



# Інтелектуальний аналіз текстів

## Робоча програма навчальної дисципліни (Силабус)

### Реквізити навчальної дисципліни

Рівень вищої освіти	<i>Другий (магістерський)</i>
Галузь знань	<i>11 Математика та статистика</i>
Спеціал	<i>113 Прикладна математика</i>
Освітня програма	<i>Наука про дані та математичне моделювання</i>
Статус дисципліни	<i>Нормативна</i>
Форма навчання	<i>очна(денна)</i>
Рік підготовки, семестр	<i>2 курс, осінній семестр</i>
Обсяг дисципліни	<i>3 кредити</i>
Семестровий контроль/ контрольні заходи	<i>Екзамен</i>
Розклад занять	<i>Лекція – 1 раз на тиждень, практичні заняття 1 раз на 2 тижні</i>
Мова викладання	<i>Українська</i>
Інформація про керівника курсу / викладачів	Лектор: канд. техн. наук, доцент, Сирота Сергій Вікторович, <a href="mailto:syrotasergiy@ill.kpi.ua">syrotasergiy@ill.kpi.ua</a> Практичні / Семінарські: канд. техн. наук, доцент, Сирота Сергій Вікторович, <a href="mailto:syrotasergiy@ill.kpi.ua">syrotasergiy@ill.kpi.ua</a>
Розміщення курсу	Дистанційний ресурс Moodle, <a href="http://moodle.eec.kpi.ua/course/view.php?id=5">http://moodle.eec.kpi.ua/course/view.php?id=5</a>

### Програма навчальної дисципліни

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

Вивчення навчальної дисципліни «Інтелектуальний аналіз текстів» дає більш глибокі знання в царині сучасних технологій автоматичного аналізу текстів природної і для вирішення інших задач професійної діяльності визначених Стандартом вищої освіти НТУУ «КПІ» для другого (магістерського) рівня вищої освіти за спеціальністю 113 Прикладна математика. **Метою** курсу є вивчення основних математичних моделей, що застосовуються для автоматичного аналізу текстів. **Предмет** навчальної дисципліни — застосування математичних методів до автоматичного аналізу текстів.

В результаті навчання студенти отримають наступні компетенції

- Здатність учитися і оволодівати сучасними знаннями.
- Здатність застосовувати знання у практичних ситуаціях.
- Здатність генерувати нові ідеї (креативність).
- Здатність до абстрактного мислення, аналізу та синтезу.
- Здатність до пошуку, оброблення та аналізу інформації з різних джерел.
- Навички у використанні інформаційних і комунікаційних технологій.
- Здатність сформулювати математичну постановку задачі, спираючись на постановку мовою предметної галузі, та обирати метод її розв'язання, що забезпечує потрібні точність і надійність результату.

#### ЗНАННЯ:

- основних положень теорії прикладної лінгвістики;

- основних методів видобутку та аналізу текстових даних для виявлення закономірностей;
- статистичних підходів, які можна загалом застосувати до довільних текстових даних будь-якої природної мови;

#### уміння:

- Використовувати статистичні методи аналізу текстів;
- Видобутку та аналізу текстових даних для виявлення закономірностей;

#### досвід:

- Видобування інформації з текстових корпусів;
- Використання спеціального програмного забезпечення для аналізу текстів

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

Вивченню дисципліни «Психологічні моделі прийняття управлінських рішень» повинне передувати вивчення дисциплін:

- «Математичний аналіз»,
- «Алгебра та геометрія»,
- «Дискретна математика»,
- «Теорія ймовірності»,
- «Математична статистика»,
- «Методи штучного інтелекту».

## **3. Зміст навчальної дисципліни**

Тема1 Огляд методів обробки природних мов

- Маркування частин мови (part-of-speech tagging)
- Синтаксичний аналіз
- Семантичний аналіз
- Двозначність
- Представлення текстів, представлення у вигляді набору слів
- Контекст слова;
- Контекстна схожість.
- Синтагматичні відношення
- Парадигматичні відношення

Тема 2. Дослідження асоціацій слів

- Тема і охоплення тема
- Термін, як тема

Тема 3. Кластеризація текстів

- Кластеризація документів, і кластеризація термінів
- Нормалізація в EM-алгоритмі
- Ієрархічна конгломеративна кластеризація, і метод K-середніх
- Пряма і непряма оцінка кластеризації

Тема 4. Методи категоризації текстів

- Класифікатор, що генерує vs класифікатор що відрізняє
- Навчальні дані
- Логістична регресія
- Носій думки, мета думки, почуття, представлення думки
- Класифікація почуттів
- Особливості, n-грам, розповсюджені зразки, і перенавчання
- Ординарна логістична регресія

- Оцінка якості прогнозу моделі

Тема 5. Аналіз думок & Контекстуальний аналіз текстів

- Передбачення, що базуються на текстах
- Контекстуальний аналіз текстів

Для зарахування лабораторного практикуму з дисципліни "Автоматизоване опрацювання текстів" слухачам необхідно:

- Обрати один із запропонованих інструментальних засобів автоматизованого опрацювання текстів.
- Ознайомитися з функціональним призначенням засобу і його можливостями.
- Встановити програмне забезпечення (програму/бібліотеку/toolkit/framework) на комп'ютері.
- За допомогою встановленого інструментального засобу провести аналіз довільно обраного тексту (текстів).
- Зробити висновки, написати відгук, викласти свої враження від використання данного засобу.
- Підготувати звіт обсягом 2-3 тис знаків за результатами виконаної роботи

### Навчальні матеріали та ресурси

#### Базова література

1. Cheng X. Z. *Text Mining and Analytics [Електронний ресурс] / Cheng Xiang Zhai // The University of Illinois at Urbana-Champaign, . – 2021. – Режим доступу до ресурсу: <https://www.coursera.org/learn/text-mining/supplement/ctYg7/welcome-to-text-mining-and-analytics>.*
2. Aggarwal C. C., Zhai C. *Mining Text Data. — Springer, 2012. — 527 с. — ISBN 9781461432234.*
3. Сирота С. В. *Автоматизоване опрацювання текстів [Електронний ресурс] / С. В. Сирота, // ЦЕО КПІ. – 2021. – Режим доступу до ресурсу: <http://moodle.eec.kpi.ua/course/view.php?id=10>*

#### Допоміжна література

1. Ananiadou S. *Text mining for biology and biomedicine / S. Ananiadou, J. McNaught. – Boston, London: Hartec House, 2006.*
2. Konchady M. *Text mining application programming / Manu Konchady. – Boston, Mass.: Charles River Media, 2008.*

#### Інформаційні ресурси

1. Центр електронної освіти КПІ <http://moodle.eec.kpi.ua/course/view.php?id=10>
2. Електронний кампус КПІ ім. Ігоря Сікорського. Матеріали з дисципліни «Алгоритми і структури даних». – Режим доступу : <http://login.kpi.ua>

### Навчальний контент

#### 4. Методика опанування навчальної дисципліни (освітнього компонента)

Далі в таблиці інформація за темами про лекції у формі календарного плану.

	Лекційні заняття
Тиждень 1-2	Тема1 Огляд методів обробки природних мов <ul style="list-style-type: none"> <li>• Маркування частин мови (part-of-speech tagging)</li> <li>• Синтаксичний аналіз</li> <li>• Семантичний аналіз</li> <li>• Двозначність</li> <li>• Представлення текстів, представлення у вигляді набору слів</li> </ul>
Тиждень 3-4	<ul style="list-style-type: none"> <li>• Контекст слова</li> <li>• Контекстна схожість</li> </ul>

	<ul style="list-style-type: none"> <li>• Синтагматичні відношення</li> <li>• Парадигматичні відношення</li> </ul>
<i>Тиждень 5-6</i>	<p>Тема 2. Дослідження асоціацій слів</p> <ul style="list-style-type: none"> <li>• Тема і охоплення тема</li> <li>• Термін, як тема</li> </ul>
<i>Тиждень 7-8</i>	<p>Тема 3 Кластеризація текстів</p> <ul style="list-style-type: none"> <li>• Кластеризація документів, і кластеризація термінів</li> <li>• Нормалізація в EM-алгоритмі</li> <li>• Ієрархічна конгломеративна кластеризація, і метод K-середніх</li> <li>• Пряма і непряма оцінка кластеризації</li> </ul>
<i>Тиждень 9-10</i>	<p>Тема 4. Методи категоризації текстів</p> <ul style="list-style-type: none"> <li>• Класифікатор, що генерує vs класифікатор що відрізняє</li> <li>• Навчальні дані</li> <li>• Логістична регресія</li> <li>• Носій думки, мета думки, почуття, представлення думки</li> <li>• Класифікація почуттів</li> <li>• Особливості, n-грам, розповсюджені зразки, і перенавчання</li> <li>• Ординарна логістична регресія</li> <li>• Оцінка якості прогнозу моделі</li> </ul>
<i>Тиждень 11</i>	<p>Тема 5. Методи категоризації текстів</p> <ul style="list-style-type: none"> <li>• Класифікатор, що генерує vs класифікатор що відрізняє</li> <li>• Навчальні дані</li> <li>• Логістична регресія</li> <li>• Класифікатор K-ближчих сусідів</li> <li>• Носій думки, мета думки, почуття, представлення думки</li> <li>• Класифікація почуттів</li> <li>• Особливості, n-грам, розповсюджені зразки, і перенавчання</li> <li>• Ординарна логістична регресія</li> <li>• Оцінка якості прогнозу моделі</li> </ul>
<i>Тиждень 12</i>	<p>Тема 5. Аналіз думок &amp; Контекстуальний аналіз текстів</p> <ul style="list-style-type: none"> <li>• Передбачення, що базуються на текстах</li> <li>• Контекстуальний аналіз текстів</li> </ul>
<i>Тиждень 13-18</i>	Здача результатів лабораторного практикуму підготовка та здача екзамену

## **5. Самостійна робота студента**

*Для засвоєння матеріалів курсу передбачені наступні види самостійної роботи*

- *Робота з рекомендованими джерелами і*
- *Виконання завдання лабораторного практикуму*

## 6. Політика навчальної дисципліни (освітнього компонента)

### Відвідування

Відвідування лекцій не обов'язкове студентам надається можливість самостійно опрацювати теми занять використовуючи дистанційний курс та рекомендовані джерела.

Поточний контроль проходить в режимі тестування і має часові обмеження.

Вимоги до слухачів курсу базуються на принципах академічної доброчесності і рівності всіх студентів. У випадку виявлення випадків запозичення без відповідних посилань об'єктів авторського права, як то: програмний код, зображення, креслення, та інший мультимедійний контент або виявлення плагіату — бали за відповідні роботи будуть анульовані і нараховані штрафні бали. Повторні порушення принципів академічної доброчесності можуть призвести до недопуску до складання заліку.

Викладачі можуть перевіряти роботи, виконані у рамках курсу, за допомогою систем виявлення плагіату.

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

На семестр пропонуються наступні види діяльності за які нараховуються бали

- Доповідь зі звітом про виконаний лабораторний практикум 40 балів
- Екзамен 60 балів

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

## 8. Додаткова інформація з дисципліни

- В разі проходження слухачами курсу:
- <https://www.coursera.org/learn/text-mining/supplement/ctYq7/welcome-to-text-mining-and-analytics> і наданні сертифікату студенту зараховується .60 балів
- Дистанційні онлайн заняття проводяться за допомогою Google Meet за запрошеннями, які публікуються в кафедральній платформі Slack та дублюються на надані студентами адреси електронної пошти;
- Офіційні звернення до викладача розглядаються через кафедральну платформу Slack або електронну адресу syrota.sergiy@iill.kpi.ua.

## Робочу програму навчальної дисципліни (силабус):

Склав доцент кафедри прикладної математики, канд. техн. наук, доц. Сирота Сергій Вікторович

Ухвалено кафедрою прикладної математики (протокол № 13 від 16.06.2022)

Погоджено Методичною комісією факультету<sup>1</sup> (протокол № 11 від 27.06.2022)

<sup>1</sup> Методичною радою університету – для загальноуніверситетських дисциплін.