

**УТВЕРЖДЕНА**

Решением Ученого совета  
АНО ВО «Центральный университет»  
«07» марта 2024 г.  
Протокол №1

**Рабочая программа дисциплины (модуля)  
«Natural Language Processing (Обработка естественного языка)»**

**Направление подготовки:** 02.04.01 Математика и компьютерные науки

**Направленность (профиль) подготовки:** Продуктовая аналитика

**Квалификация (степень) выпускника:** магистр

**Форма обучения:** очная

**Срок освоения программы:** 2 года

**Год набора:** 2024

**Москва  
2024**

## Содержание

1. Краткая характеристика дисциплины (модуля) .....	3
2. Перечень планируемых результатов обучения.....	4
3. Тематический план.....	6
4. Содержание дисциплины (модуля).....	6
5. Учебно-методическое обеспечение .....	7
6. Материально-техническое обеспечение .....	7
7. Методические и оценочные материалы .....	9

## 1. Краткая характеристика дисциплины (модуля)

Рабочая программа дисциплины (модуля) «Natural Language Processing (Обработка естественного языка)» составлена в соответствии с федеральным государственным образовательным стандартом высшего образования – магистратура по специальности 02.04.01 Математика и компьютерные науки, профиль Продуктовая аналитика, утвержденный приказом Министерства науки и высшего образования Российской Федерации № 810 от 23.08.2017 года.

Изучение дисциплины (модуля) «Natural Language Processing (Обработка естественного языка)» играет ключевую роль в создании систем, которые могут понимать и генерировать человеческий язык, что открывает новые возможности для взаимодействия между людьми и машинами. Эта область обеспечивает развитие технологий, таких как чат-боты, автоматический перевод и анализ настроений, которые становятся все более важными в различных сферах, включая бизнес, медицину и образование.

### Место дисциплины (модуля) в структуре образовательной программы

Настоящая дисциплина (модуль) включена в учебный план по программе подготовки магистратуры по направлению 02.04.01 Математика и компьютерные науки, профиль Продуктовая аналитика и входит в вариативную часть Блока 1, формируемую участниками образовательных отношений.

Дисциплина (модуль) изучается на 2 курсе в 3 семестре, доступна для прохождения при условии успешного завершения дисциплины (модуля) «Machine Learning (Машинное обучение)».

**Цель изучения дисциплины (модуля):** приобретение знаний и навыков, необходимых для разработки и применения технологий, позволяющих компьютерам анализировать, понимать и взаимодействовать с человеческим языком в различных контекстах.

### Задачи изучения дисциплины (модуля):

— формирование знаний и понимания по темам: основы предобработки текста и выделения признаков, принципы языкового моделирования и ключевые архитектуры нейронных сетей, применяемых в области NLP (RNN, Трансформер), подходы к машинному переводу, генерации текста и суммаризации, современные методы обучения и оптимизации языковых моделей (LLM), основы работы диалоговых систем и извлечения информации, применение мультимодальных моделей и технологий обработки речи;

— освоение умений: проводить предобработку и векторизацию текста, применять нейронные сети и LLM для задач NLP, обучать и оптимизировать языковые модели, разрабатывать и интегрировать NLP-приложения, оценивать и улучшать качество моделей;

— формирование навыков владения инструментами Python для NLP (HuggingFace, nltk, pytorch), методами обучения и оптимизации языковых моделей, навыками работы с мультимодальными системами и Text-to-speech.

## 2. Перечень планируемых результатов обучения

Компетенции, формируемые в результате освоения дисциплины (модуля) при проведении учебных занятий в форме контактной работы обучающихся с педагогическими работниками Университета и в форме самостоятельной работы обучающихся:

Компетенция	Содержание компетенции	Индикатор компетенции	Перечень планируемых результатов обучения по дисциплине (модулю)
УК-6.	Способен определять и реализовывать приоритеты собственной деятельности и способы ее совершенствования на основе самооценки	УК-6.1.	Знает основные методы самооценки и анализа своей деятельности, а также принципы управления временем и целеполагания
		УК-6.2	Умеет ставить реалистичные и достижимые цели, определять приоритеты в своей деятельности, а также разрабатывать и внедрять планы по совершенствованию своих навыков и компетенций на основе полученной самооценки
		УК-6.3	Имеет практический опыт применения методов самооценки в своей профессиональной деятельности, включая участие в тренингах, семинарах и проектах, направленных на развитие личной эффективности и профессионального роста
ОПК-2.	Способен создавать и исследовать новые математические модели в естественных науках, совершенствовать и разрабатывать концепции, теории и методы	ОПК-2.1.	Знает основные математические модели и методы, используемые в естественных науках, включая статистическое моделирование, дифференциальные уравнения и численные методы, а также современные подходы к исследованию и анализу данных
		ОПК-2.2	Умеет разрабатывать и адаптировать математические модели для решения конкретных проблем в естественных науках, проводить их анализ и верификацию, а также интерпретировать полученные результаты в контексте научных исследований

		ОПК-2.3	Имеет практический опыт создания и исследования математических моделей в рамках научных проектов или исследований, включая участие в публикациях, конференциях или коллаборациях, где были разработаны и апробированы новые концепции и методы
ПК-3.	Способен решать задачи профессиональной деятельности в области продуктовой аналитики, формулировать результаты анализа и выявлять последствия полученных данных для принятия обоснованных решений и оптимизации продуктов	ПК-3.1.	Знает методы и инструменты продуктовой аналитики
		ПК-3.2.	Умеет применять аналитические инструменты и программное обеспечение для обработки и визуализации данных, а также формулировать выводы на основе проведенного анализа
		ПК-3.3.	Имеет опыт работы над реальными проектами в области продуктовой аналитики, включая анализ пользовательского поведения и оптимизацию продуктов на основе полученных данных
ПК-4.	Способен публично представлять собственные и известные научные результаты	ПК-4.1.	Знает основные принципы эффективного публичного выступления, методы визуализации данных и основные требования к научным презентациям, включая структуру и содержание
		ПК-4.2.	Умеет четко и логично формулировать свои научные результаты, адаптируя их для различных аудиторий, а также использовать визуальные средства для улучшения восприятия информации
		ПК-4.3.	Имеет практический опыт участия в научных конференциях, семинарах или других мероприятиях, где успешно представлял свои и известные научные результаты, получая обратную связь и взаимодействуя с аудиторией

### 3. Тематический план

№ п/п	Наименование раздела дисциплины (модуля)	Трудоемкость, академические часы				ТКУ (текущий контроль успеваемости)
		Очная форма				
		Аудиторная работа		Контроль	Самостоятельная работа	
Лекции	Семинары (практические занятия)					
1	Основы обработки текста и векторные представления	4	4		20	Домашние задания, Тесты
2	Основные архитектуры нейронных сетей для задач NLP	4	4		20	Домашние задания, Тесты
3	Языковые модели на основе архитектуры Трансформер	4	4		20	Домашние задания, Тесты
4	Прикладные задачи NLP	3	3		20	Домашние задания, Тесты
	<i>Зачет</i>			4		
	<b>Итого:</b>	<b>15</b>	<b>15</b>	<b>4</b>	<b>80</b>	
	<b>Объем дисциплины (модуля) (в ак. ч.)</b>	<b>114</b>				
	<b>Объем дисциплины (модуля) (в зач. ед.)</b>	<b>3</b>				

### 4. Содержание дисциплины (модуля)

№п/п	Наименование раздела дисциплины (модуля)	Содержание дисциплины (модуля) по темам
1	Основы обработки текста и векторные представления	Введение в анализ текстов, базовые методы предобработки и выделения признаков. Векторные представления слов.
2	Основные архитектуры нейронных сетей для задач NLP	Языковое моделирование. Рекуррентные нейронные сети (RNN). Машинный перевод и механизм внимания.
3	Языковые модели на основе архитектуры Трансформер	Архитектура Трансформер. языковые модели на основе архитектуры кодировщик Трансформера. Генеративные языковые модели на основе архитектуры декодеровщик Трансформера (LLM, prompt tuning, RLHF). Оптимизация языковых моделей (P-tuning, LoRA, квантизация).
4	Прикладные задачи NLP	RAG системы и ранжирование. Извлечение именованных сущностей и отношений. Диалоговые системы (intent detection, slot filling). Задача суммаризации. Мультимодальные модели. NLP для кода. Text-to-speech.

## 5. Учебно-методическое обеспечение

Университет располагает полным набором лицензионного и свободно распространяемого программного обеспечения, включая продукты отечественного производства.

Каждый студент в течение всего периода обучения получает индивидуальный неограниченный доступ к электронно-библиотечной системе и электронной информационно-образовательной среде университета. Эти системы предоставляют возможность доступа к ресурсам из любой точки, где есть подключение к сети Интернет, как на территории университета, так и за его пределами.

Студентам обеспечен удаленный доступ к современным профессиональным базам данных и информационным справочным системам.

### *Основная литература:*

1. Хобсон Л. Обработка естественного языка в действии : практическое руководство / Л. Хобсон, Х. Ханнес, Х. Коул. - Санкт-Петербург : Питер, 2020. - 576 с. - (Серия «Для профессионалов»). - ISBN 978-5-4461-1371-2.

2. Николенко С., Кадури А., Архангельская Е. Глубокое обучение. — СПб.: Питер, 2018. — 480 с.: ил. — (Серия «Библиотека программиста»). — ISBN 978-5-496-02536-2.

### *Дополнительная литература:*

1. Бенгфорт, Б. Прикладной анализ текстовых данных на Python. Машинное обучение и создание приложений обработки естественного языка : практическое руководство / Б. Бенгфорт, Р. Билбро, Т. Охеда. - Санкт-Петербург : Питер, 2020. - 368 с. - (Серия «Бестселлеры O'Reilly»). - ISBN 978-5-4461-1153-4.

## 6. Материально-техническое обеспечение

Университет располагает материально-технической базой, соответствующей действующим противопожарным правилам и нормам и обеспечивающей проведение всех видов дисциплинарной и междисциплинарной подготовки, практической и научно-исследовательской работ обучающихся, предусмотренных учебным планом.

Помещения, которые представляют собой учебные аудитории для проведения занятий лекционного типа, занятий семинарского (практического) типа, групповых и индивидуальных консультаций, текущего контроля и промежуточной аттестации, а также помещения для самостоятельной работы и помещения для хранения и профилактического обслуживания учебного оборудования. Помещения укомплектованы специализированной мебелью и техническими средствами обучения, служащими для представления учебной информации большой аудитории.

Изучение дисциплины (модуля) обеспечивается в учебных аудиториях, оснащенных:

- столами и стульями;
- компьютерной техникой;
- специализированным оборудованием, включая демонстрационное оборудование.

Помещения для самостоятельной работы обучающихся, в том числе приспособленные для использования инвалидами и лицами с ограниченными возможностями здоровья, оснащены компьютерной техникой с возможностью подключения к сети «Интернет» и обеспечением доступа к в электронную информационно-образовательную среду Университета.

Обучающимся предоставляется доступ (в том числе удаленный) к ресурсам информационно-телекоммуникационной сети «Интернет», электронным ресурсам (в том числе электронным библиотечным системам, современным профессиональным базам данных и информационным справочным системам):

№	Наименование портала (издания, курса, документа)	Ссылка
1.	Научная электронная библиотека elibrary.ru библиотека	<a href="https://elibrary.ru/defaultx.asp">https://elibrary.ru/defaultx.asp</a>
2.	База данных для IT-специалистов	<a href="https://habr.com">https://habr.com</a>
3.	База данных ScienceDirect	<a href="https://www.sciencedirect.com">https://www.sciencedirect.com</a>
4.	Официальный сайт Министерства науки и высшего образования Российской Федерации	<a href="https://minobrnauki.gov.ru/">https://minobrnauki.gov.ru/</a>
5.	Федеральный портал «Российское образование»	<a href="https://www.edu.ru/">https://www.edu.ru/</a>
6.	Информационная система "Единое окно доступа к образовательным ресурсам"	<a href="http://window.edu.ru/">http://window.edu.ru/</a>
7.	Единая коллекция цифровых образовательных ресурсов	<a href="http://school-collection.edu.ru/">http://school-collection.edu.ru/</a>
8.	Федеральный центр информационно - образовательных ресурсов	<a href="http://fcior.edu.ru/">http://fcior.edu.ru/</a>

Перечень информационных технологий, используемых при осуществлении образовательного процесса по дисциплине (модулю), в том числе комплект лицензионного программного обеспечения, современные профессиональные базы данных и информационные справочные системы:

Наименование ПО	Производство	Лицензионное / свободно распространяемое
<b>Операционные системы:</b>		
Microsoft Imagine (Windows Client, Server)	зарубежное	лицензионное
<b>Браузеры:</b>		
Яндекс.Браузер	отечественное	свободно распространяемое
Google Chrome	зарубежное	свободно распространяемое
<b>Офисные приложения:</b>		
Microsoft Imagine (Visio, OneNote)	зарубежное	лицензионное
TeXstudio	зарубежное	свободно распространяемое
Adobe Acrobat Reader	зарубежное	свободно распространяемое
<b>Программное обеспечение для планирования и учета времени:</b>		
Toggle app	зарубежное	свободно распространяемое
<b>Системы управления проектами:</b>		
Microsoft Imagine (Project)	зарубежное	лицензионное
<b>Системы управления базами данных:</b>		
Microsoft Imagine (SQL Server)	зарубежное	лицензионное
<b>Системы резервного копирования (backup):</b>		
Acronis Backup Advanced for HyperV	зарубежное	лицензионное
<b>Справочно-правовые системы:</b>		
КонсультантПлюс: справочно-правовая система	отечественное	лицензионное
<b>Средства антивирусной защиты:</b>		
Kaspersky Endpoint Security для бизнеса Стандартный Russian Edition	отечественное	лицензионное
<b>Среды разработки:</b>		
Visual Studio Code	зарубежное	свободно распространяемое
Bash (Unix shell)	зарубежное	свободно распространяемое
Anaconda	зарубежное	свободно распространяемое
Robotic Operating System	зарубежное	свободно распространяемое
CopelliaSim	зарубежное	свободно распространяемое

Google Colaboratory	зарубежное	свободно распространяемое
<b>Пакеты программных средств и библиотек:</b>		
AutoPsy	зарубежное	свободно распространяемое
Interactive Disassembler (IDA)	зарубежное	свободно распространяемое
<b>Системы управления библиографической информацией:</b>		
Zotero	зарубежное	свободно распространяемое
<b>Сервисы и службы:</b>		
Bind	зарубежное	свободно распространяемое
Docker	зарубежное	свободно распространяемое

## 7. Методические и оценочные материалы

### Методические указания для обучающихся по освоению дисциплины (модуля)

В процессе изучения дисциплины (модуля) «Natural Language Processing (Обработка естественного языка)» в рамках текущего контроля успеваемости используются такие виды учебной работы, как лекции, практические занятия, домашние задания, тесты, а также различные виды самостоятельной работы обучающихся по заданию преподавателя, направленные на развитие навыков профессиональной лексики, закрепление практических профессиональных компетенций, поощрение инициатив.

*Лекция* – систематическое, последовательное, монологическое изложение преподавателем учебного материала, как правило, теоретического характера.

В процессе лекций рекомендуется вести конспект лекций: кратко и схематично фиксировать основные идеи, выводы и обобщения лекции; выделять важные мысли, ключевые слова и термины. Необходимо отметить вопросы или материалы, которые вызывают затруднения, и попытаться найти ответы в рекомендованной литературе. Если разобраться в материале не удастся, следует сформулировать вопрос и задать его преподавателю на консультации или во время семинарского (практического) занятия.

*Участие в семинаре (практическом занятии)* – активная работа студента на семинаре, его ответы на вопросы преподавателя и участие в дискуссии.

Для успешного участия в семинаре студентам рекомендуется заранее ознакомиться с темой обсуждения, прочитать необходимые материалы и подготовить вопросы. Важно активно слушать и вовлекаться в дискуссию, высказывая свои мнения и аргументируя их. При ответах на вопросы преподавателя стоит быть уверенным, четким и логичным, опираясь на изученный материал. Также полезно поддерживать диалог с однокурсниками, чтобы обогатить обсуждение и расширить свои знания.

*Домашнее задание* – набор задач по темам недели.

При работе над домашними заданиями важно внимательно ознакомиться с требованиями и сроками выполнения. Рекомендуется разбивать задания на этапы, чтобы избежать перегрузки и лучше усвоить материал. Использовать различные источники информации, включая учебники и онлайн-ресурсы, для более глубокого понимания темы.

*Тест* – особая форма проверки знаний. Проводится после освоения одной или нескольких тем и свидетельствует о качестве понимания основных понятий изучаемого материала. Тестовые задания составлены к ключевым понятиям, основным разделам, важным терминологическим категориям изучаемой дисциплины (модуля).

Для подготовки к тесту необходимо знать терминологический аппарат дисциплины (модуля), понимать смысл научных категорий и уметь их использовать в профессиональной лексике. Владение понятийным аппаратом, включённым в тестовые задания, позволяет преподавателю быстро проверить уровень понимания студентами важных методологических категорий.

*Самостоятельная работа* – работа студентов, направленная на углубленное изучение отдельных тем и вопросов учебной дисциплины (модуля).

В процессе самостоятельной работы студенты взаимодействуют с рекомендованными материалами при минимальном участии преподавателя. Задачи студента включают работу с конспектами лекций (обработка текста), повторное изучение учебных материалов планов и тезисов ответов, изучение дополнительных тем, выполнение учебно-исследовательских заданий и другое.

### **Система оценивания результатов обучения по дисциплине (модулю)**

#### **Критерии получения уровня и оценивания сформированности компетенций по дисциплине (модулю) «Natural Language Processing (Обработка естественного языка)»**

Оценивание уровня учебных достижений обучающихся по дисциплине (модулю) осуществляется в виде текущего контроля успеваемости и промежуточной аттестации.

**Промежуточная аттестация** по дисциплине (модулю) осуществляется в форме *зачета*, при этом проводится оценка компетенций, сформированных по дисциплине.

Для оценивания текущего контроля успеваемости и промежуточной аттестации используется десятибалльная шкала оценивания, которая соотносится с традиционной пятибалльной шкалой следующим образом:

Десятибалльная оценка	Пятибалльная оценка	Оценка за зачет	Общая характеристика результата обучения по дисциплине (модулю)
10	Отлично	Зачтено	Студент полностью владеет знаниями, изложенными в рабочей программе, и глубоко осмысляет дисциплину (модуль). Он самостоятельно и логически последовательно отвечает на все вопросы, акцентируя внимание на наиболее важном. Умеет анализировать, сравнивать, классифицировать, обобщать, конкретизировать и систематизировать изученный материал, выделяя ключевые моменты и устанавливая причинно-следственные связи. Четко формулирует ответы, уверенно интерпретирует результаты анализов и других исследований, а также решает сложные задачи. Студент хорошо знаком с методами исследования, необходимыми для практической деятельности, и умеет связывать теоретические аспекты дисциплины (модуля) с практическими задачами.
9	Отлично	Зачтено	
8	Отлично	Зачтено	
7	Хорошо	Зачтено	Студент обладает знаниями предмета почти в полном объеме рабочей программы и самостоятельно, логически последовательно и всесторонне отвечает на все вопросы, акцентируя внимание на наиболее значимых моментах. Он умеет
6	Хорошо	Зачтено	

Десятибалльная оценка	Пятибалльная оценка	Оценка за зачет	Общая характеристика результата обучения по дисциплине (модулю)
			анализировать, сравнивать, классифицировать, обобщать, конкретизировать и систематизировать изученный материал, выделяя его ключевые аспекты и устанавливая причинно-следственные связи. Формулирует свои ответы, уверенно интерпретирует результаты анализов и других исследований, а также решает сложные ситуационные задачи. Студент хорошо знаком с методами исследования, необходимыми для практической деятельности, и умеет связывать теоретические аспекты предмета с практическими задачами.
5	Удовлетворительно	Зачтено	Студент обладает базовыми знаниями по дисциплине (модулю), но испытывает трудности при самостоятельных ответах и использует неточные формулировки. В ходе ответов он допускает ошибки, касающиеся сути вопросов. Студент способен решать только самые простые задачи и владеет лишь минимальным набором методов исследования.
4	Удовлетворительно	Зачтено	
3	Не сдан	Не зачтено	Студент не овладел обязательным минимумом знаний по предмету и не может ответить на вопросы, даже если преподаватель задает дополнительные наводящие вопросы.
2	Не сдан	Не зачтено	
1	Не сдан	Не зачтено	

Дисциплина (модуль) «Natural Language Processing (Обработка естественного языка)» оценивается следующим образом:

Активность	Вес	Количество	Описание
<b>Накопительная оценка</b>			
Домашние задания	60%	2	Набор задач по темам недели
Тесты		7	Письменная работа с набором задач, которые нужно решить за ограниченное время
<b>Промежуточная аттестация</b>			
Зачет	40%	1	Письменная или устная работа над заданием, направленным на проверку полученных знаний и навыков по дисциплине (модулю)

**Формула расчёта итоговой оценки по дисциплине (модулю) «Natural Language Processing (Обработка естественного языка)»:**  $\langle 0,6 \times \text{накопительная оценка} (0,6 \times \text{среднее за домашние задания} + 0,4 \times \text{среднее за тесты}) + 0,4 \times \text{зачет} \rangle$ .

## Текущий контроль успеваемости обучающихся по дисциплине (модулю)

### Примерные домашние задания

#### Домашнее задание 1.

##### Задание 1.

1. Выберите текстовый корпус (например, набор новостных статей или отзывов о продуктах).
2. Выполните предобработку текста, включая:
  - Удаление пунктуации и специальных символов.
  - Приведение текста к нижнему регистру.
  - Стемминг или лемматизацию.
  - Удаление стоп-слов.
3. Реализуйте метод "мешок слов" (Bag of Words) для выделения признаков из вашего корпуса.
4. Создайте таблицу, отображающую частоту слов в вашем наборе данных.

**Ожидаемый результат:** Отчет, содержащий результаты предобработки, таблицу частоты слов и краткий анализ полученных данных.

##### Задание 2.

1. Используя библиотеку Gensim, обучите модель Word2Vec на выбранном текстовом корпусе (можно использовать корпус из первого задания).
2. Проведите анализ векторных представлений, выбрав несколько слов и найдите их ближайшие "соседи" в векторном пространстве.
3. Объясните, как векторные представления слов помогают в решении задач NLP.
4. Создайте простую языковую модель на основе N-грамм, используя ваш текстовый корпус, и оцените ее производительность на тестовом наборе данных.

**Ожидаемый результат:** Отчет с результатами обучения модели Word2Vec, списком ближайших соседей для выбранных слов и кратким анализом языковой модели.

##### Задание 3.

1. Реализуйте простую рекуррентную нейронную сеть (RNN) для задачи классификации текста (например, классификация отзывов как положительных или отрицательных). Используйте фреймворк TensorFlow или PyTorch.
2. Обучите модель на размеченном наборе данных (например, IMDb для отзывов о фильмах).
3. Сравните производительность RNN с более сложной архитектурой, такой как LSTM или GRU, и проанализируйте результаты.
4. Реализуйте механизм внимания для улучшения производительности вашей модели. Опишите, как механизм внимания помогает в контексте вашей задачи.

**Ожидаемый результат:** Отчет, содержащий результаты обучения RNN и LSTM/GRU, графики производительности и анализ влияния механизма внимания на результаты.

#### Домашнее задание 2.

##### Задание 1.

1. Изучите архитектуру Трансформер, включая механизмы внимания и позиционное кодирование. Создайте визуальную схему, иллюстрирующую основные компоненты модели.

2. Реализуйте простую языковую модель на основе кодировщика Трансформера (например, BERT) с использованием библиотеки Hugging Face Transformers. Обучите модель на задаче классификации текстов (например, классификация отзывов).

3. Проведите оценку производительности модели на тестовом наборе данных и проанализируйте, какие факторы могут влиять на результаты.

**Ожидаемый результат:** Отчет с визуальной схемой архитектуры, кодом реализации модели, результатами оценки производительности и анализом факторов.

### **Задание 2.**

1. Выберите генеративную языковую модель на основе декодера Трансформера (например, GPT-2 или GPT-3). Изучите, как работает prompt tuning и реализуйте его для вашей модели.

2. Проведите эксперимент, сравнив производительность модели с и без prompt tuning на задаче генерации текста (например, продолжение предложений).

3. Исследуйте методы оптимизации, такие как P-tuning и LoRA. Реализуйте один из этих методов для улучшения производительности вашей модели и оцените результаты.

**Ожидаемый результат:** Отчет, включающий результаты экспериментов с prompt tuning, сравнение производительности, реализацию метода оптимизации и анализ полученных данных.

### **Задание 3.**

1. Реализуйте RAG (Retrieval-Augmented Generation) систему, используя предобученные модели и набор данных для извлечения информации (например, Wikipedia). Объясните, как система использует механизмы извлечения и генерации.

2. Создайте простую диалоговую систему, реализовав задачи intent detection и slot filling. Используйте библиотеку Rasa или аналогичную для реализации.

3. Проведите тестирование вашей диалоговой системы на наборе пользовательских запросов и проанализируйте результаты, включая точность определения намерений и заполнения слотов.

**Ожидаемый результат:** Отчет, содержащий описание RAG системы, код реализации диалоговой системы, результаты тестирования и анализ производительности.

## **Примерные задания по тестам**

### **Тест 1.**

#### **Вопрос 1.**

Какие из перечисленных видов векторных представлений по построению имеют разреженную структуру векторов?

- A. Fasttext
- B. One-hot
- C. word2vec
- D. Glove

Ответ: B.

#### **Вопрос 2.**

Какие преимущества имеет FastText перед word2vec? Выберите **ВСЕ** правильные варианты:

- A. Обработка OOV-слов
- B. Существенное сокращение потребления памяти
- C. Борьба с опечатками
- D. Ускоренное обучение за счет использования хэш-таблиц

Ответ: А, С.

**Вопрос 3.**

Сколько векторов поставит в соответствие многозначному слову (например, слову "замок") модель класса word2vec?

- А. один вектор
- В. столько векторов, сколько у слова смыслов
- С. количество векторов слова – это гиперпараметр, который можно задать
- Д. модель не сойдется, вектор многозначного слова будет случайно

Ответ: А.

**Вопрос 4.**

Сколько обучаемых матриц весов в простейшем варианте архитектуры CBOW?

- А. 1
- В. 2
- С. 3
- Д. 4

Ответ: В.

**Вопрос 5.**

Примените операцию Softmax к вектору (1, 2, 0, 1). Какое максимальное значение у полученного в результате применения операции вектора? Считать  $\exp(1) = 2.7$ ,  $\exp(2) = 7.4$ .

Ответ укажите с точностью до 3-х знаков после запятой, прочие удалите без округления.

- А. 0,536
- В. 0,578
- С. 1
- Д. 7,4

Ответ: А.

**Тест 2.**

**Вопрос 1.**

Выберите **ВСЕ** верные утверждения об исходной модели Трансформер, используемой для машинного перевода:

- А. В модели Трансформер есть рекуррентные слои
- В. Модель Трансформер является моделью класса кодировщик-декодировщик
- С. В модели Трансформер есть механизм внимания
- Д. В модели Трансформер есть механизм памяти
- Е. Все утверждения верны

Ответ: В, С.

**Вопрос 2.**

Позиционные эмбединги в модели Трансформер...

- А. помогают учитывать порядок слов
- В. помогают учитывать текущий номер слоя
- С. помогают учитывать текущий номер головы модели
- Д. помогают учитывать частоту слова в большом корпусе

Ответ: А.

**Вопрос 3.**

Чем отличается слой Encoder-Decoder Attention в декодировщике от Self-Attention?

- А. Матрицы запроса (query) и ключа (key) берутся из кодировщика
- В. Матрицы запроса (query) и значения (value) берутся из кодировщика
- С. Матрицы ключа (key) и значения (value) берутся из кодировщика
- Д. Входная последовательность берется из кодировщика

Ответ: С.

**Вопрос 4.**

Сколько матриц весов будет суммарно в слоях self-attention кодировщика, если в модели 8

голов внимания и 6 слоев кодировщика?

- A. 3
- B. 8
- C. 24
- D. 25
- E. 144
- F. 150

Ответ: F.

#### Вопрос 5.

Если на вход трансформерному слою кодировщика с 8 головами внимания подается  $n$  векторов слов размерности  $m$ . Какого размера матрица получится на выходе слоя?

- A.  $[n/8, m]$
- B.  $[n, m/8]$
- C.  $[n/8, m/8]$
- D.  $[n, m]$

Ответ: D.

### Тест 3.

#### Вопрос 1.

Допустим, модель ранжирования выдала 5 топовых документов на запрос  $q$  в следующем порядке:  $[doc1, doc3, doc2, doc4, doc5]$ . При этом идеальная выдача для запроса  $q$  вместе со скорями релевантности:  $[doc1 - 3, doc2 - 2, doc5 - 1, doc3 - 0, doc4 - 0]$ . Чему будет равна метрика  $nDCG@3$  для выдачи модели? Логарифмы и результаты деления округляйте до 3 знака после запятой.

- A. 0.840
- B. 0.921
- C. 4
- D. 4.762

Ответ: A.

#### Вопрос 2.

Выберите все верные утверждения:

- A. Модель T5 имеет архитектуру кодировщик-декодировщик
- B. Модель FlanT5 – это инструктивно обученная модель T5
- C. Модель T5 предназначена исключительно для задачи суммаризации
- D. Модель BART предобучалась на задаче перевода
- E. В модели BART не используется маскированная языковая модель
- F. В модели BART больше параметров, чем в модели BERT

Ответ: A, B, F.

#### Вопрос 3.

Обратимся к алгоритму экстрактивной суммаризации TextRank. Упорядочьте его шаги в правильном порядке: а) Построение графа близости между предложениями; б) Разбиение текста на предложения; в) Вычисление меры центральности в графе предложений; г) Ранжирование предложений; д) Определение близости между предложениями; е) Векторизация предложений.

- A. a b c d e f
- B. f b e d c a
- C. b f e a c d
- D. b a f e c d

Ответ: C.

#### Вопрос 4.

В теме диалоговых систем были рассмотрены задачи intent detection и slot filling. К каким двум задачам NLP сводятся эти две задачи?

- A. Классификация текста
- B. Суммаризация
- C. Распознавание именованных сущностей
- D. Извлечение отношений

Ответ: А, С.

**Вопрос 5.**

В тексте необходимо выделить сущности: локация, персона, дата, организация. Если использовать схему IOB при разметке и модель с архитектурой BERT + полносвязный слой со скрытой размерностью векторов  $n$ , то какова будет размерность последнего слоя?

- A.  $[n, 2]$
- B.  $[n, 4]$
- C.  $[n, 5]$
- D.  $[n, 9]$

Ответ: D.

**Задания для промежуточной аттестации по дисциплине (модулю)**

№ п/п	Задание	Ответ	Компетенция
1.	Для оценки эффективности своей работы по созданию модели классификации текстов, какая метрика уравнивает Precision и Recall для оценки качества модели классификации? A) Accuracy B) ROC-AUC C) F1-score D) Recall	С	УК-6
2.	Как называется процесс приведения слова к его начальной форме (например, «бежал» → «бежать»)? Ответ дайте, одним словом, на русском или английском языках	Лемматизация/ лемматизация lemmatization/ Lemmatization/ lemmatisation/ Lemmatisation	ОПК-2
3.	Какой метод численного представления текста учитывает частоту слов и их значимость в коллекции документов? Дайте ответ одним словом (возможно, с дефисом) на русском или английском языках	TF-IDF / tfidf / tf idf / тф-идф / tfidf метод	ПК-3
4.	Как называется техника, которая улучшает устойчивость модели к новым данным за счет расширения на основе изменений тренировочного набора? Напишите ответ на русском языке, одним словом.	Аугментация/ аугментация	ПК-4