

# Data Analysis with Python

## Syllabus

Details of the academic discipline	
<b>Level of higher education</b>	<b>First (undergraduate)</b>
<b>Branch of knowledge</b>	12 Information technologies
<b>Specialty</b>	126 Information systems and technologies
<b>Educational program</b>	Integrated information systems, Information management systems and technologies, Information support of robotic systems
<b>Discipline status</b>	Academic discipline of professional and practical training (chosen by students)
<b>Form of education</b>	full-time/correspondence/distance
<b>Year of training, semester</b>	3rd year, autumn semester
<b>Scope of the discipline</b>	120 hours (36 hours – lectures, 18 hours – laboratory, 66 hours – individual work)
<b>Semester control/control measures</b>	Assessment/assessment work
<b>Lessons schedule</b>	<a href="http://rozklad.kpi.ua/Schedules/ScheduleGroupSelection.aspx">http://rozklad.kpi.ua/Schedules/ScheduleGroupSelection.aspx</a>
<b>Language of teaching</b>	Ukrainian, English
<b>Information about the course leader / teachers</b>	Lecturer: Tymofieieva Yuliia Serhiivna, <a href="mailto:yulia.s.timofeeva@gmail.com">yulia.s.timofeeva@gmail.com</a> Laboratory: Tymofieieva Yuliia Sergiivna, <a href="mailto:yulia.s.timofeeva@gmail.com">yulia.s.timofeeva@gmail.com</a>
<b>Placement of the course</b>	<a href="https://campus.kpi.ua">https://campus.kpi.ua</a>

## Program of educational discipline

### 1. Description of the educational discipline, its purpose, subject of study and learning outcomes

**Description of the discipline.** During training, students will get acquainted with the basic concepts, methods and means of preliminary training, statistical processing, visualization and data analysis, the basics of machine learning. In the laboratory classes, you will learn to use the main libraries of the Python language for working with data (NumPy, SciPy, Matplotlib, Seaborn, Scikit-learn, Pandas). The quality control of the acquired knowledge is assessed in the form of a modular test work.

**The subject of the academic discipline:** methods and means of data analysis and their implementation in the Python language.

**Interdisciplinary connections.** Basic knowledge of the disciplines: programming, algorithms and data structures, mathematical analysis, discrete mathematics, probability theory.

**The purpose of the educational discipline.** The goal of the educational discipline is to provide students with theoretical knowledge and practical abilities to apply general methods and tools for data preparation, statistical processing, visualization and analysis, using the main libraries of the Python language for working with data.

### The main tasks of the academic discipline

#### Knowledge:

- methods and means of software presentation of data;
- methods and algorithms for preliminary preparation of data for analysis;
- methods and tools of statistical processing and data visualization;
- basics of machine learning;
- the main functions of the Python language libraries for working with data.

#### Skills:

- effectively use Python language libraries for data processing, visualization and analysis;
- correctly choose methods and algorithms for data preprocessing;
- correctly choose the most informative ways of presenting data;
- use machine learning algorithms for classification, clustering, regression analysis, etc.
- effectively process, visualize and analyze data obtained as a result of own experiments and research.
- create Python programs to work with data

### 2. Pre-requisites and post-requisites of the discipline (place in the structural and logical scheme of training according to the relevant educational program)

**Prerequisites:** be able to use a computer, have basic knowledge of programming, data structures and algorithms, mathematical analysis, mathematical statistics, have experience in writing simple programs in the Python language.

**Post-requisites:** after completing the discipline, students will be able to effectively process, visualize and analyze data obtained as a result of their own experiments and research, apply basic machine learning algorithms, use NumPy, SciPy, Matplotlib libraries, Seaborn, Scikit-learn, Pandas.

### 3. Content of the academic discipline

#### Lecture classes

Chapter 1. Basic concepts of data processing

Chapter 2. Statistical data analysis

Chapter 3. Data structures of NumPy and Pandas libraries

Chapter 4. Data visualization methods

Chapter 5. Data preprocessing methods and algorithms

### Laboratory classes

1. Basic familiarity with the NumPy library.
2. Statistical data analysis.
3. Pandas data structures.
4. Data visualization using matplotlib and Seaborn.
5. Working with time series in Pandas.
6. Data preprocessing in Pandas.
7. Clustering and regression in scikit-learn.
8. Classification in scikit-learn.

### 4. Educational materials and resources

Literature:

1. Wes McKinney. Python for Data Analysis\_ Data Wrangling with Pandas, NumPy, and IPython. - O'Reilly Media, 2017. - 482 p.
2. Gayathri Rajagopalan. A Python Data Analyst's Toolkit. — apress, 2021. — 409 p.
3. Alex Campbell. Data Visualization Guide. — 2021. – 113 p.
4. Sanjeev J. Wagh. Fundamentals of Data Science. — Taylor & Francis Group, LLC, 2022. - 297
5. Avinash Navlani. Python Data Analysis. — Packt Publishing, 2021. - 463 p.
6. Joel Grus. Data Science from Scratch. - O'Reilly Media, Inc., 2019. - 513p.

## Educational content

### 5. Methods of mastering an educational discipline (educational component)

#### Lecture classes

No. z/p	The name of the topic of the lecture and a list of the main questions (a list of didactic tools, references to the literature and tasks on the SRS)
1	<p><b>Lecture 1.</b>Basic concepts of data processing</p> <p>The structure of the discipline and RSO. Concept of data, types of data, structured data, data sets and their components, main tasks of data processing, features of data processing. Classification and general overview of stages and methods of data processing. Python libraries for working with data.</p> <p>L. ( 1,3,5 )</p> <p><b>Tasks on SRS.</b>Concept of information and knowledge. Data attributes</p>
2	<p><b>Lecture 2.</b>Descriptive statistics with Python</p> <p>The main means of statistics for identifying the characteristics of a data set. Identification of statistical characteristics of data arrays using NumPy and SciPy. Hypothesis testing, general hypothesis testing algorithm, basic statistical tests.</p> <p>L. (1.5)</p> <p><b>Tasks for SRS</b>Installing libraries the Python language.</p>
3	<p><b>Lecture 3.</b>Statistical hypothesis testing with SciPy</p> <p>Hypothesis testing using the functions of the stats module of the SciPy library. Functions for one sample. Testing hypotheses about the equality of mathematical</p>

	<p>expectations. Use of one-sided alternative hypotheses. Testing hypotheses about the distribution of a random variable, stats.normaltest function.</p> <p>L. (2.5)</p> <p><b>Tasks for SRS.</b> Dispersion analysis of data.</p>
4	<p><b>Lecture 4.</b>Correlation and regression data analysis using Python.</p> <p>Correlation data analysis, corresponding NumPy and SciPy functions. Pearson and Spearman correlation coefficients. A simple linear regression model and its construction using SciPy tools.</p> <p>L. (5)</p> <p><b>Tasks for SRS.</b>Setting up the NumPy library.</p>
5	<p><b>Lecture 5.</b>NumPy library data structures</p> <p>NumPy multidimensional arrays, advantages of their use in data processing. Creating arrays, attributes of arrays, working with different dimensions. Functions for combining arrays. Arithmetic and Boolean operations with arrays.</p> <p>L. (4,5)</p> <p><b>Tasks for SRS.</b>Pandas library settings.</p>
6	<p><b>Lecture 6.</b>Pandas data structures</p> <p>The Series object, its attributes and indices. Two-dimensional DataFrame object, flexible row and column indexes. Add and remove rows and columns. Index object and its types.</p> <p>L. (3,4)</p> <p><b>Tasks for SRS.</b> Data types of the Index object</p>
7	<p><b>Lecture 7.</b>Joining and grouping data in Pandas</p> <p>Joining a DataFrame and a Series using the concat function. The merge function and its parameters, different types of merging. Data grouping using groupby. Methods of grouped objects and their use.</p> <p>L. (2,4)</p> <p><b>Tasks for SRS.</b> Preparation for the first test.</p>
8	<p><b>Lecture 8.</b>Data visualization</p> <p>The first test. Data analysis by visualization methods. Main libraries for data visualization: matplotlib, Pandas, Seaborn. Bar charts.</p> <p>L. (6)</p> <p><b>Tasks for SRS.</b> Additional graph settings in matplotlib.</p>
9	<p><b>Lecture 9.</b>Diagrams for one and two characteristics</p> <p>Histograms and their construction using Python. Scaling charts and their use to detect statistical outliers. Construction of scatter diagrams. Using heatmaps to display the correlation between quantities.</p> <p>L. (6 )</p> <p><b>Tasks for SRS.</b> Additional features of Seaborn.</p>

10	<p><b>Lecture 10.</b>Time series in Pandas</p> <p>Concept of time series. Working with time series in Pandas. Using datetime as an index in Pandas objects. Visualization of time series. Analysis of time series by means of Pandas.</p> <p>L. (4,5,6)</p> <p><b>Tasks for SRS.</b> Moving average and other indicators.</p>
11	<p><b>Lecture 11.</b>Working with files in Pandas</p> <p>Basic functions for reading files of various formats typical for representing datasets. Options for reading and writing files. Different json file orientations. Reading datasets from web pages.</p> <p>L. (4)</p> <p><b>Tasks for SRS.</b> Structure of json files.</p>
12	<p><b>Lecture 12.</b>Data pre-processing</p> <p>Basic actions during data preprocessing. Renaming data. Sorting in Pandas. Conversion of data types. Detection and processing of duplicate data. Missing data, their detection, deletion and filling.</p> <p>L. (4,5)</p> <p><b>Tasks for SRS.</b> Preparation for the second test.</p>
13	<p><b>Lecture 13.</b>Data cleaning and transformation</p> <p>The second test. Determination of erroneous and anomalous data. Data filtering. Wide and long formats and transformation between them.</p> <p>L. (4,5)</p> <p><b>Tasks for SRS.</b> Installation and configuration of the scikit-learn library</p>
14	<p><b>Lecture 14.</b>Preparation of data for training</p> <p>Main categories of machine learning models. Getting to know the scikit-learn library. Division of data into training and test samples. Scaling and centering data. Coding of categorical data. Missing data detection and filling in scikit-learn.</p> <p>L. (5)</p> <p><b>Tasks for SRS.</b> Indicator variables.</p>
15	<p><b>Lecture 15.</b>Regression and clustering in scikit-learn</p> <p>Linear regression. Model training and testing. Metrics for regression model evaluation. Hyperparameters of the regression model. Clustering problem, k-means algorithm. Metrics for evaluating the clustering model.</p> <p>L. (2,4,5)</p> <p><b>Tasks for SRS.</b>Visualization of clusters.</p>
16	<p><b>Lecture 16.</b>Classification in scikit-learn</p> <p>Binary classification and logistic regression. Metrics for evaluating the classification model. Balanced and unbalanced classes. Decision tree, ensemble methods and boosting.</p> <p>L. (1.5)</p> <p><b>Tasks for SRS.</b>Gini coefficient.</p>
17	<p><b>Lecture 17.</b>Improvement of models in scikit-learn</p>

	Finding the best hyperparameters using the GridSearch class. Cross validation. Selection, extraction, enhancement and feature creation. Analysis of principal components. L. (2.5) <b>Tasks for SRS.</b> Definition of anomalies.
18	<b>Lecture 18.</b> Recommendation systems Different types of recommendation systems. Movie recommendation system with scikit-learn. Anomaly detection problem. L. (2.5)

### Laboratory classes

No. z/p	The name of the laboratory session	Number of aud. hours
1	A basic introduction to the NumPy library	2
2	Statistical data analysis	2
3	Pandas data structures	2
4	Data visualization using matplotlib and Seaborn	2
5	Working with time series in Pandas	2
6	Data preprocessing in Pandas	2
7	Clustering and regression in scikit-learn	2
8	Classification in scikit-learn	2
9	Final lesson	2

### 6. Independent work of student

No. z/p	The name of the topic submitted for independent processing	Number of hours of SRS
1	Concept of information and knowledge. Data attributes	3
2	Installing libraries the Python language	4
3	Dispersion analysis of data	5
4	Setting up the NumPy library	5
5	Pandas library settings	4
6	Data types of the Index object	4
7	Preparation for the first part of the modular test	4
8	Additional graph settings in matplotlib	4
9	Additional features of Seaborn	5
10	Moving average and other indicators	4
11	Structure of json files	3
12	Preparation for the second part of the modular control work	4
13	Installation and configuration of the Skit-learn library	4
14	Indicator variables	3
15	Visualization of clusters	3
16	Gini coefficient	4
17	Definition of anomalies	3

### 7. Policy of academic discipline (educational component)

The system of requirements for the student:

- attending lectures and laboratory classes is a mandatory component of studying the material;
- the teacher uses his own presentation material at the lecture; uses Google Drive for teaching the material of the current lecture, additional information, tasks for laboratory work, etc.;
- questions at lectures are asked in the time allotted for this purpose;
- laboratory reports are downloaded on the eve of the defense; to defend the laboratory work, it is necessary to demonstrate the operation of the program corresponding to the task and answer the questions about the program and control questions;
- modular test papers are written in lectures without the use of aids (mobile phones, tablets, etc.); the result is downloaded in a file through a Google form to the appropriate directory of Google Drive;
- incentive points are awarded for: participation in faculty and institute olympiads in academic disciplines, participation in work competitions, preparation of reviews of scientific works, etc. The number of encouraged points is no more than 10;
- Penalty points are awarded for: late submission of laboratory work without valid reasons. The number of penalty points is no more than 10.

## 8. Types of control and rating system for evaluating learning outcomes (RSE)

A student's rating consists of the points he receives for:

1. performance and protection of laboratory work;
2. execution of modular control work;
3. incentive and penalty points.

### System of rating points and evaluation criteria

#### Laboratory tasks:

"excellent", complete answer to the question during the defense (at least 90% of the required information), complete completion of the laboratory task - 9 points;

"good", sufficiently complete answer to the question during the defense (at least 75% of the required information), complete completion of the laboratory task - 7-8 points;

"satisfactory", incomplete answer to the questions during the defense (at least 60% of the required information), minor errors in the performance of the laboratory work task - 5-6 points;

"unsatisfactory", an unsatisfactory answer and/or significant errors in the performance of the laboratory task - 0 points.

#### Modular control works:

"excellent", a complete answer (at least 90% of the required information), the task was completed without errors, the actions were justified - 14 points;

"good", a sufficiently complete answer (at least 75% of the required information), the task was completed without significant errors - 10-13 points;

"satisfactory", incomplete answer, significant errors may be present in some tasks, but at least 60% are completed correctly - 7-9 points;

"unsatisfactory", an unsatisfactory answer (incorrect performance of tasks), requires mandatory rewriting at the end of the semester - 0 points.

#### Incentive points

for the performance of creative works from the credit module (for example, participation in faculty and institute olympiads in academic disciplines, participation in work competitions, preparation of reviews of scientific works, etc.) 1-2 points, but not more than 10 in total.

#### Intersessional certification

According to the results of the educational work for the first 8 weeks, the maximum possible number of points is 41 points (3 laboratory tasks, the first part of the modular control work). At the first certification (8th week), the student receives "passed" if his current rating is not less than 21 points.

According to the results of 13 weeks of training, the maximum possible number of points is 82 points (6 laboratory tasks, a modular control work). At the second certification (14th week), the student receives "passed" if his current rating is not less than 41 points.

The maximum sum of weighted points of control measures during the semester is:

$$RD = 8 \cdot r_{l.r.} + 2 \cdot r_{mkr} + (r_z - r_{sh}) = 8 \cdot 9 + 2 \cdot 14 + (r_z - r_{sh}) = 100 + (r_z - r_{sh}),$$

where  $r_{l.r.}$  – score for the laboratory task (0...9);

$r_{mkr}$  – score for writing part of the MKR (0...14);

$r_z$  – incentive points (0...10);

$r_{zsh}$  - penalty points (0...10).

#### Test:



Students who have fulfilled all the conditions of admission to the semester certification (written all module tests, completed and defended all laboratory work) and scored the required number of points during the semester ( $RD \geq 60$ ), receive a credit score (credit) automatically in accordance with the obtained rating (table. 1 below). In this case, RD points and corresponding grades are entered into the credit and examination information.

Students who scored less than 60 points during the semester and do not have debts are required to complete the credit control work.

The final test contains five questions. Each question is worth 20 points.

Students who scored more than 60 points during the semester and fulfilled all admission conditions are given the opportunity to perform a credit test to improve their grade. In this case, the student's previous rating from the credit module is canceled and he receives an assessment taking into account the results of the credit test (Table 1 below).

**Question evaluation system:**

"excellent", complete answer (at least 90% of the required information) - 18-20 points;

"good", sufficiently complete answer (at least 75% of the required information, or minor inaccuracies) - 15-17 points;

"satisfactory", incomplete answer (at least 60% of the required information and some errors) - 12-14 points;

"unsatisfactory", unsatisfactory answer - 0-11 points.

The sum of RD points or points for credit work is transferred to the credit score according to the table:

**Table 1 — Conversion of rating points to grades on the university scale**

Scores	Rating
100-95	Perfectly
94-85	Very good
84-75	Fine
74-65	Satisfactorily
64-60	Enough
Less than 60	Unsatisfactorily
Admission conditions not met	Not allowed

## **9. Additional information on the discipline (educational component)**

### **Working program of the academic discipline (Syllabus):**

**Folded** Tymofieieva Yuliia Serhiivna

**Approved** Department of ICT (protocol No. 13 dated 15.06.2022)

**Agreed** by the methodical commission of the faculty (protocol No. 11 dated 07.07.2022)