

МИНИСТЕРСТВО ТРАНСПОРТА РОССИЙСКОЙ ФЕДЕРАЦИИ
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«РОССИЙСКИЙ УНИВЕРСИТЕТ ТРАНСПОРТА»
(РУТ (МИИТ))



Рабочая программа дисциплины (модуля),
как компонент образовательной программы
высшего образования - программы бакалавриата
по направлению подготовки
09.03.01 Информатика и вычислительная техника,
утвержденной первым проректором РУТ (МИИТ)
Тимониным В.С.

РАБОЧАЯ ПРОГРАММА ДИСЦИПЛИНЫ (МОДУЛЯ)

Анализ больших текстовых данных и текстовый поиск

Направление подготовки: 09.03.01 Информатика и вычислительная техника

Направленность (профиль): IT-сервисы и технологии обработки данных на транспорте

Форма обучения: Очная

Рабочая программа дисциплины (модуля) в виде электронного документа выгружена из единой корпоративной информационной системы управления университетом и соответствует оригиналу

Простая электронная подпись, выданная РУТ (МИИТ)
ID подписи: 937226
Подписал: руководитель образовательной программы
Проневич Ольга Борисовна
Дата: 14.06.2024

1. Общие сведения о дисциплине (модуле).

Целью освоения учебной дисциплины Анализ больших текстовых данных и текстовый поиск, программное обеспечение: Библиотеки и модули анализа данных Python является теоретическая и практическая подготовка студентов к работе с большими текстовыми данными и интеллектуальному анализу текста. Знания и компетенции, полученные в результате освоения дисциплины, помогут при автоматизированном интеллектуальном анализе больших объемов текстовой информации, что позволит успешно решать практические задачи обработки данных, возникающие в процессе профессиональной деятельности.

Задачи освоения дисциплины:

- приобретение студентами знаний о моделях и методах интеллектуального анализа текстовых данных и машинного обучения;
- развитие навыков программирования на языках, позволяющих анализировать текстовые данные;
- формирование представления о сборе, обработке и анализе данных в интерактивных средах;

2. Планируемые результаты обучения по дисциплине (модулю).

Перечень формируемых результатов освоения образовательной программы (компетенций) в результате обучения по дисциплине (модулю):

ОПК-3 - Способен решать стандартные задачи профессиональной деятельности на основе информационной и библиографической культуры с применением информационно-коммуникационных технологий и с учетом основных требований информационной безопасности;

ОПК-8 - Способен разрабатывать алгоритмы и программы, пригодные для практического применения;

ПК-1 - Способен анализировать большие данные с использованием существующей в организации методологической и технологической инфраструктуры.

Обучение по дисциплине (модулю) предполагает, что по его результатам обучающийся будет:

Знать:

- технологии, методы и инструменты развития компетенций в области анализа, хранения и обработки больших текстовых данных,
- технологии анализа больших текстовых данных и текстового, поиска,

- технологии обучения и дообучений больших текстовых моделей.

Уметь:

- использовать открытые источники информации (литература, интернет) для поиска актуальных средств статистической обработки и анализа больших текстовых данных,

- разрабатывать проектные решения, основанные на обработке текстовых данных,

- обучать модели формата text2text.

Владеть:

- навыками определения и соблюдения срока выполнения работ по разработке приложений, использующих большие текстовые данные,

- навыками программирования на языках, позволяющих анализировать текстовые данные,

- навыками суммаризации текста, определение тем текстов, поиска семантически близких текстов.

3. Объем дисциплины (модуля).

3.1. Общая трудоемкость дисциплины (модуля).

Общая трудоемкость дисциплины (модуля) составляет 3 з.е. (108 академических часа(ов)).

3.2. Объем дисциплины (модуля) в форме контактной работы обучающихся с педагогическими работниками и (или) лицами, привлекаемыми к реализации образовательной программы на иных условиях, при проведении учебных занятий:

Тип учебных занятий	Количество часов	
	Всего	Семестр №4
Контактная работа при проведении учебных занятий (всего):	80	80
В том числе:		
Занятия лекционного типа	32	32
Занятия семинарского типа	48	48

3.3. Объем дисциплины (модуля) в форме самостоятельной работы обучающихся, а также в форме контактной работы обучающихся с педагогическими работниками и (или) лицами, привлекаемыми к реализации образовательной программы на иных условиях, при проведении промежуточной аттестации составляет 28 академических часа (ов).

3.4. При обучении по индивидуальному учебному плану, в том числе при ускоренном обучении, объем дисциплины (модуля) может быть реализован полностью в форме самостоятельной работы обучающихся, а также в форме контактной работы обучающихся с педагогическими работниками и (или) лицами, привлекаемыми к реализации образовательной программы на иных условиях, при проведении промежуточной аттестации.

4. Содержание дисциплины (модуля).

4.1. Занятия лекционного типа.

№ п/п	Тематика лекционных занятий / краткое содержание
1	<p>Тема 1. Классификация текстов</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - удаление нерелевантных символов - токенизация - нейронная сеть для классификации
2	<p>Тема 2. Эмбединг в NLP</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - векторное представление слов - история эмбединга - слой эмбединга в нейронной сети - борьба с переобучение в сетях с эмбедингом
3	<p>Тема 3. Рекуррентные нейронные сети. Модели seq2seq</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - рекуррентный нейрон - простая рекуррентная сеть - проблема затухающего градиента - LSTM-сети - модели, основанные на энкодере и декодере
4	<p>Тема 4. Задача NER</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - постановка задачи NER - обученные NER-модели - подготовка данных для обучения решения задачи NER - лемматизация текста
5	<p>Тема 5. Нейронные сети в текстовой аналитике</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - слои свертки - слой эмбединга - архитектуры нейронных сетей
6	<p>Тема 6. Семантические сети</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - история семантических сетей и современное применения - граф знаний

№ п/п	Тематика лекционных занятий / краткое содержание
	<ul style="list-style-type: none"> - векторное представление слова в задача анализа семантики - word2vector - glove - формирование графа знаний
7	<p>Тема 7. Тематическое моделирование</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - задача Topic Modeling - методы тематического моделирования - LDA. Скрытое размещение Дирихле - Архитектуры нейронных сетей для больших лингвистических моделей
8	<p>Тема 8. LSA-анализ и суммаризация</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - LSA - Суммаризация и реферирование текстов
9	<p>Тема 9. Seq2Seq модели с вниманием</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - механизм внимания - архитектура энкодера с вниманием - архитектура с декодером с вниманием - особенности токенизации в моделях с вниманием
10	<p>Тема 10. Трансформеры</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - история появления - отличия от рекуррентных нейронных сетей - алгоритм работы трансформера - модель BER
11	<p>Тема 11. Большие лингвистические модели</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - история развития - предобученные LLM на русском языке - дообучение обученных LLM
12	<p>Тема 12. Настройка обученных больших лингвистических моделей</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - подключение к обученной LLM - основные гиперпараметры, влияющие на длину и содержание ответа
13	<p>Тема 13. Модели text2SQL</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - проблема задачи преобразования текста с естественного языка на искусственный язык - история моделей text2SQL - особенности обучения и использования text2SQL

4.2. Занятия семинарского типа.

Практические занятия

№ п/п	Тематика практических занятий/краткое содержание
1	<p>Тема 1. Токенизация</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - библиотеки для токенизации

№ п/п	Тематика практических занятий/краткое содержание
	<ul style="list-style-type: none"> - гиперпараметры для токенизации - проведения токенизации при различных значений гиперпараметров
2	<p>Тема 2. Обучение модели классификации текста</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - подготовка целевой метки - подготовка массива признаков - задание архитектуры нейронной сети - оценка качества обучения
3	<p>Тема 3. Обучении модели классификации текста с использование эмбединга</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - задание архитектуры нейронной сети с эмбединга-слоем - обучение модели - сравнение результатов работы моделей с эмбединга
4	<p>Тема 4. Рекуррентные нейронные сети</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - токенизация для рекуррентной сети - архитектуры нейронных сетей с рекуррентными слоями - сравнение результатов работы моделей
5	<p>Тема 5. Борьба с затухающим градиентом в рекуррентных нейронных сетях</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - LSTM - упрощенные модели LSTM - исследование зависимости времени работы LSTM от гиперпараметров
6	<p>Тема 6. Модель seq2seq без внимания</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - программирование энкодера - программирование декодера - обучение модели seq2seq - формирование конечной архитектуры модели seq2seq - формирование ответа модели
7	<p>Тема 7. Модель seq2seq с вниманием</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - программирование энкодера с вниманием - программирование декодера с вниманием - обучение модели seq2seq с вниманием - формирование конечной архитектуры модели seq2seq с вниманием - формирование ответа модели
8	<p>Тема 8. Задача NER</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - использование и анализ результатов обученных NER-моделей - обучение собственной NER-модели
9	<p>Тема 9. Анализ результатов работы нейронной сети</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - моделирование результатов работы сверточных слоев - моделирование результатов работы рекуррентных слоев в зависимости от гиперпараметров
10	<p>Тема 10. Семантические сети</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - векторное представление слов с помощью word2vector - векторное представление слов с glove - формирование графа знаний на примеры текстов о транспорте

№ п/п	Тематика практических занятий/краткое содержание
11	<p>Тема 11. Тематическое моделирование</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - пример выделения тем - обзор и применение библиотек для тематического моделирования
12	<p>Тема 12. LSA-анализ и суммаризация</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - проведения LSA анализа - Суммаризация и реферирование текстов
13	<p>Тема 13. Трансформеры</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - трансформеры open.ai - модель BERT
14	<p>Тема 14. Большие лингвистические модели</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - предобученные LLM на русском языке - дообучение обученных LLM
15	<p>Тема 15. Настройка обученных больших лингвистических моделей</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - подключение к обученной LLM - основные гиперпараметры, влияющие на длину и содержание ответа
16	<p>Тема 16. Модели text2SQL</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - анализ моделей text2SQL - разработка алгоритмов использования обученных моделей для формирование запроса на русском языке
17	<p>Тема 17. PolyAnalyst. Основы</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> -словари, их место в текстовой аналитики; -структура приложения для анализа текстовых данных
18	<p>Тема 18. Составление словарей для тематики РЖД</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - выделение тем - методы составления словарей
19	<p>Тема 19. Анализ тональности</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> -шкалы тональности текста -Роль весовых значений в анализе тональности
20	<p>Тема 20. PolyAnalyst. Узлы</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - Узлы импорта, экспорта данных, узлы текстовой аналитики; - узлы анализа данных, создание словарей
21	<p>Тема 21. импорт пакета данных в аналитическую платформу Polyanalyst;</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - Определение языка текстового массива;токенизация текста; - выделение ключевых слов;
22	<p>Тема 22. Применение текстовой аналитики и бизнес-процессы организации</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - описание бизнес-процессов с точки зрения анализа текстовых данных - методы анализа бизнес-процессов как объекта текстовых данных

№ п/п	Тематика практических занятий/краткое содержание
23	Тема 23. Графовые методы анализа текста Рассматриваемые вопросы: - Применение узла термин слов - построение графа связей терминов
24	Тема 24. Процесса анализа графа с целью выделения закономерностей и выделения основных связей Рассматриваемые вопросы: - выделение закономерностей в графе - выделение связей в графе
25	Тема 25. Создание словаря на основе графового анализа и применение его в графовом анализе с целью получения дополнительных данных и интерпретация результатов Рассматриваемые вопросы: - создание словаря на основе графа знаний - интерпретация результатов
26	Тема 26. Предварительный анализ проекта "Условия эксплуатации тепловозов и их обслуживание" Рассматриваемые вопросы: - определение инструментов анализа - проведение анализа
27	Тема 27. Применение диаграммы с целью выделения закономерностей и выделения основных связей Рассматриваемые вопросы: - выявление закономерностей - выделение основных связей
28	Тема 28. Оценка результатов анализа текстов Рассматриваемые вопросы: - применение OLAP таблицы и агрегирования данных по метрикам; - интерпретация результата

4.3. Самостоятельная работа обучающихся.

№ п/п	Вид самостоятельной работы
1	Поиск алгоритмов обработки данных в открытых источниках
2	Развертывание предобученных LLM
3	Участие в онлайн мастер классах и конференциях
4	Выполнение курсовой работы.
5	Подготовка к промежуточной аттестации
6	Подготовка к текущему контролю
7	Подготовка к практическим занятиям
8	Выполнение курсовой работы.
9	Подготовка к промежуточной аттестации.
10	Подготовка к текущему контролю.

4.4. Примерный перечень тем курсовых работ

1. Источники текстовых данных как внутри организаций так и её за пределами
2. Компьютерная лингвистика и Text Mining
3. Частотный анализ терминов в коллекции документов
4. Выделение наиболее значимых слов
5. Автоматическое извлечение наиболее важных тем
6. Кластеризация документов на основе сходства их содержания
7. Построение текстовых правил для категоризации
8. Кодирование текстовой информации с помощью Python
9. Предварительная обработка данных
10. Модуль для анализа данных pandas
11. Модуль для анализа данных scikit-learn
12. Модуль для анализа данных rummy
13. Построение модели данных
14. Введение в анализ текстов, базовые методы предобработки и выделения признаков
15. Неглубокие векторные представления слов
16. Классификация текстов
17. Разметка последовательности
18. Предобученные языковые модели.
19. Синтаксис в рамках грамматики зависимостей
20. Тематическое моделирование
21. Суммаризация и симплификация текстов
22. QA-системы, чат-боты
23. Графы знаний

5. Перечень изданий, которые рекомендуется использовать при освоении дисциплины (модуля).

№ п/п	Библиографическое описание	Место доступа
-------	----------------------------	---------------

1	<p>Язык программирования Python Г. Россум, Ф.Л.Дж. Дрейк, Д.С. Откидач Однотомное издание 2001</p>	НТБ (ЭЭ)
2	<p>Рашка, С. Python и машинное обучение: крайне необходимое пособие по новейшей предсказательной аналитике, обязательное для более глубокого понимания методологии машинного обучения : руководство / С. Рашка ; перевод с английского А. В. Логунова. — Москва : ДМК Пресс, 2017. — 418 с. — ISBN 978-5-97060-409-0</p>	https://e.lanbook.com/book/100905?category_pk=1556&publisher__fk=1028&ysclid=lw60n980di451663507
3	<p>Лю, Ю. Обучение с подкреплением на PyTorch. Сборник рецептов :</p>	https://e.lanbook.com/book/179493

	руководство / Ю. Лю ; перевод с английского А. А. Слинкина. — Москва : ДМК Пресс, 2020. — 282 с. — ISBN 978-5-97060- 853-1	
4	Крошемор, М. Алгоритмы обработки текста. 125 задач с решениями / М. Крошемор, Т. Лекрок, В. Риттер ; перевод с английского А. А. Слинкина. — Москва : ДМК Пресс, 2021. — 312 с. — ISBN 978-5-97060- 952-1	https://e.lanbook.com/book/241070

6. Перечень современных профессиональных баз данных и информационных справочных систем, которые могут использоваться при освоении дисциплины (модуля).

<https://habr.com/ru> - база знаний в виде статей, обзоров

<https://journal.tinkoff.ru/short/ai-for-all/> - база данных нейронных сетей

<https://vc.ru/services/916617-luchshie-neyroseti-bolshaya-podborka-iz-top-200-ii-generatorov-po-kategoriyam> - база данных нейронных сетей

<https://github.com/abalmumcu/bert-rest-api> - профессиональная платформа для командой работы над проектов (нейронная сеть bert)

<http://library.miit.ru/> - электронно-библиотечная система Научно-технической библиотеки МИИТ

<https://proglib.io/p/raspoznavanie-obektov-s-pomoshchyu-yolo-v3-na-tensorflow-2-0-2020-11-08> - профессиональная библиотека программистов

https://yandex.cloud/ru/blog/posts/2022/12/andrey-berger-and-yandex-cloud?utm_referrer=https%3A%2F%2Fyandex.ru%2F – библиотека профессиональных статей разработчиков Яндекс

<https://yandex.cloud/ru/blog> - библиотека профессиональных статей разработчиков Яндекс

<https://tproger.ru/translations/opencv-python-guide> - библиотека основных команд OpenCV

7. Перечень лицензионного и свободно распространяемого программного обеспечения, в том числе отечественного производства, необходимого для освоения дисциплины (модуля).

1 Операционная система семейства MicrosoftWindows

2 Пакет офисных программ MicrosoftOffice.

3 Язык программирования Python

8. Описание материально-технической базы, необходимой для осуществления образовательного процесса по дисциплине (модулю).

Компьютер преподавателя

Компьютеры студентов

экран для проектора, маркерная доска,

Проектор

9. Форма промежуточной аттестации:

Курсовая работа в 4 семестре.

Экзамен в 4 семестре.

10. Оценочные материалы.

Оценочные материалы, применяемые при проведении промежуточной аттестации, разрабатываются в соответствии с локальным нормативным актом РУТ (МИИТ).

Авторы:

доцент, к.н. Академии "Высшая инженерная школа"

Б.В. Игольников

доцент, к.н. Академии "Высшая инженерная школа"

О.Б. Проневич

Согласовано:

Директор

Б.В. Игольников

Руководитель образовательной программы

О.Б. Проневич

Председатель учебно-методической комиссии

Д.В. Паринов