Министерство науки и высшего образования Российской Федерации НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ ТОМСКИЙ ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ (НИ ТГУ)

Институт прикладной математики и компьютерных наук

УТВЕРЖДЕНО: Директор А. В. Замятин

Рабочая программа дисциплины

Обработка естественного языка

по направлению подготовки

02.03.02 Фундаментальная информатика и информационные технологии

Направленность (профиль) подготовки: Искусственный интеллект и разработка программных продуктов

> Форма обучения **Очная**

Квалификация **Бакалавр**

Год приема **2025**

СОГЛАСОВАНО: Руководитель ОП А.В. Замятин

Председатель УМК С.П. Сущенко

1. Цель и планируемые результаты освоения дисциплины

Целью освоения дисциплины является формирование следующих компетенций:

- ОПК-2. Способен применять компьютерные/суперкомпьютерные методы, современное программное обеспечение, в том числе отечественного происхождения, для решения задач профессиональной деятельности.
- ПК-1. Способен осуществлять программирование, тестирование и опытную эксплуатацию ИС с использованием технологических и функциональных стандартов, современных моделей и методов оценки качества и надежности программных средств.

Результатами освоения дисциплины являются следующие индикаторы достижения компетенций:

- ИОПК-2.1. Обладает необходимыми знаниями основных концепций современных вычислительных систем.
- ИОПК-2.2. Использует методы высокопроизводительных вычислительных технологий, современного программного обеспечения, в том числе отечественного происхождения.
- ИОПК-2.3. Использует инструментальные средства высокопроизводительных вычислений в научной и практической деятельности.
- ИПК-1.3. Кодирует на языках программирования и проводит модульное тестирование ИС.

2. Задачи освоения дисциплины

- Освоить классические методы анализа текста на естественном языке;
- Получить понимание основ векторного представления слов и применение его на практике;
- Получить понимание основ и применения на практике рекуррентной модели долгой и краткосрочной памяти;
- Получить понимание основ трансформернной модели и как применять ее на практике.

3. Место дисциплины в структуре образовательной программы

Дисциплина относится к Блоку 1 «Дисциплины (модули)».

Дисциплина относится к части образовательной программы, формируемой участниками образовательных отношений, предлагается обучающимся на выбор. Дисциплина входит в модуль «Искусственный интеллект».

4. Семестр(ы) освоения и форма(ы) промежуточной аттестации по дисциплине

Седьмой семестр, зачет

5. Входные требования для освоения дисциплины

Для успешного освоения дисциплины требуются компетенции, сформированные в ходе освоения образовательных программ предшествующего уровня образования.

Для успешного освоения дисциплины требуются результаты обучения по следующим дисциплинам: «Основы программирования», «Алгоритмы и структуры данных», «Интеллектуальные системы», «Визуализация многомерных данных», «Статистические методы машинного обучения».

6. Язык реализации

Русский

7. Объем дисциплины

Общая трудоемкость дисциплины составляет 3 з.е., 108 часов, из которых:

- -лекции: 16 ч.
- -практические занятия: 32 ч.
 - в том числе практическая подготовка: 32 ч.

Объем самостоятельной работы студента определен учебным планом.

8. Содержание дисциплины, структурированное по темам

Тема 1. Введение, история развития дисциплины, решаемые задачи, подходы, методы и инструменты

Раскрываются три основных этапа развития технологий обработки естественного языка: словарные, вероятностные и интеллектуальные алгоритмы. Даётся классификация задач. Описываются основные методы реализации алгоритмов: локальные, облачные сервисы.

Тема 2. Предварительна обработка текстовых данных

Поясняются назначение и типы предварительной обработки текста: сегментация, токенизация, лемматизация. Сравниваются лемматизация и стемминг. Поясняется роль лемматизации в построении поисковых. индексов. Пояснение недетерминированности проведения сегментации и токенизации.

Тема 3. Вероятностные алгоритмы

Приводятся основные черты вероятностных алгоритмов. Поясняется их роль в современных системах. В качестве примеров приводятся скрытые марковские модели, алгоритм Витерби, ЕМ-алгоритм. Для описания ЕМ-алгоритма поясняется назначение тематического моделирования.

Тема 4. Формальные грамматики

Определение аналитических формальных грамматик по Хомскому. Раскрытие их особенностей и принципиальных ограничений. Примеры задач, которые в настоящий момент можно решать при помощи формальных грамматик. Пояснений функций утилиты Томита-парсер.

Тема 5. Векторное представление слов

Поясняется идея замены слов точками в векторном пространстве. Приводятся примеры алгебраических операций над словами, заменёнными. точками. Определение семантической близости слов через метрики в векторном пространстве. Способы получения векторного представления. Модель Word2vec.

Тема 6. Модель Seq2seq

Пояснение преобразования последовательностей через рекуррентные ячейки. Понятия кодера и декодера. Идея долгой краткосрочной памяти. Идея дополнения кодера и декодера связью через механизм внимания.

Tema 7. Self-attention и Трансформер

Обоснование недостатков модели Seq2seq. Введение понятия Self-attention и пояснение его преимуществ. Назначение ячеек query, key и value. Описание модели Трансформер. Основные преимущества. Описание структуры кодера и декодера Трансформера.

Тема 8. BERT и GPT-3

Описание возможностей построения новых моделей на трансформере. Раздельное использование кодера и декодера. Модель BERT. Идея fine tuning. Модель GPT-3. Применение GPT-3 в практических задачах.

9. Текущий контроль по дисциплине

Текущий контроль по дисциплине проводится путем проверки выполнения практических работ и фиксируется в форме контрольной точки не менее одного раза в семестр.

Оценочные материалы текущего контроля размещены на сайте ТГУ в разделе «Информация об образовательной программе» - https://www.tsu.ru/sveden/education/eduop/.

10. Порядок проведения и критерии оценивания промежуточной аттестации

Зачет в седьмом семестре проводится в устной форме по билетам. Билет состоит из трех вопросов. Продолжительность зачета 1 час.

Оценочные материалы для проведения промежуточной аттестации размещены на сайте ТГУ в разделе «Информация об образовательной программе» - https://www.tsu.ru/sveden/education/eduop/.

11. Учебно-методическое обеспечение

- а) Электронный учебный курс по дисциплине в LMS iDo.
- б) Оценочные материалы текущего контроля и промежуточной аттестации по дисциплине.

12. Перечень учебной литературы и ресурсов сети Интернет

- а) основная литература:
- 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.
- Автоматическая обработка текстов на естественном языке и анализ данных: учеб. пособие / Большакова Е.И., Воронцов К.В., Ефремова Н.Э., Клышинский Э.С., Лукашевич Н.В., Сапин А.С. М.: Изд-во НИУ ВШЭ, 2017. 269 с.
- Введение в когнитивную лингвистику: учебное пособие. Изд. 2-е, перераб. Калининград: Изд-во БФУ им. И. Канта, 2012. 313 с.
- Николенко С., Кадурин А., Архангельская Е. Глубокое обучение. СПб.: Питер, 2018. 480 с.: ил. (Серия «Библиотека программиста»).
- Хобсон Лейн, Ханнес Хапке, Коул Ховард Обработка естественного языка в действии. СПб.: Питер, 2020. 576 с.: ил. (Серия «Для профессионалов»).

б) дополнительная литература:

- Li Deng Yang Liu Deep Learning in Natural Language Processing. ISBN 978-981-10-5209-5 https://doi.org/10.1007/978-981-10-5209-5
- Николаев И.С., Митренина О.В., Ландо Т.М. Прикладная и КОМПЬЮТЕРНАЯ ЛИНГВИСТИКА. URSS. 2017. 320 с. ISBN 978-5-9710-4633-2.
- Ян Гудфеллоу, Иошуа Бенджио, Аарон Курвилль. Глубокое обучение. Второе цветное издание, исправленное. М.: ДМК Пресс, 2018. 652 с.
 - Франсуа Шолле. Глубокое обучение на Python. СПб: Питер, 2018. 400 с.
- Daniel Jurafsky, James H. Martin. Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition. Prentice Hall, $2008.-1044\ c$

- в) ресурсы сети Интернет:
- открытые онлайн-курсы

13. Перечень информационных технологий

- а) лицензионное и свободно распространяемое программное обеспечение:
- Microsoft Visual Studio;
- публично доступные облачные технологии (Google Docs, Яндекс диск и т.п.).
- б) информационные справочные системы:

14. Материально-техническое обеспечение

— Электронный каталог Научной библиотеки ТГУ — http://chamo.lib.tsu.ru/search/query?locale=ru&theme=system
— Электронная библиотека (репозиторий) ТГУ — http://cit.ll. библиотека (репозиторий) ТГУ — http://chamo.lib.tsu.ru/search/query?locale=ru&theme=system

http://vital.lib.tsu.ru/vital/access/manager/Index

Аудитории для проведения занятий лекционного типа.

Аудитории для проведения практических занятий, индивидуальных и групповых консультаций, текущего контроля и промежуточной аттестации.

Помещения для самостоятельной работы, оснащенные компьютерной техникой и доступом к сети Интернет, в электронную информационно-образовательную среду и к информационным справочным системам.

15. Информация о разработчиках

Пожидаев Михаил Сергеевич, канд. техн. наук, доцент кафедры теоретических основ информатики

Бакланова Ольга Евгеньевна, канд. физ-мат. наук, доцент кафедры теоретических основ информатики