

**НАЦІОНАЛЬНА АКАДЕМІЯ НАУК УКРАЇНИ**  
**Центр наукових досліджень та викладання іноземних мов**



**«ЗАТВЕРДЖУЮ»**

Директор Центру наукових досліджень  
та викладання іноземних мов НАН України

к.філол.н., доцент

ЖАЛАЙ В.Я.

«21» листопада 2023 р.

**СИЛАБУС**  
**ВК 05 ДВА**

**Основи комп'ютерної лінгвістики**

підготовки здобувачів третього (освітньо-наукового) рівня вищої освіти – доктора філософії – спеціальності 035 «Філологія»

Київ 2023

<b>1. Загальна інформація про навчальну дисципліну</b>	
<b>Повна назва навчальної дисципліни</b>	Основи комп'ютерної лінгвістики
<b>Повна офіційна назва закладу вищої освіти</b>	Центр наукових досліджень та викладання іноземних мов НАН України (м. Київ)
<b>Повна назва структурного підрозділу</b>	кафедра іноземних мов
<b>Розробник</b>	Крамар Н.А., доктор філософії з філології
<b>Рівень вищої освіти</b>	Третій (освітньо-науковий) рівень НРК України – 9 рівень QF for ENEA – третій цикл, EQF for LLL – 8 рівень
<b>Семестр вивчення навчальної дисципліни</b>	1-2 семестр 3 року навчання (5-6 семестри)
<b>Обсяг навчальної дисципліни</b>	3 кредити ЄКТС / 90 годин, з яких 10 годин – лекції 8 годин – практичні заняття 72 годин – самостійна робота
<b>Мова(и) викладання</b>	англійська, українська
<b>Викладач(і)</b>	Крамар Н.А., доктор філософії з філології
<b>2. Місце навчальної дисципліни в освітній програмі</b>	
<b>Статус дисципліни</b>	вибірковий компонент, дисципліна вибору аспіранта ВК 05 ДВА
<b>Передумови для вивчення дисципліни</b>	Другий рівень вищої освіти (диплом магістра). Аспірант повинен знати: англійську мову на рівні не нижче С1 та мати навички пошуку та аналізу наукової інформації.
<b>Додаткові умови</b>	Додаткові умови відсутні.
<b>Обмеження</b>	Обмеження відсутні.
<b>3. Мета навчальної дисципліни</b>	
<b>Мета:</b> курс передбачає ознайомлення студентів з основами комп'ютерної лінгвістики та надання їм практичних навичок у застосуванні інструментів і методів аналізу тексту для різноманітних завдань, таких як екстракція даних (data mining), аналіз настроїв (sentiment analysis), тематичне моделювання, класифікація текстів тощо.	
<b>4. Зміст навчальної дисципліни</b>	

### **1-й модульний цикл**

**Тема 1. Місце комп'ютерної лінгвістики та NLP у сучасному лінгвістичному ландшафті.**

**Тема 2. Основи програмування у Python для NLP-задач.**

### **2-й модульний цикл**

**Тема 3. Попередня обробка текстових даних у Python: бібліотеки, методи, підходи.**

**Тема 4. Обробка текстових даних: від морфологічного до семантичного рівня.**

**Тема 5. Основи машинного навчання для NLP-задач.**

### **3-й модульний цикл**

**Тема 6. Векторна репрезентація слів.**

**Тема 7. Тематичне моделювання та класифікація текстових даних.**

## **5. Очікувані результати навчання навчальної дисципліни**

**У результаті вивчення навчальної дисципліни аспірант повинен знати:**

- роль комп'ютерної лінгвістики у сучасному лінгвістичному ландшафті
- історію становлення комп'ютерної лінгвістики
- основні поняття комп'ютерної лінгвістики такі як: корпус, датасет, токен, лема, стемінг, n-грам (n-gram), парсинг, синтаксичний аналіз, морфологічний аналіз, векторизація, частотний аналіз, вбудування слів (word embedding).
- особливості аналізу тексту та екстрагування інформації у Python
- особливості лематизації та стемінгу мов з різних мовних сімей
- специфіку NLP-бібліотек Python і їхню доцільність для різних завдань комп'ютерної лінгвістики
- методи word embedding (вбудування слів)
- методи векторної репрезентації слів
- основні алгоритми тематичного моделювання (LDA, LSA, PLSA)
- відмінність між класифікацією на основі правил (rule-based classification) та класифікацією на основі машинного навчання
- етичні проблеми, пов'язані з веб-скрейпінгом та аналізом чутливих даних (sensitive data)
- принципи функціонування великих мовних моделей (large language models)
- можливості доналаштування (fine-tuning) наявних мовних моделей
- основні корпуси української, англійської та інших мов та способи їх використання для цілей NLP

**вміти:**

- використовувати основні NLP-бібліотеки Python (Spacy, NLTK, Textblob, Stanza, Pandas та інші)
- створювати функції у Python
- обробляти дані різних форматів та виконувати конвертацію з одного формату в інший (txt, pdf, csv, xml тощо)
- виконувати аналіз частотності та аналіз ключових слів
- виділяти основні n-грами у текстових даних
- виконувати попередню обробку (preprocessing) даних (видалення пунктуації та нерелевантних символів, видалення стоп-слів, нормалізація правопису тощо)
- здійснювати лематизацію та стемінг мовних даних
- застосовувати SpacyMatchers та Regex для екстрагування інформації потрібного формату
- виконувати тематичне моделювання та класифікацію даних
  - етично використовувати чатботи на основі великих мовних моделей та створювати ефективні промти для них
- візуалізувати результати аналізу у вигляді діаграм та схем
- чітко викладати та презентувати результати застосування методів комп'ютерної лінгвістики в усній та письмовій формах

## **6. Роль навчальної дисципліни у досягненні програмних результатів**

Завдання навчальної дисципліни ВК 05 ДВА «Основи комп'ютерної лінгвістики» полягає у формуванні та набутті таких компетентностей: загальні компетентності: ЗК: 1, 3, 4, 5 (відповідно до переліку загальних компетентностей ОНП). Фахові компетентності: ФК 2, 4, 5, 6, 7, 8, 9, 10, 13, 14 (відповідно до переліку фахових компетентностей ОНП). Програмні результати навчання: ПРН 1.4, 2.1, 2.2, 2.3, 4.2, 4.4 (відповідно до переліку програмних результатів навчання ОНП).

### **6.1 Види навчальних занять**

Видами навчальних занять при вивченні дисципліни є лекції (Л), практичні (індивідуальні) заняття (П) та самостійна робота (С).

### **Навчально-тематичний план лекцій і практичних занять 5-6 семестри 3 року навчання**

<b>Тема</b>	<b>Назва теми</b>	<b>Кількість годин</b>
-------------	-------------------	------------------------

		Лекції	Практичні (індивідуальні)	Самостійна робота
1.	Тема 1. Місце комп'ютерної лінгвістики та NLP у сучасному лінгвістичному ландшафті.	1	1	10
2.	Тема 2. Основи програмування у Python для NLP-задач.	2	2	20
3.	Тема 3. Попередня обробка текстових даних у Python: бібліотеки, методи, підходи.	2	1	20
4.	Тема 4. Обробка текстових даних: від морфологічного до семантичного рівня.	2	1	10
5.	Тема 5. Основи машинного навчання для NLP-задач.	1	1	5
6.	Тема 6. Векторна репрезентація слів та вбудування слів (word embedding).	1	1	2
7.	Тема 7. Тематичне моделювання та класифікація текстових даних.	1	1	5
		10	8	72
Усього	90			

### Змістовий модуль 1

#### Тема 1

Л-1

Місце комп'ютерної лінгвістики та NLP у сучасному лінгвістичному ландшафті

П-1

Можливості застосування комп'ютерної лінгвістики для наукових цілей та бізнес-цілей

## **Змістовий модуль 2**

### Тема 2

Л-2

Основи програмування у Python для NLP-задач: змінні, цикли, функції

П-2

Робота з різними типами даних у Python: рядки, списки, числа, словники

Л-3

Оператори, цикли for і while, функції

П-3

Завантаження текстових даних та імпортування бібліотек у Python

### Тема 3

Л-4

Принципи попередньої обробки текстових даних у Python

Л-5

Основи роботи з регулярними виразами

П-4

Практикум з попередньої обробки тексту з використанням бібліотеки Spacy

### Тема 4

Л-6

Морфологічний аналіз текстових даних: tokenization, POS tagging, n-grams.

Л-7

Синтаксичний та семантичний аналіз текстових даних: dependency parsing, named entity recognition

П-5

Практикум з частиномовної розмітки та побудови дерева залежностей в NLTK та Spacy

## **Змістовий модуль 3**

### Тема 5

Л-8

Основи машинного навчання для NLP-задач

П-6

Тренування ML моделей				
Тема 6				
Л-9				
Векторна репрезентація слів та вбудування слів (word embedding)				
П-7				
Практикум із застосування TF-IDF				
Тема 7				
Л-10				
Тематичне моделювання та класифікація текстових даних				
П-8				
Порівняння провідних видів тематичного моделювання (LDA, NMF, BERTopic) на корпусі текстів				
<b>7.2 Види навчальної діяльності</b>				
Підготовка до лекцій.				
Підготовка до практичних занять.				
Підготовка до виконання індивідуальних завдань за модулями.				
Підготовка до іспиту.				
<b>8. Методи викладання, навчання</b>				
Модульне навчання.				
Лекції-візуалізації.				
Інтерактивні лекції.				
Flipped learning.				
Обговорення-дискусії.				
Проблемно-пошукові методи.				
Міждисциплінарне навчання.				
<b>9. Методи та критерії оцінювання</b>				
<b>9.1. Критерії оцінювання</b>				
Оцінка	Визначення	Національна шкала оцінювання	Рейтингова бальна шкала оцінювання	
ВІДМІННО	відмінне виконання лише з незначною кількістю помилок	5, 0 (відмінно)	$90 \leq RD \leq 100$	

ДОБРЕ	в цілому правильна робота з певною кількістю помилок	4, 0 (добре)	$74 \leq RD \leq 89$
ЗАДОВІЛЬНО	виконання задовольняє мінімальні критерії	3,0 (задовільно)	$60 \leq RD \leq 73$
НЕЗАДОВІЛЬНО	можливе повторне складання	2 (незадовільно)	$35 \leq RD \leq 59$
НЕЗАДОВІЛЬНО	необхідний повторний курс з навчальної дисципліни	2 (незадовільно)	$RD < 35$

### **9.2 Методи поточного формативного оцінювання**

За дисципліною передбачені такі методи поточного формативного оцінювання:  
експрес-контроль, настанови викладача в процесі виконання індивідуальних практичних завдань, перевірка, обговорення та оцінювання виконаної роботи.

### **9.3 Методи підсумкового сумативного оцінювання**

Оцінювання протягом семестру проводиться у формі усних опитувань, письмового тестування, перевірки виконання індивідуальних завдань та дискусій.

Робота повинна бути виконана самостійно.

Оцінювання аспіранта здійснюється таким чином:

1. Експрес-контроль – 30 балів
2. Практичні (індивідуальні заняття) та самостійна робота – 30 балів
3. Залік – 40 балів

Заохочувальні та штрафні бали:

1. Відсутність на лекції без поважних причин - (-) 2 бали.
  2. Відсутність на практичних (індивідуальних) заняттях без поважних причин (-) 2 бали.
  3. Підготовка публікації до друку та/або виступу на конференції (+) 10 балів
- Сума як штрафних, так і заохочувальних балів не має перевищувати 0,1R=10 балів.

Форма підсумкового контролю – іспит.

## **10. Ресурсне забезпечення навчальної дисципліни**



<p><b>10.1 Засоби навчання</b></p>	<p>Навчальний процес потребує використання мультимедійних, друкованих та електронних освітніх ресурсів.</p>
<p><b>10.2 Інформаційне та навчально методичне забезпечення</b></p>	<p><b>Основна література</b></p> <ol style="list-style-type: none"> <li>1. Bengfort, B., Bilbro, R., &amp; Ojeda, T. (2018). <i>Applied Text Analysis with Python</i>. O'Reilly Media, Inc.</li> <li>2. Bird, S., Klein, E., &amp; Loper, E. (2009). <i>Natural Language Processing with Python</i>. O'Reilly Media.</li> <li>3. Napke, H., Lane, H., &amp; Howard, C. (2019). <i>Natural Language Processing in Action</i>. Manning Publications.</li> <li>4. Hovy, D. (2012). Programming in Python for Linguists. A Gentle Introduction. Retrieved from <a href="http://www.dirkhovy.com/portfolio/papers/download/pfl_hanout.pdf">http://www.dirkhovy.com/portfolio/papers/download/pfl_hanout.pdf</a></li> </ol> <p><b>Допоміжна література</b></p> <ol style="list-style-type: none"> <li>1. Egger, R., &amp; Yu, J. (2022). A Topic Modeling Comparison Between LDA, NMF, Top2Vec, and BERTopic to Demystify Twitter Posts. <i>Frontiers in Sociology</i>, 6(7), 886498.</li> <li>2. Johnson, M. (2011). How relevant is linguistics to computational linguistics? <i>Linguistic Issues in Language Technology</i>, 6(7), 1–23.</li> <li>3. Jurafsky, D., &amp; Martin, J. H. (2008). <i>Speech and Language Processing, 2nd edition</i>. Pearson Prentice Hall.</li> <li>4. Lukeš, D., &amp; Rosa, R. (2020). An Introduction to Python for Linguists. Retrieved from <a href="https://v4py.github.io/intro.html">https://v4py.github.io/intro.html</a></li> <li>5. Manning, C., &amp; Schütze, H. (1999). <i>Foundations of Statistical Natural Language Processing</i>. MIT Press.</li> <li>6. Mitkov, R. (2009). <i>The Oxford Handbook of Computational Linguistics</i>. Oxford: Oxford University Press.</li> <li>7. Panggabean, H., &amp; Tobing, A. (2015). Computational Linguistics Application Using Python Programming. <i>IOSR Journal of Humanities and Social Science (IOSR-JHSS)</i>, 20(7), 18-30.</li> <li>8. Roth, B., &amp; Wiegand, M. (2021). Python for Linguists. <i>Computational Linguistics</i>, 47(1), 217–220.</li> <li>9. Дарчук Н. П. (2008). Комп'ютерна лінгвістика (автоматичне опрацювання тексту): підручник. К.: Видавничо-поліграфічний центр "Київський університет".</li> </ol>

	10. Жуковська, В. (2013). Вступ до корпусної лінгвістики: навчальний посібник. Житомир: Вид-во ЖДУ ім. І. Франка.
--	---

Силабус підготувала:

доктор філософії з філології,  
ст.н.с. Крамар Н.А.