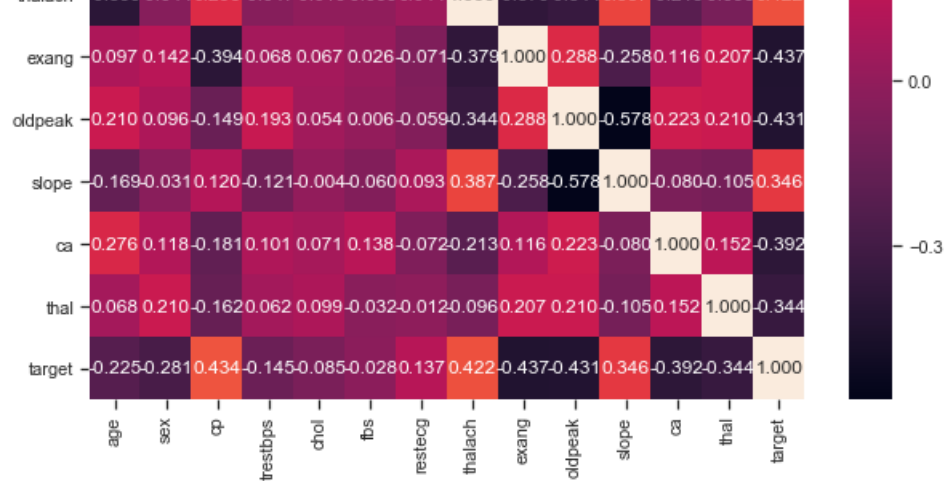
In [34]:

```
fig, ax = plt.subplots(figsize=(10,10))
sns.heatmap(data.corr(), ax=ax, annot=True, fmt='.3f')
```

Out[34]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x250ffabddd8>



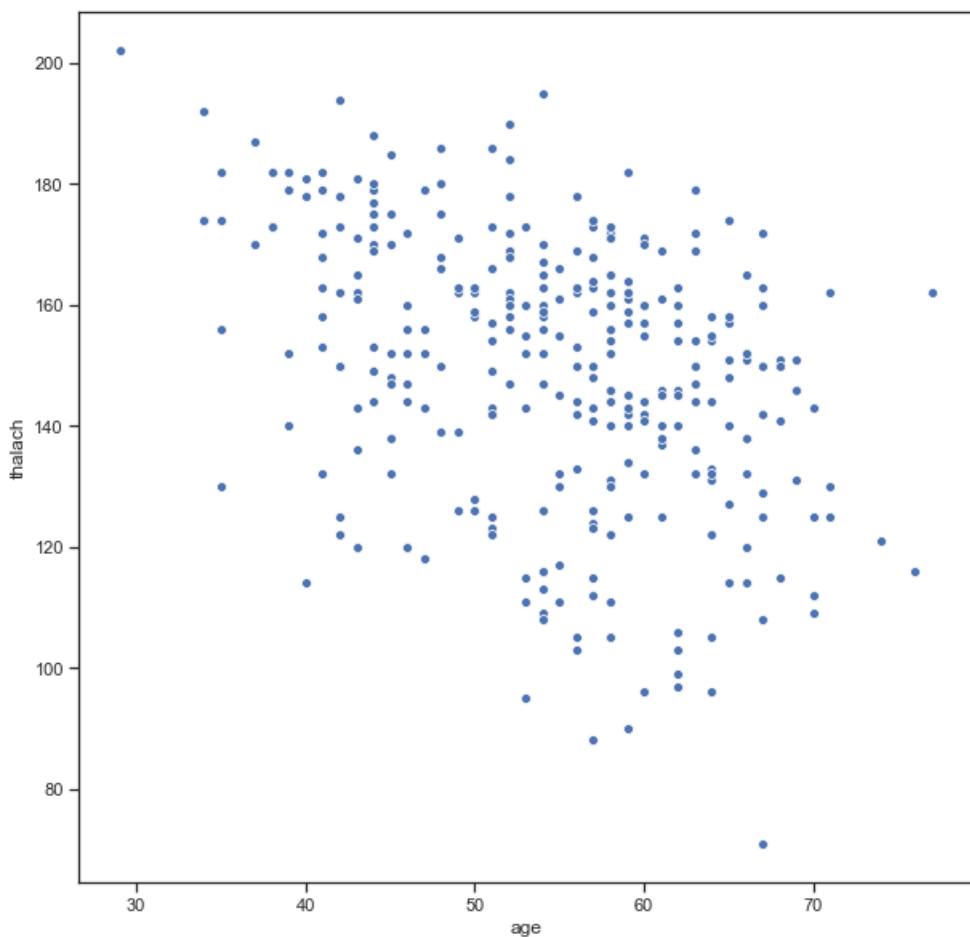


In [32]:

```
fig, ax = plt.subplots(figsize=(10,10))
sns.scatterplot(ax=ax, x='age', y='thalach', data=data)
```

Out[32]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x250ff578da0>



In [ ]: