

Минобрнауки России

**ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ БЮДЖЕТНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ
ВЫСШЕГО ОБРАЗОВАНИЯ
«ВОРОНЕЖСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ»
(ФГБОУ ВО «ВГУ»)**



УТВЕРЖДАЮ

Заведующий кафедрой
Сирота Александр Анатольевич

Кафедра технологий обработки и защиты информации

03.05.2023

РАБОЧАЯ ПРОГРАММА УЧЕБНОЙ ДИСЦИПЛИНЫ

Б1.В.02 Компьютерная лингвистика

1. Код и наименование направления подготовки/специальности:

09.04.02 Информационные системы и технологии

2. Профиль подготовки/специализация:

Системы прикладного искусственного интеллекта

3. Квалификация (степень) выпускника:

Магистратура

4. Форма обучения:

Очная

5. Кафедра, отвечающая за реализацию дисциплины:

Кафедра технологий обработки и защиты информации

6. Составители программы:

Гаршина Вероника Викторовна, к.т.н., доцент

7. Рекомендована:

№7 от 03.05.2023

8. Учебный год:

2023-2024

9. Цели и задачи учебной дисциплины:

Изучить формальные, программно-реализуемые подходы к изучению структур и закономерностей естественных языков, ознакомится с основными прикладными практическими задачами компьютерной лингвистики

Основные задачи дисциплины:

- изучить основные принципы и методов обработки естественного языка. получение навыков разработки и интеграции программных систем обработки естественного языка в программные продукты.
- знакомство с принципами проектирования программных систем, ориентированных на обработку естественно-языковых текстов.

10. Место учебной дисциплины в структуре ООП:

Дисциплина относится к вариативной части учебного плана Б1.В.

Для ее изучения требуются входные знания из курсов: математические методы в современных

информационных технологиях, искусственный интеллект, программирование и теория алгоритмов.

11. Планируемые результаты обучения по дисциплине/модулю (знания, умения, навыки), соотнесенные с планируемыми результатами освоения образовательной программы (компетенциями выпускников) и индикаторами их достижения:

Код и название компетенции	Код и название индикатора компетенции	Знания, умения, навыки
ПК-8 Способен разрабатывать новые технологии проектирования информационных систем, прогнозировать развитие информационных систем и технологий	ПК-8.1 Знает инструменты и методы моделирования бизнес-процессов, современные подходы и стандарты автоматизации организации, отраслевую документацию, основы реинжиниринга бизнес-процессов организации	Знает терминологию, базовые понятия, математические и алгоритмические методы обработки текстовой информации, этапы разработки лингвистически ориентированных программных продуктов, технологии представления и обработки текстовой информации, формальные модели представления естественного языка.
ПК-15 Способен разрабатывать и исследовать модели объектов профессиональной деятельности, предлагать и адаптировать методики решения научно-исследовательских задач, планировать и проводить исследования	ПК-15.1 Умеет обеспечивать сбор научно-технической (научной) информации, необходимой для постановки и решения задач исследования	Умеет проектировать, разрабатывать и интегрировать программные системы обработки естественного языка в программные продукты. Применять программные библиотеки и инструменты для разработки систем, ориентированных на обработку звучащей речи и текста.

12. Объем дисциплины в зачетных единицах/час:

3/108

Форма промежуточной аттестации:

Зачет с оценкой, Контрольная работа

13. Трудоемкость по видам учебной работы

Вид учебной работы	Семестр 1	Семестр 2	Всего
Аудиторные занятия	0	48	48
Лекционные занятия		32	32
Практические занятия			0
Лабораторные занятия		16	16

Вид учебной работы	Семестр 1	Семестр 2	Всего
Самостоятельная работа	0	60	60
Курсовая работа			0
Промежуточная аттестация	0	0	0
Часы на контроль			0
Всего	0	108	108

13.1. Содержание дисциплины

п/п	Наименование раздела дисциплины	Содержание раздела дисциплины	Реализация раздела дисциплины с помощью онлайн-курса, ЭУМК
1	Задачи компьютерной лингвистики в изучении естественного языка (ЕЯ).	<p>Лекционные занятия по разделу</p> <p>1. Компьютерная лингвистика: задачи, направления исследований. Проблемы моделирования естественного языка. Лингвистические ресурсы, используемые для обработки текста и речи.</p> <p>Лабораторные занятия по разделу не предусмотрены</p>	Создан электронный онлайн - курс, размещены материалы к практическим и лабораторным работам.
2	Лингвистический процессор. Алгоритмы лингвистического разбора и анализа текста. Парсеры ЕЯ-предложений	<p>Лекционные занятия по разделу</p> <p>2. Уровни текстового анализа: графематический, фонетический, морфологический, синтаксический, семантический. Основные задачи, их взаимосвязь. Лингвистический процессор - функциональная структура. Графематический анализ: выделение структурных элементов в тексте: границы предложений, слов, словари сокращений.Стемминг и лематизация. Стемминг: классификация алгоритмов стемминга. Алгоритмы стемминга для русского языка: Портера, Stemka, MyStem. Ошибки стемминга.</p> <p>3.Методы морфологического анализа, используемые в лингвистических процессорах. Морфологические словари. Программные библиотеки для работы с морфологией естественного языка.</p> <p>4.Синтаксический анализ в компьютерной лингвистике. Способы представления синтаксического разбора: синтаксическое дерево, размеченное предложение. Форматы разметки CONLLU.</p> <p>5. Классификация Хомского. Формальная модель представления синтаксиса: деревья составляющих. Грамматики составляющих. деревья зависимостей, КС - грамматики. Примеры синтаксических парсеров и инструменты их разработки и интеграции.</p> <p>Лабораторные занятия по разделу Лабораторная работа № 1, 2</p>	Создан электронный онлайн - курс, размещены материалы к практическим и лабораторным работам.

п/п	Наименование раздела дисциплины	Содержание раздела дисциплины	Реализация раздела дисциплины с помощью онлайн-курса, ЭУМК
3	Формальные методы исследования структуры ЕЯ текста и классификации полнотекстовых документов	<p>Лекционные занятия по разделу 6. Статистические методы анализа структур ЕЯ текста на морфологическом, синтаксическом, семантическом уровнях. Метод позиционных статистик. Применение статистических методов для задач атрибуции (определения авторства) и исследования естественно-языковых текстов на неизвестных языках.</p> <p>7-8. Математическая постановка задачи распознавания образов и классификации. Формальные методы определения сходства ЕЯ документов. Векторная модель. Методы определения сходства и классификации текстовых документов. Кластерный анализ текстов (рубрицирование, стилистика). Деревья принятия решений, векторные методы, Байесовский классификатор для задач компьютерной лингвистики.</p> <p>9. Модели информационного поиска: инвертированная индексация, Булева и векторная модели. Метрики оценки близости документов. Оценка качества поиска: tf-idf, точность, полнота.</p> <p>10-11. Нейросетевой подход к анализу текстов: извлечение именованных сущностей, классификация, анализ тональности высказываний.</p> <p>Лабораторные занятия по разделу Лабораторная работа № 3,4,5.</p>	Создан электронный онлайн - курс, размещены материалы к практическим и лабораторным работам.
4	Проблемы построения систем семантического анализа текстов (Text Mining)	<p>Лекционные занятия по разделу 12. Проблема семантического анализа для автоматических систем обработки текстов. DataMining и TextMining. Извлечение фактов из текстов (именованных сущностей и ключевых слов), установление взаимосвязей. Типы именованных сущностей и способы извлечения из текстов. Применение в системах обработки текстов. Проблемы разрешения омонимии, анафоры и кореферентности.</p> <p>13. Автоматические системы извлечения знаний из разнородных текстовых источников. Задачи структурирования текстовых данных. Формирование онтологии предметной области по тексту. Графовые БД. Построение семантической модели текста. Онтограф текста. Семантический анализ текстов на основе онтологии предметной области. Форматы представления, стандарты разработки, инструменты.</p> <p>Лабораторные занятия по разделу Лабораторная работа № 6, 7</p>	Создан электронный онлайн - курс, размещены материалы к практическим и лабораторным работам.

п/п	Наименование раздела дисциплины	Содержание раздела дисциплины	Реализация раздела дисциплины с помощью онлайн-курса, ЭУМК
5	Прикладные системы обработки текстов	<p>Лекционные занятия по разделу</p> <p>14. Проблемы автоматизации синтеза (генерации) текста. Методы генерации. Шаблонные системы генерации. Генерация текстов на основе БД - простой отчет, связанный отчет. ЕЯ запросы к БД. Семантические, морфологические, синтаксические проблемы синтеза текстов.</p> <p>15. Автоматическое аннотирование и реферирование: архитектура построения систем, используемые методы, прикладное использование, примеры.</p> <p>16. Вопросно-ответные системы: индексирование в информационно-поисковых системах, архитектура, способы обработки запросов, генерация различных типов ответов. Генерация диалогов в вопросно-ответных системах - чат-боты.</p> <p>Лабораторные занятия по разделу</p> <p>Лабораторная работа № 8.</p>	Создан электронный онлайн - курс, размещены материалы к практическим и лабораторным работам.

13.2. Темы (разделы) дисциплины и виды занятий

№ п/п	Наименование темы (раздела)	Лекционные занятия	Практические занятия	Лабораторные занятия	Самостоятельная работа	Всего
1	Задачи компьютерной лингвистики в изучении естественного языка (ЕЯ)	2		0	6	8
2	Лингвистический процессор. Алгоритмы лингвистического разбора и анализа текста. Парсеры ЕЯ-предложений	8		4	10	22

№ п/п	Наименование темы (раздела)	Лекционные занятия	Практические занятия	Лабораторные занятия	Самостоятельная работа	Всего
3	Формальные методы исследования структуры ЕЯ текста и классификации полнотекстовых документов	12		6	12	30
4	Проблемы построения систем семантического анализа текстов (Text Mining)	4		4	16	24
5	Прикладные системы обработки текстов	6		2	16	24
		32	0	16	60	108

14. Методические указания для обучающихся по освоению дисциплины

1) При изучении дисциплины рекомендуется использовать следующие средства:

- рекомендуемую основную и дополнительную литературу;
- методические указания и пособия;
- контрольные задания для закрепления теоретического материала;
- электронные версии учебников и методических указаний для выполнения лабораторно - практических работ (при необходимости материалы рассылаются по электронной почте).

2) Для максимального усвоения дисциплины рекомендуется проведение письменного опроса (тестирование, решение задач) студентов по материалам лекций и лабораторных работ. Подборка вопросов для тестирования осуществляется на основе изученного теоретического материала. Такой подход позволяет повысить мотивацию студентов при конспектировании лекционного материала.

3) При проведении лабораторных занятий обеспечивается максимальная степень соответствия с материалом лекционных занятий и осуществляется экспериментальная проверка методов, алгоритмов и технологий, применяемых в интеллектуальной обработке информации, излагаемых в рамках лекций.

4) При переходе на дистанционный режим обучения для создания электронных курсов, чтения лекций он-лайн и проведения лабораторно- практических занятий используется информационные ресурсы Образовательного портала "Электронный университет ВГУ (<https://edu.vsu.ru>), базирующегося на системе дистанционного обучения Moodle, развернутой в университете.

15. Перечень основной и дополнительной литературы, ресурсов интернет, необходимых

для освоения дисциплины

№ п/п	Источник
1	Боярский, К. К. Введение в компьютерную лингвистику : учебное пособие / К. К. Боярский. — Санкт-Петербург : НИУ ИТМО, 2013. — 72 с. — Текст : электронный // Лань : электронно-библиотечная система. — URL: https://e.lanbook.com/book/70822 (дата обращения: 01.05.2023). — Режим доступа: для авториз. пользователей.
2	Пентус, А. Е. Математическая теория формальных языков : учебное пособие / А. Е. Пентус, М. Р. Пентус. — 2-е изд. — Москва : ИНТУИТ, 2016. — 218 с. — ISBN 5-9556-0062-0. — Текст : электронный // Лань : электронно-библиотечная система. — URL: https://e.lanbook.com/book/100633 (дата обращения: 01.05.2023). — Режим доступа: для авториз. пользователей.
3	Ганегедара, Т. Обработка естественного языка с TensorFlow : руководство / Т. Ганегедара ; перевод с английского В. С. Яценкова. — Москва : ДМК Пресс, 2020. — 382 с. — ISBN 978-5-97060-756-5. — Текст : электронный // Лань : электронно-библиотечная система. — URL: https://e.lanbook.com/book/140584 (дата обращения: 01.05.2023). — Режим доступа: для авториз. пользователей.
4	Добров Б.В. Онтологии и тезаурусы: модели, инструменты, приложения: учебное пособие / Б.В. Добров, В.В. Иванов, Н.В. Лукашевич, В.Д. Соловьев. / - М.: Интернет-Университет Информационных Технологий; БИНОМ. Лаборатория знаний, 2009.

б) дополнительная литература:

№ п/п	Источник
1	Лукашевич Н. В. Тезаурусы в задачах информационного поиска. — Издательство МГУ имени М. В. Ломоносова, 2011
2	Маннинг К. Д., Рагхаван П., Шютце Х. Введение в информационный поиск. — Вильямс, 2011.
3	Николаев И.С., Митренина О.В., Ландо Т.М. Прикладная и компьютерная лингвистика. - URSS. 2016.
4	Кипяткова И.С., Ронжин А.Л., Крапов А.А. Автоматическая обработка разговорной русской речи. - Санкт-Петербург: ГУАП, 2013.
5	Автоматическая обработка текстов на естественном языке и компьютерная лингвистика : учеб. пособие / Большакова Е.И., Клышинский Э.С., Ландэ Д.В., Носков А.А., Пескова О.В., Ягунова Е.В. — М.: МИЭМ, 2011. — 272 с.
6	Бенгфорт Бенджамин, Билбро Ребекка, Охеда Тони Прикладной анализ текстовых данных на Python. Машинное обучение и создание приложений обработки естественного языка. — СПб.: Питер, 2019.

№ п/п	Источник
7	Моделирование распознавания рукописного текста на основе скрытых марковских моделей : монография / И. Я. Львович, Я. Е. Львович, А. П. Преображенский [и др.]. — Воронеж : ВИБТ, 2016. — 164 с. — ISBN 978-5-4446-0838-8. — Текст : электронный // Лань : электронно-библиотечная система. — URL: https://e.lanbook.com/book/157486 (дата обращения: 01.05.2023). — Режим доступа: для авториз. пользователей.
8	Маккинни, У. Python и анализ данных / У. Маккинни ; перевод с английского А. А. Слинкина. — 2-ое изд., испр. и доп. — Москва : ДМК Пресс, 2020. — 540 с. — ISBN 978-5-97060-590-5. — Текст : электронный // Лань : электронно-библиотечная система. — URL: https://e.lanbook.com/book/131721 (дата обращения: 01.05.2023). — Режим доступа: для авториз. пользователей.
9	Бонцанини, М. Анализ социальных медиа на Python. Извлекайте и анализируйте данные из всех уголков социальной паутины на Python / М. Бонцанини ; перевод с английского А. В. Логунова. — Москва : ДМК Пресс, 2018. — 288 с. — ISBN 978-5-97060-574-5. — Текст : электронный // Лань : электронно-библиотечная система. — URL: https://e.lanbook.com/book/108129 (дата обращения: 01.05.2023). — Режим доступа: для авториз. пользователей.
10	Теофили, Т. Глубокое обучение для поисковых систем : руководство / Т. Теофили ; перевод с английского Д. А. Беликова. — Москва : ДМК Пресс, 2020. — 318 с. — ISBN 978-5-97060-776-3. — Текст : электронный // Лань : электронно-библиотечная система. — URL: https://e.lanbook.com/book/140574 (дата обращения: 01.05.2023). — Режим доступа: для авториз. пользователей.

в) информационные электронно-образовательные ресурсы:

№ п/п	Источник
1	Электронный каталог Научной библиотеки Воронежского государственного университета. - (http // www.lib.vsu.ru/).
2	Образовательный портал «Электронный университет ВГУ».- (https://edu.vsu.ru/)
3	ЭБС«Университетская библиотека online» – Контракт №3010-06/23-22 от 30.12.2022 (срок предоставления с 12.01.2023 по 11.01.2024)
4	ЭБС «Консультант студента» – Лицензионный договор №3010-06/22-22 от 30.12.2022 (с дополнительным соглашением №1 от 09.01.2023) (срок предоставления с 12.01.2023 по 11.01.2024)
5	ЭБС Лань (лицензионный договор №3010-14/37-23 от 07.03.2023 (срок предоставления с 12.03.2023 по 11.03.2024)
6	Международная конференция по компьютерной лингвистике. http://www.dialog-21.ru/

№ п/п	Источник
7	Лаборатория компьютерной лингвистики Института проблем передачи информации РАН. http://proling.iitp.ru/
8	Лаборатория общей компьютерной лексикологии и лексикографии МГУ. http://www.philol.msu.ru/~lex/library.htm
9	Научно-практический журнал РЕЧЕВЫЕ ТЕХНОЛОГИИ http://speechtechnology.ru/

16. Перечень учебно-методического обеспечения для самостоятельной работы

№ п/п	Источник
1	Боярский К. К. Введение в компьютерную лингвистику. Учебное пособие. – СПб: НИУ ИТМО, 2013. – 72 с.
2	Информационные ресурсы Образовательного портала "Электронный университет ВГУ (https://edu.vsu.ru)

17. Образовательные технологии, используемые при реализации учебной дисциплины, включая дистанционные образовательные технологии (ДОТ), электронное обучение (ЭО), смешанное обучение):

Для реализации учебного процесса используются:

-
- ПО Microsoft в рамках подписки "Imagine/Azure Dev Tools for Teaching", договор №3010-16/96-18 от 29 декабря 2018г.
- ПО MATLAB Classroom ver. 7.0, 10 конкурентных бессрочных лицензий на каждый, компоненты: Matlab, Simulink, Stateflow, 1 тулбокс, N 21127/VRN3 от 30.09.2011 (за счет проекта ЕК TEMPUS/ERAMIS).
- ПО Матлаб в рамках подписки Университетская лицензия на программный комплекс для ЭВМ - MathWorks MATLAB Campus-Wide Suite по договору 3010-16/118-21 от 27.12.2021 (до 01.2025)
- ПО Редактор онтологий и фреймворк для построения баз знаний Protege. Свободно-распространяемое ПО.
- Персер русского языка ТОМИТА (Свободно-распространяемое ПО)
- Язык программирования Python, IDE Pycharm.

При проведении занятий в дистанционном режиме обучения используются технические и информационные ресурсы Образовательного портала "Электронный университет ВГУ (<https://edu.vsu.ru>), базирующегося на системе дистанционного обучения Moodle, развернутой в университете, а также другие доступные ресурсы сети Интернет.

18. Материально-техническое обеспечение дисциплины:

1) Мультимедийная лекционная аудитория (корп.1а, ауд. 290)

Учебная аудитория: компьютер преподавателя Pentium-G3420-3,2ГГц, монитор с ЖК 17", мультимедийный проектор, экран.

Система для видеоконференций Logitech ConferenceCam Group и ноутбук 15.6" FHD Lenovo V155-15API.

Специализированная мебель: доска 1 шт., столы 31 шт., стулья 64 шт.; выход в Интернет, доступ к фондам учебно-методической документации и электронным изданиям.

2) Компьютерный класс (один из №1-4 корп. 1а, ауд. № 382-385),

Учебная аудитория: специализированная мебель, персональные компьютеры на базе i3-2120-3,3ГГц, мониторы ЖК 19" (16 шт.), мультимедийный проектор, экран.

ПО: ОС Windows v.7, 8, 10, Набор утилит (архиваторы, файл-менеджеры), LibreOffice v.5-7, Foxit PDF Reader

Специализированная мебель: доска маркерная 1 шт., столы 16 шт., стулья 33 шт.; доступ к фондам учебно-методической документации и электронным изданиям, доступ к электронным библиотечным системам, выход в Интернет.

19. Оценочные средства для проведения текущей и промежуточной аттестаций

Порядок оценки освоения обучающимися учебного материала определяется содержанием следующих разделов дисциплины:

№ п/п	Разделы дисциплины (модули)	Код компетенции	Код индикатора	Оценочные средства для текущей аттестации
1	Разделы 1-5 Задачи компьютерной лингвистики в изучении естественного языка (ЕЯ). Лингвистический процессор. Алгоритмы лингвистического разбора и анализа текста. Парсеры ЕЯ-предложений. Формальные методы исследования структуры ЕЯ текста.	ПК-8	ПК-8.1	Устный опрос, собеседование. Практико-ориентированные задания по соответствующим разделам.
2	Разделы 1-5 Формальные методы классификации полнотекстовых документов. Проблемы построения систем семантического анализа текстов (Text Mining). Прикладные системы обработки текстов	ПК-15	ПК-15.1	Устный опрос, собеседование. Практико-ориентированные задания по соответствующим разделам. Контрольная работа. Лабораторные работы 1-10

Промежуточная аттестация

Форма контроля - Зачет с оценкой, Контрольная работа

Оценочные средства для промежуточной аттестации

Перечень вопросов, практико-ориентированное задание, лабораторные работы

20 Типовые оценочные средства и методические материалы, определяющие процедуры оценивания

20.1 Текущий контроль успеваемости

Текущая аттестация проводится в соответствии с Положением о текущей аттестации обучающихся по программам высшего образования Воронежского государственного университета. Текущая аттестация проводится в формах устного опроса (индивидуальный опрос, фронтальная беседа) и письменных работ (контрольные, лабораторные работы). При оценивании могут использоваться количественные или качественные шкалы оценок.

Текущий контроль успеваемости по дисциплине осуществляется с помощью следующих оценочных средств:

Устный опрос на лабораторных занятиях; Практико-ориентированное задание; Лабораторные работы.

№ п/п	Наименование оценочного средства	Представление оценочного средства в фонде	Критерии оценки
1	2	3	4
1	Устный опрос на лабораторных занятиях	Вопросы по темам/разделам дисциплины	Правильный ответ – зачтено, неправильный или принципиально неточный ответ - не зачтено
2	Практико - ориентированное задание по разделам дисциплины	Теоретические вопросы по темам \ разделам дисциплины	Шкала оценивания соответствует приведенной ниже
3	Лабораторная работа	Содержит 8 лабораторных заданий, предусматривающие разработку систем обработки текстов на основе различных алгоритмов с использованием программных средств разработки.	При успешном выполнении работ в течение семестра фиксируется возможность оценивания только теоретической части дисциплины в ходе промежуточной аттестации (зачета), в противном случае проверка задания по лабораторным работам выносится на зачет.

Пример задания для выполнения лабораторной работы

Лабораторная работа

Работа с библиотекой Natural Language Toolkit (NLTK) для решения задач

Статистической обработки текстов (английский язык)

Изучение материала на основе документации по NLTK:

<http://www.nltk.org/>

Natural Language Processing with Python – Analyzing Text with the Natural Language Toolkit

Steven Bird, Ewan Klein, and Edward Loper,

http://www.nltk.org/book_1ed/

Эксперименты проводить в <https://colab.research.google.com/>

Ход выполнения работы

Задание для выполнения:

1 – этап экспериментов (Анализ структуры отдельного корпуса (текста))

1. Выбрать корпус для экспериментов (из NLTK или загрузить свой)
2. Провести статистический анализ текста:
 - Длина текста, словарь текста, число различных слов в словаре, рассчитать параметр лексического разнообразия текста.
 - Определить число предложений, слов (провести токенизацию).
 - Убрать стоп слова (предлоги, союзы, управляющие слова) и построить частотный график встречаемости слов в тексте. Кумулятивный график частотного распределения слов.
 - Выделить частотные слова, относящиеся к одной леме (провести лематизацию)
 - На основе результатов лематизации вывести на печать слова, определяющие тематику текста (претенденты на ключевые слова). Выделить по частоте и длине.
 - Провести исследование тематической структуры текста (в каких частях текста о чем говорится) – исследовать частотное *расположение* слов в тексте - построить *график дисперсии*.
 - Распечатать ключевые слова (частотные слова), относящиеся к наиболее тематически важному разделу текста (определить по графику дисперсии). Для них построить частотный график встречаемости слов в тексте. Кумулятивный график частотного распределения слов.
 - Для ключевых слов найти им соответствующие биграммы и коллокации в тексте, оценить их частотность. Экспертным методом проверить соответствуют ли определенные словосочетания важными для уточнения тематики текста.

2 – этап экспериментов (сравнительный анализ нескольких корпусов)

1. Выбрать 2-3 корпуса для экспериментов (из NLTK или загрузить свои)
2. Провести статистический анализ этих корпусов по плану задания 1.
 - Провести сравнение по статистическим параметрам: словарь текста, число различных слов в словаре, рассчитать параметр лексического разнообразия те (насколько) текста.
 - Исследовать тематические структуры текстов (в каких частях текстов о чем говорится) – исследовать частотное *расположение* слов в тексте - построить *график дисперсии*.

Отчетность по лабораторной работе

- По проведенным исследованиям сделать отчет с заключением о статистике, стилистике, тематике исследуемых текстов.
- Отчет оформляется в Word с описанием проведенного исследования с питоновскими скриптами и комментариями в блокноте, разработанном в <https://colab.research.google.com/>.
- Дайте название своему проекту с экспериментами (щелкнув на имя блокнота в правом верхнем углу, переименуйте). Все ваши эксперименты сохраняются на вашем гугл диске.
- В отчет нужно прикрепить ссылку на созданный блокнот с экспериментами – *Поделиться* (в правом верхнем углу colab.research). Отчет выложить в ответах на задание на мудле.

Приведённые ниже задания рекомендуется использовать при проведении диагностических работ для оценки остаточных знаний по дисциплине

Компетенция ПК-8

Вопросы с выбором 1-балл

1. Что относится к лингвистическим ресурсам для разработки программного обеспечения систем компьютерной лингвистики (множественный выбор):

1. Базы словосочетаний
2. Тезаурусы
3. Онтологии
4. Текстовые корпуса
5. Базы данных
6. Компьютерные словари

Ответ: 1,2,3,4,6.

2. Поставить соответствие:

1. Графематический анализ	a) Выделение грамматической основы слова, определение частей речи, приведение слова к словарной форме.
2. Фонетический анализ	b) Выявление смысловых связей между словами и группами, извлечение семантических отношений.
3. Морфологический анализ	c) Выявление синтаксических связей между словами в предложении, построение синтаксической структуры предложения.
4. Синтаксический анализ	d) анализ звукового состава слова, позволяет вычлнить в слове звуки и определить их характеристики.
5. Семантический анализ	e) Выделение из массива данных предложений и слов (токенов).

Ответы: 1-e, 2-d, 3-a, 4-c, 5-b.

3. Поставить соответствие:

1. Токенизация	a) Процесс использует словарь и морфологический анализ, чтобы в итоге привести слово к его канонической форме (словарной форме).
2. Стемминг	b) Процесс разделения текста на предложения-компоненты или процесс разделения предложений на слова-компоненты.
3. Лемматизация	c) Процесс отсечения от слова окончаний и суффиксов, чтобы оставшаяся часть была одинаковой для всех грамматических форм слова.

Ответы: 1-b, 2-c, 3-a.

5. Какие формальные модели используются в компьютерной лингвистике? (множественный выбор)

- a. Контекстно-свободные грамматики
- b. Марковские модели

- c. Мультиагентные модели
- d. Графовые модели
- e. Автоматные модели
- f. Динамические модели

Ответы: a, b, d, e.

7. Какие утверждения верны для задания Контекстно-свободной грамматики (множественный выбор):

- a. Конечное множество A - алфавит. Его элементы называются символами. Конечные последовательности символов образуют слова в данном алфавите.
- b. Алфавит разделяется на терминальные ("окончательные") и нетерминальные ("промежуточные") символы.
- c. Среди нетерминальных символов может быть выбран один - начальный.
- d. Правила грамматики имеют вид $K \rightarrow X$, где K -нетерминальный символ, а X - слово, в которое могут входить и терминальные, и нетерминальные символы.
- e. Правила грамматики имеют вид $aKb \rightarrow aXb$, где K нетерминальный символ, окруженный как нетерминальными, так и терминальными символами a, b . X - слово, в которое могут входить и терминальные, и нетерминальные символы, окруженное нетерминальными и терминальными символами.

Ответы: a, b, c, d.

Компетенция ПК-15

Вопросы с коротким ответом 2-балла

1. Приведите название закона, отражающего эмпирическую закономерность распределения частоты слов естественного языка: "если все слова языка (или просто длинного текста) упорядочить по убыванию частоты их использования, то частота n -го слова в таком списке окажется приблизительно обратно пропорциональной его порядковому номеру n ."

Ответ: Закон Ципфа

Вопросы с развернутым ответом 3-балла

1. Опишите реализацию стеминга для русского языка на основе алгоритма Портера. Опишите шаги алгоритма.

Ответ:

Идея алгоритма: существует ограниченное количество словообразующих суффиксов, и стемминг слова происходит без использования каких-либо баз основ: только множество существующих суффиксов и вручную заданные правила.

Алгоритм состоит из пяти шагов.

На каждом шаге отсекается словообразующий суффикс и оставшаяся часть проверяется на соответствие правилам (например, для русских слов основа должна содержать не менее одной гласной). Если полученное слово удовлетворяет правилам, происходит переход на следующий шаг. Если нет — алгоритм выбирает другой суффикс для отсечения.

На первом шаге отсекается максимальный формообразующий суффикс,

на втором — буква «и»,

на третьем — словообразующий суффикс,

на четвертом — суффиксы превосходных форм, «ь» и одна из двух «н».

20.2 Промежуточная аттестация

Промежуточная аттестация может включать в себя проверку теоретических вопросов, а также, при необходимости (в случае не выполнения в течение семестра), проверку выполнения установленного перечня лабораторных заданий, позволяющих оценить уровень полученных знаний и/или практическое (ие) задание(я), позволяющее (ие) оценить степень сформированности умений и навыков.

Для оценки теоретических знаний используется перечень контрольно-измерительных материалов. Каждый контрольно-измерительный материал для проведения промежуточной аттестации включает два задания - вопросов для контроля знаний, умений и владений в рамках оценки уровня сформированности компетенции. При оценивании используется количественная шкала. Критерии оценивания представлены в приведенной ниже таблице

Для оценивания результатов обучения на экзамене используются следующие содержательные показатели (формулируется с учетом конкретных требований дисциплины)

1. знание теоретических основ учебного материала, основных определений, понятий и используемой терминологии;
2. владение навыками проведения компьютерного эксперимента, тестирования компьютерных алгоритмов обработки информации.
3. владение навыками программирования и экспериментирования с компьютерными моделями алгоритмов обработки информации в рамках выполняемых лабораторных заданий;
4. умение обосновывать свои суждения и профессиональную позицию по излагаемому вопросу;
5. умение связывать теорию с практикой, иллюстрировать ответ примерами, в том числе, собственными, умение выявлять и анализировать основные закономерности, полученные, в том числе, в ходе выполнения лабораторно-практических заданий;
6. умение проводить обоснование и представление основных теоретических и практических результатов (теорем, алгоритмов, методик) с использованием математических выкладок, блок-схем, структурных схем и стандартных описаний к ним;

Различные комбинации перечисленных показателей определяют критерии оценивания результатов обучения (сформированности компетенций) на зачете: высокий (углубленный) уровень сформированности компетенций; повышенный (продвинутый) уровень сформированности компетенций; пороговый (базовый) уровень сформированности компетенций.

Для оценивания результатов обучения на зачете с оценкой используется 4-балльная шкала: «отлично», «хорошо», «удовлетворительно», «неудовлетворительно». Для оценивания результатов обучения на зачете используется – зачтено, не зачтено по результатам тестирования.

Соотношение показателей, критериев и шкалы оценивания результатов обучения на зачете с оценкой представлено в следующей таблице.

Критерии оценивания компетенций и шкала оценок на зачете

Критерии оценивания компетенций	Уровень сформированности компетенций	Шкала оценок
---------------------------------	--------------------------------------	--------------

Обучающийся демонстрирует полное соответствие знаний, умений, навыков по приведенным критериям свободно оперирует понятийным аппаратом и приобретенными знаниями, умениями, применяет их при решении практических задач. Успешно выполнены лабораторные работы в соответствии с установленным перечнем.	Повышенный уровень	Отлично
Ответ на контрольно-измерительный материал не полностью соответствует одному из перечисленных выше показателей, но обучающийся дает правильные ответы на дополнительные вопросы. При этом обучающийся демонстрирует соответствие знаний, умений, навыков приведенным в таблицах показателям, но допускает незначительные ошибки, неточности, испытывает затруднения при решении практических задач. Успешно выполнены лабораторные работы в соответствии с установленным перечнем.	Базовый уровень	Хорошо
Обучающийся демонстрирует неполное соответствие знаний, умений, навыков приведенным в таблицах показателям, допускает значительные ошибки при решении практических задач. При этом ответ на контрольно-измерительный материал не соответствует любым двум из перечисленных показателей, обучающийся дает неполные ответы на дополнительные вопросы. Успешно выполнены лабораторные работы в соответствии с установленным перечнем.	Пороговый уровень	Удовлетворительно
Ответ на контрольно-измерительный материал не соответствует любым трем из перечисленных показателей. Обучающийся демонстрирует отрывочные, фрагментарные знания, допускает грубые ошибки. Не выполнены лабораторные работы в соответствии с установленным перечнем.	Ниже порогового уровня	Не зачтено

Пример контрольно-измерительного материала

УТВЕРЖДАЮ

Заведующий кафедрой технологий обработки и защиты информации



А.А. Сирота

__._.2023

Направление подготовки / специальность 09.04.02 Информационные системы и технологии
программа

Дисциплина Б.1.В.02 Компьютерная лингвистика

Форма обучения Очное

Вид контроля Зачет с оценкой

Вид аттестации Промежуточная

Контрольно-измерительный материал № 1

1. Стемминг и лематизация. Стемминг: классификация алгоритмов стемминга. Алгоритмы стемминга для русского языка: Портера, Stemka, MyStem. Ошибки стемминга.

2. Автоматические системы извлечения знаний из разнородных текстовых источников. Задачи структурирования текстовых данных. Извлечение именованных сущностей и отношений между ними - подходы.

Преподаватель _____ В.В.Гаршина

Примерный перечень вопросов к зачету

№	Содержание
1	Компьютерная лингвистика как междисциплинарная область. Основные задачи, решаемые компьютерной лингвистикой. Направления исследований
2	Проблемы моделирования естественного языка в компьютерной лингвистике. Лингвистические ресурсы, используемые для обработки текста и речи.
3	Классификация прикладных систем в области компьютерной лингвистики.
4	Появление компьютерной лингвистики. Этапы становления. Компьютерная лингвистика в России.
5	Уровни текстового анализа: графематический, фонетический, морфологический, синтаксический, семантический. Основные задачи, их взаимосвязь.
6	Графематический анализ: задачи, методы реализации, примеры. Графематический анализ: выделение структурных элементов в тексте: границы предложений, слов, словари сокращений.
7	Морфологический анализ. Задачи, методы реализации, примеры морфоанализаторов и инструментов разработки.
8	Стемминг и лематизация. Стемминг: классификация алгоритмов стемминга. Алгоритмы стемминга для русского языка: Портера, Stemka, MyStem. Ошибки стемминга.
9	Лематизация: используемые методы, примеры для русского языка, инструменты разработки.
10	Автоматическое выделение именованных сущностей и ключевых слов. Типы именованных сущностей и способы извлечения из текстов. Применение в системах обработки текстов/
11	Синтаксический анализ в компьютерной лингвистике. Способы представления синтаксического разбора: синтаксическое дерево, размеченное предложение. Примеры синтаксических парсеров и инструменты разработки.
12	Формальная модель представления синтаксиса: деревья составляющих. Грамматики составляющих.
13	Формальная модель представления синтаксиса: деревья зависимостей.
14	Формальная модель представления синтаксиса: КС- грамматики.

15	Проблемы автоматизации синтеза (генерации) текста. Этапы генерации (схема). Методы генерации.
16	Шаблонные системы генерации. Генерация текстов на основе БД - простой отчет, связанный отчет. ЕЯ запросы к БД.
17	Алгоритмы синтеза текста для вербализации заданного содержания. Семантические, морфологические, синтаксические проблемы синтеза текстов.
18	Автоматическое аннотирование: архитектура построения систем, используемые методы, прикладное использование, примеры действующих систем.
19	Автоматическое реферирование: архитектура построения систем, используемые методы, прикладное использование, примеры действующих систем.
20	Информационный поиск: архитектура, модели представления документов, обработка поисковых запросов, извлечение документов.
21	Модели информационного поиска: инвертированная индексация, Булева и векторная модели. Метрики оценки близости документов. Оценка качества поиска: tf-idf, точность, полнота.
22	Вопросно-ответные системы: индексирование в информационно-поисковых системах, архитектура, способы обработки запросов, генерация различных типов ответов. Генерация диалогов в вопросно-ответных системах - чат-боты.
23	DataMining и TextMining. Извлечение фактов из текстов, установление взаимосвязей. Проблемы разрешения омонимии, анафоры и кореферентности.
24	Автоматические системы извлечения знаний из разнородных текстовых источников. Задачи структурирования текстовых данных. Извлечение именованных сущностей и отношений между ними. Подходы.
25	Извлечение фактов на основе контекстно-свободных грамматик, реализуемых в Томита парсере. Примеры грамматик.