

# I СИЛАБУС НАВЧАЛЬНОЇ ДИСЦИПЛІНИ

## 1. Загальна інформація про навчальну дисципліну

Повна назва навчальної дисципліни	Методи аналізу великих даних
Повна офіційна назва закладу вищої освіти	Сумський державний університет
Повна назва структурного підрозділу	Факультет електроніки та інформаційних технологій. Кафедра прикладної математики та моделювання складних систем
Розробник(и)	Марченко Анна Вікторівна
Рівень вищої освіти	другий рівень вищої освіти, НРК – 7 рівень, QF-LLL – 7 рівень, FQ-EHEA – другий цикл
Семестр вивчення навчальної дисципліни	8 тижнів протягом 3-го семестру
Обсяг навчальної дисципліни	5 кредитів ЄКТС, 150 годин, з яких 48 години становить контактна робота з викладачем (24 годин лекцій, 24 години практичних робіт), 102 години становить самостійна робота
Мова викладання	українська

## 2. Місце навчальної дисципліни в освітній програмі

Статус дисципліни	Обов'язкова навчальна дисципліна для освітньої програми 'Прикладна математика'
Передумови для вивчення дисципліни	Системи баз даних, Алгоритми машинного навчання
Додаткові умови	Додаткові умови відсутні
Обмеження	Обмеження відсутні

## 3. Мета навчальної дисципліни

Мета дисципліни – набуття знань та практичних навичок використання методів видобування даних, які застосовуються для формування паралельних обчислень, хешування, потокової обробки даних тощо

## 4. Зміст навчальної дисципліни

Тема 1 Вступ до великих даних <i>Огляд курсу. Поняття обробки, здобуття даних. Хеш-функція. Індeksi Знайомство з тех-нологіями Hadoop та Spark</i>
Тема 2 Map Reduce <i>Розподілені файлові системи. Map Reduce (Згортка-Відображення). Алгоритми, по-будовані на MapReduce. Модель комунікативної вартості. Вступ до PySpark. Основи машинного обчислення засобами PySpark</i>
Тема 3 Пошук схожих об'єктів <i>Додатки пошуку найближчого сусіда. Розбиття документів на шингли. Сигнатури множин із збереженням схожості. Хешування документів. Метрики. Застосування хешування з урахуванням близькості</i>
Тема 4 Аналіз потоків даних <i>Потокова модель даних. Вибірка даних з потоку. Фільтрація потоків. Підрахунок різних елементів у потоці. Оцінювання моментів</i>

Тема 5 Аналіз посилань <i>PageRank. Ефективне обчислення PageRank. Тематичний PageRank. Посилальний спам. Хаби та авторитетні сторінки</i>
Тема 6 Реклама в Інтернеті <i>Онлайнкові алгоритми. Задача про паросумісність. Задача про ключові слова. Реалізація алгоритму Adwords</i>
Тема 7 Рекомендаційні системи <i>Модель рекомендаційної системи. Рекомендація на основі фільтрації вмісту. Колаборативна фільтрація. Зниження розмірності</i>
Тема 8 Аналіз графів соціальних мереж <i>Вступ до графів. Соціальні мережі як графи. Кластеризація графа соціальної мережі. Пряме знаходження спільнот. Знаходження спільнот, що перетинаються. Околиці в графах</i>
Тема 9 Зниження розмірності <i>Власні значення та власні вектори. Метод головних компонент. CUR-декомпозиція. Правильний вибір рядків та стовпців</i>
Тема 10 Машинне навчання на великих даних <i>Модель машинного навчання. Перцептрони. Метод опорних векторів. Навчання за найближчими сусідами. Порівняння методів навчання</i>

## 5. Очікувані результати навчання навчальної дисципліни

Після успішного вивчення навчальної дисципліни здобувач вищої освіти зможе:

PH1	Знати сучасні методи обробки великих даних
PH2	Визначати алгоритми обробки даних для вирішення професійних задач аналізу великих даних
PH3	Застосовувати спеціалізовані алгоритми обробки поточкових даних для аналізу швидкозмінних даних, технологій пошукових систем, методи кластеризації великогабаритних даних
PH4	Поєднувати існуючі алгоритми обробки великих даних для вирішення комплексних задач аналізу складних даних

## 6. Роль навчальної дисципліни у досягненні програмних результатів

Програмні результати навчання, досягнення яких забезпечує навчальна дисципліна.

Для спеціальності 113 Прикладна математика:

PP2	Уміти формалізувати задачі певної предметної галузі, формулювати їх математичну постановку та обирати раціональний метод вирішення; розв'язувати отримані задачі аналітичними та чисельними методами, оцінювати точність та достовірність отриманих результатів.
PP17	Володіти математичними методами обробки великих наборів даних. Вміти обирати до застосування оптимальні методи для конкретної задачі побудови моделі поведінки складної системи за існуючим набором даних та будувати графові ймовірнісні моделі для розв'язання технічних задач
PP19	Уміти оцінити на адекватність результат обробки великих масивів даних

## 7. Види навчальних занять та навчальної діяльності

### 7.1 Види навчальних занять

<b>Тема 1. Вступ до великих даних</b>
---------------------------------------

Лк1 "Вступ до великих даних" (денна) <i>Поняття обробки, здобуття даних. Хеш-функція. Індекси Знайомство з технологіями Hadoop та Spark</i>
Пр1 "Вступ до мови Scala" (денна) <i>Вступ до мови Scala та інтерфейс командного рядка</i>
<b>Тема 2. Map Reduce</b>
Лк2 "Алгоритми Map Reduce" (денна) <i>Алгоритми, побудовані на MapReduce. Модель комунікативної вартості</i>
Лк3 "Map Reduce та новий програмний стек" (денна) <i>Розподілені файлові системи. Map Reduce (Згортка-Відображення)</i>
Пр2 "Spark - швидкий обчислювальний кластер" (денна) <i>Вступ до PySpark. Основи машинного обчислення засобами PySpark</i>
<b>Тема 3. Пошук схожих об'єктів</b>
Лк4 "Пошук схожих об'єктів. Частина 1" (денна) <i>Додатки пошуку найближчого сусіда. Розбиття документів на шингли. Сигнатури множин із збереженням схожості. Хешування документів. Метрики</i>
Лк5 "Пошук схожих об'єктів. Частина 2" (денна) <i>Застосування хешування з урахуванням близькості. Корпоративна культура та організація бізнес-процесів</i>
Пр3 "Spark – надбудова SQL" (денна) <i>Виконання запитів на створення об'єктів збереження даних та вибірки даних</i>
<b>Тема 4. Аналіз потоків даних</b>
Лк6 "Аналіз потоків даних" (денна) <i>Потокова модель даних. Вибірка даних з потоку. Фільтрація потоків. Підрахунок різних елементів у потоці. Оцінювання моментів</i>
Пр4 "Spark Streaming" <i>Обробка потоків даних засобами Spark</i>
<b>Тема 5. Аналіз посилань</b>
Лк7 "PageRank" (денна) <i>PageRank. Ефективне обчислення PageRank. Тематичний PageRank. Посилальний спам. Хаби та авторитетні сторінки</i>
Пр5 "Візуалізація великих даних засобами MicroStrategy. Огляд." (денна) <i>Вступ до бізнес-аналітики. Інформаційні панелі в режимі офлайн із розширеною візуалізацією</i>
<b>Тема 6. Реклама в Інтернеті</b>
Лк8 "Реклама в Інтернеті" (денна) <i>Онлайнові алгоритми. Задача про паросумісність. Задача про ключові слова. Реалізація алгоритму Adwords</i>
Пр6 "Візуалізація великих даних засобами MicroStrategy. Побудова простої інформаційної панелі." (денна) <i>Аналіз показників для побудови світового сценарію викидів вуглецю</i>
<b>Тема 7. Рекомендаційні системи</b>
Лк9 "Рекомендаційні системи" (денна) <i>Модель рекомендаційної системи. Рекомендація на основі фільтрації змісту. Колаборативна фільтрація. Зниження розмірності</i>
Пр7 "MLLib" (денна) <i>Машинне навчання з MLLib. Створення рекомендацій фільмів</i>
Пр8 "Візуалізація великих даних засобами MicroStrategy. Побудова комплексної інформаційної панелі" (денна) <i>Вступ до державних курсів дистанційної освіти</i>
<b>Тема 8. Аналіз графів соціальних мереж</b>

Лк10 "Аналіз графів соціальних мереж" (денна) <i>Вступ до графів. Соціальні мережі як графи. Кластеризація графа соціальної мережі. Пряме знаходження спільнот. Знаходження спільнот, що перетинаються. Околиці в графах</i>
Пр9 "Аналіз графів з GraphX" (денна) <i>Дослідження веб-графів і алгоритмів графів (PageRank) за допомогою GraphX</i>
Пр10 "Візуалізація великих даних засобами MicroStrategy. Побудова звітної інформаційної панелі" (денна) <i>Інтеграція додатку MicroStrategy з Salesforce.com для створення звітної інформаційної панелі</i>
<b>Тема 9. Зниження розмірності</b>
Лк11 "Зниження розмірності" (денна) <i>Власні значення та власні вектори. Метод головних компонент. CUR-декомпозиція. Правильний вибір рядків та стовпців</i>
Пр11 "Tachyon" (денна) <i>Розгортання надійної файлової системи в пам'яті через кластер</i>
<b>Тема 10. Машинне навчання на великих даних</b>
Лк12 "Машинне навчання на великих даних" (денна) <i>Модель машинного навчання. Перцептрон. Метод опорних векторів. Навчання за найближчими сусідами. Порівняння методів навчання</i>
Пр12 "BlinkDB" (денна) <i>Використання SQL зі статистичною вибіркою</i>

## 7.2 Види навчальної діяльності

НД1	Розв'язання базових практичних завдань за допомогою онлайн-технологій
НД2	Розв'язання розширених практичних завдань за допомогою онлайн-технологій
НД3	Виконання практичних завдань
НД4	Розв'язування вправ на основі теоретичного матеріалу
НД5	Перегляд відеолекцій МВОК
НД6	Підготовка карти пам'яті
НД7	Виконання індивідуальних розрахунково-аналітичних завдань

## 8. Методи викладання, навчання

Дисципліна передбачає навчання через:

МН1	Інтерактивні лекції
МН2	Практико-орієнтоване навчання
МН3	Проблемно-пошуковий метод

## 9. Методи та критерії оцінювання

### 9.1. Критерії оцінювання

Шкала оцінювання ECTS	Визначення	Чотирибальна національна шкала оцінювання	Рейтингова бальна шкала оцінювання
-----------------------	------------	---	------------------------------------

5 (відмінно)	Відмінне виконання лише з незначною кількістю помилок	A	$90 \leq RD \leq 100$
4 (добре)	Вище середнього рівня з кількома помилками	B	$82 \leq RD < 89$
4 (добре)	Загалом правильна робота з певною кількістю помилок	C	$74 \leq RD < 81$
3 (задовільно)	Непогано, але зі значною кількістю недоліків	D	$64 \leq RD < 73$
3 (задовільно)	Виконання задовольняє мінімальні критерії	E	$60 \leq RD < 63$
2 (незадовільно)	Можливе повторне складання	FX	$35 \leq RD < 59$
2 (незадовільно)	Виконання задовольняє мінімальні критерії	F	$0 \leq RD < 34$

## 9.2 Методи поточного формативного оцінювання

МФО1	Експрес-тестування
МФО2	Оцінювання практичних завдань
МФО3	Взаємооцінювання (peer assessment)
МФО4	Оцінювання вправ на основі теорії
МФО5	Оцінювання індивідуальних розрахунково-аналітичних завдань

## 9.3 Методи підсумкового сумативного оцінювання

МСО1	Звіт за результатами виконання базових практичних завдань на основі онлайн-технологій
МСО2	Звіт за результатами виконання розширених практичних завдань на основі онлайн-технологій
МСО3	Звіт за результатами виконання практичних завдань
МСО4	Виконання вправ на основі теоретичного матеріалу
МСО5	Виконання індивідуальних розрахунково-аналітичних завдань
МСО6	Розробка карт пам'яті
МСО7	Експрес-тестування
МСО8	Поточні контрольні роботи (проміжний модульний контроль)

Контрольні заходи:

3-й семестр		100 балів
МСО1. Звіт за результатами виконання базових практичних завдань на основі онлайн-технологій		3
	3x1	3
МСО2. Звіт за результатами виконання розширених практичних завдань на основі онлайн-технологій		15
	5x3	15
МСО3. Звіт за результатами виконання практичних завдань		8
	4x2	8
МСО4. Виконання вправ на основі теоретичного матеріалу		8
	4x2	8
МСО5. Виконання індивідуальних розрахунково-аналітичних завдань		16

		16
МСО6. Розробка карт пам'яті		<b>10</b>
	10x1	10
МСО7. Експрес-тестування		<b>20</b>
	10x2	20
МСО8. Поточні контрольні роботи (проміжний модульний контроль)		<b>20</b>
		20

Контрольні заходи в особливому випадку:

## 10. Ресурсне забезпечення навчальної дисципліни

### 10.1 Засоби навчання

ЗН1	Мультимедіа, проекційна апаратура
ЗН2	Комп'ютери, комп'ютерні системи та мережі
ЗН3	Програмне забезпечення для підтримки дистанційного навчання

### 10.2 Інформаційне та навчально-методичне забезпечення

<b>Основна література</b>	
1	Юре Лесковец, Ананд Раджараман, Джеффри Д. Ульман Анализ больших наборов данных. / Пер. с англ. Слинкин А. А. – М.: ДМК Пресс, 2016. – 498 с.: ил
<b>Інформаційні ресурси в Інтернеті</b>	
2	BigData Mini Course. – UC Berkeley
3	Online Course. Mining of Massive Datasets. Jure Leskovec, Anand Rajaraman, Jeff Ullman