Разведочный анализ данных. Исследование и визуализация данных

Цель лабораторной работы: изучение различных методов визуализация данных.

Задание: Создать ноутбук, который содержит следующие разделы: Текстовое описание выбранного Вами набора данных. Основные характеристики датасета. Визуальное исследование датасета. Информация о корреляции признаков.

In [1]:

```python
import numpy as np
import pandas as pd
from sklearn import datasets
import matplotlib.pyplot as plt
import seaborn as sns

%matplotlib inline
```

In [2]:

```python
# Будем анализировать данные только на обучающей выборке
df = datasets.load_wine()
```

In [3]:

```python
print(df.DESCR)
```

.. _wine_dataset:

Wine recognition dataset
------------------------

**Data Set Characteristics:**

    :Number of Instances: 178 (50 in each of three classes)
    :Number of Attributes: 13 numeric, predictive attributes and the class
    :Attribute Information:
    - Alcohol
    - Malic acid
    - Ash
  - Alcalinity of ash
    - Magnesium
    - Total phenols
    - Flavanoids
    - Nonflavanoid phenols
    - Proanthocyanins
  - Color intensity
    - Hue
    - OD280/OD315 of diluted wines
    - Proline

    - class:
        - class_0
        - class_1
        - class_2

    :Summary Statistics:

    ============================= ==== ===== ======= =====
                                   Min  Max  Mean    SD
    ============================= ==== ===== ======= =====
    Alcohol:                      11.0 14.8  13.0   0.8
    Malic Acid:                   0.74 5.80  2.34   1.12
    Ash:                          1.36 3.23  2.36   0.27
    Alcalinity of Ash:            10.6 30.0  19.5   3.3
    Magnesium:                    70.0 162.0 99.7  14.3
    Total Phenols:                0.98 3.88  2.29   0.63
    Flavanoids:                   0.34 5.08  2.03   1.00
    Nonflavanoid Phenols:         0.13 0.66  0.36   0.12
    Proanthocyanins:              0.41 3.58  1.59   0.57
    Colour Intensity:             1.3  13.0  5.1    2.3
    Hue:                          0.48 1.71  0.96   0.23
    OD280/OD315 of diluted wines: 1.27 4.00  2.61   0.71
    Proline:                      278  1680  746    315
    ============================= ==== ===== ======= =====
```

In [4]:

```
df = pd.DataFrame(data= np.c_[df['data'], df['target']],
          columns= df['feature_names'] + ['target'])
df.head()
```

Out[4]:

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

In [5]:

```
df.dtypes
```

Out[5]:

```
alcohol             float64
malic_acid          float64
ash                 float64
alcalinity_of_ash   float64
magnesium           float64
total_phenols       float64
flavanoids          float64
```

```
flavanoids                   float64
nonflavanoid_phenols         float64
proanthocyanins              float64
color_intensity              float64
hue                          float64
od280/od315_of_diluted_wines float64
proline                      float64
target                       float64
dtype: object
```

In [6]:

```python
# Размер датасета
df.shape
```

Out[6]:

```
(178, 14)
```

In [7]:

```python
# Список колонок
df.columns
```

Out[7]:

```
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 [8]:

```python
# Проверим наличие пустых значений
# Цикл по колонкам датасета
for col in df.columns:
    # Количество пустых значений - все значения заполнены
    temp_null_count = df[df[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 [9]:

```python
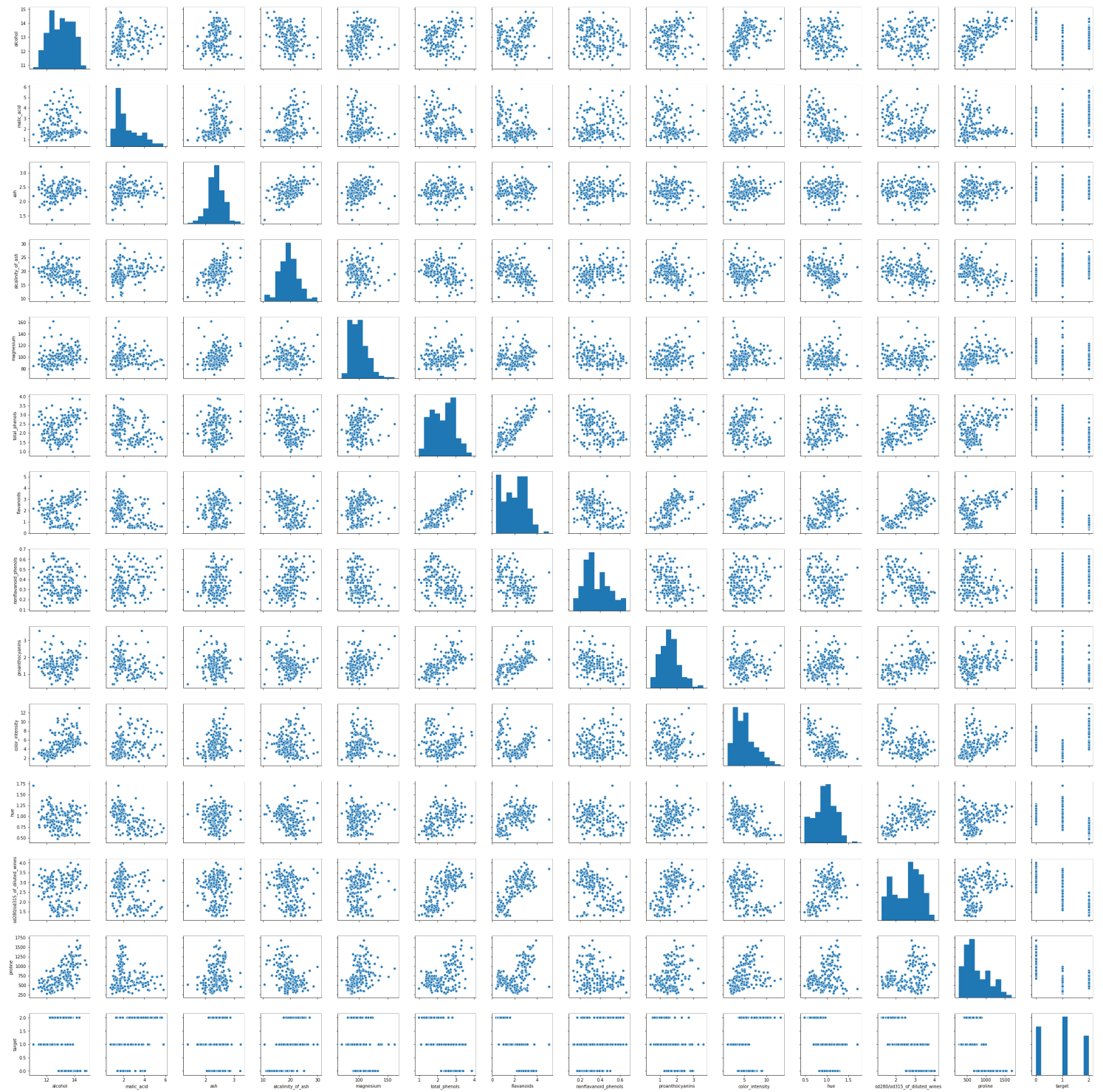# Основные статистические характеристки набора данных
df.describe()
```

Out[9]:

| | alcohol | malic_acid | ash | alcalinity_of_ash | magnesium | total_phenols | flavanoids | nonflavanoid_phenols | proanthocyanins | color_int |
|---|---|---|---|---|---|---|---|---|---|---|
| count | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.0 |
| mean | 13.000618 | 2.336348 | 2.366517 | 19.494944 | 99.741573 | 2.295112 | 2.029270 | 0.361854 | 1.590899 | 5.0 |
| std | 0.811827 | 1.117146 | 0.274344 | 3.339564 | 14.282484 | 0.625851 | 0.998859 | 0.124453 | 0.572359 | 2.3 |
| min | 11.030000 | 0.740000 | 1.360000 | 10.600000 | 70.000000 | 0.980000 | 0.340000 | 0.130000 | 0.410000 | 1.2 |
| 25% | 12.362500 | 1.602500 | 2.210000 | 17.200000 | 88.000000 | 1.742500 | 1.205000 | 0.270000 | 1.250000 | 3.2 |
| 50% | 13.050000 | 1.865000 | 2.360000 | 19.500000 | 98.000000 | 2.355000 | 2.135000 | 0.340000 | 1.555000 | 4.6 |
| 75% | 13.677500 | 3.082500 | 2.557500 | 21.500000 | 107.000000 | 2.800000 | 2.875000 | 0.437500 | 1.950000 | 6.2 |

In [10]:

```python
#Комбинация гистограмм и диаграмм рассеивания для всего набора данных.
sns.pairplot(data= df)
```

Out[10]:

<seaborn.axisgrid.PairGrid at 0x1b905944f28>



In [11]:

```python
df.corr()
```

Out[11]:

| | alcohol | malic_acid | ash | alcalinity_of_ash | magnesium | total_phenols | flavanoids | nonflavanoid_phenols | proantho |
|---|---|---|---|---|---|---|---|---|---|
| alcohol | 1.000000 | 0.094397 | 0.211545 | -0.310235 | 0.270798 | 0.289101 | 0.236815 | -0.155929 | |
| malic_acid | 0.094397 | 1.000000 | 0.164045 | 0.288500 | -0.054575 | -0.335167 | -0.411007 | 0.292977 | |

| | alcohol | malic_acid | ash | alcalinity_of_ash | magnesium | total_phenols | flavanoids | nonflavanoid_phenols | proantho |
|---|---|---|---|---|---|---|---|---|---|
| ash | 0.211545 | -0.164045 | 1.000000 | 0.443367 | 0.286587 | 0.128980 | 0.115077 | -0.186230 | |
| alcalinity_of_ash | -0.310235 | 0.288500 | 0.443367 | 1.000000 | -0.083333 | -0.321113 | -0.351370 | 0.361922 | - |
| magnesium | 0.270798 | -0.054575 | 0.286587 | -0.083333 | 1.000000 | 0.214401 | 0.195784 | -0.256294 | |
| total_phenols | 0.289101 | -0.335167 | 0.128980 | -0.321113 | 0.214401 | 1.000000 | 0.864564 | -0.449935 | |
| flavanoids | 0.236815 | -0.411007 | 0.115077 | -0.351370 | 0.195784 | 0.864564 | 1.000000 | -0.537900 | |
| nonflavanoid_phenols | -0.155929 | 0.292977 | 0.186230 | 0.361922 | -0.256294 | -0.449935 | -0.537900 | 1.000000 | - |
| proanthocyanins | 0.136698 | -0.220746 | 0.009652 | -0.197327 | 0.236441 | 0.612413 | 0.652692 | -0.365845 | |
| color_intensity | 0.546364 | 0.248985 | 0.258887 | 0.018732 | 0.199950 | -0.055136 | -0.172379 | 0.139057 | - |
| hue | -0.071747 | -0.561296 | -0.074667 | -0.273955 | 0.055398 | 0.433681 | 0.543479 | -0.262640 | |
| od280/od315_of_diluted_wines | 0.072343 | -0.368710 | 0.003911 | -0.276769 | 0.066004 | 0.699949 | 0.787194 | -0.503270 | |
| proline | 0.643720 | -0.192011 | 0.223626 | -0.440597 | 0.393351 | 0.498115 | 0.494193 | -0.311385 | |
| target | -0.328222 | 0.437776 | -0.049643 | 0.517859 | -0.209179 | -0.719163 | -0.847498 | 0.489109 | - |

In [12]:

```python
sns.heatmap(df.corr())
```

Out[12]:

<matplotlib.axes._subplots.AxesSubplot at 0x1b90a4ab9b0>