



**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)**

Курс «Технологии машинного обучения»

Отчет по лабораторной работе №1

Выполнил:  
студент группы ИУ5-64Б

Ли М.В.

Преподаватель:

Гапанюк Ю.Е

## Задание:

- Выбрать набор данных (датасет).
- Пример преобразования датасетов Scikit-learn в Pandas Dataframe.
- Для лабораторных работ не рекомендуется выбирать датасеты большого размера.
- Создать ноутбук, который содержит следующие разделы:
- Текстовое описание выбранного Вами набора данных.
- Основные характеристики датасета.
- Визуальное исследование датасета.
- Информация о корреляции признаков.
- Сформировать отчет и разместить его в своей репозитории на github.

## Выполнение:

---

### Лабораторная работа №1

---

#### Разведочный анализ данных. Исследование и визуализация данных.

##### Описание набора данных:

В качестве набора данных мы будем использовать следующий - <https://www.kaggle.com/uciml/red-wine-quality-cortez-et-al-2009> Этот набор данных содержит список красных вин с описанием их качества. И содержит в себе такие данные:

1 - fixed acidity

2 - volatile acidity

3 - citric acid

4 - residual sugar

5 - chlorides

6 - free sulfur dioxide

7 - total sulfur dioxide

8 - density

9 - pH

10 - sulphates

11 - alcohol

Output variable (based on sensory data):

12 - quality (score between 0 and 10)

## Импорт библиотек

```
B [3]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(style="whitegrid")
```

## Загрузка данных

```
B [5]: # Будем анализировать данные только на обучающей выборке
data = pd.read_csv('data/winequality-red.csv', sep=",")
```

## Основные характеристики датасета

```
B [6]: # Первые 5 строк датасета
data.head()
```

Out[6]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8	5
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8	5
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8	6
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5

```
B [4]: # Размер датасета
data.shape
```

Out[4]: (1599, 12)

```
B [5]: total_count = data.shape[0]
print('Всего строк: {}'.format(total_count))
```

Всего строк: 1599

```
B [6]: # Список колонок
data.columns
```

```
Out[6]: Index(['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
              'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
              'pH', 'sulphates', 'alcohol', 'quality'],
              dtype='object')
```

```
B [7]: # Список колонок с типами данных
data.dtypes
```

```
Out[7]: fixed acidity      float64
volatile acidity    float64
citric acid         float64
residual sugar      float64
chlorides           float64
free sulfur dioxide  float64
total sulfur dioxide float64
density            float64
pH                 float64
sulphates          float64
alcohol            float64
quality            int64
dtype: object
```

```
B [8]: # Проверим наличие пустых значений
# Цикл по колонкам датасета
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))
```

```
fixed acidity - 0
volatile acidity - 0
citric acid - 0
residual sugar - 0
chlorides - 0
free sulfur dioxide - 0
total sulfur dioxide - 0
density - 0
pH - 0
sulphates - 0
alcohol - 0
quality - 0
```

```
B [9]: # Основные статистические характеристики набора данных
data.describe()
```

```
Out[9]:
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol
count	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000
mean	8.319637	0.527821	0.270976	2.538806	0.087467	15.874922	46.467792	0.996747	3.311113	0.658149	10.422983
std	1.741096	0.179060	0.194801	1.409928	0.047065	10.460157	32.895324	0.001887	0.154386	0.169507	1.065668
min	4.600000	0.120000	0.000000	0.900000	0.012000	1.000000	6.000000	0.990070	2.740000	0.330000	8.400000
25%	7.100000	0.390000	0.090000	1.900000	0.070000	7.000000	22.000000	0.995600	3.210000	0.550000	9.500000
50%	7.900000	0.520000	0.260000	2.200000	0.079000	14.000000	38.000000	0.996750	3.310000	0.620000	10.200000
75%	9.200000	0.640000	0.420000	2.600000	0.090000	21.000000	62.000000	0.997835	3.400000	0.730000	11.100000
max	15.900000	1.580000	1.000000	15.500000	0.611000	72.000000	289.000000	1.003690	4.010000	2.000000	14.900000

```
B [10]: # Определим уникальные значения для целевого признака
data['quality'].unique()
```

```
Out[10]: array([5, 6, 7, 4, 8, 3], dtype=int64)
```

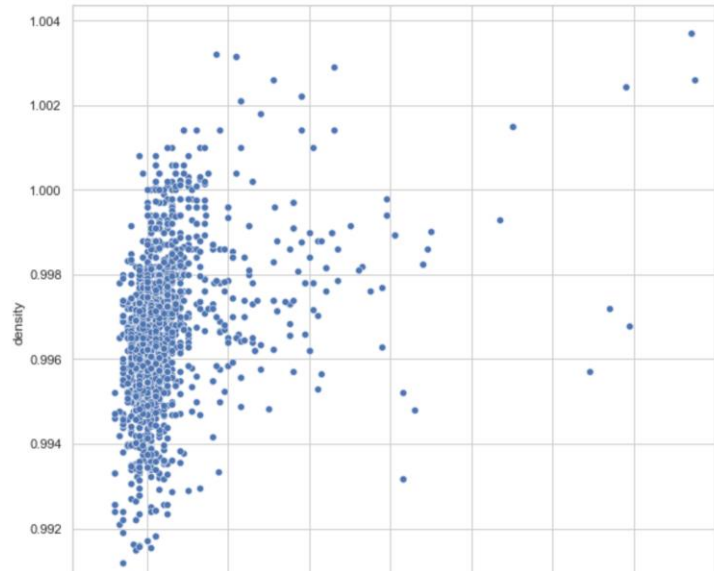
## Визуальное исследование датасета

Для визуального исследования могут быть использованы различные виды диаграмм, мы построим только некоторые варианты диаграмм, которые используются достаточно часто.

### Диаграмма рассеяния

```
In [11]: fig, ax = plt.subplots(figsize=(10,10))
sns.scatterplot(ax=ax, x='residual sugar', y='density', data=data)
```

```
Out[11]: <AxesSubplot:xlabel='residual sugar', ylabel='density'>
```

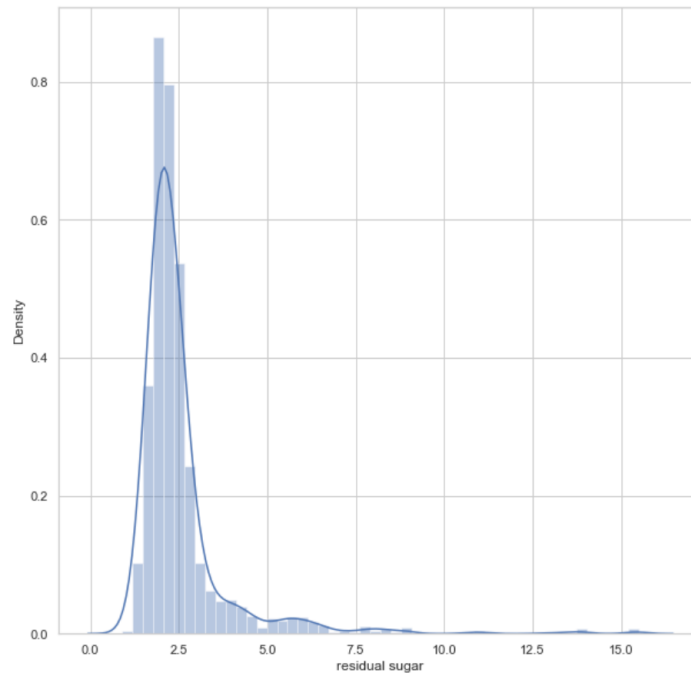


## Гистограмма

```
B [13]: In fig, ax = plt.subplots(figsize=(10,10))
sns.distplot(data['residual sugar'])
```

C:\Users\enjoy\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

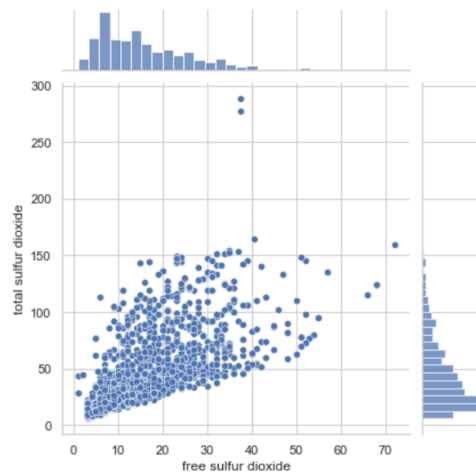
```
Out[13]: <AxesSubplot:xlabel='residual sugar', ylabel='Density'>
```



## Jointplot

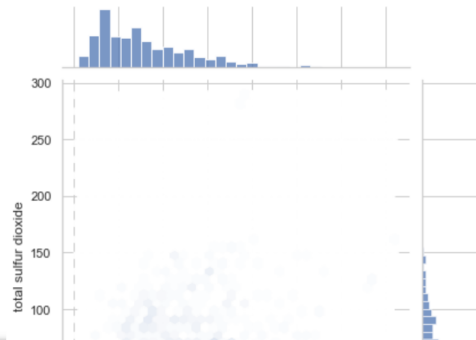
```
B [14]: In [14]: sns.jointplot(x='free sulfur dioxide', y='total sulfur dioxide', data=data)
```

```
Out[14]: <seaborn.axisgrid.JointGrid at 0x1aaa81fb490>
```



```
B [15]: In [15]: sns.jointplot(x='free sulfur dioxide', y='total sulfur dioxide', data=data, kind="hex")
```

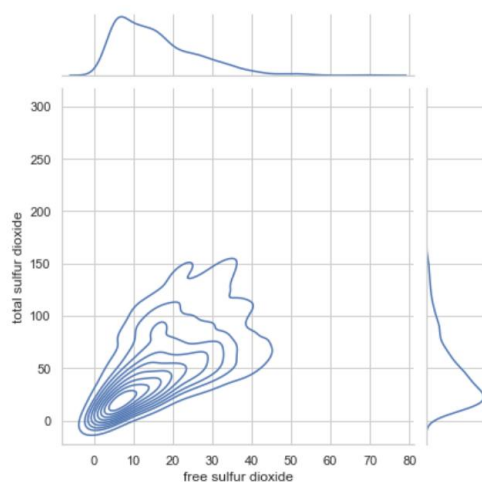
```
Out[15]: <seaborn.axisgrid.JointGrid at 0x1aaa7e15be0>
```



```
B [16]: In [16]: sns.jointplot(x='free sulfur dioxide', y='total sulfur dioxide', data=data, kind="kde")
```

```
B [16]: sns.jointplot(x='free sulfur dioxide', y='total sulfur dioxide', data=data, kind="kde")
```

```
Out[16]: <seaborn.axisgrid.JointGrid at 0x1aaa87a07c0>
```



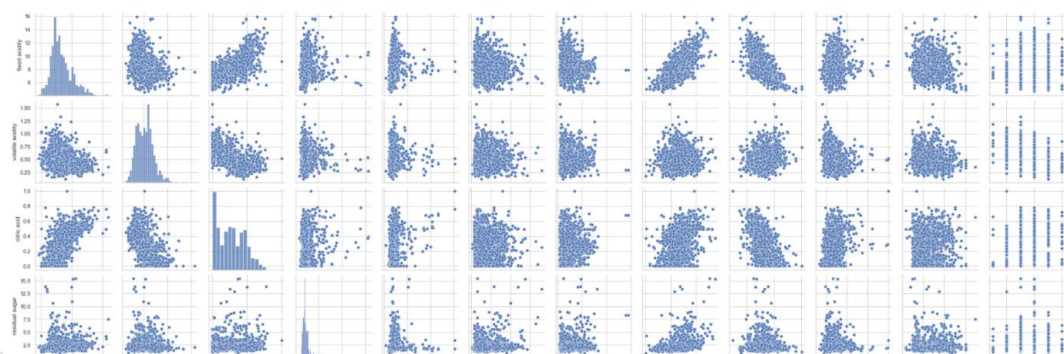
### "Парные диаграммы"

Комбинация гистограмм и диаграмм рассеивания для всего набора данных.

Выводится матрица графиков. На пересечении строки и столбца, которые соответствуют двум показателям, строится диаграмма рассеивания. В главной диагонали матрицы строятся гистограммы распределения соответствующих показателей.

```
B [17]: sns.pairplot(data)
```

```
Out[17]: <seaborn.axisgrid.PairGrid at 0x1aaa826fcd0>
```



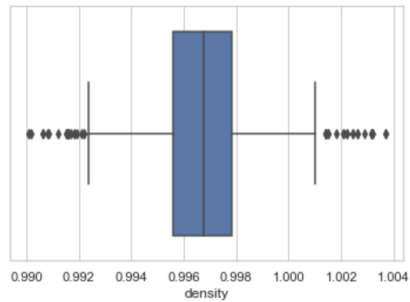


## Ящик с усами

Отображает одномерное распределение вероятности.

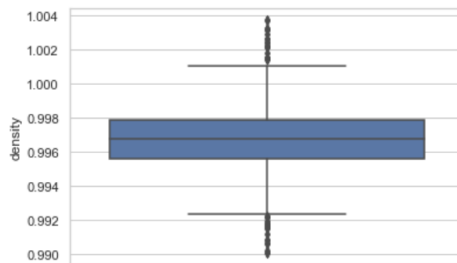
```
B [19]: sns.boxplot(x=data['density'])
```

```
Out[19]: <AxesSubplot:xlabel='density'>
```



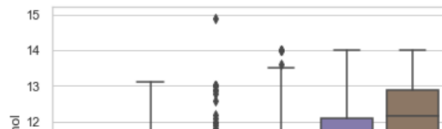
```
B [20]: sns.boxplot(y=data['density'])
```

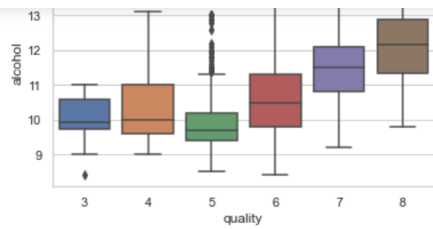
```
Out[20]: <AxesSubplot:ylabel='density'>
```



```
B [21]: sns.boxplot(x='quality', y='alcohol', data=data)
```

```
Out[21]: <AxesSubplot:xlabel='quality', ylabel='alcohol'>
```

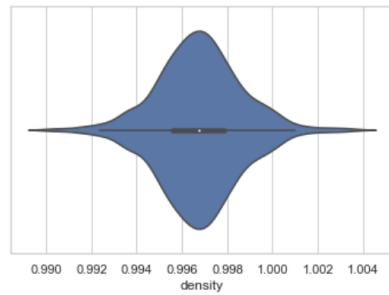




## Violin plot

B [22]: `sns.violinplot(x=data['density'])`

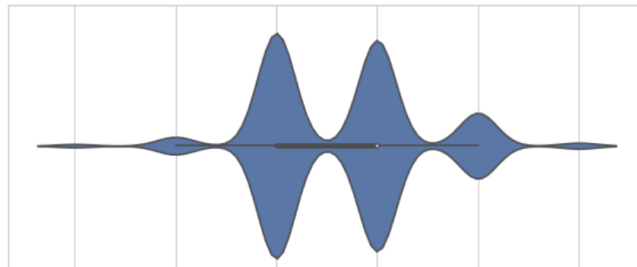
Out[22]: `<AxesSubplot:xlabel='density'>`



B [23]: `fig, ax = plt.subplots(2, 1, figsize=(10,10))`  
`sns.violinplot(ax=ax[0], x=data['quality'])`  
`sns.distplot(data['quality'], ax=ax[1])`

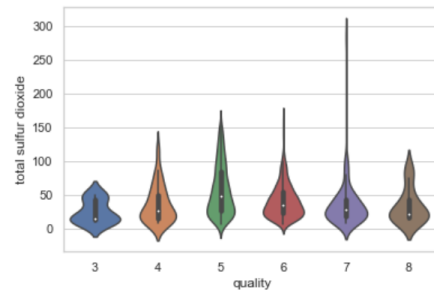
C:\Users\enjoy\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
 warnings.warn(msg, FutureWarning)

Out[23]: `<AxesSubplot:xlabel='quality', ylabel='Density'>`



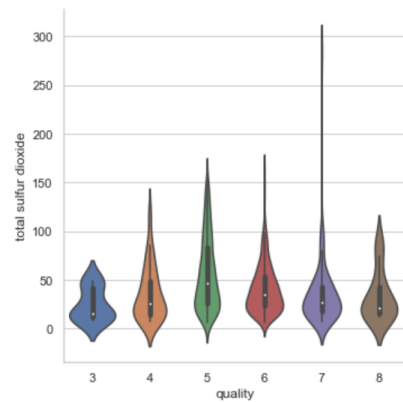
```
B [24]: sns.violinplot(x='quality', y='total sulfur dioxide', data=data)
```

```
Out[24]: <AxesSubplot: xlabel='quality', ylabel='total sulfur dioxide'>
```



```
B [25]: sns.catplot(y='total sulfur dioxide', x='quality', data=data, kind="violin", split=True)
```

```
Out[25]: <seaborn.axisgrid.FacetGrid at 0x1aab68ae310>
```



Информация о корреляции признаков

B [26]:

data.corr()

Out[26]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
fixed acidity	1.000000	-0.256131	0.671703	0.114777	0.093705	-0.153794	-0.113181	0.668047	-0.682978	0.183006	-0.061668	0.124052
volatile acidity	-0.256131	1.000000	-0.552496	0.001918	0.061298	-0.010504	0.076470	0.022026	0.234937	-0.260987	-0.202288	-0.390558
citric acid	0.671703	-0.552496	1.000000	0.143577	0.203823	-0.060978	0.035533	0.364947	-0.541904	0.312770	0.109903	0.226373
residual sugar	0.114777	0.001918	0.143577	1.000000	0.055610	0.187049	0.203028	0.355283	-0.085652	0.005527	0.042075	0.013732
chlorides	0.093705	0.061298	0.203823	0.055610	1.000000	0.005562	0.047400	0.200632	-0.265026	0.371260	-0.221141	-0.128907
free sulfur dioxide	-0.153794	-0.010504	-0.060978	0.187049	0.005562	1.000000	0.667666	-0.021946	0.070377	0.051658	-0.069408	-0.050656
total sulfur dioxide	-0.113181	0.076470	0.035533	0.203028	0.047400	0.667666	1.000000	0.071269	-0.066495	0.042947	-0.205654	-0.185100
density	0.668047	0.022026	0.364947	0.355283	0.200632	-0.021946	0.071269	1.000000	-0.341699	0.148506	-0.496180	-0.174919

B [27]:

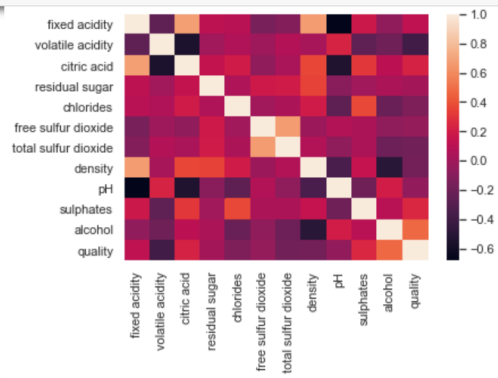
data.corr(method='pearson')

Out[27]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
fixed acidity	1.000000	-0.256131	0.671703	0.114777	0.093705	-0.153794	-0.113181	0.668047	-0.682978	0.183006	-0.061668	0.124052
volatile acidity	-0.256131	1.000000	-0.552496	0.001918	0.061298	-0.010504	0.076470	0.022026	0.234937	-0.260987	-0.202288	-0.390558
citric acid	0.671703	-0.552496	1.000000	0.143577	0.203823	-0.060978	0.035533	0.364947	-0.541904	0.312770	0.109903	0.226373
residual sugar	0.114777	0.001918	0.143577	1.000000	0.055610	0.187049	0.203028	0.355283	-0.085652	0.005527	0.042075	0.013732
chlorides	0.093705	0.061298	0.203823	0.055610	1.000000	0.005562	0.047400	0.200632	-0.265026	0.371260	-0.221141	-0.128907
free sulfur dioxide	-0.153794	-0.010504	-0.060978	0.187049	0.005562	1.000000	0.667666	-0.021946	0.070377	0.051658	-0.069408	-0.050656
total sulfur dioxide	-0.113181	0.076470	0.035533	0.203028	0.047400	0.667666	1.000000	0.071269	-0.066495	0.042947	-0.205654	-0.185100
density	0.668047	0.022026	0.364947	0.355283	0.200632	-0.021946	0.071269	1.000000	-0.341699	0.148506	-0.496180	-0.174919
pH	-0.682978	0.234937	-0.541904	-0.085652	-0.265026	0.070377	-0.066495	-0.341699	1.000000	-0.196648	0.205633	-0.057731
sulphates	0.183006	-0.260987	0.312770	0.005527	0.371260	0.051658	0.042947	0.148506	-0.196648	1.000000	0.093595	0.251397
alcohol	-0.061668	-0.202288	0.109903	0.042075	-0.221141	-0.069408	-0.205654	-0.496180	0.205633	0.093595	1.000000	0.476166
quality	0.124052	-0.390558	0.226373	0.013732	-0.128907	-0.050656	-0.185100	-0.174919	-0.057731	0.251397	0.476166	1.000000

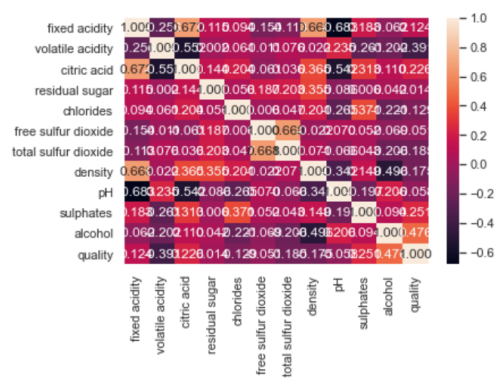
## Тепловая карта

```
B [29]: sns.heatmap(data.corr())
```



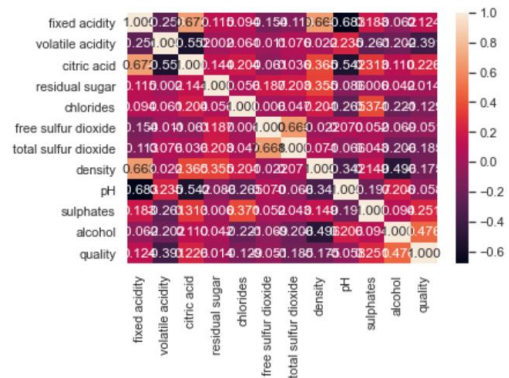
```
B [30]: # Вывод значений в ячейках
sns.heatmap(data.corr(), annot=True, fmt='.3f')
```

```
Out[30]: <AxesSubplot:>
```



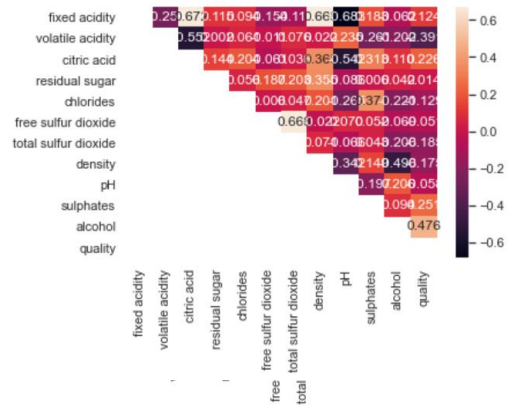
```
B [30]: # Вывод значений в ячейках
sns.heatmap(data.corr(), annot=True, fmt='.3f')
```

Out[30]: <AxesSubplot:>



```
B [31]: # Треугольный вариант матрицы
mask = np.zeros_like(data.corr(), dtype=np.bool)
# чтобы оставить нижнюю часть матрицы
# mask[np.triu_indices_from(mask)] = True
# чтобы оставить верхнюю часть матрицы
mask[np.tril_indices_from(mask)] = True
sns.heatmap(data.corr(), mask=mask, annot=True, fmt='.3f')
```

Out[31]: <AxesSubplot:>



```
B [32]: fig, ax = plt.subplots(1, 3, sharex='col', sharey='row', figsize=(15,5))
sns.heatmap(data.corr(method='pearson'), ax=ax[0], annot=True, fmt='.2f')
sns.heatmap(data.corr(method='kendall'), ax=ax[1], annot=True, fmt='.2f')
sns.heatmap(data.corr(method='spearman'), ax=ax[2], annot=True, fmt='.2f')
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

