



МЕТОДИ ДОБУВАННЯ ДАНИХ

Робоча програма навчальної дисципліни (Силабус)

Реквізити навчальної дисципліни

Рівень вищої освіти	<i>Третій (доктор філософії)</i>
Галузь знань	<i>12 Інформаційні технології</i>
Спеціальність	<i>121 Інженерія програмного забезпечення, 123 Комп'ютерна інженерія</i>
Освітня програма	<i>Комп'ютерна інженерія</i>
Статус дисципліни	<i>Вибіркова</i>
Форма навчання	<i>очна(денна)</i>
Рік підготовки, семестр	<i>2 курс, осінній семестр</i>
Обсяг дисципліни	<i>5 кредити</i>
Семестровий контроль/ контрольні заходи	<i>Залік</i>
Розклад занять	<i>Лекцій 18 (36 годин) Лабораторних 9 (18 годин)</i>
Мова викладання	<i>Українська</i>
Інформація про керівника курсу / викладачів	Лектор: д.т.н, професор, Новотарський Михайло Анатолійович novotar@gmail.com Лабораторні: д.т.н, професор, Новотарський Михайло Анатолійович novotar@gmail.com
Розміщення курсу	https://drive.google.com/drive/folders/1uRhR_udj3TYnS1mHXqfJ83_q4zNON7kh

Програма навчальної дисципліни

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

Дисципліна “Методи добування даних” спрямована на вивчення основних методів, які дозволять дослідити та визначити приховані функціональні залежності в великих наборах даних. В рамках вивчення даної дисципліни також розглядаються основні технологічні рішення у сфері науки про дані (Data Science), які забезпечують пошук прихованих невідомих раніше корисних та доступних для інтерпретації знань, необхідних для використання при прийнятті рішень у різних сферах діяльності.

Вивчення даної дисципліни майбутніми науковцями дозволить їм набути важливих компетенцій, оскільки, в зв'язку з накопиченням значних обсягів даних у різних сферах діяльності людини, виникає широкий спектр задач з добування корисних даних, на основі яких отримують нові знання, що не можуть бути одержані у інший спосіб.

Метою вивчення дисципліни «Методи добування даних» є підготовка фахівців, здатних розв'язувати комплексні проблеми в галузі дослідницько-інноваційної діяльності у сфері представлення та опрацювання даних для отримання прихованої в даних корисної інформації шляхом вивчення теоретичних та практичних положень побудови моделей прихованих в наборах даних функціональних залежностей та використання результатів аналізу даних для уточнення

наукових висновків та формування прогнозів щодо майбутніх станів об'єктів досліджень, що передбачає глибоке переосмислення наявних та створення нових цілісних знань.

Предметом дисципліни є:

- методи та способи представлення та опрацювання даних для отримання прихованої в даних корисної інформації;
- методи побудови та дослідження математичних моделей опису функціональних залежностей між параметрами даних та шуканою інформацією;
- технології адаптивних та інтелектуальних обчислень при добуванні даних для розроблення прогнозів щодо досліджуваних об'єктів та процесів.

Основні результати навчання

Здобувачі наукового ступеня доктора філософії після засвоєння навчальної дисципліни мають продемонструвати такі **компетентності**.

1. *Загальні компетентності*: здатність до абстрактного мислення, аналізу і синтезу; здатність до пошуку, оброблення та аналізу інформації з різних джерел щодо методів та технологій добування даних; здатність формування системного наукового світогляду у сфері отримання прихованої в даних корисної інформації; здатність набуття універсальних навичок усної та письмової презентації власного наукового дослідження за результатами аналізу даних, застосування сучасних інформаційних технологій у науковій діяльності.

2. *Спеціальні компетентності*: здатність ефективно застосовувати основні методи добування даних, основні статистичні методи навчання, кластерні методи машинного навчання та методи навчання без учителя при проведенні наукових досліджень; здатність інтегрувати знання з різних дисциплін, застосовувати системний підхід при проведенні досліджень; здатність аргументувати вибір методу розв'язання наукової задачі, критично оцінювати отримані результати та захищати прийняті рішення щодо вибору методів та способів обробки даних.

За результатами вивчення навчальної дисципліни «Методи добування даних» мають бути отримані такі **знання**.

1. Мати передові концептуальні та методологічні знання у сфері аналізу великих обсягів даних для отримання нових знань, а також дослідницькі навички, достатні для проведення наукових і прикладних досліджень в галузі добування даних на рівні останніх світових досягнень з комп'ютерної інженерії, IT-інфраструктур, інформаційних технологій.

2. Знати сучасні методи проведення досліджень у сфері представлення та опрацювання даних для отримання прихованої в даних корисної інформації шляхом вивчення теоретичних та практичних положень побудови моделей прихованих в наборах даних функціональних залежностей та використання результатів аналізу даних для уточнення наукових висновків та передбачення майбутніх станів об'єктів досліджень.

3. Знати і розуміти наукові і математичні положення, що лежать в основі методів опрацювання даних для отримання прихованої корисної інформації, методів побудови та дослідження математичних моделей та технологій адаптивних та інтелектуальних обчислень при добуванні даних.

Уміння, які мають бути отримані у рамках вивчення навчальної дисципліни «Методи добування даних».

1. Вміти ефективно здійснювати пошук та критичний аналіз інформації з різних джерел щодо методів та технологій добування даних.

2. Вміти розв'язувати задачі синтезу та аналізу об'єктів дослідження при опрацюванні даних для отримання прихованої корисної інформації.

3. Вміти створювати та реалізовувати математичні моделі опису функціональних залежностей, що приховані у даних;

4. Вміти застосовувати технології адаптивних та інтелектуальних обчислень при добуванні даних для передбачення майбутніх станів досліджуваних об'єктів та процесів.

Здобувачі наукового ступеня також мають бути **здатні**.

1. Застосовувати прикладні бібліотеки та програмні системи, які використовуються при машинному аналізі даних.

2. Володіти методами та технологіями програмування з використанням прикладних бібліотек та програмних систем, призначених для машинного аналізу великих наборів даних.

Таке поєднання загальних та спеціальних компетентностей, теоретичних та практичних знань, умінь та здатностей сприяє підвищенню науково-практичного рівня здобувачів наукового ступеня доктора філософії задля здійснення ними ефективних наукових досліджень.

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

Для успішного оволодіння дисципліною необхідні знання:

- основ математичного аналізу, теорії функцій та математичної статистики;
- основ функціонування операційних систем;
- основ програмування мовою Python.

Відповідно до освітньої програми необхідно попередньо оволодіти знаннями з дисциплін: “Програмування”, “Об’єктно-орієнтоване програмування”, “Системне програмування”, “Структури даних та алгоритми”, “Інженерія програмного забезпечення”, “Алгоритми та методи обчислень”, “Дискретна математика”.

Компетентності, знання та вміння, отримані в рамках вивчення даної дисципліни, можуть бути застосовані для отримання обґрунтованих результатів досліджень та підвищення наукового рівня дисертаційних робіт.

3. Зміст навчальної дисципліни

Розділ 1. Робота Pandas у середовищі Jupyter Notebook інтегратора Anaconda

Тема 1.1. Установка програмного забезпечення

Тема 1.2. Запуск Jupyter Notebook. Приклади роботи

Тема 1.3. Основні елементи інтерфейсу

Тема 1.4. Робота з бібліотекою Pandas

Розділ 2. Основні принципи роботи NumPy, SciPy, Matplotlib

Тема 2.1. Принципи роботи з NumPy

Тема 2.2. Робота з SciPy

Тема 2.3. Використання бібліотеки Matplotlib

Розділ 3. Основи бібліотеки Pandas

Тема 3.1. Структури даних Pandas

Тема 3.2. Робота з елементами Series

Тема 3.3. Робота з елементами DataFrame

Тема 3.4. Доступ до даних у структурах Pandas

Тема 3.5. Використання атрибутів для доступу до даних в Pandas

Розділ 4. Отримання випадкового набору з структур Pandas

Тема 4.1. Робота з випадковим набором даних

Тема 4.2. Додавання елементів у структури

Тема 4.3. Індексція з використанням логічних виразів

Тема 4.4. Використання isin для роботи з даними в Pandas

Тема 4.5. Робота з пропусками в даних

Розділ 5. Загальне поняття про методи аналізу та добування даних

Тема 5.1. Огляд методів та засобів.

Тема 5.2. Статистичні обмеження для аналізу даних

Тема 5.3. Основні інструменти та базові теоретичні положення

Тема 5.4. Ефект Метью

Розділ 6. Пошук подібних об'єктів

Тема 6.1. Застосування методу пошуку найближчих сусідів

Тема 6.2. Шинглінг документів

Тема 6.3. Збереження схожості скорочених наборів

Тема 6.4. Локально-чутливе хешування для документів

Розділ 7. Міри відстані

Тема 7.1. Визначення міри відстані

Тема 7.2. Відстань Евкліда

Тема 7.3. Відстань Жакарта

Тема 7.4. Косинусна відстань

Тема 7.5. Відстань редагування

Тема 7.6. Відстань Хеммінга

Розділ 8. Локально-чутливі функції

Тема 8.1. Теорія локально-чутливих функцій

Тема 8.2. LSH сімейства для інших мір відстані

Розділ 9. Застосування локально-чутливого хешування

Тема 9.1. Приклади розпізнавання об'єктів

Тема 9.2. Методи для високих ступенів подібності

Розділ 10. Аналіз потоків даних

Тема 10.1. Модель поточкових даних

Тема 10.2. Вибір даних у потоці

Тема 10.3. Фільтрація потоків

Тема 10.4. Підрахування різних елементів у потоці

Розділ 11. Оцінка моментів

Тема 11.1. Визначення моментів та практичні приклади

Тема 11.2. Підрахунки у вікні

Тема 11.3. Згасання вікон

Розділ 12. Аналіз посилань

Тема 12.1. Класифікація сторінок. PageRank

Тема 12.2. Ефективне обчислення PageRank

Розділ 13. Покращення при обчисленні Page Rank

Тема 13.1. Тематично чутливий PageRank

Тема 13.2. «Лінковий» спам

Тема 13.3. Хаби та авторитети

Розділ 14. Набори елементів, які часто зустрічаються (часті набори елементів)

Тема 14.1. Модель ринкова корзина

Тема 14.2. Ринкові корзини і алгоритм A-Priori

Розділ 15. Обробка великих наборів даних у основній пам'яті

Тема 15.1. Основні алгоритми обробки даних в пам'яті

Тема 15.2. Алгоритми з обмеженням проходів

Тема 15.3. Підрахунок частих елементів у потоці

4. Навчальні матеріали та ресурси

Базова:

1. Новотарський М.А. Методи добування даних, навч. посіб. для спеціальності 123 – комп'ютерна інженерія // <https://cloud.comsys.kpi.ua/s/YpEaokWx3HjcMmZ>
2. Новотарський М.А. Методичні вказівки до лабораторних робіт з курсу «Методи добування даних» для спеціальності 123 – комп'ютерна інженерія // <https://cloud.comsys.kpi.ua/s/YpEaokWx3HjcMmZ>
3. Han J., Kamber M. Data Mining Concepts and Techniques. – Morgan Kaufmann, 2011. –626 p.
4. Larose D.T. Data Mining Methods and Models. – NY: Wiley-IEEE Press, 2006. –344 p.
5. Rajaraman A., Ullman D.J. Mining of Massive Datasets. – Cambridge: Cambridge University Press, 2011. – 326 p.
6. Lawrence K.D., Kudyba S., Klimberg R.K. Data Mining Methods and Applications. – Auerbach Publications, 2008. –366 p.
7. Гладун А.Я., Рогушина Ю.Ф. Data Mining: пошук знань в даних. – К.: ТОВ «ВД « АДЕФ-Україна», 2016. –452 с.

Додаткова:

8. Runkowski L., Jaworski M., Duda P. Stream Data Mining: Algorithms and Their Probabilistic Properties – Springer, 2020. – 504 p.
9. Olson D.L., Delen D. Advanced Data Mining Techniques. – Springer, 2008. – 192 p.
10. Witten I.H., Frank E. Data Mining: Practical Machine Learning Tools and Techniques. – Morgan Kaufmann, 2005. – 560 p.
11. Nisbet R., Elder J., Miner G. Handbook of Statistical Analysis. Data Mining Applications. –Amsterdam: Elsevier, 2009 – 864 p.
12. Fu X., Wang L. Data Mining with Computational Intelligence. – Springer, 2005. – 287 p.
13. Черняк О.І., Захарченко П.В. Інтелектуальний аналіз даних. – К.: Знання, 2014. – 599 с.

14. Методика опанування навчальної дисципліни (освітнього компонента)

Назви розділів, тем	Кількість годин			
	Всього	У тому числі		
		Лекції	Практичні роботи	СРС
<p>1. Робота у середовищі Jupyter Notebook інтегратора Anaconda</p> <p>Тема 1.1. Установка програмного забезпечення.</p> <p>Тема 1.2. Запуск Jupyter Notebook. Приклади роботи.</p> <p>Тема 1.3. Основні елементи інтерфейсу.</p> <p>Тема 1.4. Робота з бібліотекою Pandas.</p>	16	0	6	10
<p>2. Основні принципи роботи NumPy, SciPy, Matplotlib</p> <p>Тема 2.1. Принципи роботи з NumPy.</p> <p>Тема 2.2. Робота з SciPy.</p> <p>Тема 2.3. Використання бібліотеки Matplotlib.</p>	10	2	2	6
<p>3. Основи бібліотеки Pandas</p> <p>Тема 3.1. Структури даних Pandas.</p> <p>Тема 3.2. Робота з елементами Series.</p> <p>Тема 3.3. Робота з елементами DataFrame.</p> <p>Тема 3.4. Доступ до даних у структурах Pandas.</p> <p>Тема 3.5. Використання атрибутів для доступу до даних в Pandas.</p>	12	4	4	4
<p>4. Отримання випадкового набору з структур Pandas</p> <p>Тема 4.1. Робота з випадковим набором даних.</p> <p>Тема 4.2. Додавання елементів у структури.</p> <p>Тема 4.3. Індексація з використанням логічних виразів.</p> <p>Тема 4.4. Використання isin для роботи з даними в Pandas.</p> <p>Тема 4.5. Робота з пропусками в даних.</p>	14	0	4	10
<p>5. Загальне поняття про методи аналізу та добування даних</p> <p>Тема 5.1. Огляд методів та засобів. Статистичне моделювання. Машинне навчання. Обчислювальні підходи до моделювання. Спосіб узагальнення даних. Виділення ознак.</p> <p>Тема 5.2. Статистичні обмеження для аналізу даних.</p> <p>Загальна інформативність. Принцип Бонферроні. Приклад застосування принципу Бонферроні.</p>	6	4	0	2

Назви розділів, тем	Кількість годин			
	Всього	У тому числі		
		Лекції	Практичні роботи	СРС
Тема 5.3. Основні інструменти та базові теоретичні положення. Важливість слів у документах. Хеш-функції. Індеси. Зовнішня пам'ять. Основа натуральних логарифмів. Степенева залежність. Тема 5.4. Ефект Метью. Висновки за розділом.				
6. Пошук подібних об'єктів Тема 6.1. Застосування методу пошуку найближчих сусідів. Тема 6.2. Шинглінг документів. Тема 6.3. Збереження схожості скорочених наборів. Тема 6.4. Локально-чутливе хешування для документів.	6	4	0	2
7. Міри відстані Тема 7.1. Визначення міри відстані. Тема 7.2. Відстань Евкліда. Тема 7.3. Відстань Жакарта. Тема 3.4. Косинусна відстань. Тема 3.5. Відстань редагування. Тема 3.6. Відстань Хеммінга.	8	4	0	4
8. Локально-чутливі функції Тема 8.1. Теорія локально-чутливих функцій. Тема 8.2. LSH сімейства для інших мір відстані.	4	2	0	2
9. Застосування локально-чутливого хешування Тема 9.1. Приклади розпізнавання об'єктів. Тема 9.2. Методи для високих ступенів подібності.	4	2	0	2
10. Аналіз потоків даних Тема 10.1. Модель поточкових даних. Тема 10.2. Вибір даних у потоці. Тема 10.3. Фільтрація потоків. Тема 10.4. Підрахування різних елементів у потоці.	6	2	2	2
11. Оцінка моментів Тема 11.1. Визначення моментів та практичні приклади. Тема 11.2. Підрахунки у вікні. Тема 11.3. Згасання вікон.	8	2	0	6
12. Аналіз посилань Тема 12.1. Класифікація стрінок. PageRank. Тема 12.2. Ефективне обчислення PageRank.	8	2	0	6
13. Покращення при обчисленні Page Rank Тема 13.1. Тематично чутливий PageRank. Тема 13.2. «Лінковий» спам. Тема 13.3. Хаби та авторитети.	10	2	0	8
14. Набори елементів, які часто зустрічаються (часті набори елементів) Тема 14.1. Модель ринкова корзина.	10	2	0	8

Назви розділів, тем	Кількість годин			
	Всього	У тому числі		
		Лекції	Практичні роботи	СРС
Тема 14.2. Ринкові корзини і алгоритм A-Priori.				
15. Обробка великих наборів даних у основній пам'яті				
Тема 15.1. Основні алгоритми обробки даних в пам'яті. Тема 15.2. Алгоритми з обмеженням проходів. Тема 15.3. Підрахунок частих елементів у потоці.	8	4	0	4
Залік	2	0	0	0
Всього в семестрі:	132	36	18	76

15. Самостійна робота аспіранта

Метою проведення циклу лабораторних робіт є набуття студентами необхідних практичних навичок використання методів та способів представлення та опрацювання даних для отримання прихованої інформації, методів дослідження математичних моделей опису шуканої інформації, технології добування даних для розроблення прогнозів щодо досліджуваних об'єктів та процесів.

№ з/п	Назва лабораторної роботи	Кількість ауд. годин
1	<i>Лабораторна робота № 1.</i> Створення зошита в середовищі Jupyter notebook.	2
2	<i>Лабораторна робота № 2.</i> Практичне застосування бібліотек NumPy, SciPy та Matplotlib.	2
3	<i>Лабораторна робота № 3.</i> Індексування, вибір, редагування набору даних.	2
4	<i>Лабораторна робота № 4.</i> Підсумкові функції та відображення.	3
5	<i>Лабораторна робота № 5.</i> Робота з пропущеними даними.	3
6	<i>Лабораторна робота № 6.</i> Розв'язування задачі лінійної регресії.	3
7	<i>Лабораторна робота № 7.</i> Розв'язування задачі кластеризації.	3
	Всього:	18

Політика та контроль

16. Політика навчальної дисципліни (освітнього компонента)

Під час занять з навчальної дисципліни «Методи добування даних» аспіранти повинні дотримуватись певних дисциплінарних правил:

- забороняється запізнюватись на заняття;

- при вході викладача, на знак привітання, особи, які навчаються в КПІ ім. Ігоря Сікорського повинні встати;
- не допускаються сторонні розмови або інший шум, що заважає проведенню занять;
- виходити з аудиторії під час заняття допускається лише з дозволу викладача.
- не допускається користування мобільними телефонами та іншими технічними засобами без дозволу викладача.

Лабораторні роботи здаються особисто з попередньою перевіркою теоретичних знань, які необхідні для виконання лабораторної роботи. Перевірка практичних результатів включає перевірку коду та виконання тестових завдань.

В процесі навчання викладач має право нарахувати до 5 заохочувальних балів за дострокове виконання лабораторної роботи, за проявлений творчий підхід при виконанні індивідуального завдання або за активну участь у обговоренні питань, що пов'язані з тематикою лекції або практичного заняття.

За виконання та здачу лабораторної роботи після зазначеного дедлайну, за значну кількість пропущених занять, або за порушення правил поведінки на заняттях викладач може призначити до 5 штрафних балів.

При проведенні контрольних заходів та при виконанні лабораторних робіт аспіранти повинні дотримуватися правил академічної доброчесності. При виявленні значного відсотку списування або плагіату викладач може відмовити у прийнятті даної роботи та вимагати доброчесного виконання навчального плану.

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

Види контролю з навчальної дисципліни «Методи добування даних» включають:

Поточний контроль: тестування закритими тестами.

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

Семестровий контроль: залік

Умови допуску до семестрового контролю: семестровий рейтинг більше 40 балів.

Таблиця відповідності рейтингових балів оцінкам за університетською шкалою:

<i>Кількість балів</i>	<i>Оцінка</i>
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

18. Додаткова інформація з дисципліни (освітнього компонента)

Перелік питань, які виносяться на семестровий контроль.

1. Робота з бібліотекою Pandas у середовищі Jupyter Notebook інтегратора Anaconda.
2. Порядок та основні етапи установки програмного забезпечення для вирішення задач аналізу та добування даних.
3. Реалізація практичної задачі добування даних на основі завдання, заданого викладачем.

4. Основні бібліотеки на мові Python для аналізу даних: модулі NumPy та SciPy, бібліотека Matplotlib.
5. Основні поняття про бібліотеку Pandas, структури даних Pandas.
6. Робота в Pandas з елементами Series та DataFrame.
7. Способи доступу до даних у структурах бібліотеки Pandas.
8. Використання атрибутів для доступу до даних в Pandas.
9. Отримання випадкового набору даних з структур бібліотеки Pandas.
10. Додавання елементів та індексація з використанням логічних виразів при роботі зі структурами бібліотеки Pandas.
11. Використання isin для роботи з даними в Pandas, робота з пропусками даних.
12. Поняття про методи аналізу даних.
13. Статистичні обмеження аналізу даних. Принцип Бонферроні.
14. Базові теоретичні положення та інструменти аналізу даних. Ефект Метью.
15. Метод пошуку найближчих сусідів та його застосування.
16. Збереження схожості скорочених наборів. Шинглінг документів.
17. Локально-чутливе хешування для документів.
18. Принципи формування міри відстані, відстані Евкліда, Жакарта, Хемінга. Косинусна відстань та відстань редагування.
19. Базові поняття теорії локально-чутливих функцій.
20. Сімейства локально-чутливих функцій для різних мір відстані.
21. Поняття про потокові дані. Базова модель поточкових даних.
22. Способи вибору даних в потоці.
23. Методи фільтрації потоків.
24. Способи організації підрахунку різних елементів у потоці.
26. Способи визначення моментів та практичні приклади.
27. Поняття згасання вікон та підрахунки у вікні.
28. PageRank, як інструмент для класифікації сторінок.
29. Способи ефективного обчислення для PageRank.
30. Тематично чутливий PageRank.
31. «Лінковий» спам, хаби та авторитети.
32. Модель аналізу на основі задачі про ринкову корзину
33. Використання алгоритму A-Priori для задач про ринкову корзину.
34. Методи обробки великих наборів даних у основній пам'яті.
35. Основні алгоритми обробки даних в пам'яті, алгоритм обмеження проходів.

Умова зарахування додаткових балів.

В рамках вивчення навчальної дисципліни «Методи добування даних» допускається зарахування балів, одержаних в результаті дистанційних курсів на платформі “Coursera”, за умови попереднього погодження програми даного курсу з викладачем та за умови отримання офіційного сертифікату.

Робочу програму навчальної дисципліни (силабус):

Складено , д.т.н, професор, Новотарський Михайло Анатолійович

Ухвалено кафедрою обчислювальної техніки (протокол № 18 від 25.05.2021)

Погоджено Методичною комісією ФІОТ (протокол № 10 від 14.06.2021)