



# Математическая статистика

## и основы анализа данных

### Введение

# Организационная информация

Первый семестр:

- ▶ Статистика
- ▶ Введение в анализ данных
- ▶ Введение в машинное обучение

Второй семестр:

- ▶ Машинное обучение

Лектор: Ольга Калиниченко

# Организационная информация

Профили:

- ▶ Физика;
- ▶ Биология;
- ▶ Педагогика;
- ▶ Теоретический (при наличии желающих).

Студенты кафедры инновационной педагогики выбирают профиль "педагогика".

# Организационная информация

Структура курса:

- ▶ Домашние задания на 2 недели;
- ▶ Несколько теоретических и практических задач;
- ▶ Практика на языке Python;
- ▶ Сдача заданий через телеграм-бот (подробности до конца недели);
- ▶ Гостевые лекции про реальное применение статистики  
(отдельно для направлений)

Первые две лекции вводные, первые два задания даются на одну неделю и являются отборочными.

**2 курс ФБМФ и ФМХФ:**

Результат курса можно зачесть за курс математической статистики.

Подробности на странице курса.

# Зачем это нужно?

## Физика:

- ▶ Анализ физических данных
- ▶ Оценки параметров и построение моделей физических систем
- ▶ Проверка гипотез и многое другое

## Биология:

- ▶ Проверка гипотез в экспериментах (почти в каждой статье)
- ▶ Анализ омиксных данных
- ▶ Классификация клеток, оценка уровня экспрессии генов, кластеризация микробиома, определение вторичной структуры ДНК и многое другое

## Химия:

- ▶ Нахождение химического соединения с заданными свойствами
- ▶ Оценка параметров и построение моделей реакций
- ▶ Проверка гипотез и многое другое

**Работа с погрешностями тоже основана на статистике.**

# Основные темы (весна)

1. Простые методы анализа данных и машинного обучения;
2. Точечное оценивание;

Свойства и методы поиска оценок; качество оценок и поиск наилучших оценок;  
оценки специального вида.

3. Доверительные интервалы;

4. Непараметрический подход;

Бутстреп; ядерная оценка плотности.

5. Проверка гипотез;

Критерии; p-value; множественная проверка гипотез; критерии согласия.

6. Линейная регрессия;

МНК; отбор признаков.

7. Корреляционный анализ;

8. Дисперсионный анализ.

9. Байесовский подход, метод Монте-Карло;

## Основные книги по курсу

- ▶ Лагутин М.Б., Наглядная математическая статистика;
- ▶ L. Wasserman, All of Statistics;
- ▶ Russell B. Millar, Maximum Likelihood Estimation and Inference;
- ▶ Боровков А.А., Математическая статистика.

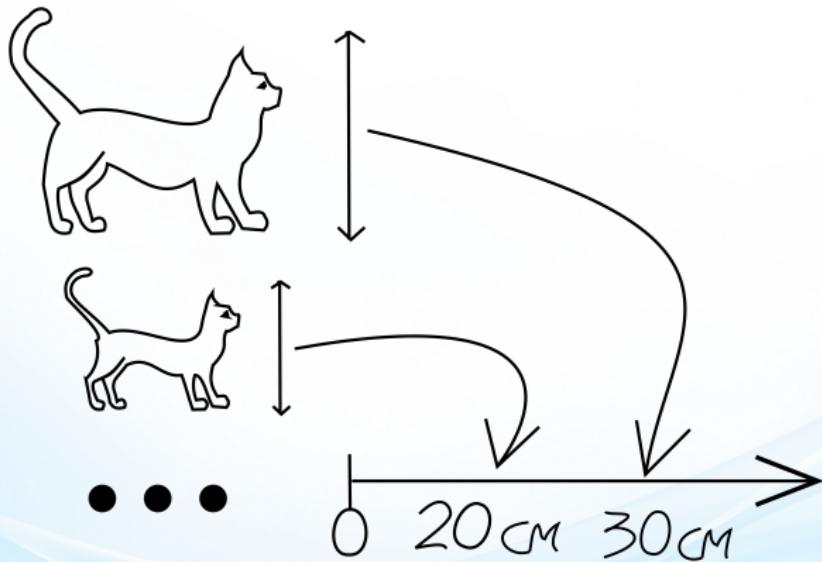


# Обзор статистики

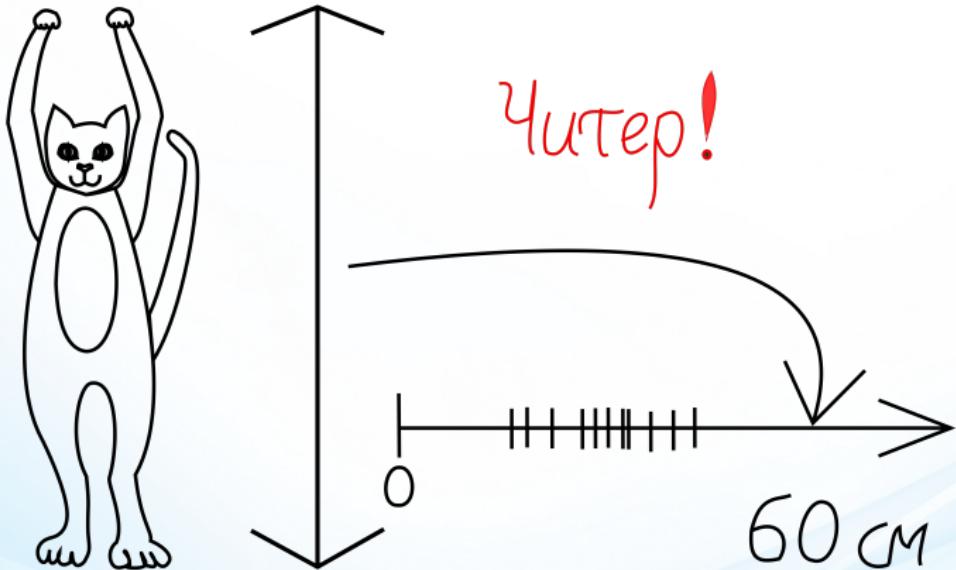
# Мурмурландия



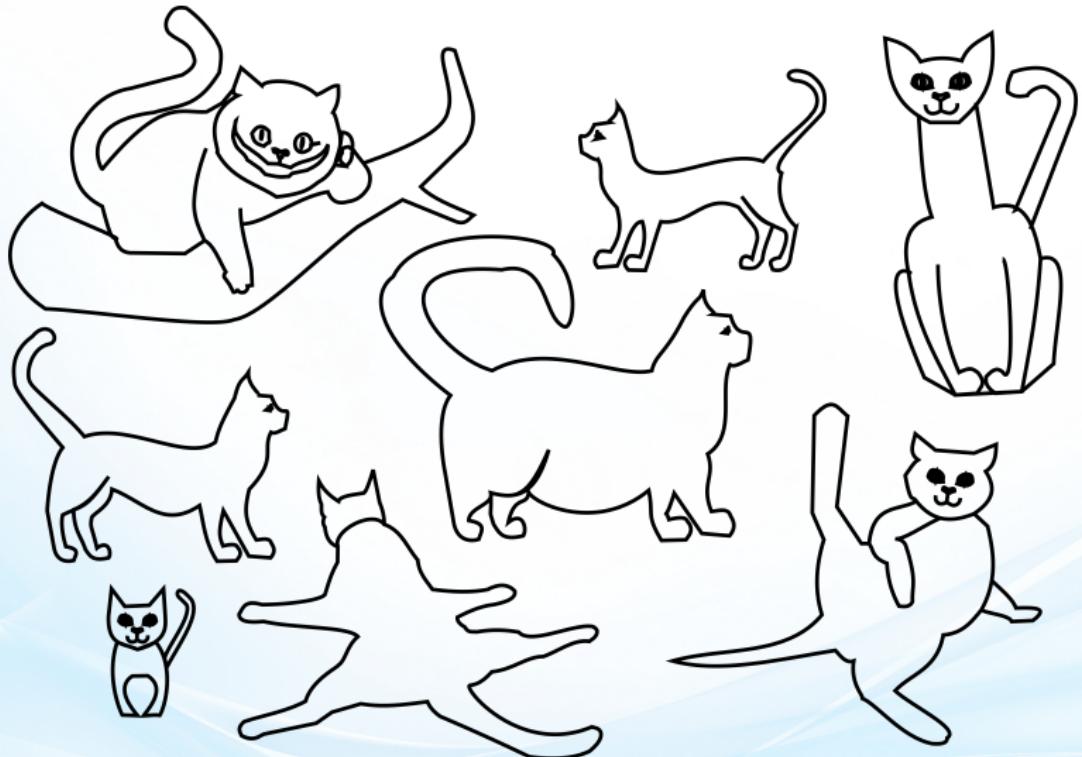
Каков средний рост котиков?



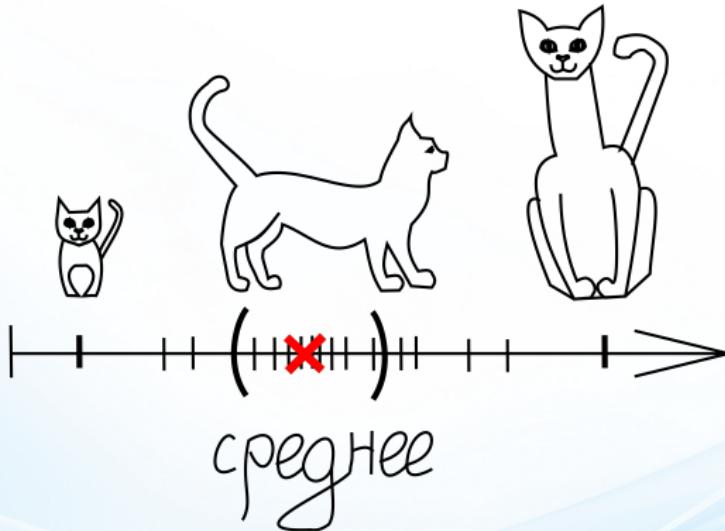
Точечное оценивание



Выбросы

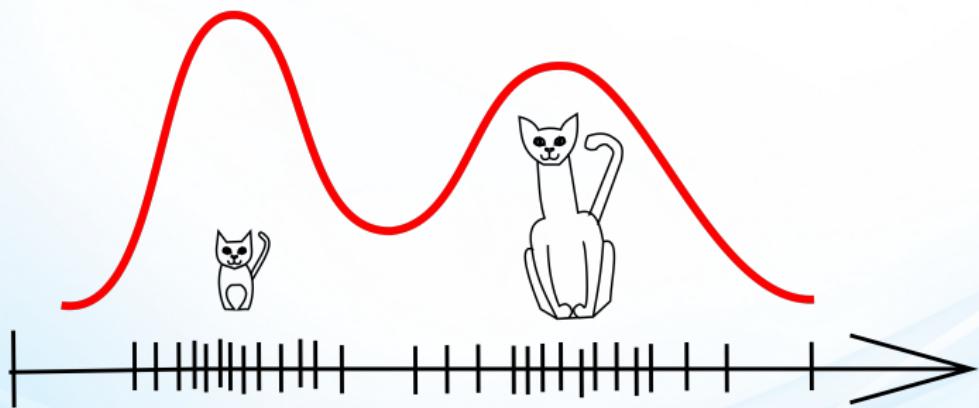


Среднее определяется неточно



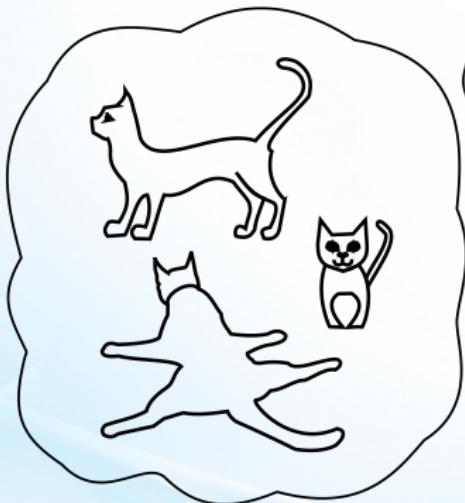
Интервальное оценивание

# Характер распределения



Непараметрическое оценивание

низкие



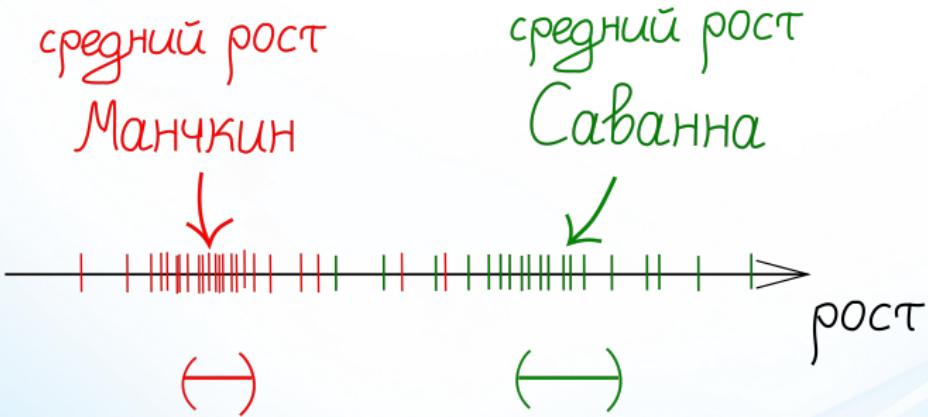
высокие





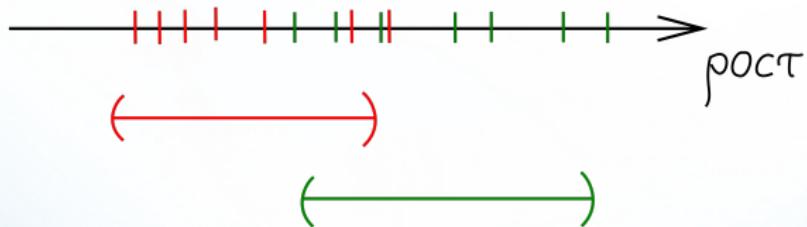
Отличается ли их средний рост?

# Собираем данные



отличается

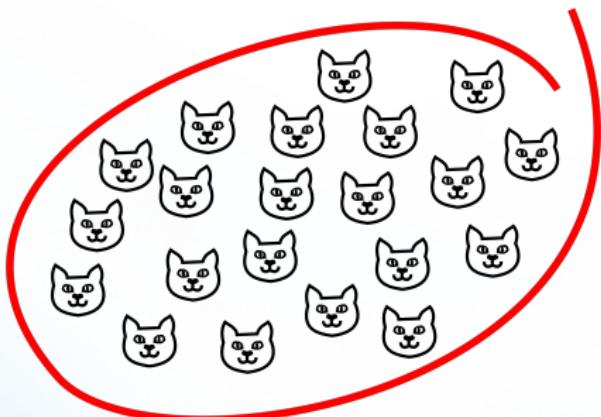
Если данных мало



непонятно

Статистические гипотезы, ANOVA

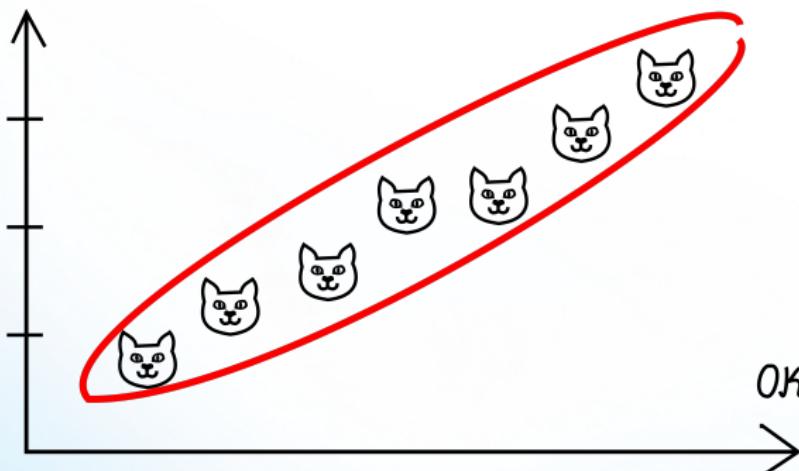
счастье

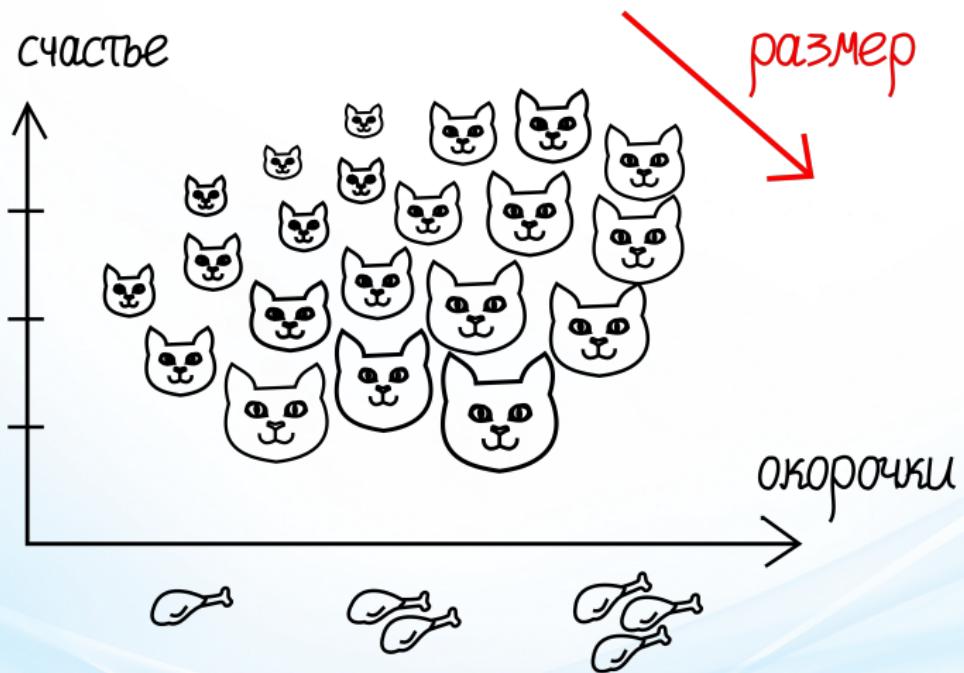


окорочки

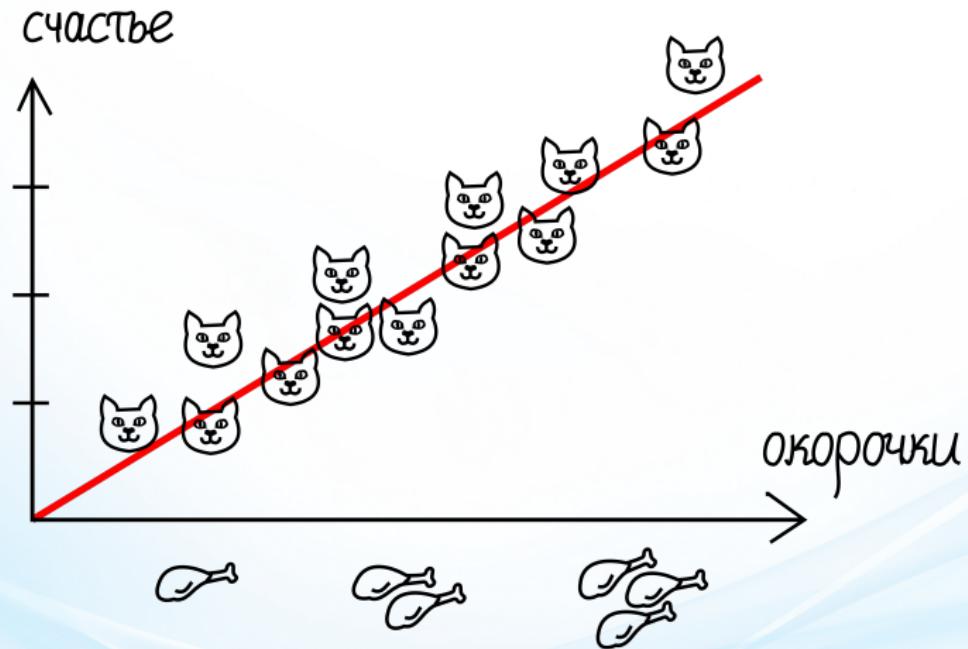


счастье





Корреляционный анализ



# Формула счастья


$$\text{счастье} = \theta_0 + \theta_1 \times \text{кал-бо} + \text{погрешности}$$

# Больше окорочков

счастье



# Формула счастья



$$= \theta_0 + \theta_1 \times \text{кал-бо} - \theta_2 \times (\text{кал-бо})^2$$

+ погрешности



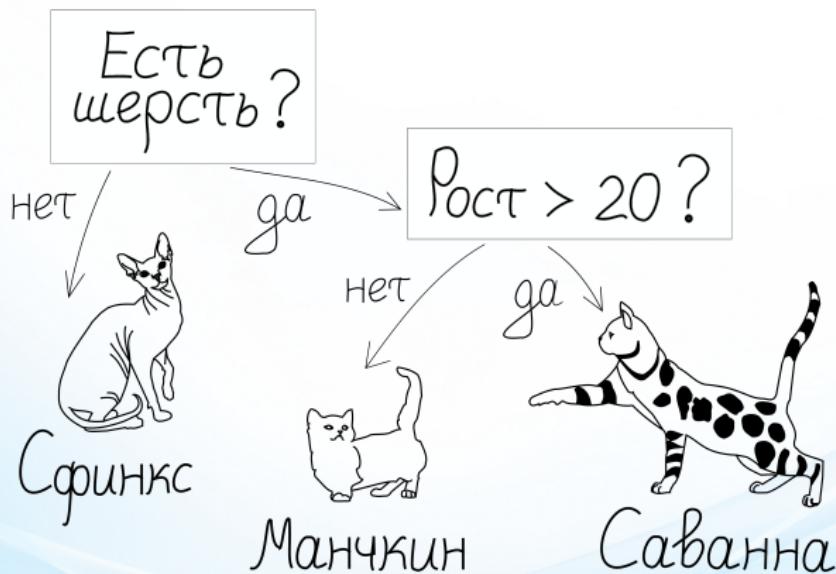
## Другие факторы

$$\begin{aligned} &= \theta_0 + \theta_1 \times \text{чили} - \theta_2 \times (\text{чили})^2 \\ &\quad + \theta_3 \times \text{шарф} \\ &\quad + \theta_4 \times \text{диван} \\ &\quad + \text{погрешности} \end{aligned}$$



Регрессионный анализ

# Классификация котиков



Классификация

# Собираем данные

котик



порода

Саванна

рост

50 см

шерсть

да



Сфинкс

30 см

нет



Манчкин

15 см

да

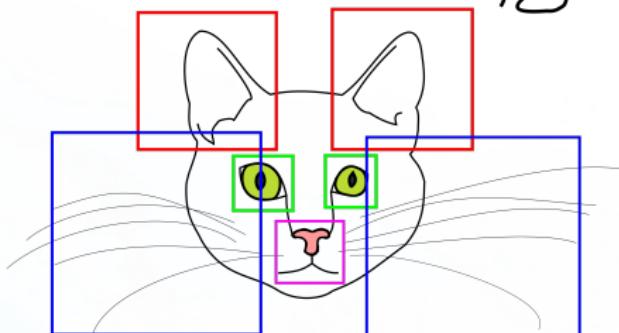


Саванна

40 см

да

# Распознавание мордочек



Нейронные сети

Художник  
курса



Гаврилова

# Примеры реальных задач

## Оценки в экспериментах

Часто в экспериментах измеряется какая-либо величина.

Пример: концентрация вещества, свечение флуоресцентного маркера, энергия частицы.

Есть данные эксперимента — измерения каких-то характеристик.

Нужно сравнить полученные распределения с предсказаниями теории или найти зависимость.

- ▶ Оценить параметры — точечные оценки
- ▶ Оценить погрешности — доверительные интервалы
- ▶ Проверить, насколько хорошо данные согласуются с теорией
- ▶ Восстановить неизвестную зависимость от параметров  
(простейший пример — МНК)

Все эти процедуры основываются на статистике.

# Проверка гипотез

Почти в каждом эксперименте надо проверять различные предположения.

Например, является трек треком нейтрино или другой частицы.

Часто гипотез не одна, а несколько.

Необходимо поддерживать небольшой "суммарный" уровень "ошибки" для того, чтобы можно было доверять выводам ученых.

**Метод решения:** множественная проверка гипотез

**Пример:** изучить роль интерлейкина-10 в подавлении воспаления.

В ходе работы неоднократно проверяются гипотезы различного рода о влиянии IL-10 на различные метаболические процессы.

*Eddie et al. Anti-inflammatory effect of IL-10 mediated by metabolic reprogramming of macrophages,*

<https://science.sciencemag.org/content/356/6337/513>



# Лекарства против респираторного заболевания

- ▶ 2 лекарства;
- ▶ несколько испытуемых.

## Каждый испытуемый

- ▶ вдыхает первое лекарство с помощью ингалятора;
- ▶ проходит упражнение беговой дорожке. Измеряется время достижения максимальной нагрузки;
- ▶ после периода восстановления эксперимент повторяется со вторым лекарством.

Отличается ли время достижения максимальной нагрузки для исследуемых лекарств?

**Метод решения:** дисперсионный анализ.

# Эффективность жаропонижающих средств

- ▶ 4 лекарства, в составе которых один и тот же активный ингридиент присутствует в разных дозировках;
- ▶ 4 группы из 15 морских свинок;
- ▶ В каждой группе фиксируется изменение температуры после введения жаропонижающего.

Есть ли различия в действии препаратов?

**Метод решения:** дисперсионный анализ.

*Bonnini S., Corain L., Marozzi M., Salmaso S. Nonparametric Hypothesis Testing - Rank and Permutation Methods with Applications in R, 2014.*

# Диагностика диабета

Около 500 пациентов, для которых известна информация:

- ▶ Пол, возраст;
- ▶ Индекс массы тела;
- ▶ Среднее кровяное давление;
- ▶ Измерения сыворотки крови;

Нужно предсказать риск развития диабета в течении года.

**Метод решения:** регрессионный анализ

## Медицина — диагностика рака кожи

Нейронную сеть обучили на 130 000 изображений,  
в которых выделено 2000 различных заболеваний кожи.

Точность диагностики 91%, как и у коллектива из 21 дерматолога из Стэнфордской медицинской школы.

<https://geektimes.ru/post/285154/>

## Предсказание генов

Задача: по последовательности генома найти места нахождения потенциальных генов, а также их смысловых частей.

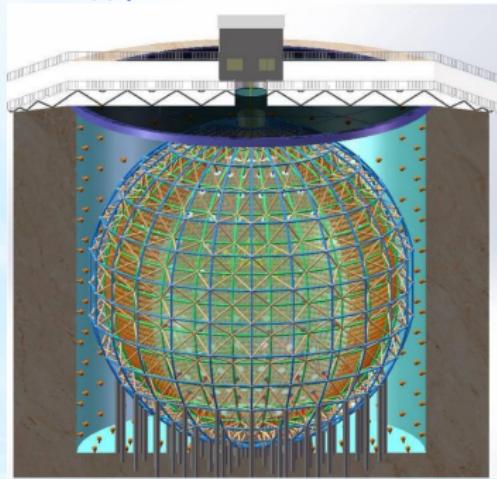
**Решение:** марковские модели, байесовский подход

# Анализ данных в физике частиц

Пример: эксперимент JUNO.

Нужно по данным количества и времена детекции нейтрино на каждом фотодетекторе определить координаты и энергию частицы.

**Метод решения:** машинное обучение.



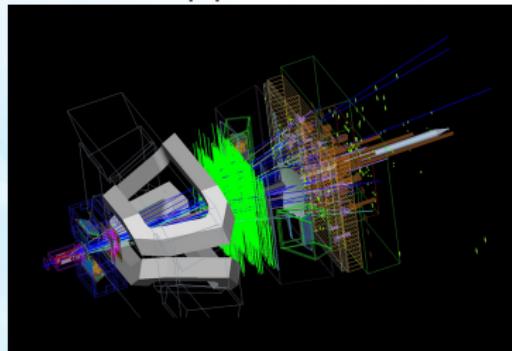
<https://neutel11.wordpress.com/2015/03/04/j-cao-juno-a-multi-purpose-neutrino-observatory/>

# Большой адронный коллайдер

Данные, полученные на детекторе LHCb, проходят через несколько триггеров, которые оставляют только потенциально интересные события.

Классификацию событий можно делать методами машинного обучения.

Используя возможности машинного обучения (CatBoost), ученым удалось повысить эффективность системы триггеров на 40-50 процентов.

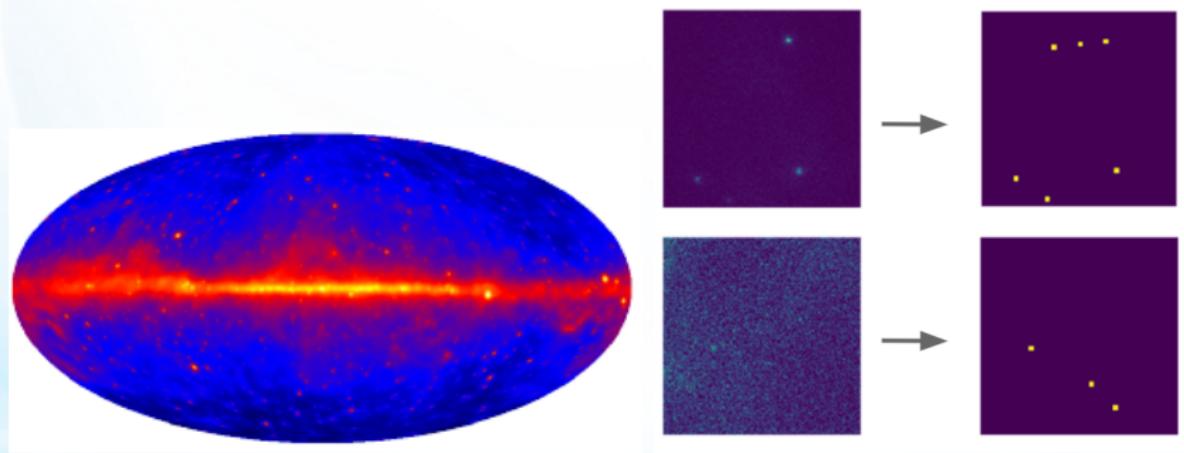


*LHCb collaboration, CERN*

<https://nplus1.ru/material/2017/10/25/cern-yandex>

# Шумоподавление изображений телескопа Ферми

Найти точечные источники на изображениях гамма-телескопа Fermi.



Метод решения: байесовский подход

Создан каталог источников и открыты "гигантские пузыри Ферми".

*The Denoised, Deconvolved, Decomposed Fermi γ-ray sky (Selig et al. 2015)*

Альтернативный подход: нейронные сети

# Аналитическая химия

Для того, чтобы новая методика количественного измерения была принята для использования, необходимо, чтобы погрешности, получаемые с помощью этой методики, были распределены нормально.

**Метод решения:** проверка гипотезы о нормальности выборки

**Сравнение воспроизводимости результатов или методик:**

Чем меньше погрешность, тем точнее методика.

Необходимо сравнить погрешности, получаемые разными методиками.

**Метод решения:** дисперсионный анализ



# Обучение с подкреплением для стимулирования химических реакций

В ходе проведения химических реакций можно во времени менять различные параметры: *температура, добавление катализаторов, и т.д.*

**Цель:** получить искомое вещество

Варианты: получить в-во с наименьшими затратами;  
получить вещество с заданными свойствами. *Например, лекарство.*

Для решения задачи может применяться обучение с подкреплением - раздел машинного обучения, с помощью которого можно "обучать" принятие решений.



# Практическая часть

# Языки для анализа данных



Слабо специализированный,  
проприетарный.



Недостаточно гибкий в современных  
задачах. Тем не менее, незаменим  
для профессионалов.



Достаточно гибкий,  
большое число библиотек.



Слишком молодой.



Jupyter Notebook — инструмент для создания документов-ноутбуков, которые сочетают в себе код и его выход, графики, текстовое описание, формулы.

Возможности:

- ▶ редактирование кода в браузере, с подсветкой синтаксиса, автоотступами и автодополнением;
- ▶ запуск кода в браузере;
- ▶ отображение результатов вычислений с медиа представлением (схемы, графики);
- ▶ работа с языком разметки Markdown и LaTeX.

Поддержка ядер: Python (по умолчанию), R, Julia, C++, Java и др.

Пример

localhost:8888 Пример

jupyter Пример (autosaved)

Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

Run

Пример ноутбука в виде ноутбука

Ячейка с кодом, в которой импортируются библиотеки:

```
In [1]: import numpy as np  
import matplotlib.pyplot as plt  
%matplotlib inline
```

Текстовая ячейка в формате **Markdown**. Например, можно сделать маркированный список

- один
- два
- три

Или написать `latex`-формулу:

$$P(A \cap B) = P(A)P(B)$$

Далее снова код и его вывод:

```
In [2]: print(np.arange(5))
```

```
[0 1 2 3 4]
```

```
In [3]: plt.figure(figsize=(5, 2))  
plt.plot([0, 1], [0, 1]);
```

1.00  
0.75  
0.50  
0.25  
0.00  
0.0 0.2 0.4 0.6 0.8 1.0

Пример

localhost:8888 Пример

jupyter Пример (autosaved)

Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

Run

## Пример ноутбука в виде ноутбука

Ячейка с кодом, в которой импортируются библиотеки:

```
In [1]: import numpy as np  
import matplotlib.pyplot as plt  
%matplotlib inline
```

Текстовая ячейка в формате **\*\*Markdown\*\***. Например, можно сделать \*маркерованный список\*

- \* один
- \* два
- \* три

Или написать `latex`-формулу:  
$$P(A \cap B) = P(A) P(B)$$

Далее снова код и его вывод:

```
In [2]: print(np.arange(5))  
[0 1 2 3 4]
```

```
In [3]: plt.figure(figsize=(5, 2))  
plt.plot([0, 1], [0, 1]);
```

A line plot showing a diagonal line from (0,0) to (1,1).



## Jupyter Lab

The screenshot shows a Jupyter Notebook interface with the following details:

- File Bar:** File, Edit, View, Run, Kernel, Tabs, Settings, Help.
- Left Sidebar:** Files, Running, Commands, Cell Tools, Tabs.
- File List:** Shows notebooks: Data.ipynb, Festa.ipynb, Julia.ipynb, Lorenz.ipynb (selected), R.ipynb, iris.csv, lightning.json, lorenz.py.
- Terminal 1:** In this Notebook we explore the Lorenz system of differential equations:
$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= \rho x - y - xz \\ \dot{z} &= -\beta z + xy\end{aligned}$$
- In [4]:** from lorenz import solve\_lorenz  
t, x\_t = solve\_lorenz(N=10)
- Output View:** Shows sliders for sigma (10.00), beta (2.67), and rho (28.00). Below the sliders is a 3D plot of the Lorenz attractor, which is a complex, fractal-like shape.
- lorenz.py:** The code for generating the attractor. It includes functions for solving the Lorenz equations and plotting them in 3D.

# Установка Jupyter

- ▶ Ubuntu и т.п. (сначала установите пакет pip3, если его нет):

```
sudo pip3 install jupyter
```

Установка Jupyter Lab:

```
sudo pip3 install jupyterlab
```

```
jupyter serverextension enable --py jupyterlab --sys-prefix
```

Запуск: `jupyter notebook` или `jupyter lab`

- ▶ В остальных случаях просто поставьте Anaconda:

<https://www.anaconda.com/download/>

- ▶ Попробовать онлайн:

<https://jupyter.org/try>

- ▶ Не получилось поставить => работа онлайн в Google Colab:

<https://colab.research.google.com/notebooks/intro.ipynb> - введение  
в Python и Colab



# Библиотеки для анализа данных

Гарантированно будем использовать:

- ▶ numpy — быстрые численные и матричные вычисления;
- ▶ scipy (stats) — статистический пакет;
- ▶ matplotlib — построение графиков;
- ▶ pandas — работа с табличными данными;
- ▶ seaborn — высокоуровневое построение графиков;
- ▶ sklearn — методы машинного обучения.

Пример установки на Ubuntu:

```
sudo pip3 install numpy
```

Также полезны:

- ▶ statsmodels — профессиональный стат. пакет;
- ▶ plotly — интерактивные графики;
- ▶ ipywidgets — виджеты.

Для освоения библиотек смотрите видеозаписи от нашей команды!



# Введение

Основная задача математической статистики

# Введение

## **Теория вероятностей**

Зная природу случайного явления,  
посчитать характеристику этого явления.

## **Математическая статистика**

По результатам экспериментальных данных  
высказать суждение о том, какова была природа этого явления.

## Классический пример

На курсе  $N$  студентов; из них  $M$  выбирает спецкурс по анализу данных.

### Задача в теории вероятностей

$P(\text{среди случайных } n \text{ чел. ровно } m \text{ слушателей спецкурса}) - ?$

Предполагается, что  $M$  известно.

### Задача в математической статистике

Среди случайных  $n$  чел. есть  $m$  слушателей спецкурса.

Оценить  $M$ .

Предполагается, что  $M$  не известно.

## Еще пример

$\xi \sim \mathcal{N}(a, \sigma^2)$  — случайная величина

### Задача в теории вероятностей

Известно, что  $a = 2.3, \sigma = 7.1$

$P(\xi \in [0, 1]) - ?$

$E\xi - ?$

### Задача в математической статистике

$x_1, \dots, x_n$  — независимые реализации случайной величины  $\xi$ .

Оценить  $a$  и  $\sigma$ .

*Вспоминаем оценки и погрешности в лабах!*

## Задача математической статистики

Пусть  $x_1, \dots, x_n$  — численные характеристики  $n$ -кратного повторения некоторого явления.

Будем их воспринимать как независимые реализации  $\xi \sim P$ .

**Задача:** по значениям  $x_1, \dots, x_n$  высказать некоторое суждение о распределении  $P$ .

**Решение:** статистический вывод или обучение.

## Основные понятия

Последовательность независимых одинаково распределенных случайных величин  $X_1, \dots, X_n$  называется **выборкой**.

Их значения  $x_1, \dots, x_n$  как числа (на конкретном исходе) называются **реализацией выборки**.

Интуитивно:  $x_i$  - различные "измерения" какой-то величины.

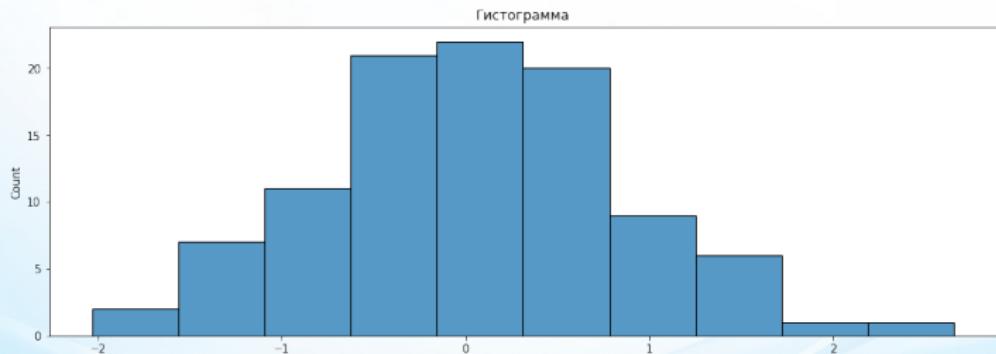
Это имеющиеся у нас данные.

Давайте посмотрим, что вообще можно делать с данными!

# Гистограмма

Пусть у нас есть реализация выборки  $x_1, \dots, x_n \in \mathbb{R}$ .

Идея: разделим всю числовую прямую на несколько "корзин" и посмотрим, сколько объектов (иксов) попало в каждую.

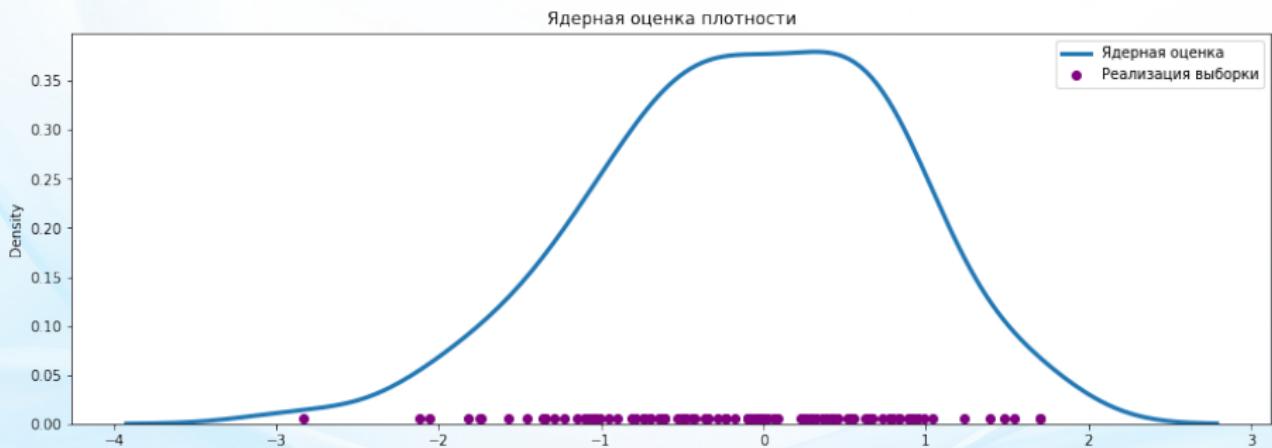


Можно построить график в виде столбиков, где высота столбика показывает, сколько объектов попало в соответствующую корзину.

Этот график по форме похож на график плотности распределения.

# Ядерная оценка плотности

Идея: как-то оценить плотность распределения.

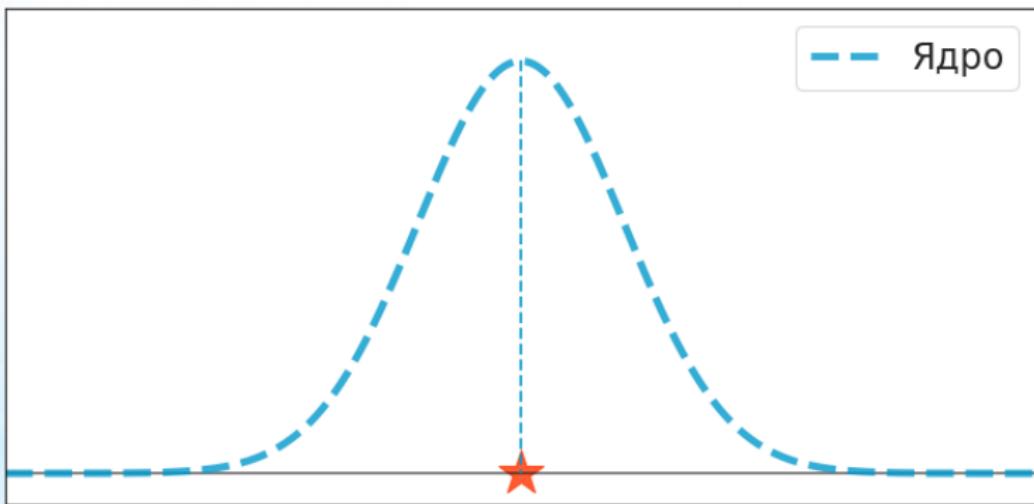


Как? Сейчас узнаем!

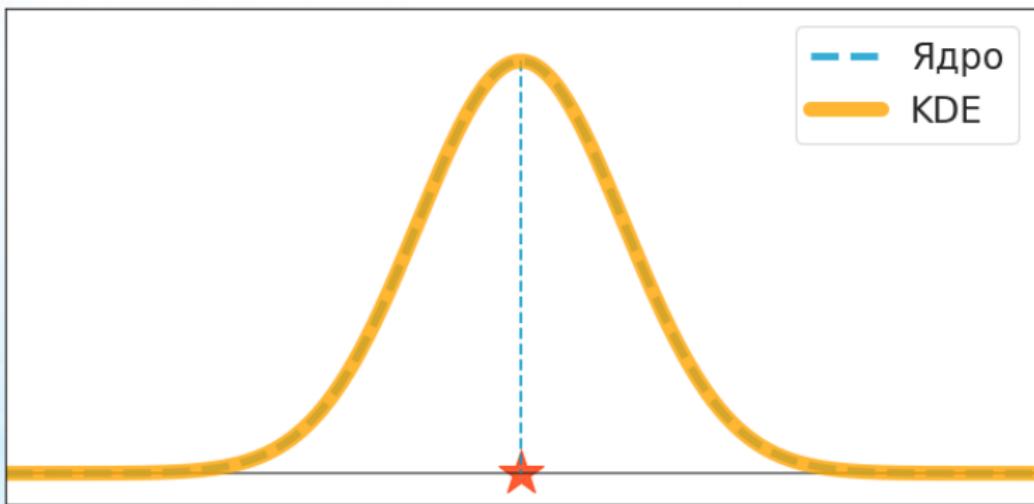
# Ядерная оценка плотности: простые примеры



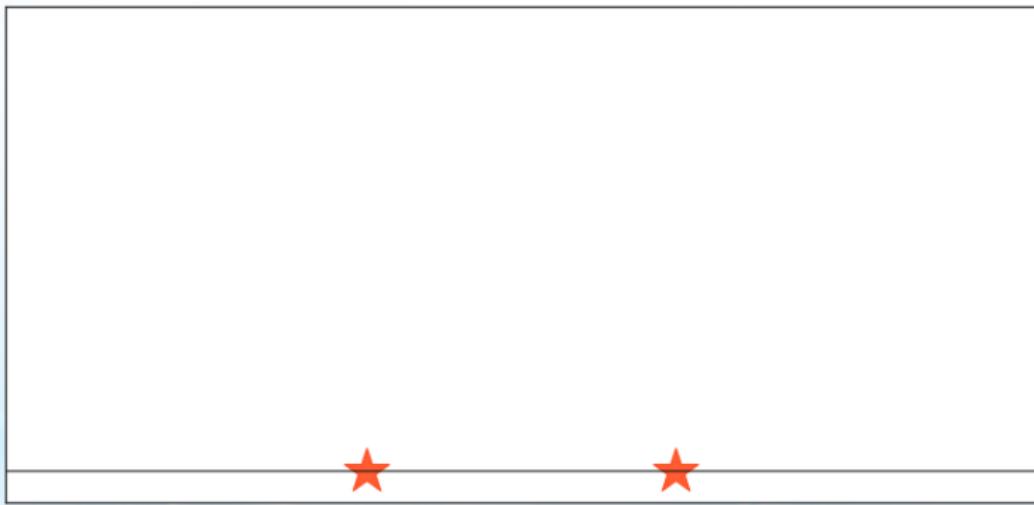
# Ядерная оценка плотности: простые примеры



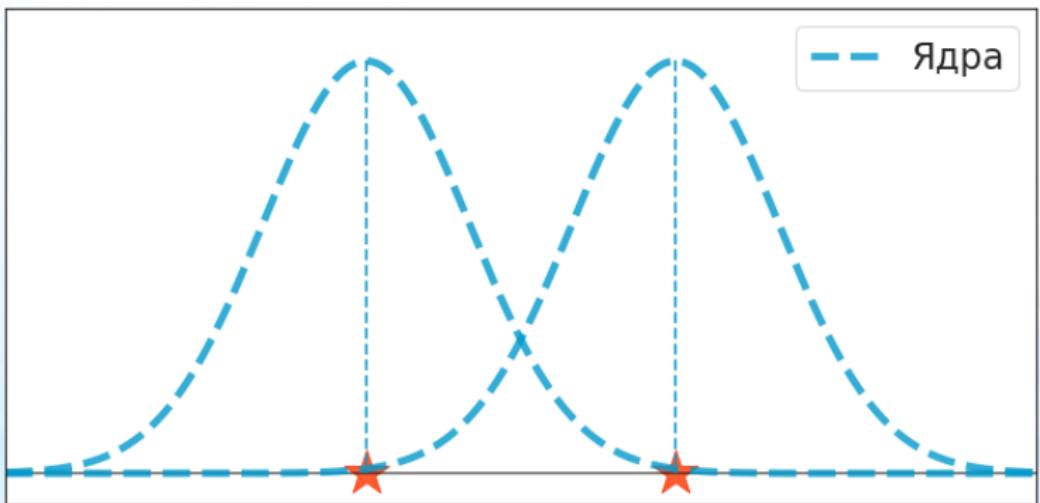
# Ядерная оценка плотности: простые примеры



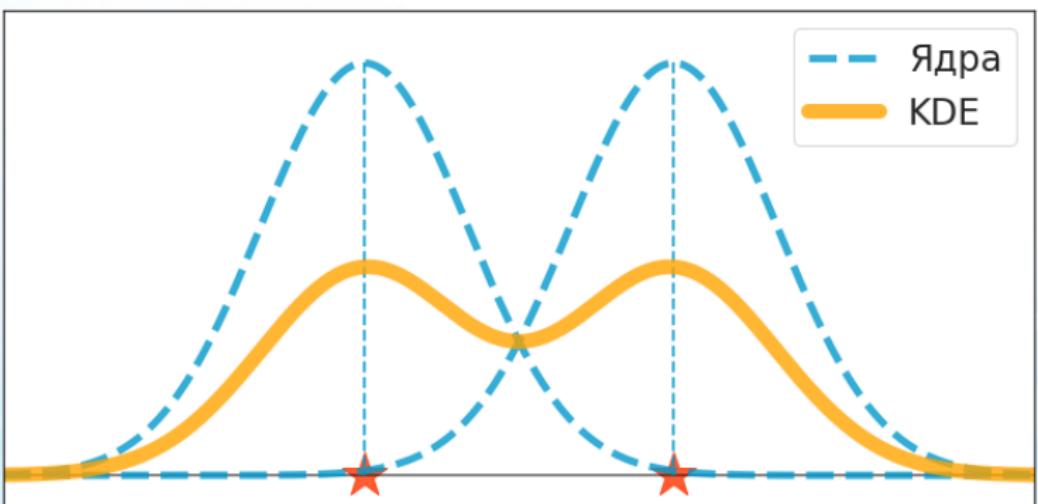
# Ядерная оценка плотности: простые примеры



## Ядерная оценка плотности: простые примеры



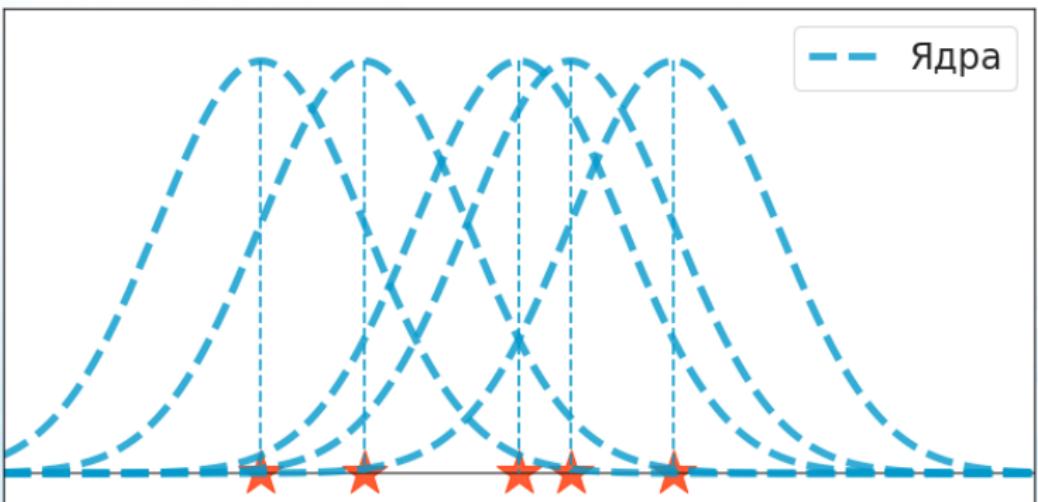
## Ядерная оценка плотности: простые примеры



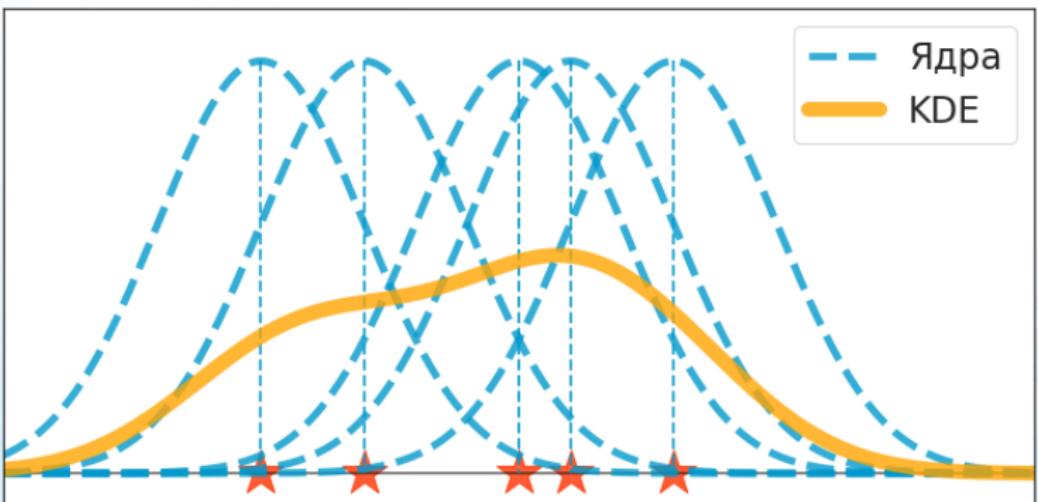
# Ядерная оценка плотности: простые примеры



## Ядерная оценка плотности: простые примеры



## Ядерная оценка плотности: простые примеры



## Определение

Пусть  $X = (X_1, \dots, X_n)$  — выборка из непрерывного распределения.

Выберем

- ▶  $q(x)$  — ядро = некоторая "базовая" симметричная плотность;
- ▶  $h > 0$  — ширина ядра, отвечающая за масштабирование.

Ядерная оценка плотности

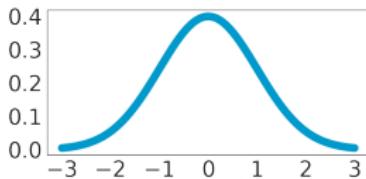
$$\hat{p}_h(x) = \frac{1}{nh} \sum_{i=1}^n q\left(\frac{x - X_i}{h}\right)$$

**Пояснение:** в каждую точку выборки поставили отмасштабированное ядро и усреднили.

# Виды ядер

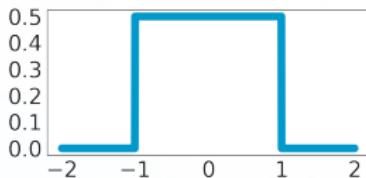
Гауссовское

$$q(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$



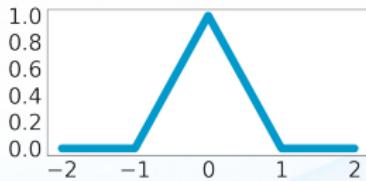
Прямоугольное

$$q(x) = \frac{1}{2} I\{|x| \leqslant 1\}$$



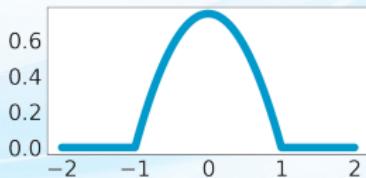
Треугольное

$$q(x) = (1 - |x|)I\{|x| \leqslant 1\}$$



Епанечникова

$$q(x) = \frac{3}{4}(1 - x^2)I\{|x| \leqslant 1\}$$



Давайте посмотрим на все это на практике!



# Статистики и оценки

## Параметрический подход к статистике

Пусть у нас есть реализация  $x_1, \dots, x_n$  выборки  $X_1, \dots, X_n$ .

Мы **предполагаем**, что элементы выборки имеют распределение из множества  $\mathcal{P}$ , где  $\mathcal{P}$  - какое-то заданное множество распределений.

Однако мы не знаем, какое именно, и хотим это как-то оценить по имеющейся реализации.

В **параметрическом** подходе мы также будем предполагать, что на множестве  $\mathcal{P}$  введен **параметр**:  $\mathcal{P} = \{P_\theta \mid \theta \in \Theta\}$ , где  $\Theta$  – множество параметров. Такое семейство  $\mathcal{P}$  называется **параметрическим**.

Задача об оценке истинного распределения сводится к задаче об оценке значения параметра.

## Статистики и оценки

Пусть  $X = (X_1, \dots, X_n)$  — выборка из неизвестного распределения  $P \in \mathcal{P}$ .

Функция от выборки  $S(X)$  называется **статистикой**.

Примеры статистик:

- ▶  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  — **выборочное среднее**;
- ▶  $\bar{X}^k = \frac{1}{n} \sum_{i=1}^n X_i^k$  — **выборочный момент  $k$ -го порядка**;
- ▶  $S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$  — **выборочная дисперсия**.

Формально: Пусть  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, \mathcal{P})$ , — вероятностно-статистическая модель, где  $\mathcal{P} = \{\mathcal{P}_\theta | \theta \in \Theta\}$  — некоторое параметрическое семейство распределений, а  $\Theta$  — множество параметров.

Если пара  $(E, \mathcal{E})$  — некоторое измеримое пространство, а  $S : \mathcal{X}^n \rightarrow E$  — измеримое отображение, то  $S$  называется **статистикой**.

# Статистики и оценки

Пусть  $X = (X_1, \dots, X_n)$  — выборка из неизвестного распределения  $P \in \mathcal{P}$ .

Функция от выборки  $S(X)$  называется **статистикой**.

Если при этом  $S$  принимает значения во множестве  $\Theta$ , то  $S(X)$  называется **точечной оценкой** или просто **оценкой** параметра  $\theta$ .

Обычно обозначается  $\hat{\theta}(X)$  или просто  $\hat{\theta}$ .

**Интуиция:** пытаемся догадаться об истинном значении параметра на основании данных.

Иногда требуется оценить значение некоторой функции  $\tau(\theta)$ . В таком случае рассматривается  $E = \tau(\theta)$ , то есть образ множества параметров.



**BGE!**