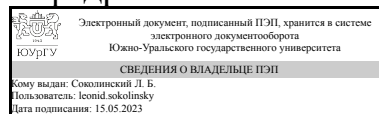


ЮЖНО-УРАЛЬСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ

УТВЕРЖДАЮ:
Заведующий выпускающей
кафедрой



Л. Б. Соколинский

РАБОЧАЯ ПРОГРАММА

дисциплины 1.Ф.М0.01 Анализ естественного языка методами искусственного интеллекта

для направления 02.04.02 Фундаментальная информатика и информационные технологии

уровень Магистратура

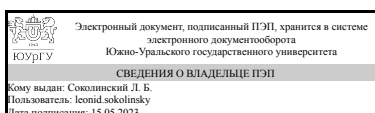
магистерская программа Машинное обучение и анализ больших данных

форма обучения очная

кафедра-разработчик Системное программирование

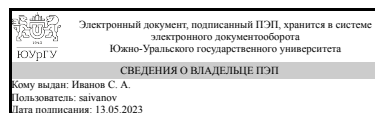
Рабочая программа составлена в соответствии с ФГОС ВО по направлению подготовки 02.04.02 Фундаментальная информатика и информационные технологии, утверждённым приказом Минобрнауки от 23.08.2017 № 811

Зав.кафедрой разработчика,
д.физ.-мат.н., проф.



Л. Б. Соколинский

Разработчик программы,
к.физ.-мат.н., доцент



С. А. Иванов

1. Цели и задачи дисциплины

Целью дисциплины является формирование базовых представлений, знаний и умений в анализе естественного языка. Основные задачи дисциплины: ознакомить студента с основными понятиями анализа и обработки текстов на естественном языке, дать понимание базовых подходов и методов при решении задач анализа естественного языка, получить практический опыт работы с различными алгоритмами машинного обучения и архитектурами искусственных нейронных сетей в рамках задач обработки естественного языка.

Краткое содержание дисциплины

Изложены наиболее важные понятия, определения и методы машинного обучения и искусственных нейронных сетей в задачах анализа естественного языка. В курс входят следующие разделы: введение в анализ естественного языка, машинное обучение и глубокие нейронные сети для решения задач анализа и обработки естественного языка, построение диалоговых систем. На практике студенты применяют навыки построения моделей машинного обучения и искусственных нейронных сетей на языке Python для решения задач морфологического анализа, классификации и кластеризации текстовых документов, анализа тональности, определения семантической близости слов, машинного перевода, построения вопросно-ответных систем, автоматического реферирования текста, построения диалоговых систем.

2. Компетенции обучающегося, формируемые в результате освоения дисциплины

Планируемые результаты освоения ОП ВО (компетенции)	Планируемые результаты обучения по дисциплине
ПК-1 Способен разрабатывать системы хранения и обработки больших данных, в том числе на основе методов искусственного интеллекта	Знает: типовые решения, библиотеки программных модулей, шаблоны, классы объектов, используемые при разработке программного обеспечения для решения задач обработки естественного языка Умеет: применять типовые решения, библиотеки программных модулей, шаблоны, классы объектов при проектировании программного обеспечения Имеет практический опыт: проектирования и реализации приложений для решения задач обработки естественного языка с использованием методов машинного обучения и нейронных сетей

3. Место дисциплины в структуре ОП ВО

Перечень предшествующих дисциплин, видов работ учебного плана	Перечень последующих дисциплин, видов работ
Глубокие нейронные сети	NoSQL-системы, Анализ и прогнозирование временных рядов методами искусственного интеллекта, Администрирование и оптимизация

Требования к «входным» знаниям, умениям, навыкам студента, необходимым при освоении данной дисциплины и приобретенным в результате освоения предшествующих дисциплин:

Дисциплина	Требования
Глубокие нейронные сети	Знает: классы задач обработки больших данных на основе методов искусственных нейронных сетей, специализированные библиотеки для создания искусственных нейронных сетей, математическую модель нейрона, технологии создания искусственных нейронных сетей, методы оптимизации, регуляризации и нормализации параметров нейронной сети и процесса ее обучения Умеет: применять современные инструментальные средства и системы программирования для разработки и обучения моделей искусственных нейронных сетей, осуществлять формализацию задачи, построение математической модели, подготовку обучающего набора данных, подбор топологии и создание искусственной нейронной сети в соответствии с поставленной задачей Имеет практический опыт: создания и обучения искусственных нейронных сетей с применением специализированных библиотек, формулирования и решения задач в области машинного обучения с использованием нейросетевого подхода

4. Объём и виды учебной работы

Общая трудоемкость дисциплины составляет 3 з.е., 108 ч., 54,5 ч. контактной работы

Вид учебной работы	Всего часов	Распределение по семестрам в часах
		Номер семестра
		2
Общая трудоёмкость дисциплины	108	108
<i>Аудиторные занятия:</i>	48	48
Лекции (Л)	16	16
Практические занятия, семинары и (или) другие виды аудиторных занятий (ПЗ)	32	32
Лабораторные работы (ЛР)	0	0
<i>Самостоятельная работа (СРС)</i>	53,5	53,5
Изучение основной и дополнительной литературы по анализу и обработке естественного языка	37,5	37,5
Подготовка к диф. зачету	16	16

Консультации и промежуточная аттестация	6,5	6,5
Вид контроля (зачет, диф.зачет, экзамен)	-	диф.зачет

5. Содержание дисциплины

№ раздела	Наименование разделов дисциплины	Объем аудиторных занятий по видам в часах			
		Всего	Л	ПЗ	ЛР
1	Введение в обработку естественного языка	6	2	4	0
2	Машинное обучение и глубокие нейронные сети для решения задач анализа и обработки естественного языка	36	12	24	0
3	Построение диалоговых систем	6	2	4	0

5.1. Лекции

№ лекции	№ раздела	Наименование или краткое содержание лекционного занятия	Кол-во часов
1	1	Введение в обработку естественного языка (NLP). Основные задачи NLP. Представления текстовых данных. Предобработка текста, лемматизация, стемминг.	2
2	2	Методы машинного обучения для классификации текстовых документов на основе частотных мер (TF-IDF). Деревья решений, наивный байесовский классификатор, логистическая регрессия в задаче классификации текстов.	2
3	2	Языковые модели. Word embeddings. Нейросетевые модели языка: word2vec, fasttext. Мера семантической близости. Классификация текстов на основе нейросетевых моделей языка.	2
4	2	Кластеризация текстовых документов. Тематическое моделирование Методы LSA, pLSA. Аддитивная регуляризация тематических моделей в BigARTM	2
5	2	Классификация текстов с помощью глубоких нейронных сетей: CNN, LSTM.	2
6	2	Задачи обработки последовательностей: машинной перевод, автоматическое реферирование (summarization), вопросно-ответные системы. Механизм внимания (attention). Архитектуры encoder-decoder-attention.	2
7	2	Transfer learning в задачах анализа текстов. Self-Attention. Архитектуры трансформеров: BERT, GPT в задачах классификации текстов, предсказания пропущенных слов, генерации текстов. Fine-tuning трансформеров.	2
8	3	Построение диалоговых систем. Архитектура диалоговых систем. Модули понимания естественного языка (NLU) и диалоговый менеджер (DM). Сложности построения диалоговых систем. Проектирование UX/UI диалоговых ассистентов в чатах и голосе. Обзор современных фреймворков для построения диалоговых систем: DeepPavlov, Rasa, Just AI Conversational Platform	2

5.2. Практические занятия, семинары

№ занятия	№ раздела	Наименование или краткое содержание практического занятия, семинара	Кол-во часов
1-2	1	Построение диалоговых систем. Архитектура диалоговых систем. Модули понимания естественного языка (NLU) и диалоговый менеджер (DM). Сложности построения диалоговых систем. Проектирование UX/UI диалоговых ассистентов в чатах и голосе. Обзор современных фреймворков для построения диалоговых систем: DeepPavlov, Rasa, Just AI Conversational	4

		Platform	
3-4	2	Классификация текстов методами машинного обучения. на основе частотных мер (TF-IDF). Деревья решений, наивный байесовский классификатор, логистическая регрессия в задаче классификации текстов.	4
5-6	2	Нейросетевые модели языка: word2vec, fasttext. Задача определения семантической близости. Классификация текстов на основе нейросетевых моделей языка.	4
7-8	2	Задача кластеризации текстовой коллекции с применением методов pLSA и фреймворка BigARTM.	4
9-10	2	Классификация текстов с помощью различных архитектур глубоких нейронных сетей: CNN, LSTM.	4
11-12	2	Реализация вопросно-ответной системы на основе нейронных сетей encoder-decoder с механизмом внимания. Решение задачи автоматического реферирования (text summarization).	4
13-14	2	Архитектуры трансформеров: BERT, GPT в задачах классификации текстов, предсказания пропущенных слов, генерации текстов.	4
15-16	3	Реализация модулей NLU и DM для чат-бота на основе одного из фреймворков диалоговых систем.	4

5.3. Лабораторные работы

Не предусмотрены

5.4. Самостоятельная работа студента

Выполнение СРС			
Подвид СРС	Список литературы (с указанием разделов, глав, страниц) / ссылка на ресурс	Семестр	Кол-во часов
Изучение основной и дополнительной литературы по анализу и обработке естественного языка	Основная литература 1, 2. Дополнительная литература 1-3.	2	37,5
Подготовка к диф. зачету	Основная литература 1, 2. Дополнительная литература 1-3	2	16

6. Фонд оценочных средств для проведения текущего контроля успеваемости, промежуточной аттестации

Контроль качества освоения образовательной программы осуществляется в соответствии с Положением о балльно-рейтинговой системе оценивания результатов учебной деятельности обучающихся.

6.1. Контрольные мероприятия (КМ)

№ КМ	Се-местр	Вид контроля	Название контрольного мероприятия	Вес	Макс. балл	Порядок начисления баллов	Учитывается в ПА
1	2	Текущий контроль	ПЗ-1. Реализация собственного POS-тэггера.	4	3	3 балла: задание выполнено полностью, 2 балла: задание выполнено полностью, но допущены	дифференцированный зачет

						незначительные ошибки, или задание выполнено более, чем 50%, 1 балла: задание выполнено полностью, но допущены серьезные ошибки, или задание выполнено менее, чем 50%, 0 баллов: задание не выполнено	
2	2	Текущий контроль	ПЗ-2. Классификация текстов методами машинного обучения	3	3	3 балла: задание выполнено полностью, 2 балла: задание выполнено полностью, но допущены незначительные ошибки, или задание выполнено более, чем 50%, 1 балла: задание выполнено полностью, но допущены серьезные ошибки, или задание выполнено менее, чем 50%, 0 баллов: задание не выполнено	дифференцированный зачет
3	2	Текущий контроль	ПЗ-3. Классификация текстов на основе нейросетевых моделей языка	4	3	3 балла: задание выполнено полностью, 2 балла: задание выполнено полностью, но допущены незначительные ошибки, или задание выполнено более, чем 50%, 1 балла: задание выполнено полностью, но допущены серьезные ошибки, или задание выполнено менее, чем 50%, 0 баллов: задание не выполнено	дифференцированный зачет
4	2	Текущий контроль	ПЗ-4. Кластеризация текстовой коллекции методами тематического моделирования	4	3	3 балла: задание выполнено полностью, 2 балла: задание выполнено полностью, но допущены незначительные ошибки, или задание выполнено более, чем 50%, 1 балла: задание выполнено полностью, но допущены серьезные ошибки, или задание выполнено менее, чем	дифференцированный зачет

						50%, 0 баллов: задание не выполнено	
5	2	Текущий контроль	ПЗ-5. Классификация текстов с помощью различных архитектур глубоких нейронных сетей	5	3	3 балла: задание выполнено полностью, 2 балла: задание выполнено полностью, но допущены незначительные ошибки, или задание выполнено более, чем 50%, 1 балла: задание выполнено полностью, но допущены серьезные ошибки, или задание выполнено менее, чем 50%, 0 баллов: задание не выполнено	дифференцированный зачет
6	2	Текущий контроль	ПЗ-6. Реализация вопросно-ответной системы. Решение задачи автоматического реферирования	6	3	3 балла: задание выполнено полностью, 2 балла: задание выполнено полностью, но допущены незначительные ошибки, или задание выполнено более, чем 50%, 1 балла: задание выполнено полностью, но допущены серьезные ошибки, или задание выполнено менее, чем 50%, 0 баллов: задание не выполнено	дифференцированный зачет
7	2	Текущий контроль	ПЗ-7. Архитектуры трансформеров: BERT, GPT в задачах классификации текстов, предсказания пропущенных слов, генерации текстов.	5	3	3 балла: задание выполнено полностью, 2 балла: задание выполнено полностью, но допущены незначительные ошибки, или задание выполнено более, чем 50%, 1 балла: задание выполнено полностью, но допущены серьезные ошибки, или задание выполнено менее, чем 50%, 0 баллов: задание не выполнено	дифференцированный зачет
8	2	Текущий контроль	ПЗ-8. Разработка диалогового агента для чат-бота	8	3	3 балла: задание выполнено полностью, 2 балла: задание выполнено полностью, но допущены	дифференцированный зачет

						незначительные ошибки, или задание выполнено более, чем 50%, 1 балла: задание выполнено полностью, но допущены серьезные ошибки, или задание выполнено менее, чем 50%, 0 баллов: задание не выполнено	
9	2	Текущий контроль	Введение в анализ естественного языка. Представления текстовых данных. Предобработка.	5	5	Компьютерный тест состоит из 5 вопросов, позволяющих оценить сформированность компетенций. На ответы отводится 7 мин. Стоимость одного вопроса - 1 балл.	дифференцированный зачет
10	2	Текущий контроль	Трансферлернинг. Обучение модели BERT. Классификация документов.	5	5	Компьютерный тест состоит из 5 вопросов, позволяющих оценить сформированность компетенций. На ответы отводится 7 мин. Стоимость одного вопроса - 1 балл.	дифференцированный зачет
11	2	Промежуточная аттестация	Итоговый тест	-	20	Компьютерный тест состоит из 20 вопросов, позволяющих оценить сформированность компетенций. На ответы отводится 1 час. 20 баллов: задание полностью выполнено без ошибок 1-19 баллов: задание выполнено частично или выполнено с ошибками 0 баллов: задание не выполнено	дифференцированный зачет

6.2. Процедура проведения, критерии оценивания

Вид промежуточной аттестации	Процедура проведения	Критерии оценивания
дифференцированный зачет	<p>При оценивании результатов учебной деятельности обучающегося по дисциплине используется балльно-рейтинговая система оценивания результатов учебной деятельности обучающихся (Положение о БРС утверждено приказом ректора от 24.05.2019 г. № 179, в редакции приказа ректора от 10.03.2022 г. № 25-13/09).</p> <p>Оценка за дисциплину формируется на основе полученных оценок за контрольно-рейтинговые мероприятия текущего контроля. Отлично: Величина рейтинга обучающегося по дисциплине 85...100 %.</p>	В соответствии с пп. 2.5, 2.6 Положения

	<p>Хорошо: Величина рейтинга обучающегося по дисциплине 75...84 %. Удовлетворительно: Величина рейтинга обучающегося по дисциплине 60...74 %. Неудовлетворительно: Величина рейтинга обучающегося по дисциплине 0...59 %. Если студент не согласен с оценкой, полученной по результатам текущего контроля, студент проходит мероприятие промежуточной аттестации в виде тестирования. Тестирование проводится в системе edu.susu.ru. Тест содержит 20 вопросов. На выполнение теста дается 60 минут. В этом случае оценка за дисциплину рассчитывается на основе полученных оценок за контрольно-рейтинговые мероприятия текущего контроля и промежуточной аттестации. Фиксация результатов учебной деятельности по дисциплине проводится в день зачёта при личном присутствии студента.</p>	
--	---	--

6.3. Паспорт фонда оценочных средств

Компетенции	Результаты обучения	№ КМ										
		1	2	3	4	5	6	7	8	9	10	11
ПК-1	Знает: типовые решения, библиотеки программных модулей, шаблоны, классы объектов, используемые при разработке программного обеспечения для решения задач обработки естественного языка	+	+	+	+	+	+	+	+	+	+	+
ПК-1	Умеет: применять типовые решения, библиотеки программных модулей, шаблоны, классы объектов при проектировании программного обеспечения	+	+	+	+	+	+	+	+	+	+	+
ПК-1	Имеет практический опыт: проектирования и реализации приложений для решения задач обработки естественного языка с использованием методов машинного обучения и нейронных сетей	+	+	+	+	+	+	+	+	+	+	+

Типовые контрольные задания по каждому мероприятию находятся в приложениях.

7. Учебно-методическое и информационное обеспечение дисциплины

Печатная учебно-методическая документация

а) *основная литература:*

Не предусмотрена

б) *дополнительная литература:*

Не предусмотрена

в) *отечественные и зарубежные журналы по дисциплине, имеющиеся в библиотеке:*

Не предусмотрены

г) *методические указания для студентов по освоению дисциплины:*

1. Гольдберг Й. Нейросетевые методы в обработке естественного языка (<https://e.lanbook.com/book/131704>) Москва : ДМК Пресс, 2019

из них: учебно-методическое обеспечение самостоятельной работы студента:

Электронная учебно-методическая документация

№	Вид литературы	Наименование ресурса в электронной форме	Библиографическое описание
1	Основная литература	Электронно-библиотечная система издательства Лань	Ганегедара, Т. Обработка естественного языка с TensorFlow : руководство / Т. Ганегедара ; перевод с английского В. С. Яценкова. — Москва : ДМК Пресс, 2020. — 382 с. https://e.lanbook.com/book/140584
2	Основная литература	Электронно-библиотечная система издательства Лань	Гольдберг, Й. Нейросетевые методы в обработке естественного языка : руководство / Й. Гольдберг ; перевод с английского А. А. Слинкина. — Москва : ДМК Пресс, 2019. — 282 с. https://e.lanbook.com/book/131704
3	Основная литература	Электронно-библиотечная система издательства Лань	Антонио, Д. Библиотека Keras – инструмент глубокого обучения. Реализация нейронных сетей с помощью библиотек Theano и TensorFlow / Д. Антонио, П. Суджит ; перевод с английского А. А. Слинкин. — Москва : ДМК Пресс, 2018. — 294 с. https://e.lanbook.com/book/111438
4	Дополнительная литература	Электронно-библиотечная система издательства Лань	Паттерсон, Д. Глубокое обучение с точки зрения практика / Д. Паттерсон, А. Гибсон. — Москва : ДМК Пресс, 2018. — 418 с. https://e.lanbook.com/book/116122
5	Дополнительная литература	Электронно-библиотечная система издательства Лань	Гудфеллоу, Я. Глубокое обучение / Я. Гудфеллоу, И. Бенджио, А. Курвилль ; перевод с английского А. А. Слинкина. — 2-е изд. — Москва : ДМК Пресс, 2018. — 652 с. https://e.lanbook.com/book/107901
6	Дополнительная литература	Электронно-библиотечная система издательства Лань	Коэльо, Л. П. Построение систем машинного обучения на языке Python / Л. П. Коэльо, В. Ричарт ; перевод с английского А. А. Слинкин. — 2-е изд. — Москва : ДМК Пресс, 2016. — 302 с. https://e.lanbook.com/book/82818

Перечень используемого программного обеспечения:

1. Python Software Foundation-Python (бессрочно)

Перечень используемых профессиональных баз данных и информационных справочных систем:

Нет

8. Материально-техническое обеспечение дисциплины

Вид занятий	№ ауд.	Основное оборудование, стенды, макеты, компьютерная техника, предустановленное программное обеспечение, используемое для различных видов занятий
Лекции	434 (36)	Компьютер, проектор
Практические занятия и семинары	114-1 (2)	Компьютерный класс, имеется выход в Интернет