

УП2015-2016

Министерство образования и науки Российской Федерации  
Федеральное государственное бюджетное образовательное учреждение  
высшего образования  
«Владимирский государственный университет  
имени Александра Григорьевича и Николая Григорьевича Столетовых»  
(ВлГУ)



УТВЕРЖДАЮ

Проректор  
по образовательной деятельности

А.А.Панфилов

2016 г.

**РАБОЧАЯ ПРОГРАММА ДИСЦИПЛИНЫ**  
**МЕТОДЫ АНАЛИЗА ДАННЫХ И ЕСТЕСТВЕННО-ЯЗЫКОВЫХ ТЕКСТОВ**

(наименование дисциплины)

Специальность 10.05.04 "Информационно-аналитические системы безопасности"  
Специализация "Автоматизация информационно-аналитической деятельности"  
Уровень высшего образования специалитет  
Форма обучения очная

Семестр	Трудоемкость зач. ед./ час.	Лекции, час.	Практич. занятия, час.	Лаборат. работы, час.	СРС, час.	Форма промежуточного контроля (экз./зачет)
8	4/144	36		36	72	Зачет с оценкой
Итого	4/144	36		36	72	Зачет с оценкой

Владимир 2016

## **1. ЦЕЛИ ОСВОЕНИЯ ДИСЦИПЛИНЫ**

**Целями освоения дисциплины** «Методы анализа данных и естественно-языковых текстов» являются обеспечение подготовки студентов в соответствии с требованиями ФГОС ВПО и учебного плана специальности 10.05.04 «Информационно-аналитические системы безопасности», ознакомление студентов с кругом задач в области автоматической обработки естественного языка (natural language processing) и компьютерной лингвистики (computational linguistics), а также с доступным программным инструментарием для решения прикладных задач обработки текста. В рамках курса рассматриваются основные понятия компьютерной лингвистики, а также существующее программное обеспечение для работы с текстами. Целями освоения дисциплины «Методы анализа данных и естественно-языковых текстов» являются:

- Изучение базовых алгоритмов анализа и интерпретации данных.
- Формирование практических навыков работы с современными пакетами прикладных программ для решения задач анализа и интерпретации данных.

При изучении курса студенты знакомятся с современными формальными методами, реализующими «восходящую» стратегию анализа: извлечение интерпретируемых зависимостей из эмпирических данных. Методы рассматриваются в рамках парадигмы интеллектуального анализа данных (data mining, knowledge discovery), являющейся важнейшим направлением современных исследований в области анализа гетерогенных данных с нечисловыми параметрами.

## **2. МЕСТО ДИСЦИПЛИНЫ В СТРУКТУРЕ ОПОП ВО СПЕЦИАЛИТЕТА**

Данная дисциплина относится к базовой части Блока Б1 (код Б1.Б.14). В учебном плане предусмотрены виды учебной деятельности, обеспечивающие синтез теоретических лекций, лабораторных работ и практических занятий.

Дисциплина изучается на 4 курсе, требования к «входным» знаниям, умениям и готовностям (пререквизитам) обучающегося определяются требованиями к уровню подготовки по специальности 10.05.04 «Информационно-аналитические системы безопасности», квалификации - специалист по курсам «Информатика», «Математическая логика и теория алгоритмов», «Структуры данных», «Базы данных и экспертные системы», «Моделирование автоматизированных информационных систем», «Методология и организация информационно-аналитической деятельности». Курс тесно взаимосвязан с другими дисциплинами, такими как «Современные платежные системы и их безопасность», «Лингвистическое обеспечение автоматизированных информационных систем», «Формализованные модели и методы решения аналитических задач» и другими.

## **3. КОМПЕТЕНЦИИ ОБУЧАЮЩЕГОСЯ, ФОРМИРУЕМЫЕ В РЕЗУЛЬТАТЕ ОСВОЕНИЯ ДИСЦИПЛИНЫ**

В результате освоения дисциплины студент должен обладать следующими профессиональными компетенциями:

ПК-2 – способностью применять методы анализа массивов данных и интерпретировать профессиональный смысл получаемых формальных результатов;

общепрофессиональными компетенциями:

ПСК-1.2 – способностью разрабатывать и применять автоматизированные технологии обработки естественно-языковых текстов и формализованных данных при решении информационно-аналитических задач;

ПСК-1.3 – способностью решать задачи анализа данных больших объемов.

В результате освоения дисциплины обучающийся должен демонстрировать следующие результаты образования:

1) **Знать:** - круг решенных и перешенных задач компьютерной лингвистики, ориентироваться в современных методах обработки текста на естественном языке; - владеть лингвистической и статистической терминологией, необходимой для чтения литературы в этой области (на русском и английском языках); - методологические основы анализа данных; -

методы снижения размерности многомерных данных; - основные свойства естественного языка как знаковой системы; - структуру естественно-языкового текста как объекта компьютерной обработки; - ограничения, накладываемые свойствами русского естественноязыкового текста на процедуры обработки; - основные типы задач по обработке текстов и основные виды автоматизированных систем, решающих эти задачи; - прикладные методы, модели и алгоритмы, применяемые в системах компьютерной обработки естественноязыковых текстов (ПК-2, ПСК-1.2, ПСК-1.3).

**2) Уметь:** - строить и анализировать частотные списки языковых единиц; - извлекать данные из текста с помощью регулярных выражений; - формулировать правила извлечения информации в терминах контекстно-свободных грамматик; - применять методы анализа массивов данных при разработке алгоритмов анализа и обработки измерительной информации; - использовать стандартную терминологию, определения и обозначения в области обработки данных; - ставить и решать практические задачи анализа данных в условиях различной полноты исходной информации; - проводить комплексный анализ данных с использованием базовых параметрических и непараметрических моделей; - применять современные автоматизированные технологии семантической обработки текстов при решении прикладных информационно-аналитических задач (ПК-2, ПСК-1.2, ПСК-1.3);

**3) Владеть:** - навыками работы с программным обеспечением для автоматического анализа текстов: морфологическими и синтаксическими анализаторами, конкордансами, системами извлечения фактов и отношений, инструментами кластеризации, классификации и тематического моделирования коллекций документов; - навыками решения формализованных математических задач анализа данных с помощью пакетов прикладных программ; - навыками работы с программными системами, реализующими автоматизированные технологии семантической обработки текстов; - методами поиска, выбора и обработки массивов документов по конкретным направлениям служебной деятельности; - навыками синтеза гуманитарного и технического знания при решении конкретных проблем автоматизации обработки текстов (ПК-2, ПСК-1.2, ПСК-1.3).

У обучаемых в процессе изучения дисциплины должны вырабатываться дополнительные компетенции, с учетом требований работодателей:

- способность применять навыки решения аналитических задач информационно-аналитической деятельности в профессиональной сфере.

#### 4. СТРУКТУРА И СОДЕРЖАНИЕ ДИСЦИПЛИНЫ

Общая трудоемкость дисциплины составляет 4 зачетных единиц, 144 часа.

№ п/п	Раздел (тема) дисциплины	Семестр	Неделя семестра	Виды учебной работы, включая самостоятельную работу студентов и трудоемкость (в часах)					Объем учебной работы, с применением интерактивных методов (в часах / %)	Формы текущего контроля успеваемости (по неделям семестра), форма промежуточной аттестации (по семестрам)	
				Лекции	Практические занятия	Лабораторные работы	Контрольные работы	СРС			
1.	Введение в анализ данных. Проблема обработки данных. Матрица данных.	8	1	2				4		1/50%	
2.	Классификация данных с использованием детерминированных моделей. Решающие поверхности и дискриминантные функции.	8	2	2		4		4		2/33%	
3.	Процедуры обучения с коррекцией ошибок: правило с фиксированным приращением, правило абсолютной коррекции, частично корректирующее правило.	8	3	2				4		1/50%	
4.	Классификация данных на основе статистических моделей. Функция потерь. Байесовская дискриминантная функция.	8	4	2		4		4		2/33%	
5.	Примеры построения статистических дискриминантных функций для различных статистических нескольких моделей данных.	8	5	2				4		1/50%	
6.	Кластер-анализ. Основные типы задач кластер-анализа. Меры подобия и функции расстояния. Выбор критерия кластеризации.	8	6	2		4		4		2/33%	Рейтинг-контроль №1
7.	Методы снижения размерностей данных. Анализ матриц исходных данных.	8	7	2				4		1/50%	
8.	Методы прогнозирования временных рядов.	8	8	2		4		4		2/33%	
9.	Системы DATA MINING. в задачах анализа и интерпретации данных.	8	9	2				4		1/50%	

№ п/п	Раздел (тема) дисциплины	Семестр	Неделя семестра	Виды учебной работы, включая самостоятельную работу студентов и трудоемкость (в часах)					Объем учебной работы, с применением интерактивных методов (в часах / %)	Формы текущего контроля успеваемости (по неделям семестра), форма промежуточной аттестации (по семестрам)	
				Лекции	Практические занятия	Лабораторные работы	Контрольные работы	СРС			
10.	Автоматическая обработка языка и компьютерная лингвистика.	8	10	2		4		4	2/33%		
11.	Инструменты морфологического анализа для русского языка. Понятия словоформа, лексема, лемма, основа.	8	11	2				4	1/50%	Рейтинг-контроль №2	
12.	Частотный анализ лексики и ключевые слова. Частотное распределение лексики в языке. Закон Ципфа. Доля нарах legomena.	8	12	2		4		4	2/33%		
13.	Локальные модели контекста. Вероятностные языковые модели	8	13	2				4	1/50%		
14.	Автоматическое определение тематики. Векторное представление текста для задач информационного поиска.	8	14	2		4		4	2/33%		
15.	ПО для кластеризации текстов. Пакеты кластеризации для R. gCLUTO. Классификация текстов.	8	15	2				4	1/50%		
16.	Извлечение мнений и оценок (Sentiment analysis).	8	16	2		4		4	2/33%		
17.	Извлечение фактов и отношений. Синтаксис и формальные языки.	8	17	2				4	1/50%		
18.	Автоматический анализ стиля. Стилометрия.	8	18	2		4		4	2/33%	Рейтинг-контроль №3	
Всего				36		36		72	27/38%	Зачет с оценкой	

### Содержание дисциплины «Методы анализа данных и естественно-языковых текстов»

**Тема 1.** Введение в анализ данных. Проблема обработки данных. Матрица данных. Гипотезы компактности и скрытых факторов. Структура матрицы данных и задачи обработки. Матрица объект-объект и признак-признак. Расстояние и близость. Измерение признаков. Отношения и их представление. Основные проблемы измерений. Основные типы шкал. Проблема адекватности. Основные задачи анализа и интерпретации данных.

**Тема 2.** Классификация данных с использованием детерминированных моделей. Решающие поверхности и дискриминантные функции. Линейные дискриминантные функции классификатор по минимуму расстояния. Линейная разделимость. Кусочно-линейные дискриминантные функции. Нелинейные дискриминантные функции. Фи-машины. Потенциальные функции как дискриминантные функции. Пространство весов.

**Тема 3.** Процедуры обучения с коррекцией ошибок: правило с фиксированным приращением, правило абсолютной коррекции, частично корректирующее правило. Обобщенные градиентные методы. Персептронный критерий. Процедуры обучения на основе минимальной среднеквадратичной ошибки: псевдоинверсный метод, метод Хо-Кашьпа.

**Тема 4.** Классификация данных на основе статистических моделей. Функция потерь. Байесовская дискриминантная функция. Принятие решения по максимуму правдоподобия. Оптимальная дискриминантная функция для нормально распределенных образов. Дискриминантная функция Фишера. Множественный дискриминантный анализ. Пошаговый дискриминантный анализ. Ошибки классификации.

**Тема 5.** Примеры построения статистических дискриминантных функций для различных статистических нескольких моделей данных. Обучение для статистических дискриминантных функций. Оценки максимального правдоподобия, байесовские оценки. Непараметрическое оценивание. Парзеновские окна, метод непараметрического оценивания на основе К-ближайшего соседства.

**Тема 6.** Кластер-анализ. Основные типы задач кластер-анализа. Меры подобия и функции расстояния. Выбор критерия кластеризации. Кластерные методы, основанные на евклидовой метрике. Иерархическая кластеризация. Метод К-внутригрупповых средних. Использование методов теории графов в задачах кластеризации. Кластеризация на основе анализа плотностей вероятностей.

**Тема 7.** Методы снижения размерностей данных. Анализ матриц исходных данных. Метод главных компонент. Корреляционная матрица и ее основные свойства. Собственные векторы и собственные числа корреляционной матрицы. Приведение корреляционной матрицы к диагональной форме. Геометрическая интерпретация главных компонент на плоскости. Модели факторного анализа. Оценка факторных нагрузок методом максимального правдоподобия и центроидным методом. Вращение факторов и их интерпретация. Использование кластеризации признаков для снижения размерности. Многомерное шкалирование (МИ). Метрический и неметрический подход к МИ. Методы ортогонального проектирования. Нелинейные методы МИ. Многомерное шкалирование неметрических данных. Многомерные развертки.

**Тема 8.** Методы прогнозирования временных рядов. Классификация методов прогнозирования. Оценивание трендов. Методы скользящего среднего. Экспоненциальное сглаживание. Регрессионный анализ и прогнозирование. Линейные параметрические модели временных рядов. Методы оценивания моделей авторегрессии, скользящего среднего и смешанных моделей. Сезонные модели. Прогнозирование на основе параметрических моделей. Прогнозирование с использованием нейронных сетей.

**Тема 9.** Системы DATA MINING в задачах анализа и интерпретации данных. Понятие об интеллектуальных системах анализа и интерпретации данных. DATA MINING - системы извлечения новых знаний из данных. Типы систем DATA MINING - предметно-ориентированные аналитические системы, статистические пакеты, нейронные сети, деревья решений, обнаружение логических закономерностей, генетические алгоритмы, системы визуализации многомерных данных

**Тема 10.** Автоматическая обработка языка и компьютерная лингвистика. Задачи автоматической обработки текста в научных исследованиях. Основные задачи компьютерной лингвистики и история развития автоматической обработки языка. Иерархия языковых уровней и стандартный цикл обработки текста (графематика — морфология — синтаксис — семантика). Основные задачи автоматической обработки текста: токенизация и нормализация текста; сегментация на предложения; стемминг; лемматизация и частеречные теги; снятие омонимии; парсинг — поверхностный и полный; кореференция и разрешение анафоры. Задачи высокоуровневого анализа: извлечение фактов и отношений, анализ оценок (sentiment analysis).

**Тема 11.** Инструменты морфологического анализа для русского языка. Понятия словоформа, лексема, лемма, основа. Стемминг. Алгоритм Портера. Stemka. Лемматизация. Словарный метод — грамматический словарь Зализняка. mystem. АОТ. ruymorphy. Грамматическая омонимия. Разметка частей речи. TreeTagger. TnT.

**Тема 12.** Частотный анализ лексики и ключевые слова. Частотное распределение лексики в языке. Закон Ципфа. Доля нарек legomena. Скорость роста словаря. Коэффициент лексического разнообразия (type/token ratio). Распределение лексики в текстах коллекции. Взвешенная частотность. TF-IDF. Прочие меры лексической дисперсии. Мера отклонения пропорций DP и DPnorm. Извлечение ключевых слов. Метод контрастного корпуса. Отношение правдоподобия. Диахронический анализ лексической частотности.

**Тема 13.** Локальные модели контекста. Вероятностные языковые модели. Понятие N-грамм. Буквенные и словарные n-граммы. Контекстное окно. Применения N-грамм в автоматической обработке языка. Роль биграмм и триграмм. Определение языка по письменности. Языковые модели. Цепь Маркова. Коллокации. Формальные определения и лингвистический смысл коллокаций. Меры ассоциации. Коэффициент взаимной информации (MI). T-score. Отношение правдоподобия (log-likelihood). Статистические тесты ассоциации: хи-квадрат и Fisher exact test. Выделение коллокаций по синтаксическому шаблону. Разрывные коллокации.

**Тема 14.** Автоматическое определение тематики. Векторное представление текста для задач информационного поиска. Открытые и закрытые классы слов. Стоп-слова. Динамические списки стоп слов. Порог отсечения по частотности и DF. Диагностивная семантика. Совместная встречаемость и семантическая близость.

Кластеризация текстов. Задачи и область применения кластерных методов. Виды кластеризации: плоские, агglomerативные, нечеткие. Меры близости: евклидово расстояние, косинусная мера. Популярные алгоритмы кластеризации: k-средних, DBCLUST, спектральные алгоритмы.

**Тема 15.** ПО для кластеризации текстов. Пакеты кластеризации для R. gCLUTO. Классификация текстов. Машинное обучение с учителем и без учителя в задачах классификации текстов. Популярные алгоритмы классификации: наивный байесовский метод, метод опорных векторов, деревья принятия решений. ПО для классификации текстов. SVMLight. Пространственное моделирование семантических отношений (word space). Латентный семантический анализ. Вероятностный латентно семантический анализ. Тематическое моделирование. Метод латентного размещения Дирихле. ПО для латентного семантического анализа и тематического моделирования. Mallet. sTMT.

**Тема 16.** Извлечение мнений и оценок (Sentiment analysis). Область применения методов извлечения мнений и оценок. Типы оценочных текстов: позитивный, негативный, нейтральный. Оценочные шкалы. Классификация документов по оценке. Извлечение оценочных предложений и фрагментов. Определение предмета оценки. Методы извлечения оценок. Словарные методы. Машинное обучение. Комбинирование источников. Проблемы и ограничения методов извлечения оценок.

**Тема 17.** Извлечение фактов и отношений. Синтаксис и формальные языки. Иерархия грамматик Хомского. Регулярные грамматики. Контекстно-свободные грамматики. Основные понятия: терминал, нетерминал, правило. Форма записи Бакуса-Наура. Текст и дискурс. Методы сегментации текста с обучающей выборкой и без. Понятие связности текста. Автоматическое определение отношений связности. Коммуникативная структура текста. Понятия тема, рема, информационный статус. Теория риторической структуры. Риторические отношения. Анализ нарративной структуры. Разрешение анафоры и анализ кореференции. Методы извлечения информации из текстов на естественном языке. Словарные методы. Синтаксические шаблоны. Распознавание именованных сущностей. Извлечение отношений. Извлечение ключевых слов текста.

**Тема 18.** Автоматический анализ стиля. Стилометрия. Автоматическое определение авторства: краткая история и обзор методов. Формальные и лингвистические черты для стилистического анализа. Автоматическое определение жанровой принадлежности текста.

## **5. ОБРАЗОВАТЕЛЬНЫЕ ТЕХНОЛОГИИ**

Изучение дисциплины предполагает не только запоминание и понимание, но и анализ, синтез, рефлексию, формирует универсальные умения и навыки, являющиеся основой становления специалиста по специальности 10.05.04 «Информационно-аналитические системы безопасности».

Для реализации компетентностного подхода предлагается интегрировать в учебный процесс интерактивные образовательные технологии, включая информационные и коммуникационные технологии (ИКТ), при осуществлении различных видов учебной работы:

- учебную дискуссию;
- электронные средства обучения (слайд-лекции, электронные тренажеры, компьютерные тесты);
- дистанционные (сетевые) технологии.

Как традиционные, так и лекции инновационного характера могут сопровождаться компьютерными слайдами или слайд-лекциями. Основное требование к слайд-лекции – применение динамических эффектов (анимированных объектов), функциональным назначением которых является наглядно-образное представление информации, сложной для понимания и осмысливания студентами, а также интенсификация и диверсификация учебного процесса.

Удельный вес занятий, проводимых в интерактивных формах, определяется главной целью ОПОП специальности 10.05.04 «Информационно-аналитические системы безопасности», особенностью контингента обучающихся и содержанием конкретных дисциплин, и в целом, в учебном процессе, они составляют не менее 30 процентов аудиторных занятий.

Занятия лекционного типа для соответствующих групп студентов согласно требованиям стандарта высшего образования не могут составлять более 55 процентов аудиторных занятий. Программа дисциплины соответствует данным требованиям.

Таким образом, применение интерактивных образовательных технологий придает инновационный характер практически всем видам учебных занятий, включая лекционные. При этом делается акцент на развитие самостоятельного, продуктивного мышления, основанного на диалогических дидактических приемах, субъектной позиции обучающегося в образовательном процессе. Тем самым создаются условия для реализации компетентностного подхода при изучении данной дисциплины.

## **6. ОЦЕНОЧНЫЕ СРЕДСТВА ДЛЯ ТЕКУЩЕГО КОНТРОЛЯ УСПЕВАЕМОСТИ, ПРОМЕЖУТОЧНОЙ АТТЕСТАЦИИ ПО ИТОГАМ ОСВОЕНИЯ ДИСЦИПЛИНЫ И УЧЕБНО-МЕТОДИЧЕСКОЕ ОБЕСПЕЧЕНИЕ САМОСТОЯТЕЛЬНОЙ РАБОТЫ СТУДЕНТОВ**

Для текущего контроля успеваемости предлагается использование рейтинговой системы оценки, которая носит интегрированный характер и учитывает успешность студента в различных видах учебной деятельности, степень сформированности у студента общекультурных и профессиональных компетенций.

Примерный перечень заданий для текущих контрольных мероприятий:

### **Текущий контроль**

#### **Вопросы рейтинг-контроля №1**

- Введение в анализ данных. Проблема обработки данных. Матрица данных. Гипотезы компактности и скрытых факторов.
- Структура матрицы данных и задачи обработки. Матрица объект-объект и признак-признак.
- Расстояние и близость. Измерение признаков. Отношения и их представление.
- Основные проблемы измерений. Основные типы шкал. Проблема адекватности. Основные задачи анализа и интерпретации данных.
- Классификация данных с использованием детерминированных моделей.
- Решающие поверхности и дискриминантные функции. Линейные дискриминантные функции классификатор по минимуму расстояния.

- Линейная разделимость. Кусочно-линейные дискриминантные функции.
- Нелинейные дискриминантные функции.
- Фи-машины. Потенциальные функции как дискриминантные функции. Пространство весов.
- Процедуры обучения с коррекцией ошибок: правило с фиксированным приращением, правило абсолютной коррекции, частично корректирующее правило.
- Обобщенные градиентные методы. Персепtronный критерий. Процедуры обучения на основе минимальной среднеквадратичной ошибки: псевдоинверсный метод, метод Хо-Кашпана.
- Классификация данных на основе статистических моделей.
- Функция потерь. Байесовская дискриминантная функция.
- Принятие решения по максимуму правдоподобия. Оптимальная дискриминантная функция для нормально распределенных образов.
- Дискриминантная функция Фишера. Множественный дискриминантный анализ.
- Пошаговый дискриминантный анализ. Ошибки классификации.
- Примеры построения статистических дискриминантных функций для различных статистических нескольких моделей данных.
- Обучение для статистических дискриминантных функций.
- Оценки максимального правдоподобия, байесовские оценки.
- Непараметрическое оценивание. Парзеновские окна, метод непараметрического оценивания на основе К-ближайшего соседства.
- Кластер-анализ. Основные типы задач кластер-анализа.
- Меры подобия и функции расстояния. Выбор критерия кластеризации. Кластерные методы, основанные на евклидовой метрике.
- Иерархическая кластеризация. Метод К-внутригрупповых средних.
- Использование методов теории графов в задачах кластеризации.
- Кластеризация на основе анализа плотностей вероятностей.
- Методы снижения размерностей данных. Анализ матриц исходных данных.
- Метод главных компонент. Корреляционная матрица и ее основные свойства.
- Собственные векторы и собственные числа корреляционной матрицы. Приведение корреляционной матрицы к диагональной форме.
- Геометрическая интерпретация главных компонент на плоскости.
- Модели факторного анализа. Оценка факторных нагрузок методом максимального правдоподобия и центроидным методом.

### **Вопросы рейтинг-контроля №2**

- Вращение факторов и их интерпретация.
- Использование кластеризации признаков для снижения размерности.
- Многомерное шкалирование (МИ). Метрический и неметрический подход к МИ.
- Методы ортогонального проектирования.
- Нелинейные методы МИ. Многомерное шкалирование неметрических данных. Многомерные развертки.
- Методы прогнозирования временных рядов.
- Классификация методов прогнозирования. Оценивание трендов.
- Методы скользящего среднего. Экспоненциальное сглаживание.
- Регрессионный анализ и прогнозирование. Линейные параметрические модели временных рядов.
- Методы оценивания моделей авторегрессии, скользящего среднего и смешанных моделей.
- Сезонные модели. Прогнозирование на основе параметрических моделей. Прогнозирование с использованием нейронных сетей.
- Системы DATA MINING, в задачах анализа и интерпретации данных.
- Понятие об интеллектуальных системах анализа и интерпретации данных. DATA MINING - системы извлечения новых знаний из данных. Типы систем DATA MINING - предметно-ориентированные аналитические системы, статистические пакеты, нейронные сети, деревья

решений, обнаружение логических закономерностей, генетические алгоритмы, системы визуализации многомерных данных

- Автоматическая обработка языка и компьютерная лингвистика. Задачи автоматической обработки текста в научных исследованиях.
- Основные задачи компьютерной лингвистики и история развития автоматической обработки языка.
- Иерархия языковых уровней и стандартный цикл обработки текста (графематика — морфология — синтаксис — семантика).
- Основные задачи автоматической обработки текста: токенизация и нормализация текста; сегментация на предложения; стемминг; лемматизация и частеречные теги; снятие омонимии; парсинг — поверхностный и полный; кореференция и разрешение анафоры.
- Задачи высокоуровневого анализа: извлечение фактов и отношений, анализ оценок (sentiment analysis).
- Инструменты морфологического анализа для русского языка. Понятия словоформа, лексема, лемма, основа. Стемминг. Алгоритм Портера. Stemka. Лемматизация.
- Словарный метод — грамматический словарь Зализняка, mystem. АОТ. pymorphy. Грамматическая омонимия. Разметка частей речи. TreeTagger. TnT.
- Частотный анализ лексики и ключевые слова. Частотное распределение лексики в языке.
- Закон Ципфа. Доля  $h_{\alpha}$  из  $legomena$ . Скорость роста словаря. Коэффициент лексического разнообразия (type/token ratio).
- Распределение лексики в текстах коллекции. Взвешенная частотность. TF-IDF. Прочие меры лексической дисперсии. Мера отклонения пропорций DP и DPnorm.
- Извлечение ключевых слов. Метод контрастного корпуса. Отношение правдоподобия. Диахронический анализ лексической частотности.
- Локальные модели контекста. Вероятностные языковые модели. Понятие N-грамм.
- Буквенные и словарные n-граммы. Контекстное окно. Применения N-грамм в автоматической обработке языка. Роль биграмм и триграмм.

### Вопросы рейтинг-контроля №3

- Буквенные и словарные n-граммы. Контекстное окно. Применения N-грамм в автоматической обработке языка. Роль биграмм и триграмм.
- Определение языка по письменности. Языковые модели. Цепь Маркова. Коллокации.
- Формальные определения и лингвистический смысл коллокаций. Меры ассоциации.
- Коэффициент взаимной информации (MI). T-score. Отношение правдоподобия (log-likelihood).
- Статистические тесты ассоциации: хи-квадрат и Fisher exact test.
- Выделение коллокаций по синтаксическому шаблону. Разрывные коллокации.
- Автоматическое определение тематики. Векторное представление текста для задач информационного поиска.
- Открытые и закрытые классы слов. Стоп-слова. Динамические списки стоп слов. Порог отсечения по частотности и DF. Дистрибутивная семантика. Совместная встречаемость и семантическая близость.
- Кластеризация текстов. Задачи и область применения кластерных методов.
- Виды кластеризации: плоские, агglomerативные, нечеткие. Меры близости: евклидово расстояние, косинусная мера. Популярные алгоритмы кластеризации: k-средних, DBCLUST, спектральные алгоритмы.
- ПО для кластеризации текстов. Пакеты кластеризации для R. gCLUTO.
- Классификация текстов. Машинальное обучение с учителем и без учителя в задачах классификации текстов.
- Популярные алгоритмы классификации: наивный байесовский метод, метод опорных векторов, деревья принятия решений.
- ПО для классификации текстов. SVMLight. Пространственное моделирование семантических отношений (word space). Латентный семантический анализ.

- Вероятностный латентно семантический анализ. Тематическое моделирование. Метод латентного размещения Дирихле. ПО для латентного семантического анализа и тематического моделирования. Mallet, sTMT.
- Извлечение мнений и оценок (Sentiment analysis). Область применения методов извлечения мнений и оценок.
- Типы оценочных текстов: позитивный, негативный, нейтральный. Оценочные шкалы. Классификация документов по оценке. Извлечение оценочных предложений и фрагментов.
- Определение предмета оценки. Методы извлечения оценок. Словарные методы. Машинное обучение. Комбинирование источников. Проблемы и ограничения методов извлечения оценок.
- Извлечение фактов и отношений. Синтаксис и формальные языки.
- Иерархия грамматик Хомского. Регулярные грамматики. Контекстно-свободные грамматики. Основные понятия: терминал, нетерминал, правило. Форма записи Бакуса-Наура.
- Текст и дискурс. Методы сегментации текста с обучающей выборкой и без. Понятие связности текста.
- Автоматическое определение отношений связности. Коммуникативная структура текста. Понятия тема, рема, информационный статус.
- Теория риторической структуры. Риторические отношения. Анализ нарративной структуры. Разрешение анафоры и анализ кореференции.
- Методы извлечения информации из текстов на естественном языке. Словарные методы. Синтаксические шаблоны.
- Распознавание именованных сущностей. Извлечение отношений. Извлечение ключевых слов текста.
- Автоматический анализ стиля. Стилометрия.
- Автоматическое определение авторства: краткая история и обзор методов.
- Формальные и лингвистические черты для стилистического анализа. Автоматическое определение жанровой принадлежности текста.

**Перечень вопросов к зачету с оценкой (промежуточной аттестации по итогам освоения дисциплины):**

1. Введение в анализ данных. Проблема обработки данных. Матрица данных. Гипотезы компактности и скрытых факторов.
2. Структура матрицы данных и задачи обработки. Матрица объект-объект и признак-признак.
3. Расстояние и близость. Измерение признаков. Отношения и их представление.
4. Основные проблемы измерений. Основные типы шкал. Проблема адекватности. Основные задачи анализа и интерпретации данных.
5. Классификация данных с использованием детерминированных моделей.
6. Решающие поверхности и дискриминантные функции. Линейные дискриминантные функции классификатора по минимуму расстояния.
7. Линейная разделимость. Кусочно-линейные дискриминантные функции.
8. Нелинейные дискриминантные функции.
9. Фи-машины. Потенциальные функции как дискриминантные функции. Пространство весов.
10. Процедуры обучения с коррекцией ошибок: правило с фиксированным приращением, правило абсолютной коррекции, частично корректирующее правило.
11. Обобщенные градиентные методы. Персепtronный критерий. Процедуры обучения на основе минимальной среднеквадратичной ошибки: псевдоинверсный метод, метод Хо-Кашьпа.
12. Классификация данных на основе статистических моделей.
13. Функция потерь. Байесовская дискриминантная функция.
14. Принятие решения по максимуму правдоподобия. Оптимальная дискриминантная функция для нормально распределенных образов.
15. Дискриминантная функция Фишера. Множественный дискриминантный анализ.
16. Пошаговый дискриминантный анализ. Ошибки классификации.
17. Примеры построения статистических дискриминантных функций для различных статистических нескольких моделей данных.
18. Обучение для статистических дискриминантных функций.

19. Оценки максимального правдоподобия, байесовские оценки.
20. Непараметрическое оценивание. Парзеновские окна, метод непараметрического оценивания на основе К-ближайшего соседства.
21. Кластер-анализ. Основные типы задач кластер-анализа.
22. Меры подобия и функции расстояния. Выбор критерия кластеризации. Кластерные методы, основанные на евклидовой метрике.
23. Иерархическая кластеризация. Метод К-внутригрупповых средних.
24. Использование методов теории графов в задачах кластеризации.
25. Кластеризация на основе анализа плотностей вероятностей.
26. Методы снижения размерностей данных. Анализ матриц исходных данных.
27. Метод главных компонент. Корреляционная матрица и ее основные свойства.
28. Собственные векторы и собственные числа корреляционной матрицы. Приведение корреляционной матрицы к диагональной форме.
29. Геометрическая интерпретация главных компонент на плоскости.
30. Модели факторного анализа. Оценка факторных нагрузок методом максимального правдоподобия и центроидным методом.
31. Вращение факторов и их интерпретация.
32. Использование кластеризации признаков для снижения размерности.
33. Многомерное шкалирование (МИ). Метрический и неметрический подход к МИ.
34. Методы ортогонального проектирования.
35. Нелинейные методы МИ. Многомерное шкалирование неметрических данных. Многомерные развертки.
36. Методы прогнозирования временных рядов.
37. Классификация методов прогнозирования. Оценивание трендов.
38. Методы скользящего среднего. Экспоненциальное сглаживание.
39. Регрессионный анализ и прогнозирование. Линейные параметрические модели временных рядов.
40. Методы оценивания моделей авторегрессии, скользящего среднего и смешанных моделей.
41. Сезонные модели. Прогнозирование на основе параметрических моделей. Прогнозирование с использованием нейронных сетей.
42. Системы DATA MINING. в задачах анализа и интерпретации данных.
43. Понятие об интеллектуальных системах анализа и интерпретации данных. DATA MINING - системы извлечения новых знаний из данных. Типы систем DATA MINING -предметно-ориентированные аналитические системы, статистические пакеты, нейронные сети, деревья решений, обнаружение логических закономерностей, генетические алгоритмы, системы визуализации многомерных данных
44. Автоматическая обработка языка и компьютерная лингвистика. Задачи автоматической обработки текста в научных исследованиях.
45. Основные задачи компьютерной лингвистики и история развития автоматической обработки языка.
46. Иерархия языковых уровней и стандартный цикл обработки текста (графематика — морфология — синтаксис — семантика).
47. Основные задачи автоматической обработки текста: токенизация и нормализация текста; сегментация на предложения; стемминг; лемматизация и частеречные теги; снятие омонимии; парсинг — поверхностный и полный; кореференция и разрешение анафоры.
48. Задачи высокогоуровневого анализа: извлечение фактов и отношений, анализ оценок (sentiment analysis).
49. Инструменты морфологического анализа для русского языка. Понятия словоформа, лексема, лемма, основа. Стемминг. Алгоритм Портера. Stemka. Лемматизация.
50. Словарный метод — грамматический словарь Зализняка. mystem. АОТ. rutmorph. Грамматическая омонимия. Разметка частей речи. TreeTagger. TnT.
51. Частотный анализ лексики и ключевые слова. Частотное распределение лексики в языке.
52. Закон Ципфа. Доля hapax legomena. Скорость роста словаря. Коэффициент лексического разнообразия (type/token ratio).

53. Распределение лексики в текстах коллекции. Взвешенная частотность. TF-IDF. Прочие меры лексической дисперсии. Мера отклонения пропорций DP и DPnorm.
54. Извлечение ключевых слов. Метод контрастного корпуса. Отношение правдоподобия. Диахронический анализ лексической частотности.
55. Локальные модели контекста. Вероятностные языковые модели. Понятие N-грамм.
56. Буквенные и словарные n-граммы. Контекстное окно. Применения N-грамм в автоматической обработке языка. Роль биграмм и триграмм.
57. Определение языка по письменности. Языковые модели. Цепь Маркова. Коллокации.
58. Формальные определения и лингвистический смысл коллокаций. Меры ассоциации.
59. Коэффициент взаимной информации (MI). T-score. Отношение правдоподобия (log-likelihood).
60. Статистические тесты ассоциации: хи-квадрат и Fisher exact test.
61. Выделение коллокаций по синтаксическому шаблону. Разрывные коллокации.
62. Автоматическое определение тематики. Векторное представление текста для задач информационного поиска.
63. Открытые и закрытые классы слов. Стоп-слова. Динамические списки стоп слов. Порог отсечения по частотности и DF. Дистрибутивная семантика. Совместная встречаемость и семантическая близость.
64. Кластеризация текстов. Задачи и область применения кластерных методов.
65. Виды кластеризации: плоские, агglomerативные, нечеткие. Меры близости: евклидово расстояние, косинусная мера. Популярные алгоритмы кластеризации: k-средних, DBCLUST, спектральные алгоритмы.
66. ПО для кластеризации текстов. Пакеты кластеризации для R. gCLUTO.
67. Классификация текстов. Машинное обучение с учителем и без учителя в задачах классификации текстов.
68. Популярные алгоритмы классификации: наивный байесовский метод, метод опорных векторов, деревья принятия решений.
69. ПО для классификации текстов. SVMLight. Пространственное моделирование семантических отношений (word space). Латентный семантический анализ.
70. Вероятностный латентно семантический анализ. Тематическое моделирование. Метод латентного размещения Дирихле. ПО для латентного семантического анализа и тематического моделирования. Mallet, sTMT.
71. Извлечение мнений и оценок (Sentiment analysis). Область применения методов извлечения мнений и оценок.
72. Типы оценочных текстов: позитивный, негативный, нейтральный. Оценочные шкалы. Классификация документов по оценке. Извлечение оценочных предложений и фрагментов.
73. Определение предмета оценки. Методы извлечения оценок. Словарные методы. Машинное обучение. Комбинирование источников. Проблемы и ограничения методов извлечения оценок.
74. Извлечение фактов и отношений. Синтаксис и формальные языки.
75. Иерархия грамматик Хомского. Регулярные грамматики. Контекстно-свободные грамматики. Основные понятия: терминал, нетерминал, правило. Форма записи Бакуса-Наура.
76. Текст и дискурс. Методы сегментации текста с обучающей выборкой и без. Понятие связности текста.
77. Автоматическое определение отношений связности. Коммуникативная структура текста. Понятия тема, рема, информационный статус.
78. Теория риторической структуры. Риторические отношения. Анализ нарративной структуры. Разрешение анафоры и анализ кореференции.
79. Методы извлечения информации из текстов на естественном языке. Словарные методы. Синтаксические шаблоны.
80. Распознавание именованных сущностей. Извлечение отношений. Извлечение ключевых слов текста.
81. Автоматический анализ стиля. Стилометрия.
82. Автоматическое определение авторства: краткая история и обзор методов.

83. Формальные и лингвистические черты для стилистического анализа. Автоматическое определение жанровой принадлежности текста.

**Перечень тем лабораторных работ:**

**Лабораторная работа №1.** Предварительный анализ данных с использованием специализированного программного обеспечения (по выбору преподавателя).

**Лабораторная работа №2.** Изучение методов дискриминантного анализа с использованием специализированного программного обеспечения (по выбору преподавателя).

**Лабораторная работа №3.** Изучение методов кластер-анализа с использованием специализированного программного обеспечения (по выбору преподавателя).

**Лабораторная работа №4.** Изучение методов факторного-анализа с использованием специализированного программного обеспечения (по выбору преподавателя).

**Лабораторная работа №5.** Классификация данных и изучение методов снижения размерности данных с использованием специализированного программного обеспечения (по выбору преподавателя).

**Лабораторная работа №6.** Изучение методов прогнозирования временных рядов с использованием специализированного программного обеспечения (по выбору преподавателя).

**Вопросы и задания для самостоятельной работы студентов:**

- Классификация данных с использованием детерминированных моделей.
- Классификация данных на основе статистических моделей.
- Кластер-анализ данных.
- Методы снижения размерностей данных.
- Методы прогнозирования временных рядов.
- Системы DATA MINING в задачах анализа и интерпретации данных.
- Современные пакеты прикладных программ для решения задач обработки экспериментальных данных.
- Частотный анализ лексики и ключевые слова.
- Локальные модели контекста. Вероятностные языковые модели.
- Автоматическое определение тематики при исследовании текстов.
- Извлечение мнений и оценок при исследовании текстов.
- Извлечение фактов и отношений при исследовании текстов.
- Автоматический анализ стиля при исследовании текстов.
- Основы обработки неструктурированных (текстовых) данных в корпоративных информационных системах (ERP, АСУП и др.) и современных веб-приложениях.
- Системы корпоративного поиска ESR (Enterprise Search and Retrieval) и их компоненты.
- Информационные и программные средства лингвистического обеспечения. Организация лингвистического обеспечения в АСОИУ.
- Информационно-поисковые языки. Системы метаданных.
- Классификационные, вербальные, фактографические языки.
- Лингвистические процессы. Системы автоматической обработки текста. Лингвистические базы данных.
- Методы и решения в системах организации знаний: автоформализация, формализация, лексикографическое (словарное) и логико-интуитивное описание, организация, анализ и извлечение знаний.
- Онтологии: языки, инструменты создания, структура.
- Статистические методы Text Mining.
- Законы распределения частот слов. Закон Ципфа. Распределение Мандельброта. Закон Бредфорда. Формирование ядра релевантных текстов.
- Проектирование компонентов лингвистического обеспечения АСОИУ.
- Разработка информационных компонентов для систем электронного документооборота (СЭД): проектирование словарей ключевых понятий; предметных, именных указателей;

тематических словарей (по группам документов предприятия); электронных картотек товаров и услуг; информационно-поисковых индексов.

- Естественно-языковые интерфейсы. Алгоритмы морфологического анализа и лемматизации. Синтаксический и семантический анализ.

## 7. УЧЕБНО-МЕТОДИЧЕСКОЕ И ИНФОРМАЦИОННОЕ ОБЕСПЕЧЕНИЕ ДИСЦИПЛИНЫ

### **а) Основная литература:**

1. Статистический анализ данных, моделирование и исследование вероятностных закономерностей. Компьютерный подход / Б.Ю. Лемешко, С.Б. Лемешко, С.Н. Постовалов и др. - М.: НИЦ ИНФРА-М, 2015. - 890 с. Режим доступа: <http://znanium.com/catalog.php?bookinfo=515227>
2. Численный вероятностный анализ неопределенных данных/ДобронецБ.С., ПоповаО.А. - Краснояр.: СФУ, 2014. - 168 с.: ISBN 978-5-7638-3093-4 Режим доступа: <http://znanium.com/catalog.php?bookinfo=549444>
3. Информационные системы предприятия: Учебное пособие / А.О. Варфоломеева, А.В. Коряковский, В.П. Романов. - М.: НИЦ ИНФРА-М, 2013. - 283 с. Режим доступа: <http://znanium.com/catalog.php?bookinfo=344985>

### **б) Дополнительная литература:**

1. Интеллектуальный анализ данных и систем управления бизнес-правилами в телекоммуникациях: Монография / Р.Р. Вейнберг. - М.: НИЦ ИНФРА-М, 2015. - 173 с. Режим доступа: <http://znanium.com/catalog.php?bookinfo=520998>
2. Прикладные методы анализа статистических данных: учеб. пособие / Горяннова Е.Р., Панков А.Р., Платонов Е.Н. - М. : ИД Высшей школы экономики, 2012. <http://www.studentlibrary.ru/book/ISBN9785759808664.html> 310с.
3. Методы многомерного анализа статистических данных : учеб. пособие/ В.М. Симчера. - М. : Финансы и статистика, 2008. - <http://www.studentlibrary.ru/book/ISBN9785279031849.html> 400 с.

### **в) Периодические издания:**

1. Журнал «Вопросы защиты информации». Режим доступа: [http://i-vimi.ru/editions/detail.php?SECTION\\_ID=155/](http://i-vimi.ru/editions/detail.php?SECTION_ID=155/);
2. Журнал "Information Security/Информационная безопасность". Режим доступа: <http://www.itsec.ru/insec-about.php>.
3. Ежемесячный теоретический и прикладной научно-технический журнал «Информационные технологии». Режим доступа <http://novtex.ru/IT/>.

### **г) Программное обеспечение и Интернет-ресурсы:**

<http://www.dialog-21.ru/>— Диалог.Международная конференция по компьютерной лингвистике.

<http://nlpub.ru>— Каталог лингвистических ресурсов для обработки русского языка.

<http://www.regular-expressions.info>— The Premier website about Regular Expressions.

<http://sentiment.christopherpotts.net/>— Sentiment symposium tutorial.

<http://www.aclweb.org/anthology/>— ACL Anthology

A Digital Archive of Research Papers in Computational Linguistics.

**Программные средства.** Для успешного освоения дисциплины, студент использует следующие программные средства: - Программа построения частотных словарей.

<http://alingva.ru/index.php/lingvosoft/12-ngramfrequency>; - mystem. Морфологический анализатор для русского языка, <http://company.yandex.ru/technologies/mystem/>- LSA. Латентно-семантический анализ текстовых данных.

<http://alingva.ru/index.php/lingvosoft/17--lsa>; - Tomita-парсер. Инструмент для извлечения структурированных данных из текста на естественном языке. <http://api.yandex.ru/tomita/>; - Модуль Perl Text::NSP. N-gram statistics and association measures.

<http://search.cpan.org/dist/Text-NSP/lib/Text/NSP/Measures.pm>; - Stanford Topic Modeling Toolbox; - <http://nlp.stanford.edu/software/tmt/tmt-0.4/>

Тестовые массивы текстов для обработки публикуются на сайте:  
<http://maslinsky.spb.ru/courses/cmfa2013/>

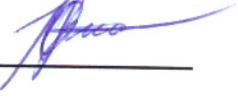
## **8. МАТЕРИАЛЬНО-ТЕХНИЧЕСКОЕ ОБЕСПЕЧЕНИЕ ДИСЦИПЛИНЫ (МОДУЛЯ)**

ауд. 408-2, Лекционная аудитория, количество студенческих мест – 50, площадь 60 м<sup>2</sup>, оснащение: мультимедийное оборудование (интерактивная доска Hitachi FX-77WD, проектор BenQ MX 503 DLP 2700ANSI XGA), ноутбук Lenovo Idea Pad B5045

ауд. 427а-2, лаборатория сетевых технологий, количество студенческих мест – 14, площадь 36 м<sup>2</sup>, оснащение: компьютерный класс с 8 рабочими станциями Core 2 Duo E8400 с выходом в Internet, 3 маршрутизатора Cisco 2800 Series, 6 маршрутизаторов Cisco 2621, 6 коммутаторов Cisco Catalyst 2960 Series, 3 коммутатора Cisco Catalyst 2950 Series, коммутатор Cisco Catalyst Express 500 Series, проектор BenQ MP 620 P, экран настенный рулонный. Лицензионное программное обеспечение: операционная система Windows 7 Профессиональная, офисный пакет приложений Microsoft Office Профессиональный плюс 2007, бесплатно распространяемое программное обеспечение: линейка интегрированных сред разработки Visual Studio Express 2012, программный продукт виртуализации Oracle VM VirtualBox 5.0.4, симулятор сети передачи данных Cisco Packet Tracer 7.0, интегрированная среда разработки программного обеспечения IntelliJ IDEA Community Edition 15.0.3.

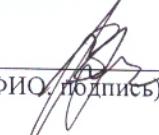
ауд. 427б-2, УНЦ «Комплексная защита объектов информатизации», количество студенческих мест – 15, площадь 52 м<sup>2</sup>, оснащение: компьютерный класс с 7 рабочими станциями Alliance Optima P4 с выходом в Internet, коммутатор D-Link DGS-1100-16 мультимедийный комплект (проектор Toshiba TLP X200, экран настенный рулонный), прибор ST-031Р «Пиранья-Р» многофункциональный поисковый, прибор «Улан-2» поисковый, вибраакустический генератор шума «Соната АВ 1М», имитатор работы средств нелегального съема информации, работающих по радиоканалу «Шиповник», анализатор спектра «GoodWill GSP-827», индикатор поля «SEL SP-75 Black Hunter», устройство блокирования работы систем мобильной связи «Мозайка-3», устройство защиты телефонных переговоров от прослушивания «Прокрут 2000», диктофон Edic MINI Hunter, локатор «Родник-2К» нелинейный, комплекс проведения акустических и вибраакустических измерений «Спрут мини-А», видеорегистратор цифровой Best DVR-405, генератор Шума «Гном-3», учебно-исследовательский комплекс «Сверхширокополосные беспроводные сенсорные сети» (Nano Xaos), сканирующий приемник «Icom IC-R1500», анализатор сетей Wi-Fi Fluke AirCheck с активной антенной. Лицензионное программное обеспечение: Windows 8 Профессиональная, офисный пакет приложений Microsoft Office Профессиональный плюс 2010, бесплатно распространяемое программное обеспечение: линейка интегрированных сред разработки Visual Studio Express 2012, инструмент имитационного моделирования AnyLogic 7.2.0 Personal Learning Edition, интегрированная среда разработки программного обеспечения IntelliJ IDEA Community Edition 14.1.4.

Программа составлена в соответствии с требованиями ФГОС ВО по специальности  
10.05.04 "Информационно-аналитические системы безопасности", специализация  
«автоматизация информационно-аналитической деятельности»

Рабочую программу составил доцент кафедры ИЗИ к.т.н. Монахов Ю.М.  
(ФИО, подпись) 

Рецензент  
(представитель работодателя) Заместитель руководителя РАЦ ООО «ИнфоЦентр»

к.т.н. Вертилевский Н.В.

(место работы, должность, ФИО, подпись) 

Программа рассмотрена и одобрена на заседании кафедры ИЗИ

Протокол № 7 от 28.12.16 года

Заведующий кафедрой д.т.н., профессор

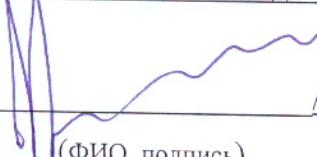
 /M.Yu. Монахов/

(ФИО, подпись)

Рабочая программа рассмотрена и одобрена на заседании учебно-методической комиссии по специальности 10.05.04 "Информационно-аналитические системы безопасности", специализация «автоматизация информационно-аналитической деятельности»

Протокол № 4 от 28.12.16 года

Председатель комиссии д.т.н., профессор

 /M.Yu. Монахов/

(ФИО, подпись)

### ЛИСТ ПЕРЕУТВЕРЖДЕНИЯ РАБОЧЕЙ ПРОГРАММЫ ДИСЦИПЛИНЫ (МОДУЛЯ)

Рабочая программа одобрена на 2017/18 учебный год

Протокол заседания кафедры № 1 от 28.08.17 года

Заведующий кафедрой д.т.н., профессор

 /M.Yu. Монахов/

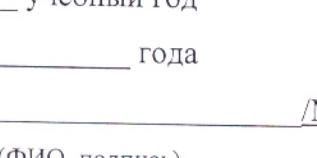
(ФИО, подпись)

### ЛИСТ ПЕРЕУТВЕРЖДЕНИЯ РАБОЧЕЙ ПРОГРАММЫ ДИСЦИПЛИНЫ (МОДУЛЯ)

Рабочая программа одобрена на \_\_\_\_\_ учебный год

Протокол заседания кафедры № \_\_\_\_\_ от \_\_\_\_\_ года

Заведующий кафедрой д.т.н., профессор

 /M.Yu. Монахов/

(ФИО, подпись)

**Министерство образования и науки Российской Федерации**  
Федеральное государственное бюджетное образовательное учреждение  
высшего образования  
**«Владимирский государственный университет**  
**имени Александра Григорьевича и Николая Григорьевича Столетовых»**  
**(ВлГУ)**

Институт \_\_\_\_\_

Кафедра \_\_\_\_\_

Актуализированная  
рабочая программа  
рассмотрена и одобрена  
на заседании кафедры  
протокол № \_\_\_\_ от \_\_\_\_ 20 \_\_\_\_ г.  
Заведующий кафедрой  
\_\_\_\_\_  
(подпись, ФИО)

**Актуализация рабочей программы дисциплины**

---

(наименование дисциплины)

Направление подготовки

Профиль/программа подготовки

Уровень высшего образования

Форма обучения

Владимир 20 \_\_\_\_

Рабочая программа учебной дисциплины актуализирована в части рекомендуемой литературы.

Актуализация выполнена: \_\_\_\_\_  
(подпись, должность, ФИО)

а) основная литература: \_\_\_\_\_

б) дополнительная литература: \_\_\_\_\_

в) периодические издания: \_\_\_\_\_

в) интернет-ресурсы: \_\_\_\_\_