

Документ подписан простой электронной подписью  
Информация о владельце:  
ФИО: Смирнов Сергей Николаевич  
Должность: врио ректора  
Дата подписания: 17.10.2024 16:18:59  
Уникальный программный ключ:  
69e375c64f7e975d4e8830e7b4fcc2ad1bf35f08

Министерство науки и высшего образования  
Российской Федерации  
Федеральное государственное бюджетное  
образовательное учреждение высшего образования  
«Тверской государственный университет»



Рабочая программа дисциплины (с аннотацией)  
**Основы компьютерной лингвистики**

Направление подготовки  
03.02 — Прикладная математика и информатика

Профиль подготовки  
ИСКУССТВЕННЫЙ ИНТЕЛЛЕКТ И АНАЛИЗ ДАННЫХ

для студентов 4 курса  
ФОРМА ОБУЧЕНИЯ — очная

Составитель(и):  
• к.ф.-м.н. доц. Карлов Б.Н.

Тверь — 2023

## **I. Аннотация**

### **1. Цель и задачи дисциплины:**

Цель курса — ознакомить студентов с основными понятиями компьютерной лингвистики, с различными способами задания языков, с возможностью применения ЭВМ для обработки естественных языков.

### **2. Место дисциплины в структуре ООП**

Дисциплина входит в раздел «Элективные дисциплины» части, формируемой участниками образовательных отношений, блока 1.

**Предварительные знания и навыки.** Знание курсов «Математическая логика и теория алгоритмов», «Теория автоматов и формальных языков», «Теория вероятностей и математическая статистика».

**Дальнейшее использование.** Полученные знания используются для итоговой государственной аттестации, прохождении практики, а также в дальнейшей трудовой деятельности выпускников.

### **3. Объем дисциплины: 3 зачетных единицы, 108 академических часа, в том числе:**

**контактная аудиторная работа:** лекции 30 часов, практические занятия 30 часов;

**контактная внеаудиторная работа:** контроль самостоятельной работы 0 часов, в том числе курсовая работа 0 часов;

**самостоятельная работа:** 48 часов, в том числе контроль 0 часов.

### **4. Перечень планируемых результатов обучения по дисциплине, соотнесенных с планируемыми результатами освоения образовательной программы**

Планируемые результаты освоения образовательной программы (формируемые компетенции)	Планируемые результаты обучения по дисциплине
ПК-3, Способен осуществлять концептуальное моделирование проблемной области и проводить формализацию представления знаний в системах искусственного интеллекта	ПК-3.1, Разрабатывает концептуальную модель проблемной области системы искусственного интеллекта
ПК-5, Способен использовать инструментальные средства для решения задач машинного обучения	ПК-5.1, Осуществляет оценку и выбор инструментальных средств для решения поставленной задачи

ПК-7 Способен осуществлять сбор и подготовку данных для систем искусственного интеллекта

ПК-7.1 Осуществляет поиск данных в открытых источниках, специализированных библиотеках и репозиториях

**5. Форма промежуточной аттестации и семестр прохождения:** зачёт в 7 семестре.

**6. Язык преподавания:** русский.

**II. Содержание дисциплины, структурированное по темам (разделам) с указанием отведенного на них количества академических часов и видов учебных занятий**

Учебная программа – наименование разделов и тем	Всего (час.)	Контактная работа (час.)					Самостоя- тельная работа, в том числе Контроль (час.)
		Лекции		Практиче- ские заня- тия		Контроль сам. раб., в т.ч. курсовая ра- бота	
		Всего	в т.ч. прак- тическая подготовка	Всего	в т.ч. прак- тическая		
1	2	3	4	5	6	7	8
Лексический анализ	24	6		6/0			12
Статистические мето- ды обработки языка	28	8		8/0			12
Синтаксический ана- лиз	28	8		8/0			12
Применение нейрон- ных сетей для обра- ботки языка	28	8		8/0			12
<b>Итого</b>	<b>108</b>	<b>30</b>	<b>0</b>	<b>30/0</b>	<b>0</b>	<b>0</b>	<b>48</b>

### Учебная программа дисциплины

#### 1. Лексический анализ.

- 1) Конечные автоматы и конечные преобразователи.
- 2) Алгоритм моделирования недетерминированного конечного автомата.
- 3) Префиксные деревья. Представление множества слов в виде префиксного дерева.
- 4) Применение конечных преобразователей для морфологического анализа слов.

- 5) Алгоритм Портера.
- 6) Исправление орфографических ошибок. Редакционное расстояние. Алгоритм Вагнера-Фишера.
2. Статистические методы обработки языка.
  - 1) Условные вероятности. Формула Байеса.
  - 2) N-граммы. Сглаживание Лапласа, Уиттена-Белла, Гуда-Тьюринга.
  - 3) Применение N-грамм для исправления орфографических ошибок с учётом контекста.
  - 4) Цепи Маркова, скрытые марковские модели. Алгоритм Витерби.
  - 5) Определение частей речи на основе правил и на основе скрытых марковских моделей.
  - 6) Модели word2vec и GloVe.
  - 7) Оценка качества моделей. Энтропия, перекрёстная энтропия.
3. Синтаксический анализ.
  - 1) Порождающие грамматики. Иерархия Хомского.
  - 2) Контекстно-зависимые и контекстно-свободные грамматики. Деревья вывода.
  - 3) Эквивалентность контекстно-зависимых грамматик и линейно-ограниченных автоматов.
  - 4) Классические категориальные грамматики. Эквивалентность КС-грамматик и классических категориальных грамматик.
  - 5) Системы составляющих и деревья зависимостей. Связь деревьев зависимостей и систем составляющих.
  - 6) Синтаксический анализ на основе КС-грамматик. Алгоритмы Кока-Янгера-Касами и Эрли.
  - 7) Слабо-контекстные грамматики. Множественные контекстно-свободные грамматики.
  - 8) Головные грамматики, линейные индексные грамматики, комбинаторные категориальные грамматики, ТАГ-грамматики.
  - 9) Категориальные грамматики зависимостей (КГЗ). Алгоритм анализа КГЗ.
  - 10) Основы  $\lambda$ -исчисления.  $\lambda$ -термы,  $\beta$ -редукция, нормальная форма. Формулировка теоремы Чёрча-Россера.
  - 11) Применения комбинаторных категориальных грамматик для представления семантики в виде  $\lambda$ -термов.
4. Применение нейронных сетей для обработки языка.
  - 1) Искусственные нейроны. Искусственные нейронные сети. Функции активации.
  - 2) Метод обратного распространения ошибки.
  - 3) Рекуррентные нейронные сети.
  - 4) Архитектура LSTM.
  - 5) Использование нейронных сетей в компьютерной лингвистике: модели языка, машинный перевод.

### III. Образовательные технологии

Учебная программа – наименование разделов и тем	Вид занятия	Образовательные технологии
Лексический анализ	лекции, практические занятия	изложение теоретического материала, решение задач
Статистические методы обработки языка	лекции, практические занятия	изложение теоретического материала, решение задач
Синтаксический анализ	лекции, практические занятия	изложение теоретического материала, решение задач
Применение нейронных сетей для обработки языка	лекции, практические занятия	изложение теоретического материала, решение задач

### IV. Оценочные материалы для проведения текущей и промежуточной аттестации

#### Типовые контрольные задания и/или критерии для проверки индикатора ПК-3.1

Требования к обучающемуся	Типовые контрольные задания для оценки знаний, умений, навыков	Показатели и критерии оценивания, шкала оценивания
Знать понятия порождающей грамматики, системы составляющих, дерева зависимостей, связь систем составляющих и деревьев зависимостей с грамматиками	<p>Примеры вопросов к зачёту:</p> <ul style="list-style-type: none"> <li>• Порождающие грамматики. Иерархия Хомского.</li> <li>• Контекстно-зависимые и контекстно-свободные грамматики. Деревья вывода.</li> <li>• Эквивалентность контекстно-зависимых грамматик и линейно-ограниченных автоматов.</li> <li>• Классические категориальные грамматики. Эквивалентность КС-грамматик и классических категориальных грамматик.</li> <li>• Системы составляющих и деревья зависимостей. Связь деревьев зависимостей и систем составляющих.</li> <li>• Синтаксический анализ на основе КС-грамматик. Алгоритмы Кока-Янгера-Касами и Эрли.</li> <li>• Слабо-контекстные грамматики.</li> </ul>	оценка 3 — знает определения основных понятий, оценка 4 — кроме того знает основные свойства порождающих грамматик, систем составляющих и деревьев зависимостей, оценка 5 — кроме того знает алгоритмы анализа

Требования к обучающемуся	Типовые контрольные задания для оценки знаний, умений, навыков	Показатели и критерии оценивания, шкала оценивания
	<p>Множественные контекстно-свободные грамматики.</p> <ul style="list-style-type: none"> <li>• Головные грамматики, линейные индексные грамматики, комбинаторные категориальные грамматики, ТАГ-грамматики.</li> <li>• Категориальные грамматики зависимостей (КГЗ). Алгоритм анализа КГЗ.</li> <li>• Основы <math>\lambda</math>-исчисления. <math>\lambda</math>-термы, <math>\beta</math>-редукция, нормальная форма. Формулировка теоремы Чёрча-Россера.</li> <li>• Применения комбинаторных категориальных грамматик для представления семантики в виде <math>\lambda</math>-термов.</li> </ul>	
<p>Уметь строить различные грамматики и автоматы по описанию языка</p>	<p>Примеры задач для контрольных работ:</p> <ul style="list-style-type: none"> <li>• Постройте детерминированный конечный преобразователь для выполнения следующей операции. На вход подаётся слово <math>w\\$</math> в алфавите <math>\{a,b,c,\\$\}</math>, символ <math>\\$</math> помечает конец входного слова. Требуется вставить символ <math>c</math> после каждого блока символов <math>a</math> нечётной длины.</li> <li>• Постройте ТАГ-грамматику для языка <math display="block">L = \{ a^i b^j c^k : 0 &lt; i &lt; j &lt; k \}.</math> </li> <li>• Постройте комбинаторную категориальную грамматику для языка <math display="block">L = \{ w * w^{-1} * w : w \in \{ 0, 1 \}^* \}.</math> </li> <li>• Постройте контекстно-свободную грамматику для языка <math display="block">L = \{ a^i b^j c^k d^l : i &lt; k \text{ или } j \geq l \}.</math> </li> <li>• Постройте классическую категориальную грамматику для языка <math display="block">L = \{ a^i b^j c^k : i, j, k \geq 0, i = j \text{ или } j = k \}.</math> </li> <li>• Постройте категориальную грамматику зависимостей для языка</li> </ul>	<p>оценка 3 — умеет строить конечные автоматы и конечные преобразователи по описанию языка или отношения, оценка 4 — кроме того умеет строить контекстно-свободные и классические категориальные грамматики по описанию языка, оценка 5 — кроме того умеет слабо контекстные грамматики по описанию языка</p>

Требования к обучающемуся	Типовые контрольные задания для оценки знаний, умений, навыков	Показатели и критерии оценивания, шкала оценивания
	<p><math>L = \{ w \in \{ a, b, c, d \}^+ :  w _a +  w _c \geq  w _b +  w _d \text{ и все символы } a \text{ стоят правее } c \}</math>.</p> <p>Примеры тем для самостоятельной работы:</p> <ul style="list-style-type: none"> <li>• Напишите программу, реализующую алгоритм Портера для английского и русского языков.</li> <li>• Напишите программу, которая по контекстно-свободной грамматике и входному слову определяет, выводимо ли слово в грамматике. Реализуйте алгоритмы Кока-Янгера-Касами и Эрли.</li> <li>• Напишите программу, которая по классической категориальной грамматике и предложению на русском языке строит все его размеченные деревья зависимостей.</li> </ul>	
<p>Уметь строить формальное представление синтаксиса и семантики предложений на естественных языках</p>	<p>Примеры задач для контрольных работ:</p> <ul style="list-style-type: none"> <li>• Постройте размеченную иерархизованную систему составляющих для предложения «К востоку от боровых озёр лежат громадные мещёрские болота — мшары».</li> <li>• Постройте размеченное дерево зависимостей для предложения «Для точной диагностики заболеваний внутренних органов человека рентген незаменим».</li> <li>• Дано предложение «Мастер внимательно осматривал станок». Его словам сопоставлены следующие категории и <math>\lambda</math>-термы:  — мастер — NP : M  — внимательно — C : V  — осматривал — <math>((S \setminus NP) \setminus C) / NP</math>  : <math>\lambda x y z. Ozx y</math>  — станок — NP : C</li> </ul>	<p>оценка 3 — умеет строить системы составляющих и деревья зависимостей для предложений, оценка 4 — кроме того умеет строить размеченные иерархизованные системы составляющих и размеченные деревья зависимостей, оценка 5 — кроме того умеет представлять семантику в виде <math>\lambda</math>-термов</p>

Требования к обучающемуся	Типовые контрольные задания для оценки знаний, умений, навыков	Показатели и критерии оценивания, шкала оценивания
	Сократите категории до S и упростите получающийся $\lambda$ -терм. На первом шаге примените к первой категории правило ( $> T$ ), а на втором шаге примените к третьей и четвертой категориям правило ( $>$ ).	

### Типовые контрольные задания и/или критерии для проверки индикатора ПК-5.1

Требования к обучающемуся	Типовые контрольные задания для оценки знаний, умений, навыков	Показатели и критерии оценивания, шкала оценивания
Знать основные методы лексического анализа текста	<p>Примеры вопросов к зачёту:</p> <ul style="list-style-type: none"> <li>• Конечные автоматы и конечные преобразователи.</li> <li>• Алгоритм моделирования недетерминированного конечного автомата.</li> <li>• Префиксные деревья. Представление множества слов в виде префиксного дерева.</li> <li>• Применение конечных преобразователей для морфологического анализа слов.</li> <li>• Алгоритм Портера.</li> <li>• Исправление орфографических ошибок. Редакционное расстояние. Алгоритм Вагнера-Фишера.</li> </ul>	оценка 3 — знает понятия конечного автомата и конечного преобразователя, оценка 4 — кроме того знает основные методы морфологического анализа слов, оценка 5 — кроме того знает алгоритмы исправления орфографических ошибок
Уметь использовать статистические методы обработки текста	<p>Примеры задач для контрольных работ:</p> <ul style="list-style-type: none"> <li>• Найдите частоты биграмм для предложения «Ворон к ворону летит, ворон ворону кричит». Примените к получившимся частотам сглаживание Уиттена-Белла. Найдите вероятность предложения «Летит к ворону ворон».</li> <li>• Предположим, что в предложении написано слово «стор» с опечаткой.</li> </ul>	оценка 3 — умеет использовать формулу Байеса, оценка 4 — кроме того умеет использовать N-граммы, оценка 5 — кроме того умеет выполнять сглаживание различными спосо-



Требования к обучающемуся	Типовые контрольные задания для оценки знаний, умений, навыков	Показатели и критерии оценивания, шкала оценивания
	<p>Словарь содержит правильные слова «стар», «стог», «стон», «сток», «стоп», «тор», «сто», «сор», «створ». Пусть частоты этих слов в корпусе равны 0,2%, 0,1%, 0,3%, 0,5%, 0,1%, 1%, 0,2%, 0,3%, 0,3%. Пусть вероятность замены символа равна 1%, удаления – 2%, вставки – 3%. Найдите наиболее вероятное правильное написание слова.</p> <p>Примеры тем для самостоятельной работы:</p> <ul style="list-style-type: none"> <li>• Напишите программу, которая по корпусу русских текстов находит частоты N-грамм с использованием сглаживания Гуда-Тьюринга и после этого по заданному предложению вычисляет его вероятность.</li> </ul>	бами

### Типовые контрольные задания и/или критерии для проверки индикатора ПК-7.1

Требования к обучающемуся	Типовые контрольные задания для оценки знаний, умений, навыков	Показатели и критерии оценивания, шкала оценивания
Знать основные статистические методы анализа текста	<p>Примеры вопросов к зачёту:</p> <ul style="list-style-type: none"> <li>• Условные вероятности. Формула Байеса.</li> <li>• N-граммы. Сглаживание Лапласа, Уиттена-Белла, Гуда-Тьюринга.</li> <li>• Применение N-грамм для исправления орфографических ошибок с учётом контекста.</li> <li>• Цепи Маркова, скрытые марковские модели. Алгоритм Витерби.</li> <li>• Определение частей речи на основе правил и на основе скрытых марковских моделей.</li> <li>• Модели word2vec и GloVe.</li> </ul>	оценка 3 — знает понятие N-граммы, сглаживания, скрытой марковской модели, моделей word2vec и GloVe, оценка 4 — кроме того знает алгоритмы, использующие эти модели, оценка 5 — кроме того знает методы оценки качества

Требования к обучающемуся	Типовые контрольные задания для оценки знаний, умений, навыков	Показатели и критерии оценивания, шкала оценивания
	<ul style="list-style-type: none"> <li>• Оценка качества моделей. Энтропия, перекрёстная энтропия.</li> </ul>	моделей
Знать основные применения нейронных сетей для обработки текстов	<p>Примеры вопросов к зачёту:</p> <ul style="list-style-type: none"> <li>• Искусственные нейроны. Искусственные нейронные сети. Функции активации.</li> <li>• Метод обратного распространения ошибки.</li> <li>• Рекуррентные нейронные сети.</li> <li>• Архитектура LSTM.</li> <li>• Использование нейронных сетей в компьютерной лингвистике: модели языка, машинный перевод.</li> </ul>	оценка 3 — знает простейшие архитектуры нейронных сетей и их применения для обработки текстов, оценка 4 — кроме того знает основные методы обучения, оценка 5 — кроме того знает более сложные архитектуры

## V. Учебно-методическое и информационное обеспечение дисциплины

### 1. Рекомендованная литература

#### а) Основная литература

[1] Волосатова, Т.М. Информатика и лингвистика [Электронный ресурс]: Учебное пособие/Волосатова Т.М., Чичварин Н.В. — Электрон. дан. — М.: НИЦ ИНФРА-М, 2016. — 196 с.: 60x90 1/16. — (Высшее образование: Бакалавриат) (Переплёт 7БЦ) ISBN 978-5-16-010977-0 — Режим доступа: <https://znanium.com/catalog/document?id=422587>

[2] Марченков, С.С. Конечные автоматы [Электронный ресурс]: учеб. пособие — Электрон. дан. — Москва: Физматлит, 2008. — 56 с. — Режим доступа: <https://e.lanbook.com/book/59510>. — Загл. с экрана.

[3] Короткова, М.А. Задачник по курсу "Математическая лингвистика и теория автоматов": учебное пособие для вузов [Электронный ресурс]: учеб. пособие / М.А. Короткова, Е.Е. Трифонова. — Электрон. дан. — Москва: НИЯУ МИФИ, 2012. — 92 с. — Режим доступа: <https://e.lanbook.com/book/75843>. — Загл. с экрана.

#### б) Дополнительная литература

[4] Федосеева, Л.И. Основы теории конечных автоматов и формальных языков [Электронный ресурс]: учеб. пособие / Л.И. Федосеева, Р.М. Адилов, М.Н. Шмокин. — Электрон. дан. — Пенза: ПензГТУ, 2013. — 136 с. — Ре-

жим доступа: <https://e.lanbook.com/book/62703>. — Загл. с экрана.

## 2. Программное обеспечение

Наименование помещений	Программное обеспечение
Ауд. 201а (компьютерная лаборатория ПМиК) (170002, Тверская обл., г. Тверь, пер. Садовый, д. 35)	Перечень программного обеспечения (со свободными лицензиями): Linux Kubuntu, KDE, TeXLive, TeXStudio, LibreOffice, GIMP, Gwenview, ImageMagick, Okular, Skanlite, Google Chrome, KDE Connect, Konversation, KRDC, KTorrent, Thunderbird, Elisa, VLC media player, PulseAudio, KAppTemplate, KDevelop, pgAdmin4, PostgreSQL, Qt, QtCreator, R, RStudio, Visual Studio Code, Perl, Python, Ruby, clang, clang++, gcc, g++, nasm, flex, bison, Maxima, Octave, Dolphin, HTop, Konsole, KSystemLog, Xterm, Ark, Kate, Kcalc, Krusader, Spectacle, Vim.

## 3. Современные профессиональные базы данных и информационные справочные системы

- [1] ЭБС «ZNANIUM.COM» <http://www.znanium.com>
- [2] ЭБС «Университетская библиотека онлайн» <https://biblioclub.ru>
- [3] ЭБС IPRbooks <http://www.iprbookshop.ru>
- [4] ЭБС «Лань» <http://e.lanbook.com>
- [5] ЭБС «Юрайт» <https://urait.ru>
- [6] ЭБС ТвГУ <http://megapro.tversu.ru/megapro/Web>
- [7] Научная электронная библиотека eLIBRARY.RU (подписка на журналы) [https://elibrary.ru/projects/subscription/rus\\_titles\\_open.asp](https://elibrary.ru/projects/subscription/rus_titles_open.asp)
- [8] Репозиторий ТвГУ <http://eprints.tversu.ru>

## 4. Перечень ресурсов информационно-телекоммуникационной сети «Интернет», необходимых для освоения дисциплины

- [1] Natural Language Processing, <http://www.learnerstv.org/Free-Computer-Science-Video-lectures-ltv676-Page1.htm>
- [2] Natural Language Processing with Deep Learning, <http://web.stanford.edu/class/cs224n/>
- [3] Московский центр непрерывного математического образования, <http://www.mccme.ru/>

## VI. Методические материалы для обучающихся по освоению дисциплины

Важной составляющей данного раздела РПД являются требования к

рейтинг-контролю с указанием баллов, распределенных между модулями и видами работы обучающихся.

Максимальная сумма баллов по учебной дисциплине, заканчивающейся зачетом, по итогам семестра составляет 100 баллов (50 баллов - 1-й модуль и 50 баллов - 2-й модуль).

Студенту, набравшему 40 баллов и выше по итогам работы в семестре, в экзаменационной ведомости и зачетной книжке выставляется оценка «зачтено». Студент, набравший до 39 баллов включительно, сдает зачет.

Распределение баллов по модулям устанавливается преподавателем и может корректироваться.

### Примеры задач для подготовки к контрольным работам

1. Постройте недетерминированный конечный автомат с 11 состояниями для языка

$L = \{w \in \{a,b\}^* : |w| \geq 4 \text{ и первые две буквы слова не равны последним двум}\}$   
и обоснуйте его правильность. С помощью алгоритма моделирования проверьте, какие из слов *abb*, *baaba*, *bbab*, *abbabba* распознаются этим автоматом.

2. С помощью алгоритма Вагнера-Фишера найдите редакционное расстояние между словами «алгоритм» и «логарифм».
3. Дана КС-грамматика с начальным нетерминалом E и следующими правилами:

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \varepsilon \\ F &\rightarrow c \mid (E) \end{aligned}$$

С помощью алгоритма Эрли определите, какие из следующих слов выводимы в этой грамматике, и постройте для них деревья вывода:  $(c+c)^*c$ ,  $c+c+c+c$ ,  $((c))$ ,  $(c+c)$ ,  $c*c**c$ ,  $c*c$ .

4. Постройте комбинаторную категориальную грамматику для языка

$$L = \{w\#w^{-1} : w - \text{правильное скобочное слово}\}$$

и обоснуйте ее правильность.

5. Постройте КЗ-грамматику для языка  $L = \{a^n : n \text{ — составное число}\}$ .
6. Постройте размеченную иерархизованную систему составляющих и размеченное дерево зависимостей для предложения «Тощая торговка вяленой воблой торчала среди ящиков».
7. Обобщенной КЗ-грамматикой (ОКЗ-грамматикой) называется порождающая грамматика  $G = (N, \Sigma, P, S)$ , правила которой имеют вид  $\xi A \eta \rightarrow \xi \alpha \eta$ , где  $\xi, \eta, \alpha \in (\Sigma \cup N)^*$ ,  $A \in N$ . Докажите, что любой рекурсивно перечислимый язык порождается некоторой ОКЗ-грамматикой.
8. Докажите, что класс языков типа 0 замкнут относительно тасовки:

$$\text{TAC}(L_1, L_2) = \{x_1 y_1 \dots x_n y_n : x_i, y_j \in \Sigma^*, x_1 \dots x_n = x, y_1 \dots y_n = y, x \in L_1, y \in L_2\}.$$

## Требования к рейтинг контролю (7 семестр)

**Контрольная работа 1.** Темы: конечные преобразователи, статистические методы обработки текста, КЗ-грамматики. Пример задания:

1. Постройте детерминированный конечный преобразователь для выполнения следующей операции. На вход подаётся слово  $w\$$  в алфавите  $\{a,b,c,\$\}$ , символ  $\$$  помечает конец входного слова. Требуется удалить те символы  $a$ , которые стоят между двумя символами  $b$ . Например, из слова  $acbabaabscab\$$  должно получиться  $acbbaabscab\$$ .
2. Найдите частоты биграмм для предложения «Без устали, без устали смотрю, смотрю в окно». Примените к получившимся частотам сглаживание Уиттена-Белла. Найдите вероятность предложения «Смотрю в окно без устали».
3. Постройте неукорачивающую грамматику, порождающую следующий язык:

$$L = \{a^n b^m a^n b^m : m, n > 0\}.$$

За решение каждой задачи выставляется максимум 10 баллов.

**Самостоятельная работа 1.** Темы: конечные преобразователи, статистические методы обработки текста, КЗ-грамматики. Пример задания:

Напишите программу, которая исправляет орфографические ошибки в отдельных словах с использованием формулы Байеса.

За решение задачи выставляется максимум 10 баллов.

**Контрольная работа 2.** Темы: слабо контекстные грамматики, категориальные грамматики, системы составляющих, деревья зависимостей. Пример задания:

1. Постройте ТАГ-грамматику для языка

$$L = \{a^i b^j c^{i+2} : j > i > 0\}$$

2. Постройте категориальную грамматику зависимостей для языка

$$L = \{a^i b^{2i} w : i > 0, w \in \{c, d\}^*, |w|_c < |w|_a\}$$

3. Постройте 3-множественную КС-грамматику для языка

$$L = \{www^{-1}w^{-1} : w \in \{a, b\}^*\}$$

4. Постройте размеченную иерархизованную систему составляющих и размеченное дерево зависимостей для предложения «К востоку от боровых озёр лежат громадные мещёрские болота — мшары».

5. Для следующего предложения сократите в указанном порядке категории до S и приведите получающийся  $\lambda$ -терм к нормальной форме.

Мастер	внимательно	осматривал	станок
<u>NP : M</u>	C : B	(((S \ NP) \ C) / NP) :	NP : C
(>T)		<u><math>\lambda xyz.Ozxy</math></u>	
		(>)	

За решение каждой задачи выставляется максимум 10 баллов.

**Самостоятельная работа 2.** Темы: слабо контекстные грамматики, системы

составляющих, деревья зависимостей. Пример задания:

Напишите программу, которая по системе составляющих строит согласованное с ней дерево зависимостей.

За решение задачи выставляется максимум 10 баллов.

**Общая сумма.** В сумме за все задачи выставляет не более 100 баллов.

### Вопросы к зачёту

1. Лексический анализ.
  - 1) Конечные автоматы и конечные преобразователи.
  - 2) Алгоритм моделирования недетерминированного конечного автомата.
  - 3) Префиксные деревья. Представление множества слов в виде префиксного дерева.
  - 4) Применение конечных преобразователей для морфологического анализа слов.
  - 5) Алгоритм Портера.
  - 6) Исправление орфографических ошибок. Редакционное расстояние. Алгоритм Вагнера-Фишера.
2. Статистические методы обработки языка.
  - 1) Условные вероятности. Формула Байеса.
  - 2) N-граммы. Сглаживание Лапласа, Уиттена-Белла, Гуда-Тьюринга.
  - 3) Применение N-грамм для исправления орфографических ошибок с учётом контекста.
  - 4) Цепи Маркова, скрытые марковские модели. Алгоритм Витерби.
  - 5) Определение частей речи на основе правил и на основе скрытых марковских моделей.
  - 6) Модели word2vec и GloVe.
  - 7) Оценка качества моделей. Энтропия, перекрёстная энтропия.
3. Синтаксический анализ.
  - 1) Порождающие грамматики. Иерархия Хомского.
  - 2) Контекстно-зависимые и контекстно-свободные грамматики. Деревья вывода.
  - 3) Эквивалентность контекстно-зависимых грамматик и линейно-ограниченных автоматов.
  - 4) Классические категориальные грамматики. Эквивалентность КС-грамматик и классических категориальных грамматик.
  - 5) Системы составляющих и деревья зависимостей. Связь деревьев зависимостей и систем составляющих.
  - 6) Синтаксический анализ на основе КС-грамматик. Алгоритмы Кока-Янгера-Касами и Эрли.
  - 7) Слабо-контекстные грамматики. Множественные контекстно-свободные грамматики.
  - 8) Головные грамматики, линейные индексные грамматики, комбинаторные категориальные грамматики, ТАГ-грамматики.
  - 9) Категориальные грамматики зависимостей (КГЗ). Алгоритм анализа

КГЗ.

- 10) Основы  $\lambda$ -исчисления.  $\lambda$ -термы,  $\beta$ -редукция, нормальная форма. Формулировка теоремы Чёрча-Россера.
  - 11) Применения комбинаторных категориальных грамматик для представления семантики в виде  $\lambda$ -термов.
4. Применение нейронных сетей для обработки языка.
- 1) Искусственные нейроны. Искусственные нейронные сети. Функции активации.
  - 2) Метод обратного распространения ошибки.
  - 3) Рекуррентные нейронные сети.
  - 4) Архитектура LSTM.
  - 5) Использование нейронных сетей в компьютерной лингвистике: модели языка, машинный перевод.

## **VII. Материально-техническая база, необходимая для осуществления образовательного процесса по дисциплине**

### **Для аудиторной работы**

Наименование помещений	Материально-техническое оснащение помещений
Ауд. 308 (170002, Тверская обл., г. Тверь, пер. Садовый, д. 35)	Набор учебной мебели, экран проектор.

### **Для самостоятельной работы**

Наименование помещений	Материально-техническое оснащение помещений
Ауд. 201а (компьютерная лаборатория ПМиК) (170002, Тверская обл., г. Тверь, пер. Садовый, д. 35)	Набор учебной мебели, доска маркерная, компьютер, сервер (системный блок), концентратор сетевой.

## **VIII. Сведения об обновлении рабочей программы дисциплины**

№ п/п	Обновленный раздел рабочей программы дисциплины	Описание внесённых изменений	Дата и протокол заседания кафедры, утвердившего изменения
1	11. 2) Программное обеспечение	Внесены изменения в список ПО	От 24.08.2023 года, протокол

			№ 1 ученого совета факультета
2	V. 1) Рекомендуемая литература	Обновление ссылок на литературу	От 24.08.2023 года, протокол № 1 ученого совета факультета
3			
4			