



# BIG DATA AND ANALYTICS

## Syllabus

### Details of the academic discipline

Level of higher education	<i>Second (master's)</i>
Branch of knowledge	<i>F Information technologies</i>
Specialty	<i>F2 Software engineering</i>
Educational program	<i>Software Engineering of Multimedia and Information Retrieval Systems</i>
Discipline status	<i>Selective</i>
Form of education	<i>full-time</i>
Year of training, semester	<i>2nd year of training, 3rd semester</i>
Scope of the discipline	<i>Lectures: 36 hours, laboratory work: 18 hours, independent work: 66 hours.</i>
Semester control/ control measures	<i>Final test, modular test, calendar control</i>
Lessons schedule	<i>According to the schedule of the current academic year (<a href="http://roz.kpi.ua/">http://roz.kpi.ua/</a>)</i>
Language of teaching	<i>English</i>
Information about the head of the course / teachers	<i>Lecturer: Ph.D., associate professor, Liubov Oleshchenko, oleshchenkoliubov@gmail.com Laboratory work: Ph.D., associate professor, Liubov Oleshchenko, oleshchenkoliubov@gmail.com</i>
Placement of the course	<i>Google classroom: Access is given to registered students.</i>

### Program of educational discipline

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

*The study of the discipline "Big data and analytics" allows students to develop the competencies necessary for solving practical problems of professional and scientific activities related to the analysis of big data.*

***The purpose** of studying the discipline "Big data and analytics" is to form students' abilities to apply software methods of big data analytics to solve professional and scientific problems.*

***The subject** of the discipline "Big data and analytics" is software methods and technologies of big data analytics. After mastering the discipline "Big data and analytics", **the learning outcomes** are:*

#### *knowledge:*

- big data storage and support technologies;*
- the Hadoop stack for big data and the Hadoop Distributed File System (HDFS);*
- components of the Hadoop stack, including the YARN resource and task management system, HDFS, and the MapReduce programming model;*
- methods of parallel programming using Spark, built-in Spark libraries;*
- data placement strategies, using Zookeeper capabilities;*
- Spark Streaming and Sliding Window Analytics, Spark GraphX & Graph Analytics;*
- machine learning algorithms using Map Reduce for Big Data Analytics.*

#### *skills:*

- use big data analytics technologies using the capabilities of the Python programming languages.*

**experience:**

- *apply acquired knowledge for further scientific and professional activities.*

*The study of the discipline "Big data and analytics" contributes to the formation of higher education students who study under the educational program "Software Engineering of Multimedia and Information Retrieval Systems" competencies necessary for solving practical problems of professional activity related to big data analytics:*

- *Ability to abstract thinking, analysis and synthesis.*
- *Ability to conduct research at an appropriate level.*
- *Ability to develop and implement scientific and/or applied projects in the field of software engineering.*

*The study of the discipline "Big data and analytics" contributes to the formation of higher education students studying under the educational program "Software Engineering of Multimedia and Information Retrieval Systems" competencies necessary for solving practical problems of professional activities related to big data analytics:*

**GC01** *Ability to abstract thinking, analysis and synthesis.*

**GC03** *Ability to conduct research at the appropriate level.*

**PC02** *Ability to develop and implement scientific and / or applied projects in the field of software engineering.*

## **2. Pre-requisites and post-requisites of the discipline**

### **(the place of the course in the scheme of studies in accordance with curriculum)**

*The successful study of the discipline "Big data and analytics" is preceded by the study of the regulatory disciplines "Programming and "Databases" of the curriculum for bachelors in the specialty F2 "Software engineering".*

*The discipline "Big data and analytics" ensures the implementation of course projects and master's theses in the specialty F2 "Software engineering".*

## **3. Content of the academic discipline**

*The discipline "Big data and analytics" involves the study of the following topics:*

*Topic 1. Architectural models of big data.*

*Topic 2. Software methods of big data analytics.*

*Modular test.*

*Final test.*

## **4. Educational materials and resources**

### **Basic literature:**

1. *Big Data Tutorial.* <https://www.edureka.co/blog/big-data-tutorial>
2. *Big Data & Analytics.* [https://www.tutorialspoint.com/big\\_data\\_tutorials.htm](https://www.tutorialspoint.com/big_data_tutorials.htm)

### **Additional literature:**

1. *MapReduce. Tutorial. Electronic resource. Access mode:* [https://hadoop.apache.org/docs/r1.2.1/mapred\\_tutorial.html](https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html)
2. *Python programming language. Electronic resource. Access mode:* <https://www.python.org/>
3. *Running ZooKeeper in Production. Electronic resource. Access mode:* <https://docs.confluent.io/platform/current/zookeeper/deployment.html>
4. *Running ZooKeeper, A Distributed System Coordinator. Electronic resource. Access mode:* <https://kubernetes.io/docs/tutorials/stateful-application/zookeeper/>

5. *Hadoop YARN Architecture. Electronic resource. Access mode:*  
<https://www.geeksforgeeks.org/hadoop-yarn-architecture/>
6. *Understanding YARN architecture and features. Electronic resource. Access mode:*  
[https://docs.cloudera.com/HDPDocuments/HDP3/HDP-3.0.0/data-operating-system/content/apache\\_yarn.html](https://docs.cloudera.com/HDPDocuments/HDP3/HDP-3.0.0/data-operating-system/content/apache_yarn.html)
7. *Machine Learning Algorithm K-means using Map Reduce for Big Data Analytics. Electronic resource. Access mode:*  
<http://www.nitttrc.edu.in/nptel/courses/video/106104189/lec25.pdf>
8. *Big Data Predictive Analytics. Electronic resource. Access mode:*  
<https://www.predictiveanalyticstoday.com/big-data-analytics-and-predictive-analytics/>
9. *Page Rank Algorithm and Implementation. Electronic resource. Access mode:*  
<https://www.geeksforgeeks.org/page-rank-algorithm-implementation/>
10. *PageRank Algorithm in Big Data. Electronic resource. Access mode:*  
<http://www.nitttrc.edu.in/nptel/courses/video/106104189/lec31.pdf>
11. *GraphX Programming Guide. Electronic resource. Access mode:*  
<https://spark.apache.org/docs/latest/graphx-programming-guide.html>

*The materials are freely available on the Internet.*

## Educational content

### 5. Methodology

No.	Type of training session	Description of the lesson	Hours
<i>Topic 1. Architectural models of big data.</i>			
1	<i>Lecture 1. Big data and its role in the modern world. Problems of storage and analytics of big data.</i>	<i>Growth of data in today's world. Growth of Internet of Things devices. Definition of big data and its classification. Structured and unstructured data. Cloud computing. Fog computing.</i>	2
2	<i>Lecture 2. Big Data Infrastructure. Distributed Data and Analytics.</i>	<i>Big data infrastructure. Distributed data and its analytics.</i>	2
3	<i>Lecture 3. The life cycle of big data analytics.</i>	<i>Big data analytics life cycle. Software methods for extracting, transforming, and loading big data.</i>	2
4	<i>Laboratory work 1. Big Data Analytics and Machine Learning.</i>	<i>Lesson 1. Data analysis and preparation. Loading the test.csv dataset (50,000 records). Primary analysis: checking data types, detecting missing values, processing categorical variables, scaling numeric features.</i>	2
		<i>Lesson 2. Building a machine learning model. Create and train a credit score classification model (Logistic Regression, Random Forest, XGBoost). Model evaluation: accuracy, recall, precision, F1-measure.</i>	2
		<i>Lesson 3. Building an intelligent assessment system. Integrate the model into a credit risk assessment system. Classify users by credit score (high, medium, low). Visualize classification results.</i>	2

5	Lecture 4. Web scraping, web crawling, indexing, and web automation technologies .	Technologies and software methods for web scraping, web crawling, indexing and web automation. Machine-readable data and APIs. HTML parsers. DOM analysis. Semantic annotation recognition. Web page analyzers using computer vision.	2
6	Laboratory work 2. Analysis of financial tweets users for stocks traded on NYSE, NASDAQ and SNP.	Lesson 1. Preprocessing and analysis of tweets. Downloading the stockerbot-export.csv file (28,000 tweets). Data structure analysis: text cleaning (removing stop words, symbols, hashtags, mentions), normalization, tokenization, lemmatization. Studying the distribution of tweets by companies and symbols.	2
		Lesson 2. Building a financial sentiment classifier. Labeling or using existing labels for sentiment analysis (positive / negative / neutral sentiment). Text vectorization (TF-IDF or Embeddings). Model training (Logistic Regression, SVM, Naive Bayes or LSTM). Model evaluation by metrics (accuracy, F1-score).	2
		Lesson 3. Analysis of results and tracking sentiment. Analysis of model results. Visualization of sentiment by company (graphs, diagrams, comparisons). Building an interface or script to monitor current sentiment regarding individual companies or cryptocurrencies.	2
7	Lecture 5. Virtualization technologies for big data analytics.	The role of virtualization for big data analytics. Virtualization technologies. Layers of abstraction. Hypervisors.	2
8	Lecture 6. Container technology for executing program code on the server.	Container technology for executing software code on the server. Big data engineering. SaaS, PaaS and IaaS.	2
9	Lecture 7. Hadoop technologies for big data.	Hadoop technologies for big data. Scalability with big data.	2
10	Lecture 8. MapReduce software model and framework. HDFS capabilities.	MapReduce programming model and framework. Data storage and processing in distributed file systems. Distributed databases. Hadoop Distributed File System (HDFS).	2
11	Lecture 9. Kafka distributed streaming platform. Advantages of Cassandra distributed DBMS.	Data ingestion problem. Kafka distributed streaming platform. Advantages of Cassandra distributed DBMS.	2
12	Lecture 10. Apache Spark platform and its capabilities.	The problem of computational function. Spark technology. Comparison of Spark and MapReduce. Capabilities of Spark MLlib, Spark Streaming GraphX modules.	2
13	Lecture 11. Apache Flink capabilities for streaming and batch data processing.	Apache Flink architecture. Stateful functions. Flink ML capabilities.	2

<i>Topic 2. Software methods for big data analytics.</i>			
14	<i>Lecture 12. Distributed configuration and service management system for distributed applications ZooKeeper.</i>	<i>Running ZooKeeper in production. Hardware. Memory. CPU and GPU. Configuration options.</i>	2
15	<i>Laboratory work 3. Distributed Big Data Computing Using Spark Cluster.</i>	<i>Lesson 1. Installing and configuring Spark. Installing Apache Spark on a local machine (or in a Google Colab environment with PySpark). Checking its functionality. Connecting to a Spark session and configuring a Spark cluster.</i>	2
		<i>Lesson 2. Loading and processing big data. Import a large dataset (CSV). Read data using Spark DataFrame API. Primary processing: cleaning, filtering, aggregations.</i>	2
		<i>Lesson 3. Distributed analytics. Performing analytical operations in distributed mode: grouping, counting, sorting, calculating statistics. Plotting or saving results. Analyzing the advantages of distributed processing compared to traditional processing.</i>	2
16	<i>Lecture 13. Using Hadoop YARN for Big Data Analytics.</i>	<i>Hadoop YARN Architecture. Map Reduce in Hadoop. Using Hadoop YARN for Big Data Analytics.</i>	2
17	<i>Lecture 14. Software methods for clustering big data.</i>	<i>Software methods for clustering big data . Using the K-means and Map Reduce algorithms.</i>	2
18	<i>Lecture 15. Predictive analytics of big data.</i>	<i>Predictive analytics of big data, basic methods and technologies . Regression analysis of big data.</i>	2
19	<i>Lecture 16. Using Big Data cloud providers .</i>	<i>Overview of Big Data cloud providers. AWS Big Data platform.</i>	2
20	<i>Lecture 17. Using the PageRank algorithm for big data analytics.</i>	<i>Page Rank Algorithm and its Software Implementation for Big Data Analytics . Property Graph. Example Property Graph. Graph Operators. Aggregated Operator List. Property Operators. Structural Operators. Join Operators . Neighborhood Aggregation. Aggregate Messages (aggregateMessages). Map Reduce Triplets Transition Guide (Legacy). Computer Degree Information. Neighbor Gathering. Caching and Decaching. Pregel API. Graph Constructors.</i>	2
21	<i>Lecture 18. Conclusion lesson.</i>	<i>Review of the studied material. Modular control work.</i>	2

## 6. Self-study

The discipline "Big data and analytics" is based on independent preparation for classroom classes on theoretical topics.

<i>No</i>	<i>The name of the topic submitted for independent processing</i>	<i>Hours of study</i>	<i>References</i>
1	<i>Preparation for the lecture 1</i>	2	[1]
2	<i>Preparation for lecture 2</i>	2	[1], [1 add.]
3	<i>Preparation for the lecture 3</i>	2	[1]
4	<i>Preparation for the laboratory work 1</i>	4	[2]
5	<i>Preparation for the lecture 4</i>	2	[1]
6	<i>Preparation for the laboratory work 2</i>	4	[2]
7	<i>Preparation for the lecture 5</i>	2	[1]
8	<i>Preparation for the lecture 6</i>	2	[1]
9	<i>Preparation for the lecture 7</i>	2	[1]
10	<i>Preparation for the lecture 8</i>	2	[1]
11	<i>Preparation for the lecture 9</i>	2	[1]
12	<i>Preparation for lecture 10</i>	2	[1]
13	<i>Preparation for lecture 11</i>	2	[1]
14	<i>Preparation for lecture 12</i>	2	[1], [1 add.]
15	<i>Preparation for the laboratory work 3</i>	4	[2]
16	<i>Preparation for lecture 13</i>	2	[3 add.], [4 add.]
17	<i>Preparation for lecture 14</i>	2	[5 add.], [6 add.]
18	<i>Preparation for lecture 15</i>	2	[7 add.]
19	<i>Preparation for lecture 16</i>	2	[8 add.]
20	<i>Preparation for lecture 17</i>	2	[9 add.], [10 add.], [11 (add.)]
21	<i>Preparation for modular test</i>	20	[1]

### 7. Course policy

*Forms of organizing the educational process, types of training sessions and assessment of learning outcomes are regulated by the Regulations on the Organization of the Educational Process at the National Technical University of Ukraine "Igor Sikorsky Kyiv Polytechnic Institute".*

**Class attendance policy.** Attendance at lecture classes is mandatory. Attendance at laboratory classes may be occasional and required for the defense of laboratory work. The presence or absence of a student at a class is not assessed by awarding or deducting points. If a student cannot attend classes, he or she is still responsible for studying the theoretical material and completing practical assignments.

**Policy on ethical norms in the classroom:** discipline; compliance with subordination; honesty; responsibility; respect for those present, turning off phones.

**Policy on assessing learning outcomes.** The policy on assessing learning outcomes is regulated by the Regulations on the system of assessing learning outcomes at Igor Sikorsky Kyiv Polytechnic Institute. According to the Regulations, each grade is given in accordance with the criteria developed by the teacher and announced to students in advance. If a student fails to complete all four laboratory tests, he/she will not be allowed to take the test. Failure to pass the current control measure (modular test) without good reason is assessed as 0 points.

The policy and principles of academic integrity are regulated by the norms set of the Code of Honor of the National Technical University of Ukraine "Igor Sikorsky Kyiv Polytechnic Institute" ([https://kpi.ua/files/honorcode\\_en.pdf](https://kpi.ua/files/honorcode_en.pdf)); Regulations on the Organization of the Educational Process, Regulations on the System for Preventing Academic Plagiarism, Regulations on the Commission on Ethics and Academic Integrity. Plagiarism and other forms of violation of the principles of academic integrity are unacceptable. The student must complete all laboratory practical tasks independently using open sources of information and acquired knowledge and skills.

Plagiarism and other forms of violation of the principles of academic integrity are unacceptable. All work for current and semester tests must be completed independently by the student using open sources of information and the acquired knowledge and skills.

All works that violate the principles of academic integrity (the program code does not match the assignment option, the identity of the program code among different works, etc.) are evaluated at 0 points. To gain access to the test, the student must independently complete the laboratory work (without changing the current rating). In the case of semester control work that violates the principles of academic integrity, the semester control report is marked "Eliminate".

Policy on appealing the results of assessment of control measures. According to the "Regulations on resolving conflict situations at Igor Sikorsky Kyiv Polytechnic Institute" ([https://osvita.kpi.ua/sites/default/files/downloads/regulations\\_resolving\\_conflict\\_situations\\_2020.pdf](https://osvita.kpi.ua/sites/default/files/downloads/regulations_resolving_conflict_situations_2020.pdf)), students have the right to appeal the results of control measures with arguments, explaining which criterion they disagree with according to the assessment. A student can raise any issue related to the procedure of control measures and expect that it will be considered in accordance with predetermined procedures.

Policy on assigning incentive points. According to the "Regulations on the system of assessing learning outcomes at Igor Sikorsky Kyiv Polytechnic Institute", incentive points are not included in the main RSO scale, and their sum cannot exceed 10 points.

Incentive points are awarded for a creative approach in performing laboratory work (the maximum number of points for all work is 10 points), as well as for participation in scientific projects and conferences related to the topic of this course.

## 8. Monitoring and grading policy

During the semester, students perform **3 laboratory works**. The maximum number of points for each laboratory work: 20 points.

Points are awarded for:

- quality of performance of the laboratory work: 0-8 points;
- answer during the defense of the laboratory work: 0-8 points;
- timely submission of work for defense: 0-4 points.

Performance evaluation criteria:

- 7-8 points – the work is done qualitatively, in full;
- 5-6 points - the work is done qualitatively, in full, but has shortcomings;
- 3-4 points - the work is done qualitatively, but not in full, has flaws;
- 1-2 points – the work is not done well, not in full, has flaws;
- 0 points – the work is incomplete or contains significant errors.

Answer evaluation criteria:

- 7-8 points – the answer is complete, well-argued;
- 5-6 points – the answer is incomplete, but well argued;
- 3-4 points – there are minor errors in the answer;
- 1-2 points – there are significant errors in the answer;
- 0 points - there is no answer or the answer is incorrect.

Criteria for evaluating the timeliness of work submission for defense:

- 4 points – the work is presented for defense no later than the specified deadline;
- 0-3 points – the work is submitted for defense later than the specified deadline.

The maximum number of points for performing and defending computer practicals:  
20 points × 3 comp. practice = 60 points.

The assignment for **the modular test** consists of 5 questions - 3 theoretical and 2 practical.

The answer to each theoretical/practical question is evaluated by 8 points.

Evaluation criteria for each theoretical/practical question of the modular test:

- 7-8 points – the answer is correct, complete, well-argued;
- 5-6 points – the answer is correct, but incomplete or poorly argued;
- 3-4 points – there are minor errors in the answer;
- 1-2 points – there are significant errors in the answer;
- 0 points - there is no answer or the answer is incorrect.

The maximum number of points for a modular test:

8 points × 3 theoretical questions + 8 points × 2 practical questions = 40 points.

The rating scale for the discipline is equal to:

$R = R_{\text{laboratory work}} + R_{\text{modular test}} = R_S = 60 \text{ points} + 40 \text{ points} = 100 \text{ points}.$

Calendar control: is carried out twice a semester as a monitoring of the current state of fulfillment of the syllabus requirements.

At the first attestation, the student receives "passed" if his current rating is at least 50% of the maximum number of points that the student can receive before the first certification (20 points).

At the second attestation, the student receives "passed" if his current rating is at least 50% of the maximum number of points that the student can receive before the second certification (30 points).

Semester control: **Final test**.

Conditions for admission to semester control:

with a semester rating of not less than 60% (60 points) and the enrollment of all the works of the computer workshop.

Completion and defense of a computer workshop is a necessary condition for admission to the final test.

Table of correspondence of rating points to grades on the university scale:

Score	Grade
100-95	Excellent
94-85	Very good
84-75	Good
74-65	Satisfactory
64-60	Sufficient
Below 60	Fail
Course requirements are not met	Not Graded

## 9. Additional information about the course

The availability of a certificate of completion of a similar course in big data analytics is evaluated as 20 points (if the course topic corresponds to the subject of the laboratory practicum). Writing research articles or participating in conferences/projects in the relevant field of this discipline can also be evaluated as an additional 5 points.

For example, a laboratory practicum may be automatically credited with 20 points upon submission in Classroom of a Coursera certificate corresponding to the practicum's subject. Certificates of completion of the following or similar courses may be used to credit Computer Workshop 1 (20 points):

**Introduction to Big Data with Spark and Hadoop (IBM)** — a course with practical lab exercises on Hadoop, Spark, HDFS, Hive, MapReduce, SparkSQL in a Jupyter/Docker environment.

<https://www.coursera.org/learn/introduction-to-big-data-with-spark-hadoop>

**Machine Learning With Big Data** — overview of machine learning methods for big data using Spark, with practical tasks for scaling ML models.

<https://www.coursera.org/learn/big-data-machine-learning>

**Google Data Analytics Professional Certificate (Google)** — a series of 8 courses covering data analysis to practical projects.

<https://www.coursera.org/professional-certificates/google-data-analytics>

Certificates of completion of the following or similar courses may be used to credit Computer Workshop 2 (20 points):

**NLP: Twitter Sentiment Analysis** — a practical project implementing a Naive Bayes-based classifier for sentiment analysis on large datasets of tweets, including financial messages.

<https://www.coursera.org/projects/twitter-sentiment-analysis>

Certificates of completion of the following or similar courses may be used to credit Computer Workshop 3 (20 points):

**Introduction to Big Data with Spark and Hadoop** — covering both Hadoop and distributed processing with Spark.

<https://www.coursera.org/learn/introduction-to-big-data-with-spark-hadoop>

**Scalable Machine Learning on Big Data using Apache Spark** — course on using Spark for scalable ML tasks on clusters, SparkSQL, SparkML Pipelines.

<https://www.coursera.org/learn/machine-learning-big-data-apache-spark>

*Big Data Analysis with Scala and Spark* — distributed computing with Scala and Spark, including the MapReduce paradigm.

<https://www.coursera.org/learn/scala-spark-big-data>

*Big Data Specialization (UC San Diego)* — specialization covering Hadoop, Spark, MapReduce, and Hive.

<https://www.coursera.org/specializations/big-data>

etc.

*The student must inform the instructor about a completed or planned course and clarify the results and prospects for crediting the acquired learning outcomes obtained through non-formal or informal education.*

### **Syllabus of the course**

**Is designed by teacher** PhD, Associate Professor, Liubov Oleshchenko

**Adopted by Computer Systems Software Department** (protocol № 3, 29.09.2025)

**Approved by the Faculty Board of Methodology** (protocol № 2, 16.10.2025)