



Побудова масштабованих систем обробки даних у реальному часі

Робоча програма навчальної дисципліни (Силабус)

Реквізити навчальної дисципліни

Рівень вищої освіти	<i>Перший (бакалаврський)</i>
Галузь знань	<i>12 Інформаційні технології</i>
Спеціальність	<i>121 Інженерія програмного забезпечення</i>
Освітня програма	<i>Інженерія програмного забезпечення інтелектуальних кібер-фізичних систем в енергетиці</i>
Статус дисципліни	<i>Нормативна</i>
Форма навчання	<i>Очна (денна)</i>
Рік підготовки, семестр	<i>4 курс, весняний семестр</i>
Обсяг дисципліни	<i>4,5 кредитів ЄКТС/135 год (36 год. лекції, 18 год. практичні заняття, 81 год. самостійна робота)</i>
Семестровий контроль/ контрольні заходи	<i>Екзамен, модульна контрольна робота, календарний контроль</i>
Розклад занять	<i>Згідно розкладу на весняний семестр поточного навчального року (rozklad.kpi.ua)</i>
Мова викладання	<i>Українська</i>
Інформація про керівника курсу / викладачів	<i>Лектор: старший викладач, Колумбет В.П., kvplinux@gmail.com, тел. 093-405-35-39 Практичні: старший викладач, Колумбет В.П., kvplinux@gmail.com, тел. 093-405-35-39</i>
Розміщення курсу	<i>Кампус</i>

Програма навчальної дисципліни

1. Опис навчальної дисципліни, її мета, предмет вивчення та результати навчання

Вивчення дисципліни «Побудова масштабованих систем обробки даних у реальному часі» дозволяє сформуванню у здобувачів освіти компетенції, необхідні для розв'язання практичних задач професійної та наукової діяльності, пов'язаної з підготовкою, аналізом великих даних у масштабованих системах обробки даних.

Метою вивчення дисципліни «Побудова масштабованих систем обробки даних у реальному часі» є формування у студентів навичку застосовувати програмні методи та засоби для побудови масштабованих систем, обробки, аналізу великих даних та прийняття управлінських рішень в різних галузях науки та бізнесу.

Предметом дисципліни «Побудова масштабованих систем обробки даних у реальному часі» є загальні принципи та підходи до побудови масштабованих систем, оброблення великих даних,

можливості мов програмування Python та R для побудови систем, оброблення, аналізу та візуалізації великих даних.

Вивчення дисципліни «Побудова масштабованих систем обробки даних у реальному часі» сприяє формуванню у студентів фахової компетентності (ФК) за освітньою програмою:

ФК08 - Здатність застосовувати фундаментальні і міждисциплінарні знання для успішного розв'язання завдань інженерії програмного забезпечення.

Вивчення дисципліни «Побудова масштабованих систем обробки даних у реальному часі» сприяє формуванню у студентів наступних програмних результатів навчання (ПРН) за освітньою програмою:

ПРН01 - Аналізувати, цілеспрямовано шукати і вибирати необхідні для вирішення професійних завдань інформаційно-довідникові ресурси і знання з урахуванням сучасних досягнень науки і техніки.

ПРН07 - Знати і застосовувати на практиці фундаментальні концепції, парадигми і основні принципи функціонування мовних, інструментальних і обчислювальних засобів інженерії програмного забезпечення.

ПРН18 - Знати та вміти застосовувати інформаційні технології обробки, зберігання та передачі даних.

2. Пререквізити та постреквізити дисципліни (місце в структурно-логічній схемі навчання за відповідною освітньою програмою)

Успішному вивченню дисципліни «Побудова масштабованих систем обробки даних у реальному часі» знання отримані при вивченні дисциплін «Програмування» та «Бази даних» навчального плану підготовки бакалаврів за спеціальністю 121 Інженерія програмного забезпечення.

Отримані при засвоєнні дисципліни «Побудова масштабованих систем обробки даних у реальному часі» теоретичні знання та практичні уміння сприяють успішному виконанню курсових проєктів, стартапів та дипломної роботи бакалаврів.

3. Зміст навчальної дисципліни

Дисципліна «Побудова масштабованих систем обробки даних у реальному часі» передбачає вивчення таких тем:

Тема 1. Джерела та типи великих даних. Підготовка та аналіз даних з використанням мови Python.

Тема 2. Архітектурні моделі BigData. Можливості мови R.

Модульна контрольна робота.

Екзамен.

4. Навчальні матеріали та ресурси

Основна література:

1. BIG DATA: Інноваційні можливості підвищення прибутковості агробізнесу / [Електронний ресурс] – Режим доступу: <http://www.agrobusiness.com.ua/ideii-i-trendy/8383-bigdata-innovatsiini-mozhlyvostipidvyschennia-prybutkovostiagrobiznesu.html>
2. IoT Analytics Platform / Електронний ресурс. Режим доступу: <https://blog.codecentric.de/en/2016/07/iot-analytics-platform/>
3. IoT Fundamentals: Big Data & Analytics / Електронний ресурс. Режим доступу: <https://www.netacad.com/courses/iot/big-data-analytics>

4. Akka // Електронний ресурс. Режим доступу: <https://akka.io/>
5. A brief introduction to two data processing architectures — Lambda and Kappa for Big Data / Електронний ресурс. Режим доступу: <https://towardsdatascience.com/a-briefintroduction-to-two-data-processing-architectures-lambda-and-kappa-for-big-data4f35c28005bb>
6. What Is Apache Hadoop? / [Електронний ресурс] – Режим доступу: <http://hadoop.apache.org/>
7. Apache Kafka. A distributed streaming platform / [Електронний ресурс] – Режим доступу: <https://kafka.apache.org/>
8. Apache Cassandra / [Електронний ресурс] – Режим доступу: <http://cassandra.apache.org/>
9. Apache Spark. A fast and general engine for large-scale data processing / [Електронний ресурс] – Режим доступу: <https://spark.apache.org>
10. IoT-Analyse-Plattform: Floating Bus Data / Електронний ресурс. Режим доступу: https://www.youtube.com/watch?v=VYxc-3ZRRL4&ab_channel=codacentricAG.

Додаткова література:

1. MapReduce. Tutorial. Електронний ресурс. Режим доступу: https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html
2. Pandas. Електронний ресурс. Режим доступу: https://www.w3schools.com/python/pandas/pandas_intro.asp
3. Pandas documentation: веб-сайт. URL: <https://pandas.pydata.org/docs/>
4. Parcel documentation: веб-сайт. URL: <https://parceljs.org/docs/>
5. Python documentation: веб-сайт. URL: <https://docs.python.org/3/>
6. Slatkin B. *Effective Python: 90 Specific Ways to Write Better Python*. Addison-Wesley Professional, 2019. 480 p.
7. *Programming with Databases – Python*. Електронний ресурс. Режим доступу: <https://swcarpentry.github.io/sql-novice-survey/10-prog/index.html>
8. Мова програмування R. Електронний ресурс. Режим доступу: <https://coderlessons.com/tutorials/mashinnoe-obuchenie/r-programmirovaniye/r-programmirovaniye>

Матеріали знаходяться у вільному доступі в Інтернеті.

Навчальний контент

5. Методика опанування навчальної дисципліни (освітнього компонента)

№ з/п	Тип навчального заняття	Опис навчального заняття
<i>Тема 1. Джерела та типи великих даних. Підготовка та аналіз даних з використанням мови Python.</i>		
1	<i>Лекція 1. Джерела великих даних. Інтернет Речей. Визначення BigData.</i>	<i>Розвиток Інтернету Речей та збільшення обсягу інформації. Використання платформи Kaggle та DrivenData для роботи з великими обсягами даних. Аналіз визначення великих даних. Наведення прикладів використання великих обсягів даних у реальних сценаріях. Розгляд відкритих та приватних даних, структурованих та неструктурованих. Вивчення хмарних та туманних обчислень. Розгляд рухомих та нерухомих даних. Огляд інфраструктури для обробки великих обсягів даних. Аналіз розподілених даних та їх обробка.</i>

2	<p>Комп'ютерний практикум 1. Аналіз джерел відкритих даних.</p> <p>Завантаження та даних збереження в форматі csv.</p>	<p>Мета: аналіз джерел відкритих даних через Open Government Partnership та платформи, що надають відкриті дані, для можливостей збереження та візуалізації інформації. Використання web-сайтів www.knoema.com та www.gartinder.org для вивчення власності персональних даних у випадках їх відсутності на локальному рівні та обмежень електронних таблиць під час завантаження інформації.</p>
3	<p>Лекція 2. Розроблення програмного забезпечення для аналізу web-сайтів, які надають відкриті дані за допомогою Python Pandas.</p> <p>Відкриті дані, їх формати та засоби обробки.</p>	<p>Можливості інструментів для обробки даних. Роль мови програмування Python у проведенні аналізу інформації. Порівняння традиційного підходу до обробки великих обсягів даних і нового покоління аналітики. Загальний цикл життя процесу аналізу даних. Огляд відкритих даних, їх різноманітних форматів та інструментів для їх обробки. Вивчення технік web-скрапінгу та процесів вилучення, трансформації і завантаження даних.</p>
4	<p>Комп'ютерний практикум 2. Аналіз та візуалізація даних у Python.</p>	<p>Мета: продемонструвати володіння життєвим циклом аналізу даних через використання наданого набору інформації та відзначених інструментів. Здійснити імпорт пакетів Python, необхідних для обробки даного набору, та використати інструменти Python та Jupyter для підготовки даних до аналізу. Провести аналіз та побудувати графіки на основі отриманих результатів.</p>
5	<p>Лекція 3. Форматування даних про час та дату, читання та запис файлів в Python.</p>	<p>Опрацювання даних часу та дати у мові програмування Python. Робота з читанням та записом файлів у середовищі Python. Взаємодія з зовнішніми програмами в рамках Python.</p>
6	<p>Лекція 4. Програмування Python та SQLite.</p> <p>Призначення утиліти csvsql.</p>	<p>Робота Python з SQLite. Основні операції SQL. Призначення утиліти csvsql. Метод execute().</p>
7	<p>Лекція 5. Процедура імпорту даних із файлів у Pandas.</p> <p>Імпорт даних з Інтернету за допомогою Pandas.</p> <p>Засоби для кореляційного аналізу в Pandas.</p>	<p>Використання статистичних методів у великоаналітичних завданнях та їх реалізація за допомогою бібліотеки Pandas у Python. Здійснення імпорту даних з різних джерел, таких як файли та Інтернет, у середовище Pandas. Використання Pandas для проведення описової статистики та інструментів для кореляційного аналізу даних.</p>
8	<p>Комп'ютерний практикум 3. Кореляційний аналіз у Python.</p>	<p>Мета: демонстрація вмінь проведення кореляційного аналізу даних з використанням вказаного набору та інструментів Python для розрахунку ступеня кореляції. Необхідно налаштувати набір даних, визначити, чи існує кореляція між змінними в даному наборі, використовуючи Python для обчислення кореляції між двома наборами змінних, та створити візуальне представлення результатів аналізу.</p>

9	Лекція 6. Оброблення відсутніх даних. Перетворення типів даних та маніпулювання дата фреймами.	Перетворення типів даних. Оброблення відсутніх даних. Маніпулювання дата фреймами у Python.
10	Лекція 7. Регресійний аналіз даних в Python.	Дослідження регресійного аналізу та його різновидів. Використання регресійного аналізу у здійсненні аналізу інформації.
11	Комп'ютерний практикум 4. Побудова лінійної регресії в Python.	Мета: освоїти концепції лінійної регресії та роботи з даними для прогнозування в середовищі Python. Аналіз запропонованих даних щодо обсягів продажів, та побудова лінійної регресії для прогнозування річного чистого обсягу продажів на основі кількості магазинів у даному районі.
12	Лекція 8. Помилки в аналізі даних та прогностичній аналітиці. Оцінка помилок регресії засобами Python.	Недоліки у виконанні аналізу даних та прогностичній аналітиці. Визначення та оцінка помилок у регресійному аналізі за допомогою інструментів Python. Роль бібліотеки scikit-learn у цьому контексті.
13	Лекція 9. Алгоритми класифікації даних. Застосування класифікацій.	Труднощі у сфері класифікації. Методи класифікації. Графічне зображення результатів класифікації. Використання та перевірка правильності класифікацій. Огляд моделі класифікатора на основі дерева рішень.
14	Лекція 10. Бібліотека Pyplot. Засіб Plotly. Типи візуалізації даних. Використання бібліотек Folium та Leaflet.js для побудови карт. Візуалізація аномалій.	Бібліотека Pyplot. Засіб Plotly. Різноманітність методів візуалізації даних. Графічне відображення аномалій. Використання Folium та Leaflet.js для створення картографічних візуалізацій.
15	Комп'ютерний практикум 5.	Мета: освоїти бібліотеку Pyplot та Засіб Plotly, моделей класифікатора на основі дерева рішень.
<i>Тема 2. Архітектурні моделі BigData. Можливості мови R.</i>		
16	Лекція 11. Аналіз даних в R. Фактори, списки, фрейми та дії над ними.	Еволюція мови програмування R. Функціональні можливості R. Робота з об'єктами, пакетами та функціями. Операції з векторами, матрицями у мові R. Робота із факторами, списками, фреймами та їх обробка.
17	Комп'ютерний практикум 6.	Мета: освоїти функціональні можливості R, роботу з об'єктами, пакетами та функціями.
18	Лекція 12. Експорт, імпорт та оброблення даних в R.	Використання R для аналізу часових рядів. Експорт та імпорт даних в R. Оброблення даних в R.
19	Лекція 13. Основні інструменти аналізу та візуалізації даних в R.	Функція "plot" і її параметри. Типи графіків в R. Управління загальними параметрами - аргументами графічних функцій.

20	Комп'ютерний практикум 7. Аналіз та візуалізація даних в R.	Мета: вивчити функціональність мови програмування R для проведення аналізу та візуалізації даних. Використати бібліотеку <code>dplyr</code> у R для оптимізації та трансформації даних, а також використати бібліотеку <code>ggplot2</code> для візуалізації отриманих результатів.
21	Лекція 14. Архітектурні моделі Big Data. Технології віртуалізації. Гіпервізори. Контейнерна технологія виконання програмного коду на сервері. SaaS, PaaS і IaaS.	Моделі архітектури в області Big Data Engineering. Віртуалізаційні технології та рівні абстракції. Гіпервізори. Використання контейнерної технології для виконання програмного коду на сервері. Процес інжинірингу даних.
22	Лекція 15. Технології Hadoop Big Data. Розподілена обробка MapReduce. HDFS.	Масштабування за рахунок використання великих обсягів даних. Методи зберігання та оброблення інформації у розподілених файлових системах. Використання розподілених баз даних. Використання розподіленої файлової системи Hadoop (HDFS).
23	Комп'ютерний практикум 8.	Мета: освоїти технології Hadoop, використання контейнерних технологій.
24	Лекція 16. Розподілена потокова платформа Kafka. Переваги Cassandra.	Проблема прийому даних. Переваги Cassandra . Розподілена потокова платформа Kafka.
25	Лекція 17. Платформа Apache Spark. Lambda та Карра архітектури оброблення великих даних.	Проблема обчислювальної функції та її вирішення за допомогою технології Spark. Порівняння Spark та MapReduce. Використання Spark та <code>sparklyr</code> для обробки великих обсягів даних в мові програмування R. Архітектура Lambda та її переваги та недоліки. Також, архітектура Карра та аналіз її переваг і недоліків.
26	Комп'ютерний практикум 9. Розподілені обчислення даних з використанням Spark кластера та мови R.	Мета: встановити Spark на власному комп'ютері та провести розподілені обчислення для конкретного набору даних, використовуючи Spark-кластер та мову програмування R.
27	Лекція 18. Підсумкове лекційне заняття. Екзамен.	Повторення вивченого матеріалу. Екзаменаційна контрольна робота.

6. Самостійна робота студента

Дисципліна «Побудова масштабованих систем обробки даних у реальному часі» ґрунтується на самостійній підготовці до аудиторних занять на теоретичні та практичні теми.

№ з/п	Назва теми, що виноситься на самостійне опрацювання	Кількість годин
1	Підготовка до лекції 1	2
2	Підготовка до комп'ютерного практикуму 1	2

3	<i>Підготовка до лекції 2</i>	2
4	<i>Підготовка до комп'ютерного практикуму 2</i>	2
5	<i>Підготовка до лекції 3</i>	2
6	<i>Підготовка до лекції 4</i>	2
7	<i>Підготовка до лекції 5</i>	2
8	<i>Підготовка до комп'ютерного практикуму 3</i>	2
9	<i>Підготовка до лекції 6</i>	2
10	<i>Підготовка до лекції 7</i>	2
11	<i>Підготовка до комп'ютерного практикуму 4</i>	2
12	<i>Підготовка до лекції 8</i>	2
13	<i>Підготовка до лекції 9</i>	2
14	<i>Підготовка до лекції 10</i>	2
15	<i>Підготовка до комп'ютерного практикуму 5</i>	
16	<i>Підготовка до лекції 11</i>	2
17	<i>Підготовка до лекції 12</i>	2
18	<i>Підготовка до комп'ютерного практикуму 6</i>	
19	<i>Підготовка до лекції 13</i>	2
20	<i>Підготовка до комп'ютерного практикуму 7</i>	2
21	<i>Підготовка до лекції 14</i>	2
22	<i>Підготовка до лекції 15</i>	2
23	<i>Підготовка до лекції 16</i>	2
25	<i>Підготовка до комп'ютерного практикуму 8</i>	
26	<i>Підготовка до лекції 17</i>	2
27	<i>Підготовка до комп'ютерного практикуму 9</i>	2
28	<i>Підготовка до модульної контрольної роботи</i>	27

Політика та контроль

7. Політика навчальної дисципліни (освітнього компонента)

Відвідування лекційних занять є обов'язковим.

- Відвідування занять комп'ютерного практикуму може бути епізодичним та за потреби захисту робіт комп'ютерного практикуму.*

- *Правила поведінки на заняттях: активність, повага до присутніх, відключення телефонів.*
- *Дотримання політики академічної доброчесності.*

8. Види контролю та рейтингова система оцінювання результатів навчання (PCO)

Протягом семестру студенти виконують 9 комп'ютерних практикумів. Максимальна кількість балів за кожний комп'ютерний практикум: 7 балів.

Бали нараховуються за:

- *якість виконання комп'ютерного практикуму: 0-4 бали;*
- *відповідь під час захисту комп'ютерного практикуму: 0-4 бали; - своєчасне представлення роботи до захисту: 0-2 бали.*

Критерії оцінювання якості виконання:

- 3-4 бали – робота виконана якісно, в повному обсязі;*
- 1-2 бали – робота виконана якісно, в повному обсязі, але має недоліки; 0 балів – робота виконана не в повному обсязі, або містить суттєві помилки.*

Критерії оцінювання відповіді:

- 3-4 бали – відповідь повна, добре аргументована;*
- 1-2 бал – у відповіді є суттєві помилки;*
- 0 балів – немає відповіді або відповідь невірна.*

Критерії оцінювання своєчасності представлення роботи до захисту:

- 1-2 бали – робота представлена до захисту не пізніше вказаного терміну; 0 балів – робота представлена до захисту пізніше вказаного терміну.*

Максимальна кількість балів за виконання та захист комп'ютерних практикумів: 7 балів × 9 комп. практ. = 63 бали.

*Завдання на **модульну контрольну роботу** складається з 7 питань – 6 теоретичних та 1 практичного. Відповідь на кожне теоретичне запитання оцінюється 5 балами, практичне - 7. Критерії оцінювання кожного теоретичного/практичного запитання модульної контрольної роботи:*

- 5(7) балів – відповідь вірна, повна, добре аргументована;*
- 3-4(5-6) балів – відповідь вірна, але неповна або погано аргументована;*
- 2(3-4) бали – у відповіді є незначні помилки;*
- 1(1-2) бал – у відповіді є суттєві помилки;*
- 0 балів – немає відповіді або відповідь невірна.*

Максимальна кількість балів за модульну контрольну роботу:

5 балів × 5 теоретичні запитання + 7 балів × 1 практичні запитання = 37 балів.

Рейтингова шкала з дисципліни дорівнює:

$R = R_C = 63 \text{ балів} + 37 \text{ балів} = 100 \text{ балів.}$

За описом: $R = R_{\text{комп.практ}} + R_{\text{МКР}} = 63 + 37 \text{ балів} = 100 \text{ балів}$

Календарний контроль: провадиться двічі на семестр як моніторинг поточного стану виконання вимог силабусу.

На першій атестації (8-й тиждень) студент отримує «зараховано», якщо його поточний рейтинг не менше 50 % від максимальної кількості балів, яку може отримати студент до першої атестації (20 балів).

На другій атестації (14-й тиждень) студент отримує «зараховано», якщо його поточний рейтинг не менше 50 % від максимальної кількості балів, яку може отримати студент до другої атестації (30 балів).

Семестровий контроль: **екзамен**.

Умови допуску до семестрового контролю:

При семестровому рейтингу (r_c) не менше 60 % (60 балів) та зарахуванні усіх робіт комп'ютерного практикуму.

Необхідною умовою допуску до екзамену є виконання і захист комп'ютерного практикуму.

Таблиця відповідності рейтингових балів оцінкам за університетською шкалою:

Кількість балів	Оцінка
100-95	Відмінно
94-85	Дуже добре
84-75	Добре
74-65	Задовільно
64-60	Достатньо
Менше 60	Незадовільно
Не виконані умови допуску	Не допущено

1. Додаткова інформація з дисципліни (освітнього компонента)

Сертифікати проходження онлайн курсів «IoT Fundamentals: Big Data & Analytics» та «Programming Essentials in Python» дозволяють зарахувати модульну контрольну роботу згідно отриманої загальної оцінки в Мережевій академії (Cisco Networking Academy), написання статей або участь у конференціях та проєктах за відповідною тематикою також оцінюється як додаткові 5 балів.

Складено старшим викладачем ІПЗЕ, Колумбетом Вадимом Петровичем

Ухвалено кафедрою інженерії програмного забезпечення в енергетиці НН ІАТЕ (протокол № 34 від 10.05.2024 р.)

Погоджено Методичною комісією НН ІАТЕ (протокол № 9 від 31.05.2024 р.)