

Московский государственный технический университет имени Н.Э.Баумана

Кафедра «Системы обработки информации и управления»

Рубежный контроль №1
по дисциплине
«Методы машинного обучения»
на тему
«Методы обработки данных»

Выполнил:
Студент ИУ5-24М
Гаврилюк А.Г.

Москва, 2020

Задача:

Для заданного набора данных постройте основные графики, входящие в этап разведочного анализа данных с использованием библиотек Matplotlib и Seaborn. В случае наличия пропусков в данных удалите строки или колонки, содержащие пропуски. Какие графики Вы построили и почему? Какие выводы о наборе данных Вы можете сделать на основании построенных графиков? Проведите корреляционный анализ. Сделайте выводы о возможности построения моделей машинного обучения и о возможном вкладе признаков в модель.

Набор данных:

https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_wine.html#sklearn.datasets.load_wine

In [17]:

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

In [18]:

```
from sklearn.datasets import load_wine
X, y = load_wine(return_X_y=True)
print(X.shape)
```

(178, 13)

In []:

```
def make_dataframe(ds_function):
    ds = ds_function()
    df = pd.DataFrame(data= np.c_[ds['data'], ds['target']],
                      columns= list(ds['feature_names']) + ['target'])
    return df
```

In [19]:

```
data = make_dataframe(load_wine)
data.head()
```

Out[19]:

	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	total_phenols	flavanoids	nonflavanoid_phenols	proanthocyanins	color_inte
0	14.23	1.71	2.43	15.6	127.0	2.80	3.06	0.28	2.29	
1	13.20	1.78	2.14	11.2	100.0	2.65	2.76	0.26	1.28	
2	13.16	2.36	2.67	18.6	101.0	2.80	3.24	0.30	2.81	
3	14.37	1.95	2.50	16.8	113.0	3.85	3.49	0.24	2.18	
4	13.24	2.59	2.87	21.0	118.0	2.80	2.69	0.39	1.82	

Количество пустых значений в колонках

In [6]:

```
for col in data.columns:
    temp_null_count = data[data[col].isnull()].shape[0]
    print('{} - {}'.format(col, temp_null_count))
```

```
alcohol - 0
malic_acid - 0
ash - 0
alcalinity_of_ash - 0
magnesium - 0
total_phenols - 0
flavanoids - 0
nonflavanoid_phenols - 0
proanthocyanins - 0
color_inte - 0
```

```
nonflavanoid_phenols - 0
proanthocyanins - 0
color_intensity - 0
hue - 0
od280/od315_of_diluted_wines - 0
proline - 0
target - 0
```

In [7]:

```
data.describe()
```

Out[7]:

	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	total_phenols	flavanoids	nonflavanoid_phenols	proanthocya
count	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000
mean	13.000618	2.336348	2.366517	19.494944	99.741573	2.295112	2.029270	0.361854	1.590000
std	0.811827	1.117146	0.274344	3.339564	14.282484	0.625851	0.998859	0.124453	0.570000
min	11.030000	0.740000	1.360000	10.600000	70.000000	0.980000	0.340000	0.130000	0.410000
25%	12.362500	1.602500	2.210000	17.200000	88.000000	1.742500	1.205000	0.270000	1.250000
50%	13.050000	1.865000	2.360000	19.500000	98.000000	2.355000	2.135000	0.340000	1.550000
75%	13.677500	3.082500	2.557500	21.500000	107.000000	2.800000	2.875000	0.437500	1.950000
max	14.830000	5.800000	3.230000	30.000000	162.000000	3.880000	5.080000	0.660000	3.580000

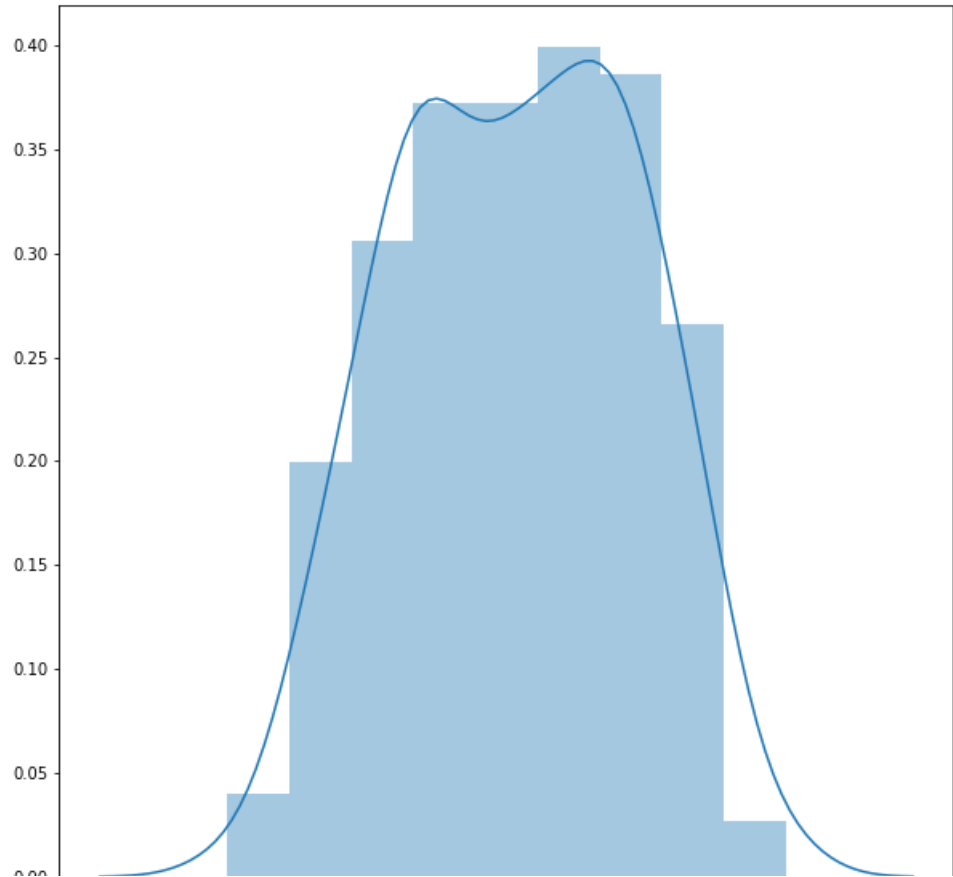
Распределение значений целевого признака

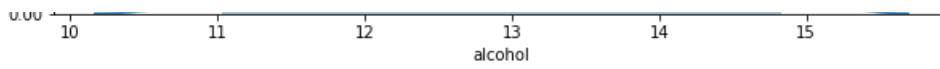
In [12]:

```
fig, ax = plt.subplots(figsize=(10,10))
sns.distplot(data['alcohol'])
```

Out[12]:

<matplotlib.axes._subplots.AxesSubplot at 0x12fcd0cd0>





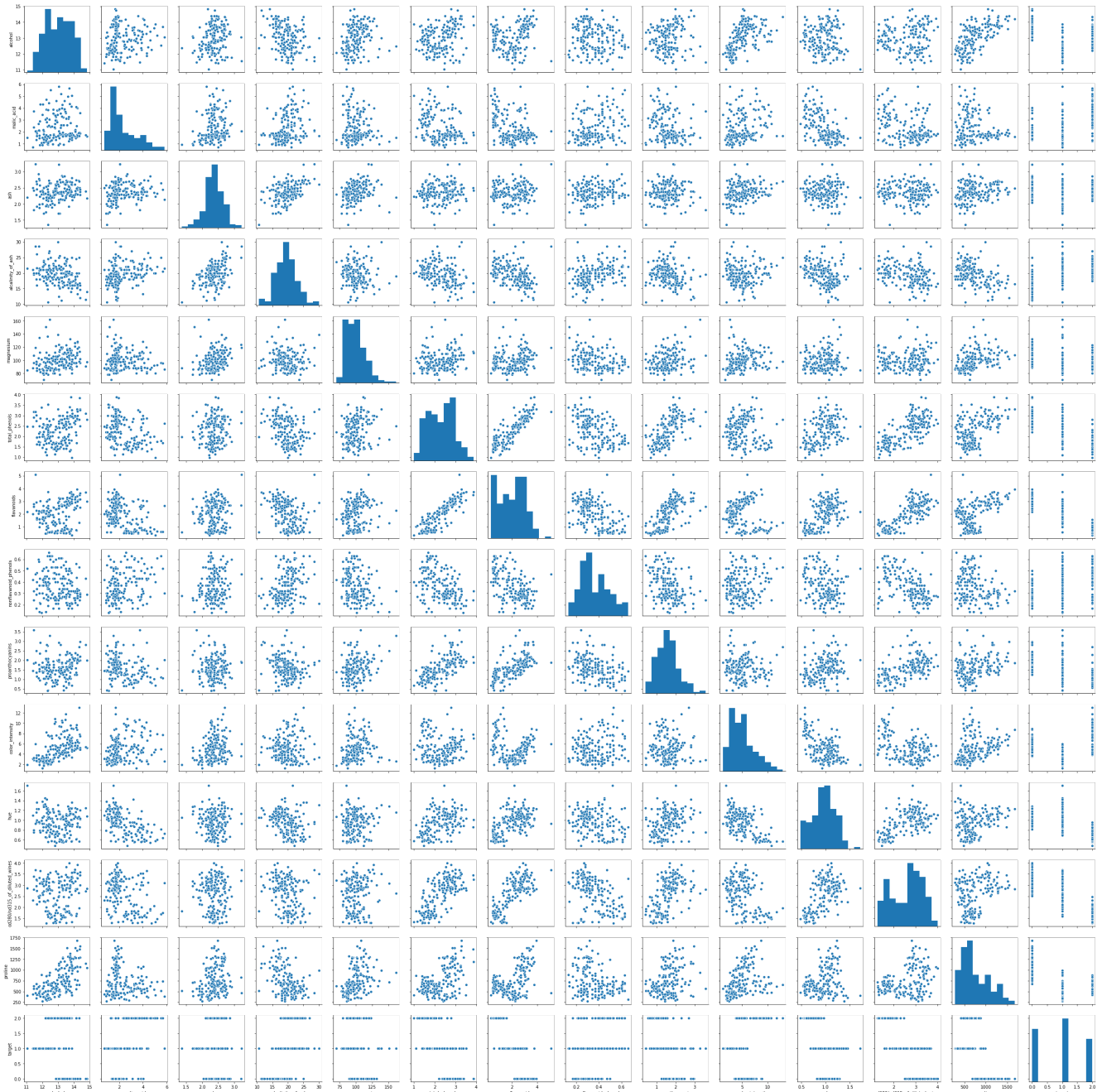
Парные диаграммы

In [9]:

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

Out[9]:

<seaborn.axisgrid.PairGrid at 0x129995bd0>



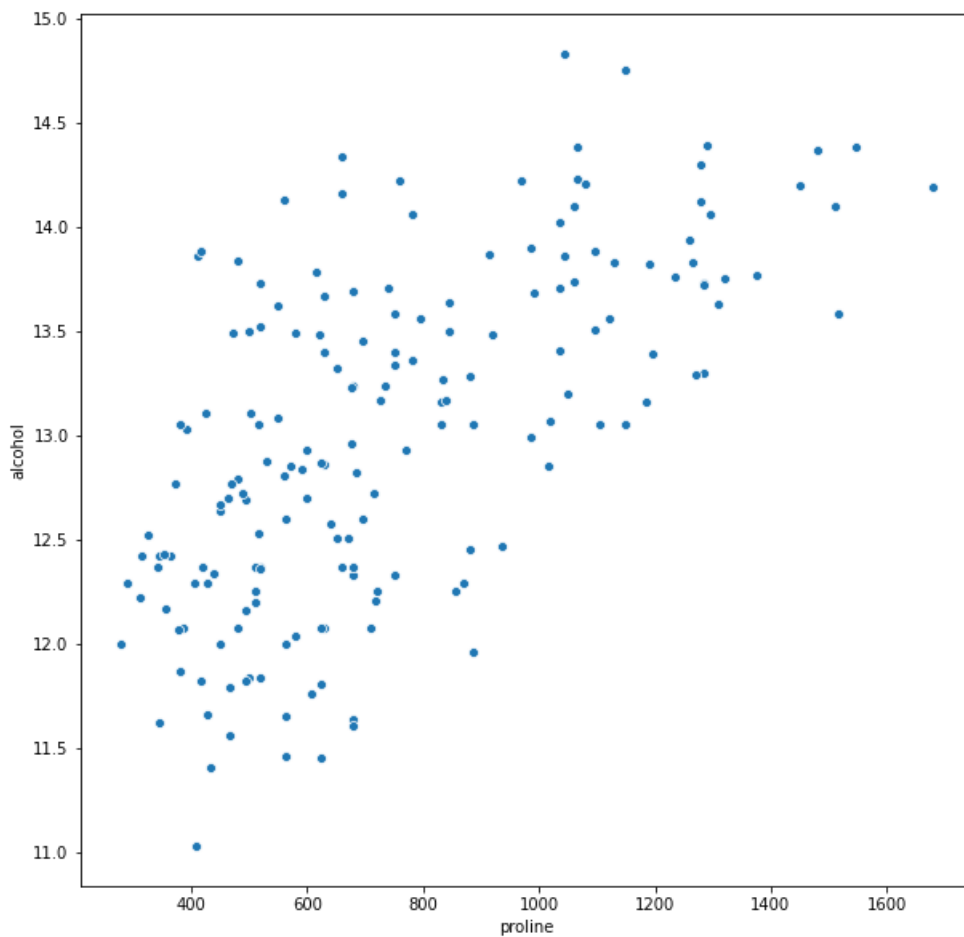
Находим почти линейную зависимость

In [14]:

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

Out[14]:

<matplotlib.axes._subplots.AxesSubplot at 0x131bfa090>



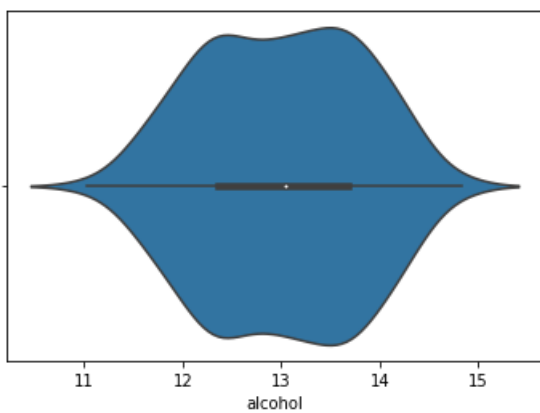
Violin plot

In [13]:

```
sns.violinplot(x=data['alcohol'])
```

Out[13]:

<matplotlib.axes._subplots.AxesSubplot at 0x131baacd0>



Корреляционная матрица

In [22]:

```
data.corr()
```

Out[22]:

	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	total_phenols	flavanoids	nonflavanoid_phe
alcohol	1.000000	0.094397	0.211545	-0.310235	0.270798	0.289101	0.236815	-0.15
malic_acid	0.094397	1.000000	0.164045	0.288500	-0.054575	-0.335167	-0.411007	0.29
ash	0.211545	0.164045	1.000000	0.443367	0.286587	0.128980	0.115077	0.18
alcalinity_of_ash	0.310235	0.288500	0.443367	1.000000	-0.083333	-0.321113	-0.351370	0.36
magnesium	0.270798	-0.054575	0.286587	-0.083333	1.000000	0.214401	0.195784	-0.25
total_phenols	0.289101	-0.335167	0.128980	-0.321113	0.214401	1.000000	0.864564	-0.44
flavanoids	0.236815	-0.411007	0.115077	-0.351370	0.195784	0.864564	1.000000	-0.53
nonflavanoid_phenols	0.155929	0.292977	0.186230	0.361922	-0.256294	-0.449935	-0.537900	1.00
proanthocyanins	0.136698	-0.220746	0.009652	-0.197327	0.236441	0.612413	0.652692	-0.36
color_intensity	0.546364	0.248985	0.258887	0.018732	0.199950	-0.055136	-0.172379	0.13
hue	0.071747	-0.561296	0.074667	-0.273955	0.055398	0.433681	0.543479	-0.26
od280/od315_of_diluted_wines	0.072343	-0.368710	0.003911	-0.276769	0.066004	0.699949	0.787194	-0.50
proline	0.643720	-0.192011	0.223626	-0.440597	0.393351	0.498115	0.494193	-0.31
target	0.328222	0.437776	0.049643	0.517859	-0.209179	-0.719163	-0.847498	0.48

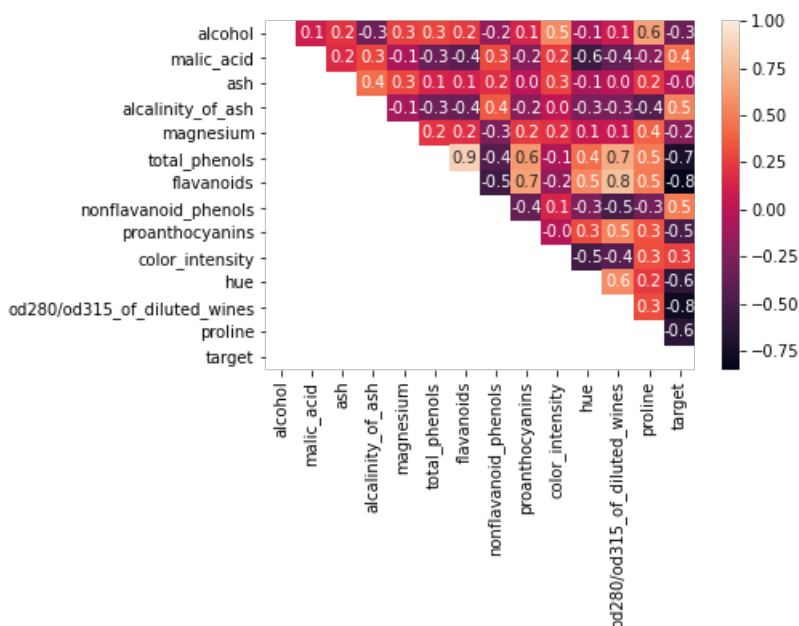
Матрица корреляций по Пирсону

In [23]:

```
mask = np.zeros_like(data.corr(), dtype=np.bool)
mask[np.tril_indices_from(mask)] = True
sns.heatmap(data.corr(), mask=mask, annot=True, fmt='.1f')
```

Out[23]:

<matplotlib.axes._subplots.AxesSubplot at 0x137de4ed0>



In []: