

**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 магистратуры  
по направлению подготовки  
09.04.01 Информатика и вычислительная техника,  
утвержденной первым проректором РУТ (МИИТ)  
Тимониным В.С.

**РАБОЧАЯ ПРОГРАММА ДИСЦИПЛИНЫ (МОДУЛЯ)**

**Информационный поиск и анализ текстов**

Направление подготовки: 09.04.01 Информатика и вычислительная техника

Направленность (профиль): Искусственный интеллект и предиктивная аналитика в транспортных системах

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: 5665  
Подписал: заведующий кафедрой Нутович Вероника  
Евгеньевна  
Дата: 01.09.2024

## 1. Общие сведения о дисциплине (модуле).

Цель дисциплины «Информационный поиск и анализ текстов» заключается в освоении принципов интеллектуального анализа текста в задачах информационного поиска и извлечения информации.

В рамках дисциплины формируются знания о компьютерной лингвистике, способах предварительной обработки и преобразования языка, определении связи между словами, построении моделей классификации текста и структуре и особенностях построения поисковых систем.

На лабораторных занятиях у обучающихся формируются навыки работы с современными библиотеками интеллектуального анализа и обработки естественного языка Apache OpenNLP и FreeLing языка программирования Java и поисковых систем Elasticsearch, Apache Solr или Apache Lucene.

## 2. Планируемые результаты обучения по дисциплине (модулю).

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

**ПК-3** - Способен спроектировать, разработать, обучить, оценить и развернуть модели искусственного интеллекта в соответствии с методологией MLOps.

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

### **Знать:**

- принципы текстовой классификации;
- отличительные особенности полнотекстового, релевантного и нечетного поиска;
- отличительные особенности задач классификации для информационного поиска и извлечения информации.

### **Уметь:**

- проводить автоматизированную обработку естественного языка;
- представлять текстовую информацию в векторной плоскости представления слов и связанных выражений;
- оценивать релевантность результатов поиска на ранжированных и неранжированных данных;

### **Владеть:**

- навыком предварительной обработки и построения семантических связей между словами;

- навыком построения и настройки поисковых систем;
- навыком интегрирования и развертывания среды разработки и обучения модели искусственного интеллекта для информационного поиска.

### 3. Объем дисциплины (модуля).

#### 3.1. Общая трудоемкость дисциплины (модуля).

Общая трудоемкость дисциплины (модуля) составляет 5 з.е. (180 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	Сем. №2
Контактная работа при проведении учебных занятий (всего):	32	32
В том числе:		
Занятия лекционного типа	16	16
Занятия семинарского типа	16	16

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 148 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	<b>Анализ текста.</b> Рассматриваемые вопросы: <ul style="list-style-type: none"> <li>- представления языка – фонетика, морфология, синтаксис, семантика;</li> <li>- морфологический и синтаксический анализ;</li> <li>- обзор направлений анализа текста: классификация, машинный перевод, , извлечение ключевых слов (NER), информационный поиск, извлечение информации.</li> </ul>
2	<b>Задачи информационного поиска.</b> Рассматриваемые вопросы: <ul style="list-style-type: none"> <li>- классификация и фильтрация документов;</li> <li>- кластеризация полнотекстовых документов;</li> <li>- контекстно-зависимый анализ;</li> <li>- этапы информационного поиска.</li> </ul>
3	<b>Информационный поиск.</b> Рассматриваемые вопросы: <ul style="list-style-type: none"> <li>- полнотекстовый поиск и словари (тезариусы);</li> <li>- булевский поиск, релевантный и нечеткий поиск.</li> </ul>
4	<b>Методы нечеткого поиска.</b> Рассматриваемые вопросы: <ul style="list-style-type: none"> <li>- расширение множества признаков;</li> <li>- фрагментное индексирование (N-граммы);</li> <li>- хэширование по сигнатуре;</li> <li>- метрические деревья: аннотированное суффиксное дерево (АСД).</li> </ul>
5	<b>Индексирование и ранжирование текстов.</b> Рассматриваемые вопросы: <ul style="list-style-type: none"> <li>- инвертированный индекс;</li> <li>- латентно-семантическим индексирование (LSI);</li> <li>- вероятностный латентно-семантический анализ: латентное размещение Дирихле.</li> </ul>
6	<b>Векторизация текста.</b> Рассматриваемые вопросы: <ul style="list-style-type: none"> <li>- построение векторного набора слов: быстрое кодирование (One-Hot Encoding), вещественные векторы;</li> <li>- методы «мешка слов»: Bag of Words(BoW), Continuous Bag of Words (CBow), Deep CBow;</li> <li>- модель скип-грамм (Skip Gram);</li> <li>- статистические меры: частота слова (TF), TF/IDF;</li> <li>- семантические вложения (эмбединги);</li> </ul>
7	<b>Оценка эффективности информационного поиска.</b> Рассматриваемые вопросы: <ul style="list-style-type: none"> <li>- релевантность и оценка релевантности ;</li> <li>- критерии оценки эффективности;</li> <li>- оценка неранжированных и ранжированных результатов поиска.</li> </ul>
8	<b>Поисковые системы.</b> Рассматриваемые вопросы: <ul style="list-style-type: none"> <li>- типы информационных систем;</li> <li>- устройство и принцип работы поисковых систем;</li> <li>- понятие семантического ядра.</li> </ul>

#### 4.2. Занятия семинарского типа.

## Лабораторные работы

№ п/п	Наименование лабораторных работ / краткое содержание
1	<b>Предварительная обработка текста.</b> В результате выполнения лабораторных работ студент производит настройку среды разработки и осваивает механизмы обнаружения и выделения предложений в наборе текстовых данных с помощью библиотек OpenNLP и FreeLing языка программирования Java.
2	<b>Сегментация слов.</b> В результате выполнения лабораторных работ студент осваивает механизмы выделения связанных слов (токенов) на наборе текстовых данных с помощью библиотек OpenNLP и FreeLing языка программирования Java.
3	<b>Маркировка частей речи.</b> В результате выполнения лабораторных работ студент осваивает механизмы выделения частей речи на наборе текстовых данных с помощью библиотек OpenNLP и FreeLing языка программирования Java.
4	<b>Лемматизация слов.</b> В результате выполнения лабораторных работ студент осваивает механизмы объединения словоформ на наборе текстовых данных с помощью библиотек OpenNLP и FreeLing языка программирования Java.
5	<b>Кластеризация частей речи.</b> В результате выполнения лабораторных работ студент осваивает механизмы разделения частей речи на наборе текстовых данных с помощью библиотек OpenNLP и FreeLing языка программирования Java.
6	<b>Распознавание именованных сущностей.</b> В результате выполнения лабораторных работ студент осваивает механизмы поиска и выделения именованных объектов на наборе текстовых данных с помощью библиотек OpenNLP и FreeLing языка программирования Java.
7	<b>Полнотекстовый поиск</b> В результате выполнения лабораторных работ студент осваивает механизмы полнотекстового поиска на наборе текстовых данных с помощью библиотек Elasticsearch, Apache Solr или Apache Lucene языка программирования Java.
8	<b>Поиск по нескольким полям</b> В результате выполнения лабораторных работ студент осваивает механизмы информационного поиска по нескольким полям (терминам) на наборе текстовых данных с помощью библиотек Elasticsearch, Apache Solr или Apache Lucene языка программирования Java.

### 4.3. Самостоятельная работа обучающихся.

№ п/п	Вид самостоятельной работы
1	Изучение рекомендованной литературы.
2	Подготовка к лабораторным работам.
3	Выполнение курсового проекта.
4	Подготовка к промежуточной аттестации.
5	Подготовка к текущему контролю.

### 4.4. Примерный перечень тем курсовых проектов

1. Разработать систему релевантного информационного поиска для предметной области «Новостная лента»
2. Разработать систему релевантного информационного поиска для

предметной области «Библиотека»

3. Разработать систему релевантного информационного поиска для предметной области «Аренда жилья»

4. Разработать систему релевантного информационного поиска для предметной области «Онлайн магазин»

5. Разработать систему релевантного информационного поиска для предметной области «Электронная почта»

6. Разработать систему релевантного информационного поиска для предметной области «Веб-учебник»

7. Разработать систему релевантного информационного поиска для предметной области «Онлайн энциклопедия»

8. Разработать систему релевантного информационного поиска для предметной области «Сайт отзывов»

9. Разработать систему релевантного информационного поиска для предметной области «Поиск вакансий»

10. Разработать систему релевантного информационного поиска для предметной области «Одноклассники»

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

№ п/п	Библиографическое описание	Место доступа
1	Бенджамин Бенгфорт. Прикладной анализ текстовых данных на Python. Машинное обучение и создание приложений обработки естественного языка. – Санкт-Петербург: Питер, 2021. – 368 с. – ISBN 978-5-4461-1153-4.	<a href="https://ibooks.ru/bookshelf/365298/reading">https://ibooks.ru/bookshelf/365298/reading</a> (дата обращения: 1.11.2022). – Текст : электронный
2	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.	<a href="https://e.lanbook.com/book/140584">https://e.lanbook.com/book/140584</a> (дата обращения: 1.11.2022). – Текст : электронный
3	Даг, Т. Релевантный поиск с использованием Elasticsearch и Solr / Т. Даг, Б. Джон ; перевод с английского А. Н. Киселев. — Москва : ДМК Пресс, 2018. — 408 с. — ISBN 978-5-97060-592-9.	<a href="https://e.lanbook.com/book/111439">https://e.lanbook.com/book/111439</a> (дата обращения: 1.11.2022). – Текст : электронный
4	Бринк Х. Машинное обучение / Х. Бринк, Д. Ричардс, М. Феверолф. – Санкт-Петербург :	<a href="https://ibooks.ru/bookshelf/355472/reading">https://ibooks.ru/bookshelf/355472/reading</a> (дата обращения: 1.11.2022). – Текст :

	Питер, 2017. – 336 с. – ISBN 978-5-496-02989-6	электронный
5	Ингерсолл, Г. С. Обработка неструктурированных текстов. Поиск, организация и манипулирование / Г. С. Ингерсолл, Т. С. Мортон, Э. Л. Фэррис. — Москва : ДМК Пресс, 2015. — 414 с. — ISBN 978-5-97060-144-0.	<a href="https://e.lanbook.com/book/73069">https://e.lanbook.com/book/73069</a> (дата обращения: 1.11.2022). – Текст : электронный

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

Электронно-библиотечная система Научно-технической библиотеки РУТ(МИИТ) (<http://library.miit.ru/>)

Электронно-библиотечная система издательства «Лань» (<http://e.lanbook.com/>)

Электронно-библиотечная система [ibooks.ru](http://ibooks.ru) (<http://ibooks.ru/>)

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

Прикладное программное обеспечение

Браузер Microsoft Internet Explorer или его аналоги

Пакет офисных программ Microsoft Office или его аналоги

Среда разработки PyCharm Community Edition

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

Учебные аудитории для проведения учебных занятий, оснащенные компьютерной техникой и наборами демонстрационного оборудования.

Для практических занятий – наличие персональных компьютеров вычислительного класса.

9. Форма промежуточной аттестации:

Курсовой проект во 2 семестре.

Экзамен во 2 семестре.

10. Оценочные материалы.

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

Авторы:

заведующий кафедрой, доцент, к.н.  
кафедры «Цифровые технологии  
управления транспортными  
процессами»

В.Е. Нутович

старший преподаватель кафедры  
«Цифровые технологии управления  
транспортными процессами»

Е.А. Заманов

Согласовано:

Заведующий кафедрой ЦТУТП

В.Е. Нутович

Председатель учебно-методической  
комиссии

Н.А. Андриянова