

Министерство образования Республики Беларусь  
Учреждение образования  
«Брестский государственный технический университет»  
Кафедра ИИТ

Лабораторная работа №1  
По дисциплине: «ОМО»  
Тема: «Знакомство с анализом данных: предварительная обработка и визуализация»

Выполнил:  
Студенты 3-го курса  
Группы АС-65  
Осовец М. М.  
Проверил:  
Крощенко А. А.

**Цель работы:** Получить практические навыки работы с данными с использованием библиотек Pandas для манипуляции и Matplotlib для визуализации. Научиться выполнять основные шаги предварительной обработки данных, такие как очистка, нормализация и работа с различными типами признаков.

### Вариант 3

Выборка Iris. Классический набор данных для классификации, содержащий измерения длины и ширины чашелистиков и лепестков для трех видов ирисов.

Задачи:

**1. Загрузите данные и проверьте, есть ли в них пропущенные значения.**

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler

# Загружаем CSV (обязательно должен лежать рядом в папке)
df = pd.read_csv('iris.csv')

# Первые 5 строк
df.head()
```

	sepal.length	sepal.width	petal.length	petal.width	variety
0	5.1	3.5	1.4	0.2	Setosa
1	4.9	3.0	1.4	0.2	Setosa
2	4.7	3.2	1.3	0.2	Setosa
3	4.6	3.1	1.5	0.2	Setosa
4	5.0	3.6	1.4	0.2	Setosa

Проверка на пропуски

```
df.isnull().sum()
```

```
sepal.length    0
sepal.width     0
petal.length    0
petal.width     0
variety         0
dtype: int64
```

## 2. Выведите количество образцов каждого вида ириса.

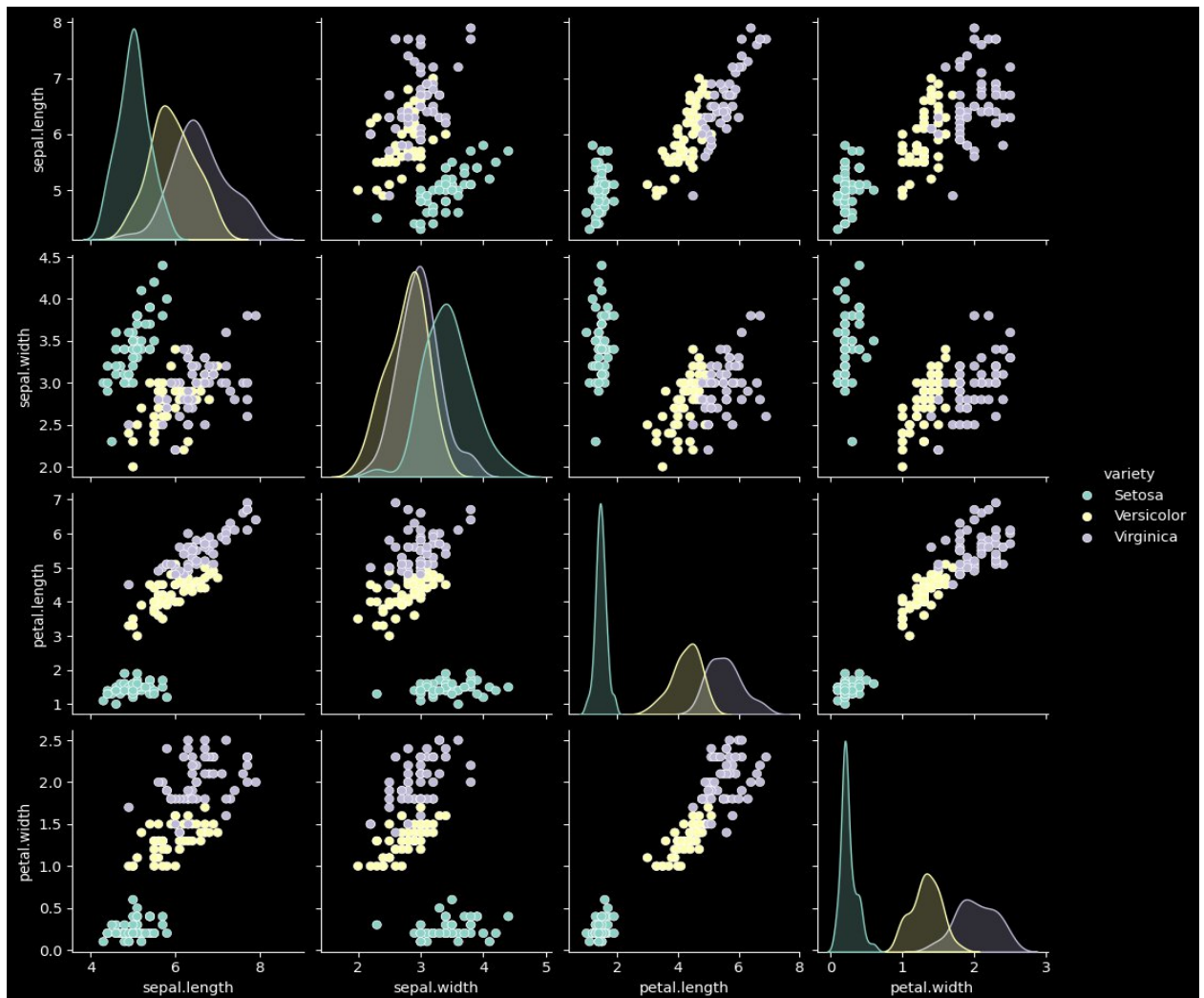
```
df['variety'].value_counts()
```

```
variety
Setosa      50
Versicolor 50
Virginica   50
Name: count, dtype: int64
```

## 3. Постройте парные диаграммы рассеяния (pair plot) для всех признаков, чтобы визуально оценить их разделимость.

```
sns.pairplot(df, hue='variety')
```

```
plt.show()
```



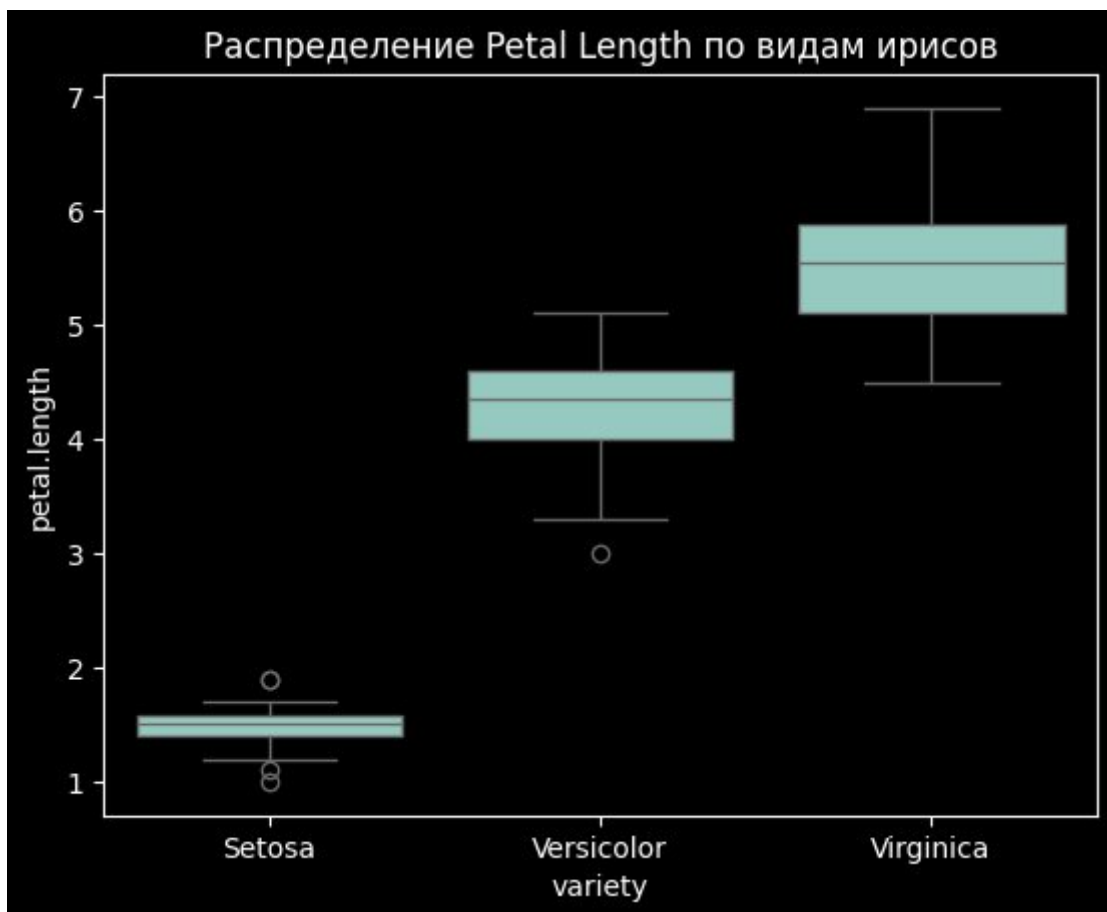
**4. Для каждого вида ириса рассчитайте среднее значение по каждому из четырех признаков.**

```
df.groupby('variety').mean()
```

	sepal.length	sepal.width	petal.length	petal.width
variety				
Setosa	5.006	3.428	1.462	0.246
Versicolor	5.936	2.770	4.260	1.326
Virginica	6.588	2.974	5.552	2.026

**5. Создайте "ящик с усами" (box plot) для признака Petal Length (cm), чтобы сравнить его распределение по разным видам ирисов.**

```
sns.boxplot(x='variety', y='petal.length', data=df)
plt.title('Распределение Petal Length по видам ирисов')
plt.show()
```



**6. Стандартизируйте данные (приведите к нулевому среднему и единичному стандартному отклонению).**

```
features = df.drop(columns=['variety'])
```

```

scaler = StandardScaler()

scaled_features = scaler.fit_transform(features)

df_scaled = pd.DataFrame(scaled_features,
columns=features.columns)

df_scaled['variety']=df['variety']

print("ДО стандартизации:")
print(features.describe())
print("\n" + "="*50 + "\n")

print("После стандартизации")
print(df_scaled.drop(columns=['variety']).describe())

```

ДО стандартизации:

	sepal.length	sepal.width	petal.length	petal.width
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.057333	3.758000	1.199333
std	0.828066	0.435866	1.765298	0.762238
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

=====

После стандартизации

	sepal.length	sepal.width	petal.length	petal.width
count	1.500000e+02	1.500000e+02	1.500000e+02	1.500000e+02
mean	-4.736952e-16	-7.815970e-16	-4.263256e-16	-4.736952e-16
std	1.003350e+00	1.003350e+00	1.003350e+00	1.003350e+00
min	-1.870024e+00	-2.433947e+00	-1.567576e+00	-1.447076e+00
25%	-9.006812e-01	-5.923730e-01	-1.226552e+00	-1.183812e+00
50%	-5.250608e-02	-1.319795e-01	3.364776e-01	1.325097e-01
75%	6.745011e-01	5.586108e-01	7.627583e-01	7.906707e-01
max	2.492019e+00	3.090775e+00	1.785832e+00	1.712096e+00

**Вывод:**

Получили практические навыки работы с данными с использованием библиотек Pandas для манипуляции и Matplotlib для визуализации. Научились выполнять основные шаги предварительной обработки данных, такие как очистка, нормализация и работа с различными типами признаков.