

**ГОУ ВПО Российско-Армянский (Славянский)  
университет**



**УЧЕБНО-МЕТОДИЧЕСКИЙ КОМПЛЕКС ДИСЦИПЛИНЫ**

**Наименование дисциплины: Б1.В.01 Статистический анализ данных**

**Автор: Костанян Артавазд Арамович, к.ф.м.н., ст. преподаватель**

**Направление подготовки: 11.04.04 Электроника и наноэлектроника**  
**Наименование образовательной программы: Квантовая и оптическая электроника**

**Согласовано:**

Заведующий Кафедрой общей физики и квантовых наноструктур

Айрапетян Д.Б.



(подпись)

# 1. АННОТАЦИЯ

## 1.1. Краткое описание содержания данной дисциплины;

Данный курс знакомит с базовыми методами статистического анализа данных: классические линейные методы регрессии – логистическая и линейная регрессии, классификации: логистическая регрессия и линейный дискриминантный анализ, древовидные модели, модели без учителя и т.д. Центральная задача анализа данных заключается в выборе наилучшего статистического метода для данной задачи. В течении курса студенты знакомятся также с методами перекрестной проверки, бутстрапом, ознакомятся с нелинейными методами. В рамках курса исследуются также древовидные методы, включая группировку, случайные леса, метод опорных векторов — набор подходов для выполнения как линейной, так и нелинейной классификации.

## 1.2. Трудоемкость в академических кредитах и часах, формы итогового контроля (экзамен/зачет);

2 з.е. (72ч.), зачет

## 1.3. Взаимосвязь дисциплины с другими дисциплинами учебного плана специальности (направления)

Методы машинного обучения в материаловедении, Элементы квантовой и оптической электроники, Компьютерные технологии в физике, Квантоворазмерные системы нанoeлектроники.

## 1.4. Результаты освоения программы дисциплины:

<b>Код компетенции</b> (в соответствии рабочим с учебным планом)	<b>Наименование компетенции</b> (в соответствии рабочим с учебным планом)	<b>Код индикатора достижения компетенций</b> (в соответствии рабочим с учебным планом)	<b>Наименование индикатора достижений компетенций</b> (в соответствии рабочим с учебным планом)
УК-1.	Способен осуществлять критический анализ проблемных ситуаций на основе системного подхода, выработать стратегию действий	УК-1.1	"Знает методы анализа проблемной ситуации как системы, выявляя ее составляющие и связи между ними; знает способы

		УК-1.2	определения пробелы в информации, необходимой для решения проблемной ситуации, и проектирования процессов по их устранению." "Умеет критически оценивать надежность источников информации, работать с противоречивой информацией из разных источников; Разрабатывать и содержательно аргументировать стратегию решения проблемной ситуации на основе системного и междисциплинарных подходов"
		УК-1.3	Владеет навыками использования логико-методологического инструментария для критической оценки современных концепций философского и социального характера в своей предметной области.
УК-2.	Способен управлять проектом на всех этапах его жизненного цикла	УК-2.1 УК-2.2 УК-2.3	Знает, как формулировать цели, задачи, значимости, ожидаемые результатов проекта. Умеет определять потребности в ресурсах для реализации проекта; Разрабатывать план реализации проекта. Владеет навыками контроля реализации проекта; навыками оценки эффективности реализации проекта и разработки плана действий по его корректировке.
ПК-2	Способен разрабатывать эффективные алгоритмы решения сформулированных задач с использованием современных языков	ПК-2.1	Знает методы разработки эффективных алгоритмов решения научно-исследовательских задач

	программирования и обеспечивать их программную реализацию	ПК-2.2  ПК-2.3	Умеет использовать алгоритмы решения исследовательских задач с использованием современных языков программирования Владеет навыками разработки стратегии и методологии исследования изделий микро- и нанoeлектроники
--	---	----------------------	--

## 2. УЧЕБНАЯ ПРОГРАММА

### 2.1. Цели и задачи дисциплины

**Целью** изучения дисциплины «Статистический анализ данных» является формирование знаний и навыков в области применения технологий для анализа данных. Она включает в себя разработку, оценку и адаптацию статистических моделей в предметной среде основы анализа данных, например, метод ближайших соседей, байесовские классификаторы, метод опорных векторов, древовидные методы, случайный лес, градиентный бустинг, а также комплекс прикладных программ для реализации данных технологий.

**Задачами** дисциплины являются:

*приобрести*

- умение применять методы статистического анализа для решения практических задач по обработке данных
- умение анализировать задачи статистического анализа данных и осуществлять взвешенный выбор того или иного решения

*ознакомить*

- студентов с компьютерными технологиями обработки многомерных неструктурированных массивов разнородной статистической информации

*научить*

- интерпретировать полученные результаты построенных статистических моделей.

В результате изучения данного курса студенты получают знания о принципах применения статистического анализа данных, о его методах анализа и инжиниринга,

приобретут навыки и умения применения методов решающих деревьев, случайного леса, и т.п..

**2.2. Трудоемкость дисциплины и виды учебной работы (в академических часах и зачетных единицах) (удалить строки, которые не будут применены в рамках дисциплины)**

Виды учебной работы	Всего, в акад. часах	Распределение по семестрам
		3 сем
1	2	3
<b>1. Общая трудоемкость изучения дисциплины по семестрам, в т. ч.:</b>	<b>72</b>	<b>72</b>
1.1. Аудиторные занятия, в т. ч.:	<b>16</b>	<b>16</b>
1.1.1. Лекции	<b>10</b>	<b>10</b>
1.1.2. Практические занятия, в т. ч.	<b>6</b>	<b>6</b>
1.2. Самостоятельная работа, в т. ч.:	<b>56</b>	<b>56</b>
1.3. Консультации		
Итоговый контроль (Экзамен, Зачет, диф. зачет - указать)	<b>Зачет</b>	<b>Зачет</b>

**2.3. Содержание дисциплины**

**2.3.1. Тематический план и трудоемкость аудиторных занятий (модули, разделы дисциплины и виды занятий) по рабочему учебному плану**

Разделы и темы дисциплины	Всего (ак. часов)	Лекции( ак. часов)	Практ. Занятия (ак. часов)
1	2=3+4	3	4
Тема 1. Введение	1	2	
Тема 2. Задача предмета анализа данных	1		
Тема 3. Линейная регрессия	2	2	2
Тема 4. Классификации	2		
Тема 5. Методы повторной выборки	2	2	2
Тема 6. Регуляризация линейной модели	2		
Тема 7. Нелинейные модели	1	2	
Тема 8. Древовидные модели	1		
Тема 9. Модели без учителя	2	2	2
Тема 10. Многократное тестирование	2		
<b>ИТОГО</b>	<b>16</b>	<b>10</b>	<b>6</b>

## 2.3.2. Краткое содержание разделов дисциплины в виде тематического плана

### **Тема 1. Введение**

Что такое статистический анализ данных? Зачем и как оценивать функцию предсказания?

### **Тема 2. Задача предмета**

Компромисс между точностью прогнозирования и интерпретируемостью модели. Контролируемая и неконтролируемая модели. Проблемы регрессии и классификации. Оценка точности модели. Измерение качества соответствия. Компромисс между смещением и дисперсией.

### **Тема 3. Линейная регрессия**

Простая линейная регрессия. Оценка коэффициентов. Оценка точности коэффициентов. Оценка точности модели. Множественная линейная регрессия. Оценка коэффициентов регрессии. Другие соображения по регрессионной модели. Качественные предикторы. Расширения линейной модели. Сравнение линейной регрессии с K-ближайшими соседями

### **Тема 4. Классификации**

Обзор классификации. Логистическая регрессия. Оценка коэффициентов регрессии. Множественная логистическая регрессия. Мультиномиальная логистическая регрессия. Генеративные модели для классификации. Линейный дискриминантный анализ. Квадратичный дискриминантный анализ. Наивный байесовский алгоритм. Аналитическое сравнение. Обобщенные линейные модели. Обобщенные линейные модели в большей общности

### **Тема 5. Методы повторной выборки**

Перекрестная проверка. Подход к набору проверки. Перекрестная проверка с исключением одного. Перекрестная проверка k-кратности. Компромисс смещения и дисперсии для перекрестной проверки k-кратности. Перекрестная проверка в задачах классификации. Бутстрап.

### **Тема 6. Регуляризация линейной модели**

Выбор подмножества. Пошаговый выбор. Выбор оптимальной модели. Методы сжатия. Регрессия гребня. Лассо. Выбор параметра настройки. Методы сокращения размерности. Регрессия главных компонент. Частичные наименьшие квадраты. Данные больших размерностей. Регрессия в больших размерностях. Интерпретация результатов в больших размерностях.

## **Тема 7. Нелинейные модели**

Полиномиальная регрессия. Ступенчатые функции. Базисные функции. Сплайны регрессии. Кусочные полиномы. Ограничения и сплайны. Представление базиса сплайна. Выбор количества и расположения узлов. Сравнение с полиномиальной регрессией. Сглаживающие сплайны. Обзор сглаживающих сплайнов. Выбор параметра сглаживания. Локальная регрессия. Обобщенные аддитивные модели (GAM).

## **Тема 8. Древоподобные модели**

Основы деревьев решений. Деревья регрессии. Деревья классификации. Деревья против линейных моделей. Преимущества и недостатки деревьев. Бэггинг. Случайные леса. Бустинг. Деревья аддитивной регрессии Байеса. Краткое изложение методов ансамбля деревьев

## **Тема 9. Модели без учителя**

Проблема модели без учителя. Анализ главных компонент. Что такое главные компоненты? Объяснение доли дисперсии. Подробнее о PCA. Другие применения главных компонент. Пропущенные значения и заполнение матриц. Методы кластеризации. Кластеризация методом k-средних. Иерархическая кластеризация.

## **Тема 10. Многократное тестирование**

Краткий обзор проверки гипотез. Проверка гипотезы. Ошибки типа I и типа II. Проблема множественного тестирования. Частота ошибок по семействам. Что такое частота ошибок по семействам? Подходы к контролю частоты ошибок по семействам.

### **2.3.3. Краткое содержание практических занятий**

Разработка и представление результатов, на основе изученных методов в формате ноутбуков (Jupyter Notebook) разработанных средствами языка Python.

1. Практическое занятие 1: простая линейная регрессия, множественная линейная регрессия, многомерное качество соответствия, качественные предикторы.
2. Практическое занятие 2: логистическая регрессия, линейный дискриминантный анализ, квадратичный дискриминантный анализ, наивный байесовский алгоритм, классификация k-NN, регрессия Пуассона
3. Практическое занятие 3: приближение набора проверки, перекрестная проверка. перекрестная проверка k-кратности, бутстрап.

4. Практическое занятие 4: методы выбор подмножеств: пошаговый выбор, выбор оптимального подмножества, регрессия гребня и лассо, регрессия главных компонент, частичные наименьшие квадраты
5. Практическое занятие 5: полиномиальная регрессия и ступенчатые функции, сплайны, сглаживающие сплайны, обобщенные аддитивные модели, GAM регрессии, тест ANOVA и аддитивные модели
6. Практическое занятие 6: деревья решений, деревья регрессии и классификации, бэггинг, бустинг, деревья аддитивной регрессии Байеса
7. Практическое занятие 7: анализ главных компонент, кластеризация методом k-средних, иерархическая кластеризация.
8. Практическое занятие 8: проверка гипотезы, ошибки типа I и типа II, частота ошибок по семействам, частота ложных открытий

#### 2.3.4. Материально-техническое обеспечение дисциплины

Компьютерная аудитория, проектор

#### 2.4. Модульная структура дисциплины с распределением весов по формам контролей

Формы контролей	Вес формы (форм) текущего контроля в результирующей оценке текущего контроля (по модулям)		Вес формы промежуточного контроля в итоговой оценке промежуточного контроля		Вес итоговой оценки промежуточного контроля в результирующей оценке промежуточных контролей		Вес итоговой оценки промежуточного контроля в результирующей оценке промежуточных контролей (семестровой оценке)	Вес результирующей оценки промежуточных контролей и оценки итогового контроля в результирующей оценке итогового контроля
	M1 <sup>1</sup>	M2	M1	M2	M1	M2		
<b>Вид учебной работы/контроля</b>	M1 <sup>1</sup>	M2	M1	M2	M1	M2		
Контрольная работа <i>(при наличии)</i>			0.5	0.5				
Устный опрос <i>(при наличии)</i>								
Тест <i>(при наличии)</i>								
Лабораторные работы <i>(при наличии)</i>	0.5	0.5						
Письменные домашние задания <i>(при наличии)</i>								
Реферат <i>(при наличии)</i>								
Эссе <i>(при наличии)</i>								
Проект <i>(при наличии)</i>								
Решение задач	0.5	0.5						

<sup>1</sup> Учебный Модуль

Веса результирующих оценок текущих контролей в итоговых оценках промежуточных контролей					0.5	0.5		
Веса оценок промежуточных контролей в итоговых оценках промежуточных контролей								
Вес итоговой оценки 1-го промежуточного контроля в результирующей оценке промежуточных контролей							0.5	
Вес итоговой оценки 2-го промежуточного контроля в результирующей оценке промежуточных контролей							0.5	
Вес результирующей оценки промежуточных контролей в результирующей оценке итогового контроля								0.5
<b>Вес итогового контроля (Экзамен/зачет)</b> в результирующей оценке итогового контроля								0.5
	$\Sigma = 1$							

### 3. Теоретический блок

#### 3.1. Материалы по теоретической части курса

3.1.1. *Учебник*: G.James, D.Witten, T.Hastie, R.Tibshirani, J.Taylor, An Introduction to Statistical Learning With Applications in Python, Springer International Publishing (2023), <https://www.statlearning.com>

3.1.2. *Краткие конспекты лекций (слайды)*: <https://www.statlearning.com/resources-python>

### 4. Фонды оценочных средств

#### 4.1. Планы практических занятий

1. Решения регрессионных задач с использованием простой и множественной линейной регрессий
2. Решение задач классификации: логистическая регрессия, линейный дискриминантный анализ, наивный байесовский алгоритм, классификация k-NN
3. Задача приближения набора проверки, перекрестная проверка, перекрестная проверка k-кратности, бутстрап

4. Имплементация метода выбора подмножеств: пошаговый выбор, выбор оптимального подмножества, регрессия гребня и лассо, регрессия главных компонент, частичные наименьшие квадраты
5. Задача регрессии с использованием полиномиальной регрессии и ступенчатых функций, сплайны, сглаживающие сплайны, имплементация и ознокомление с тестом ANOVA
6. Классификация и регрессия с деревом регрессии и классификации, бэггинг, бустинг, деревья аддитивной регрессии Байеса
7. Реализация анализа главных компонент, кластеризация методом k-средних, иерархическая кластеризация
8. Проверка гипотезы, оценка ошибок типа I и типа II, частота ошибок по семействам

#### 4.2. Материалы по практической части курса

- 4.2.1 Задачи из учебника: G.James, D.Witten, T.Hastie, R.Tibshirani, J.Taylor, An Introduction to Statistical Learning With Applications in Python, Springer International Publishing (2023), <https://www.statlearning.com>
- 4.2.2 Jupyter ноутбуки на Python из ресурсов к задачку: [https://github.com/intro-stat-learning/ISLP\\_labs/blob/stable/Ch02-statlearn-lab.ipynb](https://github.com/intro-stat-learning/ISLP_labs/blob/stable/Ch02-statlearn-lab.ipynb)

#### 10.4. Вопросы и задания для самостоятельной работы студентов

1. Что такое распределение Стьюдента и t-критерий Стьюдента?
2. Что такое F-распределение и F-критерий?
3. Что такое проверка гипотез и как она работает (нулевая гипотеза, альтернативная гипотеза, p-значения и т. д.)?
4. Упражнение 8 на странице 129 из учебника.
5. Упражнение 9 на странице 129 из учебника.
6. Что такое классификатор Байеса и граница принятия решений Байеса?
7. Доказать, что представление логистической функции и логит-представление для модели логистической регрессии эквивалентны.
8. Доказать, что при предположении, что наблюдения в k-м классе взяты из распределения  $N(\mu_k, \sigma^2)$ , байесовский классификатор относит наблюдение к классу, для которого дискриминантная функция максимальна.
9. Упражнение 13 на странице 196 из учебника.

10. Упражнение 14 на странице 197 из учебника.
11. Что такое k-кратная кросс-валидация?
12. Упражнение 5 на странице 225 из учебника.
13. Упражнение 6 на странице 226 из учебника.
14. Упражнение 8 и 9 на странице 286 из учебника.
15. Упражнение 6, 7 и 8 на странице 327 из учебника.

#### **10.5. Перечень вопросов для зачета**

1. Точность прогнозирования и интерпретируемость статистической модели. Оценка точности модели, измерение качества соответствия.
2. Простая и множественная линейные регрессии. Оценка коэффициентов регрессии, точности модели. Качественные предикторы. Выбросы (outliers).
3. Логистическая регрессия. Множественная и мультиномиальная логистические регрессии.
4. Линейный и квадратичный дискриминантные анализы.
5. Наивный байесовский алгоритм. Сравнение методов классификации.
6. Перекрестная проверка. Подход к набору проверки.
7. Перекрестная проверка с исключением одного. Перекрестная проверка k-кратности.
8. Перекрестная проверка в задачах классификации. Бутстрап.
9. Выбор подмножества. Пошаговый выбор.
10. Выбор оптимальной модели. Методы сжатия.
11. Регрессия гребня. Лассо. Выбор параметра настройки.
12. Методы сокращения размерности. Регрессия главных компонент. Частичные наименьшие квадраты.
13. Полиномиальная регрессия. Ступенчатые функции.
14. Базисные функции. Сплайны регрессии. Кусочные полиномы.
15. Локальная регрессия. Обобщенные аддитивные модели. GAM для задач регрессии и классификации
16. Основы деревьев решений. Деревья регрессии. Деревья классификации. Деревья против линейных моделей. Бэггинг. Случайные леса.
17. Бустинг. Деревья аддитивной регрессии Байеса.

18. Прогнозирование временных рядов. Краткое описание RNN.
19. Анализ главных компонент. Что такое главные компоненты?
20. PCA и другие применения главных компонент. Пропущенные значения и заполнение матриц.
21. Методы кластеризации. Кластеризация методом k-средних. Иерархическая кластеризация.
22. Проверка гипотезы. Ошибки типа I и типа II. Проблема множественного тестирования. Частота ошибок по семействам.

#### 10.6. Образцы зачетных билетов

**ГОУ ВПО РОССИЙСКО-АРМЯНСКИЙ УНИВЕРСИТЕТ**  
**ИНЖЕНЕРНО-ФИЗИЧЕСКИЙ ИНСТИТУТ**  
**Кафедра общей физики и квантовых наноструктур**

**Направление: Электроника и наноэлектроника**  
**Дисциплина: Статистический анализ данных**  
**(магистратура II курс, I семестр)**

**Экзаменационный билет № \*\***

1. Регрессия главных компонент. Частичные наименьшие квадраты.
2. Методы кластеризации. Кластеризация методом k-средних.
3. Бустинг. Деревья аддитивной регрессии Байеса.

**ГОУ ВПО РОССИЙСКО-АРМЯНСКИЙ УНИВЕРСИТЕТ**  
**ИНЖЕНЕРНО-ФИЗИЧЕСКИЙ ИНСТИТУТ**  
**Кафедра общей физики и квантовых наноструктур**

**Направление: Электроника и наноэлектроника**  
**Дисциплина: Статистический анализ данных**  
**(магистратура II курс, I семестр)**

**Экзаменационный билет № \*\*\***

1. Логистическая регрессия. Множественная и мультиномиальная логистические регрессии.
2. Деревья классификации. Бэггинг.
3. PCA и другие применения главных компонент. Пропущенные значения и заполнение матриц.

2025г.

---

## **11. Методический блок**

### **11.4. Методика преподавания**

Лекционные занятия споследующими практическими заданиями, выполняемые как на вычислительных средствах в аудитории, так и дома на персональных компьютерах