

Московский государственный технический университет имени Н. Э. Баумана

Факультет «Информатика и системы управления»

Кафедра ИУ5

Отчёт по

лабораторной работе № 1

«Технологии машинного обучения»

Подготовил:

Кан Андрей Дмитриевич

Группа ИУ5-54Б

Подпись_____

Дата_____

Москва
2021г.

Цель лабораторной работы: изучение различных методов визуализация данных.

Краткое описание. Построение основных графиков, входящих в этап разведочного анализа данных.

Рекомендуемые инструментальные средства можно посмотреть [здесь](#).

Задание:

- Выбрать набор данных (датасет). Вы можете найти список свободно распространяемых датасетов [здесь](#).
- Для первой лабораторной работы рекомендуется использовать датасет без пропусков в данных, например из [Scikit-learn](#).
- Пример преобразования датасетов Scikit-learn в Pandas Dataframe можно посмотреть [здесь](#).

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.
- Сформировать отчет и разместить его в своем репозитории на github.

Текст программы и экранные формы выполнения:

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

Загрузка данных

```
In [4]: wine = load_wine()
```

```
In [5]: type(wine)
```

```
Out[5]: sklearn.utils.Bunch
```

```
In [6]: for x in wine:  
         print(x)
```

```
data  
target  
frame  
target_names  
DESCR  
feature_names
```

```
In [7]: wine['feature_names']
```

```
Out[7]: ['alcohol',  
         'malic_acid',  
         'ash',  
         'alcalinity_of_ash',  
         'magnesium',  
         'total_phenols',  
         'flavanoids',  
         'nonflavanoid_phenols',  
         'proanthocyanins',  
         'color_intensity',  
         'hue',  
         'od280/od315_of_diluted_wines',  
         'proline']
```

```
In [8]: data1 = pd.DataFrame(data= np.c_[wine['data'], wine['target']],  
                             columns= wine['feature_names'] + ['target'])
```

2) Основные характеристики датасета

```
In [10]: # первые 5 строк датасета  
data1.head()
```

```
Out[10]:
```

	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	total_phenols	flavanoids	nonflavanoid_phenols	proanthocyanins	color_intensity	hue	od280/od315_of_diluted_wines	target
0	14.23	1.71	2.43	15.6	127.0	2.80	3.06	0.28	2.29	5.64	1.04	1.01	1
1	13.20	1.78	2.14	11.2	100.0	2.65	2.76	0.26	1.28	4.38	1.05	1.04	1
2	13.16	2.36	2.67	18.6	101.0	2.80	3.24	0.30	2.81	5.68	1.03	1.01	1
3	14.37	1.95	2.50	16.8	113.0	3.85	3.49	0.24	2.18	7.80	0.86	1.01	1
4	13.24	2.59	2.87	21.0	118.0	2.80	2.69	0.39	1.82	4.32	1.04	1.01	1

```
In [11]: # Размер датасета - 178 строк, 14 колонок  
data1.shape
```

```
Out[11]: (178, 14)
```

```
In [12]: # Список колонок  
data1.columns
```

```
Out[12]: Index(['alcohol', 'malic_acid', 'ash', 'alcalinity_of_ash', 'magnesium',  
               'total_phenols', 'flavanoids', 'nonflavanoid_phenols',  
               'proanthocyanins', 'color_intensity', 'hue',  
               'od280/od315_of_diluted_wines', 'proline', 'target'],  
              dtype='object')
```

```
In [14]: # Список колонок с типами данных  
data1.dtypes
```

```
In [15]: # Проверим наличие пустых значений
# Цикл по колонкам датасета
for col in data1.columns:
    # Количество пустых значений - все значения заполнены
    temp_null_count = data1[data1[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_intensity - 0
hue - 0
od280/od315_of_diluted_wines - 0
proline - 0
target - 0
```

```
In [16]: # Основные статистические характеристики набора данных
data1.describe()
```

```
Out[16]:
```

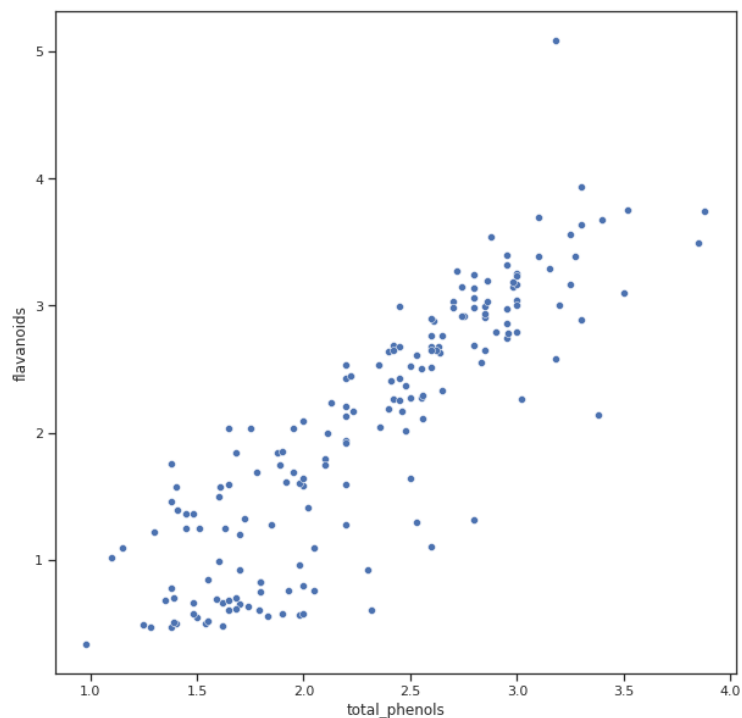
	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	total_phenols	flavanoids	nonflavanoid_phenols	proanthocyanins	color_intensity
count	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000	178.000000
mean	13.000618	2.336348	2.366517	19.494944	99.741573	2.295112	2.029270	0.361854	1.590899	5.058090
std	0.811827	1.117146	0.274344	3.339564	14.282484	0.625851	0.998859	0.124453	0.572359	2.318286
min	11.030000	0.740000	1.360000	10.600000	70.000000	0.980000	0.340000	0.130000	0.410000	1.280000
25%	12.362500	1.602500	2.210000	17.200000	88.000000	1.742500	1.205000	0.270000	1.250000	3.220000
50%	13.050000	1.865000	2.360000	19.500000	98.000000	2.355000	2.135000	0.340000	1.555000	4.690000
75%	13.677500	3.082500	2.557500	21.500000	107.000000	2.800000	2.875000	0.437500	1.950000	6.200000
max	14.830000	5.800000	3.230000	30.000000	162.000000	3.880000	5.080000	0.660000	3.580000	13.000000

3) Визуальное исследование датасета

Диаграмма рассеяния

```
In [60]: fig, ax = plt.subplots(figsize=(10,10))
sns.scatterplot(ax=ax, x='total_phenols', y='flavanoids', data=data1)

Out[60]: <AxesSubplot:xlabel='total_phenols', ylabel='flavanoids'>
```

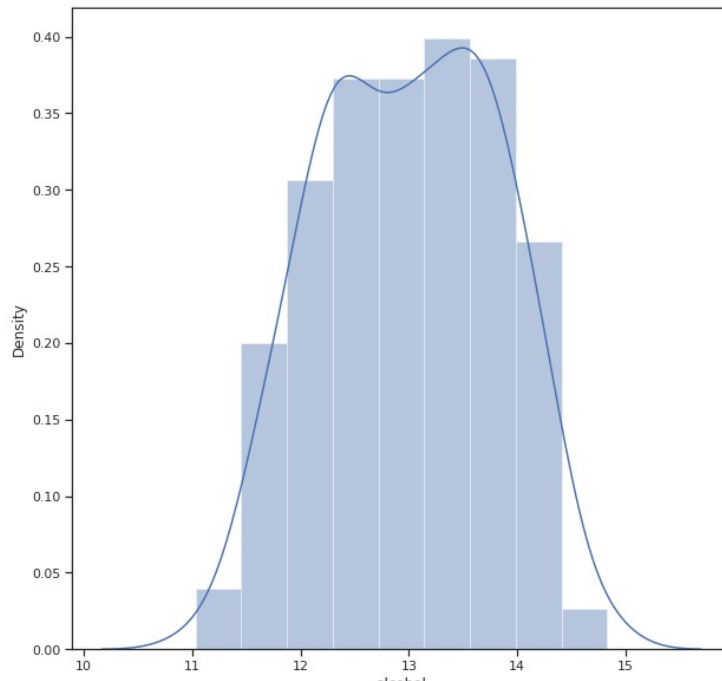


Гистограмма

```
In [27]: fig, ax = plt.subplots(figsize=(10,10))  
sns.distplot(data1['alcohol'])
```

/home/ripperonik/anaconda3/lib/python3.8/site-packages/seaborn/distributions.py:2557: 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)

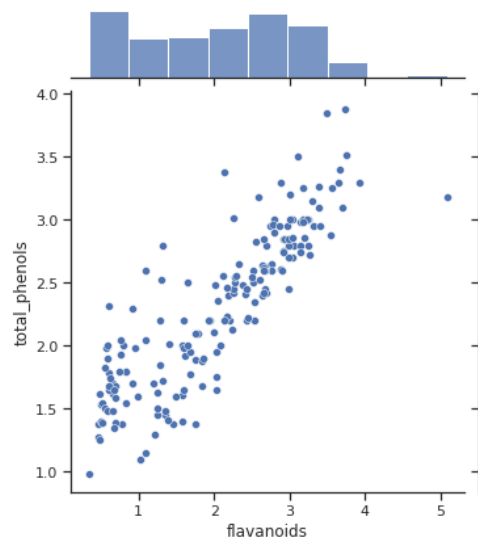
```
Out[27]: <AxesSubplot:xlabel='alcohol', ylabel='Density'>
```



Jointplot

```
In [57]: sns.jointplot(x='flavanoids', y='total_phenols', data=data1)
```

```
Out[57]: <seaborn.axisgrid.JointGrid at 0x7f538b434d90>
```



4) Информация о корреляции признаков

In [50]: data1.corr()

Out[50]:

	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	total_phenols	flavanoids	nonflavanoid_phenols	proanthocyanins
alcohol	1.000000	0.094397	0.211545	-0.310235	0.270798	0.289101	0.236815	-0.155929	0.136681
malic_acid	0.094397	1.000000	0.164045	0.288500	-0.054575	-0.335167	-0.411007	0.292977	-0.220771
ash	0.211545	0.164045	1.000000	0.443367	0.286587	0.128980	0.115077	0.186230	0.009681
alcalinity_of_ash	-0.310235	0.288500	0.443367	1.000000	-0.083333	-0.321113	-0.351370	0.361922	-0.197327
magnesium	0.270798	-0.054575	0.286587	-0.083333	1.000000	0.214401	0.195784	-0.256294	0.236441
total_phenols	0.289101	-0.335167	0.128980	-0.321113	0.214401	1.000000	0.864564	-0.449935	0.612413
flavanoids	0.236815	-0.411007	0.115077	-0.351370	0.195784	0.864564	1.000000	-0.537900	0.652692
nonflavanoid_phenols	-0.155929	0.292977	0.186230	0.361922	-0.256294	-0.449935	-0.537900	1.000000	-0.365845
proanthocyanins	0.136681	-0.220771	0.009681	-0.197327	0.236441	0.612413	0.652692	-0.365845	1.000000
color_intensity	0.546364	0.248985	0.258887	0.018732	0.199950	-0.055136	-0.172379	0.139057	-0.025211
hue	-0.071747	-0.561296	-0.074667	-0.273955	0.055398	0.433681	0.543479	-0.262640	0.295558
od280/od315_of_diluted_wines	0.072343	-0.368710	0.003911	-0.276769	0.066004	0.699949	0.787194	-0.503270	0.519001
proline	0.643720	-0.192011	0.223626	-0.440597	0.393351	0.498115	0.494193	-0.311385	0.330411
target	-0.328222	0.437776	-0.049643	0.517859	-0.209179	-0.719163	-0.847498	0.489109	-0.499111

С целевым признаком коррелируют nonflavanoid_phenols, malic_acid, alcalinity_of_ash. Также существует отрицательная корреляция с total_phenols, flavanoids, hue, proline. Это значит, что рост целевой компоненты приводит к уменьшению других компонентов.

In [53]: sns.heatmap(data1.corr())

Out[53]: <AxesSubplot:>

