



# Введение в АД





О курсе





# Команда Физтех.Статистики





## Проводимые нами учебные курсы

- ▶ Введение в анализ данных;
- ▶ DS-поток;
- ▶ Phystech@DataScience (в т.ч. статистика на ЛФИ);
- ▶ Мат. статистика на ФБМФ и ФМХФ (семинары);
- ▶ Прикладная статистика на кафедрах в магистратуре;
- ▶ Машинное обучение в ШАД (частично);
- ▶ АБ-тестирование в ШАД.



## Организационная информация

Telegram-бот

@miptstats\_ds24\_bot



Код регистрации **NUMpy@2024!**

Сайт команды [miptstats.github.io](https://miptstats.github.io)

Почта [mipt.stats@yandex.ru](mailto:mipt.stats@yandex.ru)



## Цели курса

1. **Дать представление об анализе данных;**
2. Обучить базовым инструментам анализа данных;
3. Рассказать о практическом смысле объектов теории вероятностей;
4. Помочь определиться с кафедрой.

# Яндекс

В курсе предполагается участие кафедры анализа данных:

- ▶ Гостевая лекция от Яндекса;
- ▶ Работа по курсу учитывается при отборе на кафедру.



## План занятий (О — обязательная, Ф — факультативная)

Дата		Тема	Лектор	Дедлайн по ДЗ
03.02	О	<i>Введение</i>	 Никита Волков	16.02
10.02	О	<i>Библиотеки Питона. Онлайн, посещ. необ.</i>	 Алексей Горбулев	
17.02	О	<i>Линейная регрессия. Общий ML-пайплайн.</i>	 Никита Волков	23.02
24.02		<i>Занятия нет, делаем домашки</i>		01.03
02.03	О	<i>Введение в нейросети</i>	 Лидия Троешестова	8.03
09.03	Ф	<i>Компьютерное зрение и генеративные модели</i>	 Лидия Троешестова	15.03



## План занятий (О — обязательная, Ф — факультативная)

Дата		Тема	Лектор	Дедлайн по ДЗ
16.03	Ф	<i>Обработка текстов</i>	 Кирилл Овчаренко	29.03
23.03	О	<i>Гостевая лекция от Яндекса</i>		
30.03	О	<i>Теория вероятностей на практике</i>	 Алексей Горбулев	05.04
06.04	Ф	<i>Байесовские классификаторы</i>	 Никита Волков	12.04
13.04	Ф	<i>Кластеризация и пониж. размерности</i>	 Никита Волков	19.04
20.04	Ф	<i>Доп. лекция</i>		26.04



# Система оценивания обязательной части

**Официальное название:** Введение в анализ данных

**Правила:**

1. если  $L > 25\%$  и ( $B < 50\%$  или  $C < 20\%$ ),  
то  $O = 3 + 4 * L$ ;
2. если  $L > 25\%$  и  $B \geq 50\%$  и  $C \geq 20\%$ ,  
то  $O = \max(3 + 4 * L, 3 + B + 7 * LC)$ .
3. если  $L \leq 25\%$ ,  
то  $O = 1 + 3 * T$ .

**Обозначения:**

- ▶  $L$  — доля выполнения легких заданий (их немного);
- ▶  $C$  — доля выполнения сложных заданий (их много);
- ▶  $T$  — доля выполнения тестов;
- ▶  $LC$  — доля выполнения всех заданий (кроме тестов);
- ▶  $B$  — доля правильных ответов на вопросы в боте на занятии;
- ▶  $O$  — итоговая оценка, округляется вверх.



# Система оценивания обязательной части

**Официальное название:** Введение в анализ данных

**Правила:**

1. если  $L > 25\%$  и ( $B < 50\%$  или  $C < 20\%$ ),  
то  $O = 3 + 4 * L$ ;
2. если  $L > 25\%$  и  $B \geq 50\%$  и  $C \geq 20\%$ ,  
то  $O = \max(3 + 4 * L, 3 + B + 7 * LC)$ .
3. если  $L \leq 25\%$ ,  
то  $O = 1 + 3 * T$ .

**Следствия:** для получения оценки

- ▶ уд(3) достаточно выполнить 33.4% тестов.
- ▶ хор(5) достаточно выполнить 25.1% легких заданий.
- ▶ отл(8) достаточно выполнить 50.1% всех заданий и ответить на 50% вопросов.

**Если списать:** все участники списывания сдают устный зачет.





# Система оценивания факультативной части

**Официальное название:** Введение в анализ данных: доп. главы

**Правила:**

1.  $O = 9*Ф + 2*В$

**Обозначения:**

- ▶ Ф — доля выполнения домашних заданий;
- ▶ В — доля правильных ответов на вопросы в боте на занятии;
- ▶ О — итоговая оценка, округляется вверх.

**Как записаться?**

Просто ходить на занятия и сдавать домашние задания.

В конце семестра разошлем форму, в которой вы отметитесь, что хотите поставить оценку в ведомость.



## Для кого наш курс?

### **Обязательная часть курса:**

- ▶ ПМИ, гр. 251 — курс обязателен.
- ▶ ПМФ, ИВТ, гр. 252, 253 — курс факультативен.
- ▶ Студенты других физтех-школ также могут сдавать курс.

*Для записи достаточно зарегистрироваться в боте.*

*Для проставления оценки нужно в конце семестра взять ведомость.*

### **Факультативная часть курса:**

- ▶ Факультативно для всех.

В силу ограниченных возможностей проверяющих при большом количестве желающих возможность сдавать курс может быть ограничена.



## Правила комфорта

- ▶ Постарайтесь задавать вопросы на занятии в тот момент, когда это актуально, не перебивая на полуслове.  
Другой вопрос лучше задать в перерыве или после занятия.
- ▶ Цените труд проверяющих :)  
В каком из случаев проверяющему больше захочется пойти навстречу автору вопроса?
  - ▶ *"Объясните вашу претензию, почему вы мне сняли баллы, я же все сделал, я не согласен"*
  - ▶ *"Добрый день! По такой-то задаче вы написали ..., но я считаю ..., потому что ..., и у меня в работе написано ..."*



## Особенности проверки домашних заданий

Иногда к нам приходят отзывы:

*"Проверяющие проверяют рандомно,  
за одну и ту же ошибку разным студентам  
сняли разное количество баллов."*

Пусть на курсе 300 человек,  
каждый проверяющий проверяет 30 работ.

*Сколько результатов проверок нужно ему посмотреть для сравнения?*

*Сколько всего таких сравнений?*

*Оцените количество времени, которое нужно затратить.*



# Особенности проверки домашних заданий

## Как мы решаем проблему?

- ▶ Общие критерии для всех проверяющих в табличке, с общим текстом и количеством баллов.  
Проверяющему достаточно поставить галочку.
- ▶ Сравнение частоты применения критерия между проверяющими с помощью статистического t-test'a.
- ▶ Сравнение среднего балла между проверяющими с помощью статистического t-test'a.
- ▶ Студенты 4 курса DS-потока разрабатывают ML-модель, которая ищет похожие комментарии проверяющих.

Мы стараемся, но мы не волшебники :)

Если вы заметили несправедливость проверки, пожалуйста, напишите нам. Мы посмотрим и при необходимости поправим. И учтем это для совершенствования наших методов проверки.



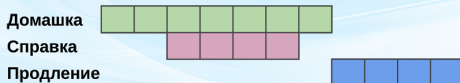
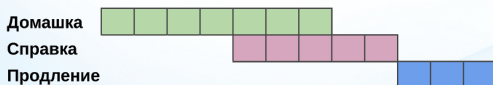
# Переносы дедлайнов по уважительным причинам

## Уважительные причины

- ▶ Медицинская справка с подписью и печатью.
- ▶ Приказ по институту об освобождении.

## На сколько можно перенести

На количество дней пересечения интервала выполнения задания и датам по справке от даты дедлайна или окончания справки.





## Тех. ассистент



Алина Крищук

tg: @angeliyya

- ▶ Организация проверки домашних заданий
- ▶ Перенос дедлайнов по уважительным причинам
- ▶ Разные технические вопросы



# DS-поток

Программа 3-4 курсов  
*Продвинутый анализ данных*





## DS-поток

Семестр	DS-поток	Основной поток ПМИ
<b>5</b>	Математическая статистика	Математическая статистика
	Машинное обучение	Машинное обучение
	<i>Практика</i>	<i>Практика по мат. статистике</i>
	Основы прикладной статистики	Курс по выбору x 2
	Курс по выбору	
<b>6</b>	Дискр. случ. процессы и временные ряды	Случайные процессы
	Глубокое обучение и его приложения	Вычислительная математика
	Прикладная статистика и анализ данных	Параллельные и распределенные вычисления
	<i>Практика</i>	Курс по выбору
<b>Кафедра АД</b>	Курс ШАД	Методы прикладной статистики
<b>7</b>	Байесовский подход в анализе данных	Курс по выбору x 2
	<i>Практика</i>	
<b>8</b>	Прикладные задачи машинного обучения	Курс по выбору x 2
	<i>Практика</i>	

### Примечания.

1. По одним и тем же курсам лекторы разные.
2. По факту практика не является отдельным курсом.
3. В процессе обучения перейти в DS-поток невозможно.



# DS-поток

Семестр	DS-поток	Основной поток ИВТ
<b>5</b>	Математическая статистика	Математическая статистика
	Машинное обучение	Машинное обучение
	<i>Практика</i>	<i>Практика по мат. статистике</i>
	Основы прикладной статистики	Курс по выбору x 2
	Курс по выбору	
<b>6</b>	Дискр. случ. процессы и временные ряды	Машинное обучение, ч. 2
	Глубокое обучение и его приложения	Вычислительная математика
	Прикладная статистика и анализ данных	Современные компьютерные сети
	<i>Практика</i>	Курс по выбору
<b>Кафедра АД</b>	Курс ШАД	Методы прикладной статистики
<b>7</b>	Байесовский подход в анализе данных	Курс по выбору x 2
	<i>Практика</i>	
<b>8</b>	Прикладные задачи машинного обучения	Курс по выбору x 2
	<i>Практика</i>	

## Примечания.

1. По одним и тем же курсам лекторы разные.
2. По факту практика не является отдельным курсов.
3. В процессе обучения перейти в DS-поток невозможно.



## DS-поток

Семестр	DS-поток	Основной поток ПМФ
<b>5</b>	Математическая статистика	Математическая статистика
	Машинное обучение	Вычислительная математика
	<i>Практика</i>	Курс по выбору
	Основы прикладной статистики	
<b>6</b>	Дискр. случ. процессы и временные ряды	Случайные процессы
	Глубокое обучение и его приложения	Квантовая механика, ч. 2
	Прикладная статистика и анализ данных	Курс по выбору x 2
	<i>Практика</i>	
<b>Кафедра АД</b>	Курс ШАД	Методы прикладной статистики
<b>7</b>	Байесовский подход в анализе данных	Машинное обучение
	<i>Практика</i>	Курс по выбору x 2
<b>8</b>	Прикладные задачи машинного обучения	Курс по выбору x 2
	<i>Практика</i>	

### Примечания.

1. По одним и тем же курсам лекторы разные.
2. По факту практика не является отдельным курсов.
3. В процессе обучения перейти в DS-поток невозможно.



# DS-поток

## Чему мы учим

1. В меру глубокое математическое понимание статистики и машинного обучения.
2. Применение математических моделей на реальных данных, в том числе на реальных задачах.
3. Умение составлять полноценные выводы.

## Реальная практика на DS-потоке

1. Реальные примеры из практики;
2. Соревнования на Kaggle;
3. Гостевые лекторы, применяющие анализ данных на практике;
4. Разбор статей на тему анализа данных;
5. Бонусы за участие в хакатонах, соревнованиях по анализу данных и прочую активность;



# Отбор на DS-поток

## Что будет учитываться?

### 1. Необходимое условие:

оценка не менее отл(8) по обяз. части курса

"Введение в анализ данных" и не менее хор(7) по факульт. части.

2. Работа в семестре по курсу, грамотное оформление ДЗ.

3. Оценка по теории вер., в меньшей степени — другие предметы.

## Что нужно делать для отбора?

1. Трудиться в течение семестра.

2. В июне подать заявку.

3. Ждать. Результаты летом.

DS-поток адаптирован для ПМИ. Студенты других напр. могут попасть на DS-поток, но возможна доп. нагрузка и проблемы с расписанием.



## Другие образовательные направления по АД

### **Школа анализа данных Яндекса (ШАД)**

Двухгодичная программа дополнительного образования, специализирующаяся на анализе данных, разработке ML-моделей, создании систем хранения и обработки больших данных и др.. Независима по отношению к обучению в МФТИ.  
Подробнее: [shad.yandex.ru](http://shad.yandex.ru)

### **Кафедра анализа данных (Яндекс) в рамках ФПМИ**

- ▶ По 2-3 курса в семестр в основном из курсов ШАДа.
- ▶ Написание и защита диплома.

### **Другие кафедры в рамках ФПМИ**

- ▶ По 2-3 курса в семестр.
- ▶ Написание и защита диплома.



## Куда пойти?

2 курс

Введение в анализ данных

3-4 курсы

DS-поток

Кафедра  
анализа данных

Школа анализа  
данных (ШАД)

1. Если анализ данных интересен, то хорошее решение:  
**DS-поток + кафедра анализа данных.**
2. Если выбираете кафедру анализа данных,  
то кафедра рекомендует пойти на DS-поток.
3. **По результатам нашего курса кафедра анализа данных может зачесть тех. собеседование при отборе на кафедру.**



# Бонусы при отборе на кафедру АД

## Этапы отбора:

1. Контест (математика, алгоритмы)
2. Техническое собеседование (математика, алгоритмы)
3. Мотивационное собеседование

**Техническое собеседование** засчитывается автоматом, если

1. Оценки ОТЛ за обе части курса Введение в АД
2. Среднее оценок по Введение в АД и ср. балла не менее 4 из 5

**Если пройти на DS-поток:**

1. Если оценки ОТЛ или ХОР за Введение в АД и по DS-поток, то *техническое собеседование засчитывается*
2. Если попал в топ-5 в рейтинге до ДЗ на DS-потоке, то *контест и техническое собеседование засчитываются*





Что такое анализ данных?



## Кому нужен анализ данных

1. *"Я математик, практическое применение не интересует".*

Скорее всего АД не нужен,  
но часто математики им начинают интересоваться.

2. *"Я математик, но хочу применять свои знания на практике".*

**АД для вас, ждем вас на DS-потоке**

3. *"Я программист, и хочу писать только код".*

Скорее всего АД в подробностях не нужен,  
но стоит понимать, чем занимаются коллеги-аналитики.

4. *"Я программист, но хочу глубоко разбираться в тонкостях математических методов".*

**АД для вас, ждем вас на DS-потоке**



Так что, анализ данных  
это математика  
или программирование?

Давайте разбираться...



# Посмотрим на лекции по статистике

Статистика, прикладной поток 12. Контроль FWER и FDR. Критерии согласия

3. Метод Ньютона

Выходные параметры с  $\alpha = \frac{\alpha}{m}$  (Bonferroni)

Точнее. Если  $\forall i: P(X_i) \sim U(0, 1)$ , тогда  $\text{FWER}(\beta, \alpha) \approx \alpha$ .

2. Если  $n$  не очень велико, то  $\text{FWER}(\beta, \alpha) \approx \alpha$ .

1. Если  $n$  очень велико, то  $\text{FWER}(\beta, \alpha) \approx \alpha$ .

$\text{FWER}(\beta, \alpha) = P(\bigcup_{i=1}^m \{P_i < \alpha_i\}) \leq \sum_{i=1}^m P(P_i < \alpha_i) = m \cdot \alpha = \alpha$ .

$\text{FWER}(\beta, \alpha) \approx \alpha$ .

Среднеквадратичная ошибка

$\hat{P}_0$  - правая

$\hat{P}_0$  - левая

Статистика, прикладной поток 11. Учебные Практическая значимость. Множественная проверка

### Критерии (напоминание)

Часто критерий имеет вид  $S = \{T(x) \geq c_\alpha\}$ , где  $T(X)$  — статистика критерия.

$\alpha$  выбирается ДО эксперимента,  $c_\alpha$  вычисляется из условия  $P_0(T(X) > c_\alpha) \leq \alpha$ .

$S = \{T(x) > c_\alpha\}$      $S = \{T(x) < c_\alpha\}$      $S = \{|T(x)| > c_\alpha\}$

$Pr(t)$

$c_\alpha$

$c_\alpha$

$c_\alpha$      $c_\alpha$

Статистика, прикладной поток 5. Бутстрап. Ядерные оценки плотности. Проверка статистической гипотезы

### Метод бутстрапа

**Этап 2.**

Процедуру генерации выборки повторить  $B$  раз:  $X_i^* = (X_{i1}^*, \dots, X_{in}^*)$ , где  $1 \leq i \leq B$ .

Далее по каждой выборке посчитаем значение статистики  $T$ , получив выборку значений  $T_i^* = T(X_i^*), \dots, T_B^* = T(X_B^*)$ .

**Этап 3.**

Полученную выборку использовать для аппроксимации значения оценки, которая называется бутстрапной оценкой.

Например, бутстрапная оценка дисперсии имеет вид

$$\hat{\sigma}_{boot}^2 = \frac{1}{B} \sum_{b=1}^B T_b^{*2} - \left( \frac{1}{B} \sum_{b=1}^B T_b^* \right)^2$$

Статистика, прикладной поток 15. Оптимальные оценки. Эквивариантные оценки. Доказательства теорем

Тогда  $\hat{F}_n(x) - F(x) \leq \hat{F}_n(U_{(1/n)} - \epsilon) - F(U_{(1/n)} - \epsilon) =$

$\hat{F}_n(U_{(1/n)} - \epsilon) - F(U_{(1/n)} - \epsilon) + \frac{1}{N} \sum_{i=1}^N \frac{F(U_{(i/n)} - \epsilon) - F(U_{(i-1/n)} - \epsilon)}{U_{(i/n)} - U_{(i-1/n)}} \leq$

$\hat{F}_n(U_{(1/n)} - \epsilon) - F(U_{(1/n)} - \epsilon) + \frac{1}{N} \sum_{i=1}^N \frac{F(U_{(i/n)} - \epsilon) - F(U_{(i-1/n)} - \epsilon)}{U_{(i/n)} - U_{(i-1/n)}} \leq$

Аналогично

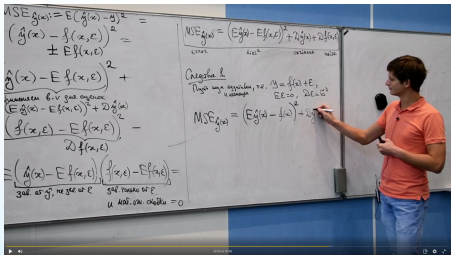
$\hat{F}_n(x) - F(x) \geq \hat{F}_n(U_{(1/n)}) - F(U_{(1/n)}) - \frac{1}{N} \sum_{i=1}^N \frac{F(U_{(i/n)}) - F(U_{(i-1/n)})}{U_{(i/n)} - U_{(i-1/n)}} \geq$

Плюс  $x$  — промежуточные

$|\hat{F}_n(x) - F(x)|$



# Посмотрим на лекции по машинному обучению



## Двуслойная нейронная сеть

### Матричное представление

Пусть  $x = (x_1, x_2, \dots, x_d)^T$  — элемент выборки,  
 $u = (u_1, u_2, \dots, u_H)^T$  — выход I слоя,  $y = (y_1, y_2, \dots, y_M)^T$  — выход II слоя,  
 $W_1 = (w_{1j})_{jH}$  — м-ца весов I слоя,  $W_2 = (w_{2m})_{mM}$  — м-ца весов II слоя,  
 $b_1 = (b_{1j})_{jH}$  — в-р сдвигов I слоя,  $b_2 = (b_{2m})_{mM}$  — в-р сдвигов II слоя.

Тогда работу двуслойной нейронной сети можно представить как:

- $u = \sigma_1(x^T W_1 + b_1)$
- $y = \sigma_2(u^T W_2 + b_2) = \sigma_2(\sigma_1(x^T W_1 + b_1)^T W_2 + b_2)$



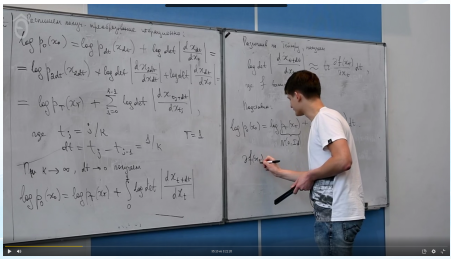
## Classifier guidance

Чтобы явно включить информацию о классе в процесс диффузии, можно обучить классификатор  $f_\theta(y|x, t)$  на зашумленном изображении  $x_t$  и использовать градиенты  $\nabla_x \log f_\theta(y|x_t)$ , чтобы направлять процесс диффузионного семплирования к информации об условии  $y$  путем изменения прогнозирования шума. Можно вывести, что  $\nabla_x \log q(x_t) = -\frac{1}{\sqrt{1-\alpha_t}} \epsilon_\theta(x_t, t)$ .

$$\nabla_x \log q(x_t, y) = \nabla_x \log q(x_t) + \nabla_x \log q(y|x_t) =$$

$$-\frac{1}{\sqrt{1-\alpha_t}} \epsilon_\theta(x_t, t) + \nabla_x \log f_\theta(y|x_t) =$$

$$-\frac{1}{\sqrt{1-\alpha_t}} (\epsilon_\theta(x_t, t) - \sqrt{1-\alpha_t} \nabla_x \log f_\theta(y|x_t))$$





# Посмотрим на научные статьи

## McInnes, L, Healy, J, *UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction*, ArXiv e-prints 1802.03426, 2018

The choice of non-expansive maps in Definition 2 is due to Spivak note that it closely mirrors the work of Carlsson and Memoli in [11] logical methods for clustering as applied to finite metric spaces. The significant since pure isometries are too strict and do not provide la Hom-sets.

In [13] Spivak constructs a pair of adjoint functors,  $\mathbf{Real}$  and  $\mathbf{Sing}$  the categories  $\mathbf{sFuzz}$  and  $\mathbf{EPMet}$ . These functors are the natural e the classical realization and singular set functors from algebraic top functor  $\mathbf{Real}$  is defined in terms of standard fuzzy simplices  $\Delta_{\leq a}^n$  as

$$\mathbf{Real}(\Delta_{\leq a}^n) \triangleq \left\{ (t_0, \dots, t_n) \in \mathbb{R}^{n+1} \mid \sum_{i=0}^n t_i = -\log(a), t_i \geq 0 \right\}$$

similarly to the classical realization functor  $|\cdot|$ . The metric on  $\mathbf{Real}$  simply inherited from  $\mathbb{R}^{n+1}$ . A morphism  $\Delta_{\leq a}^n \rightarrow \Delta_{\leq b}^n$  exists only if is determined by a  $\Delta$  morphism  $\sigma : [n] \rightarrow [m]$ . The action of  $\mathbf{Real}$  morphism is given by the map

$$(x_0, x_1, \dots, x_n) \mapsto \frac{\log(b)}{\log(a)} \left( \sum_{i \in \sigma^{-1}(0)} x_{i_0}, \sum_{i \in \sigma^{-1}(1)} x_{i_1}, \dots, \sum_{i \in \sigma^{-1}(m)} x_{i_m} \right)$$

Such a map is clearly non-expansive since  $0 \leq a \leq b \leq 1$  implies that  $\log(b)/\log(a) \leq 1$ .

We then extend this to a general simplicial set  $X$  via colimits, defining

$$\mathbf{Real}(X) \triangleq \mathop{\mathrm{colim}}_{\Delta_{\leq a}^n \rightarrow X} \mathbf{Real}(\Delta_{\leq a}^n)$$

Since the functor  $\mathbf{Real}$  preserves colimits, it follows that there exists a right adjoint functor. Again, analogously to the classical case, we find the right adjoint denoted  $\mathbf{Sing}$ , is defined for an extended pseudo metric space  $Y$  in terms of its action on the category  $\Delta \times I$ :

$$\mathbf{Sing}(Y) : ([n], [0, a]) \mapsto \mathop{\mathrm{hom}}_{\mathbf{EPMet}}(\mathbf{Real}(\Delta_{\leq a}^n), Y)$$

For our case we are only interested in finite metric spaces. To correspond with this we consider the subcategory of bounded fuzzy simplicial sets  $\mathbf{Fin-sFuzz}$ . We therefore use the analogous adjoint pair  $\mathbf{FinReal}$  and  $\mathbf{FinSing}$ . Formally we define the finite fuzzy realization functor as follows:

```

Algorithm 2 Constructing a local fuzzy simplicial set
function LOCALFUZZYSIMPLICIALSET( $X, x, n$ )
  knn, knn-dists  $\leftarrow$  APPROXNEARESTNEIGHBORS( $X, x, n$ )
   $\rho \leftarrow$  knn-dists[1]  $\triangleright$  Distance to nearest neighbor
   $\sigma \leftarrow$  SMOOTHKNNDIST(knn-dists,  $n, \rho$ )  $\triangleright$  Smooth approximator to knn-distance
  fs-set0  $\leftarrow X$ 
  fs-set1  $\leftarrow$  {([ $x, y$ ], 0) |  $y \in X$ }
  for all  $y \in$  knn do
     $d_{x,y} \leftarrow$  max{0, dist( $x, y$ ) -  $\rho$ }/ $\sigma$ 
    fs-set1  $\leftarrow$  fs-set1  $\cup$  {([ $x, y$ ], exp(- $d_{x,y}$ ))}
  return fs-set
  
```

```

Algorithm 3 Compute the normalizing factor for distances  $\sigma$ 
function SMOOTHKNNDIST(knn-dists,  $n, \rho$ )
  Binary search for  $\sigma$  such that  $\sum_{i=1}^n \exp(-(knn-dists_i - \rho)/\sigma) = \log_2(n)$ 
  return  $\sigma$ 
  
```



Figure 8: Visualization of the full 3 million word vectors from the GoogleNews dataset as embedded by UMAP.

is contained in  $U$ , then  $g$  is constant in  $B$  and hence  $\sqrt{\det(g)}$  is constant can be brought outside the integral. Thus, the volume of  $B$  is

$$\sqrt{\det(g)} \int_B dx^1 \wedge \dots \wedge dx^n = \sqrt{\det(g)} \frac{\pi^{n/2} r^n}{\Gamma(n/2 + 1)}$$

where  $r$  is the radius of the ball in the ambient  $\mathbb{R}^n$ . If we fix the volume of the ball to be  $\frac{\pi^{n/2}}{\Gamma(n/2+1)}$  we arrive at the requirement that

$$\det(g) = \frac{1}{r^{2n}}$$

since  $g$  is assumed to be diagonal with constant entries we can solve for  $g_{ij}$  as

$$g_{ij} = \begin{cases} \frac{1}{r^2} & \text{if } i = j, \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

geodesic distance on  $\mathcal{M}$  under  $g$  from  $p$  to  $q$  (where  $p, q \in B$ ) is defined as

$$\inf_{c \in C} \int_a^b \sqrt{g(\dot{c}(t), \dot{c}(t))} dt$$

where  $C$  is the class of smooth curves  $c$  on  $\mathcal{M}$  such that  $c(a) = p$  and  $c(b) = q$ , and  $\dot{c}$  denotes the first derivative of  $c$  on  $\mathcal{M}$ . Given that  $g$  is as defined in (2) we see that this can be simplified to

$$\begin{aligned} & \frac{1}{r} \inf_{r \in C} \int_a^b \sqrt{\dot{c}(t), \dot{c}(t)} dt \\ & \rightarrow \frac{1}{r} \inf_{r \in C} \int_a^b \|\dot{c}(t)\| dt \\ & = \frac{1}{r} d_{\mathbb{R}^n}(p, q). \end{aligned} \quad (3)$$

□

### B Proof that $\mathbf{FinReal}$ and $\mathbf{FinSing}$ are adjoint

**Theorem 2.** The functors  $\mathbf{FinReal} : \mathbf{Fin-sFuzz} \rightarrow \mathbf{FinEPMet}$  and  $\mathbf{FinSing} : \mathbf{FinEPMet} \rightarrow \mathbf{Fin-sFuzz}$  form an adjunction with  $\mathbf{FinReal}$  the left adjoint and  $\mathbf{FinSing}$  the right adjoint.



# Посмотрим на научные статьи

## Diederik P Kingma, Max Welling: *Auto-Encoding Variational Bayes*,

### ArXiv 1312.6114, 2014

#### 2.2 The variational bound

The marginal likelihood is composed of a sum over the marginal likelihoods of individual datapoints  $\log p_{\theta}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)})$ , which can each be rewritten as:

$$\log p_{\theta}(\mathbf{x}^{(i)}) = D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})) + \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) \quad (1)$$

The first RHS term is the KL divergence of the approximate from the true posterior. Since this KL-divergence is non-negative, the second RHS term  $\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$  is called the (variational) lower bound on the marginal likelihood of datapoint  $i$ , and can be written as:

$$\log p_{\theta}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} [-\log q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)}) + \log p_{\theta}(\mathbf{x}, \mathbf{z})] \quad (2)$$

which can also be written as:

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})) + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})}$$

We want to differentiate and optimize the lower bound  $\mathcal{L}(\theta, \phi; \mathbf{x})$  parameters  $\phi$  and generative parameters  $\theta$ . However, the gradient is a bit problematic. The usual (naïve) Monte Carlo gradient estimator:  $\nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [f(\mathbf{z})] = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [f(\mathbf{z}) \nabla_{\phi} \log q_{\phi}(\mathbf{z}|\mathbf{x})] \approx \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}^{(l)}) \nabla_{\phi} \log q_{\phi}(\mathbf{z}|\mathbf{x})$ . This gradient estimator exhibits high variance and is impractical for our purposes.

#### 2.3 The SGVB estimator and AEVB algorithm

In this section we introduce a practical estimator of the lower bound parameters. We assume an approximate posterior in the form  $q_{\phi}(\mathbf{z}|\mathbf{x})$  technique can be applied to the case  $q_{\phi}(\mathbf{z}|\mathbf{x})$ , i.e. where we do not use variational Bayesian method for inferring a posterior over the parameters.

Under certain mild conditions outlined in section 2.4 for a chosen approximate posterior  $\bar{\mathbf{z}} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$  using a different (an auxiliary) noise variable:

$$\bar{\mathbf{z}} = g_{\phi}(\epsilon, \mathbf{x}) \quad \text{with} \quad \epsilon \sim p(\epsilon)$$

See section 2.4 for general strategies for choosing such an approximate posterior  $q_{\phi}(\mathbf{z}|\mathbf{x})$ . We can now form Monte Carlo estimates of expectation  $\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}$  as follows:

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} [f(\mathbf{z})] = \mathbb{E}_{p(\epsilon)} [f(g_{\phi}(\epsilon, \mathbf{x}^{(i)}))] \approx \frac{1}{L} \sum_{l=1}^L f(g_{\phi}(\epsilon^{(l)}, \mathbf{x}^{(i)}))$$

We apply this technique to the variational lower bound (eq. 2), yielding our generic Stochastic Gradient Variational Bayes (SGVB) estimator  $\tilde{\mathcal{L}}^A(\theta, \phi; \mathbf{x}^{(i)}) \approx \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$ :

$$\tilde{\mathcal{L}}^A(\theta, \phi; \mathbf{x}^{(i)}) = \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{x}^{(i)}, \mathbf{z}^{(i,l)}) - \log q_{\phi}(\mathbf{z}^{(i,l)}|\mathbf{x}^{(i)})$$

$\mathcal{L}(\theta, \phi; \mathbf{x})$  is the variational lower bound of the marginal likelihood of datapoint  $i$ :

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = \int q_{\phi}(\mathbf{z}|\mathbf{x}) (\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})) dz \quad (16)$$

The expectations on the RHS of eqs (14) and (16) can obviously be written as a sum of three separate expectations, of which the second and third component can sometimes be analytically solved, e.g. when both  $p_{\theta}(\mathbf{z})$  and  $q_{\phi}(\mathbf{z}|\mathbf{x})$  are Gaussian. For generality we will here assume that each of these expectations is intractable.

Under certain mild conditions outlined in section (see paper) for chosen approximate posteriors parameterize conditional samples  $\bar{\mathbf{z}} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$  as

$$\bar{\mathbf{z}} = g_{\phi}(\epsilon, \mathbf{x}) \quad \text{with} \quad \epsilon \sim p(\epsilon) \quad (17)$$

) and a function  $g_{\phi}(\epsilon, \mathbf{x})$  such that the following holds:

$$\mathbb{E}_{p(\epsilon)} [f(g_{\phi}(\epsilon, \mathbf{x}))] = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [f(\mathbf{z})] \quad (18)$$

approximate posterior  $q_{\phi}(\mathbf{z})$ :

$$\bar{\theta} = h_{\phi}(\zeta) \quad \text{with} \quad \zeta \sim p(\zeta) \quad (19)$$

e, choose a prior  $p(\zeta)$  and a function  $h_{\phi}(\zeta)$  such that the following

$$p(\zeta) (\log p_{\theta}(\mathbf{X}) + \log p_{\alpha}(\theta) - \log q_{\phi}(\theta)) \Big|_{\theta=h_{\phi}(\zeta)} d\zeta \quad (20)$$

we introduce a shorthand notation  $f_{\phi}(\mathbf{x}, \mathbf{z}, \theta)$ :

$$f_{\phi}(\mathbf{x}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p_{\alpha}(\theta) - \log q_{\phi}(\theta) \quad (21)$$

(18), the Monte Carlo estimate of the variational lower bound, given

$$\tilde{\mathcal{L}}(\theta; \mathbf{X}) \approx \frac{1}{L} \sum_{i=1}^L f_{\phi}(\mathbf{x}^{(i)}, g_{\phi}(\epsilon^{(i)}, \mathbf{x}^{(i)}), h_{\phi}(\zeta^{(i)})) \quad (22)$$

where  $\epsilon^{(i)} \sim p(\epsilon)$  and  $\zeta^{(i)} \sim p(\zeta)$ . The estimator only depends on samples from  $p(\epsilon)$  and  $p(\zeta)$  which are obviously not influenced by  $\phi$ , therefore the estimator can be differentiated w.r.t.  $\phi$ .

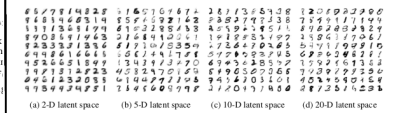
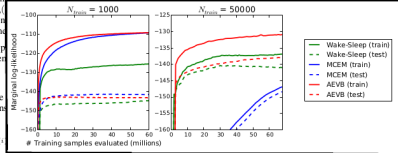


Figure 5: Random samples from learned generative models of MNIST for different dimensionalities of latent space.





# Посмотрим на примеры использования библиотек

## Variational Autoencoders

<https://github.com/pyro-ppl/pyro/blob/dev/examples/vae/vae.py>

```

1  # Copyright (c) 2017-2019 Uber Technologies, Inc.
2  # SPDX-License-Identifier: Apache-2.0
3
4  import argparse
5
6  import numpy as np
7  import torch
8  import torch.nn as nn
9  import visdom
10
11 import pyro
12 import pyro.distributions as dist
13 from pyro.infer import SVI, JitTrace_ELBO, Trace_ELBO
14 from pyro.optim import Adam
15 from utils.mnist_cached import MNISTCached as MNIST
16 from utils.mnist_cached import setup_data_loaders
17 from utils.vae_plots import mnist_test_tsne, plot_tik, plot_vae_samp
18
19 # define the PyTorch module that parameterizes the
20 # diagonal gaussian distribution q(z|x)
21 class Encoder(nn.Module):
22     def __init__(self, z_dim, hidden_dim):
23         super().__init__()
24         # setup the three linear transformations used
25         self.fc1 = nn.Linear(784, hidden_dim)
26         self.fc21 = nn.Linear(hidden_dim, z_dim)
27         self.fc22 = nn.Linear(hidden_dim, z_dim)
28         # setup the non-linearities
29         self.softplus = nn.Softplus()
30
31     def forward(self, x):
32         # setup the forward computation on the image x
33         # first shape the mini-batch to have pixels in the rightmost
34         x = x.reshape(-1, 784)
35         # then compute the hidden units
36         hidden = self.softplus(self.fc1(x))
37         # then return a mean vector and a (positive) square root over
38         # each of size batch_size x z_dim
39         z_loc = self.fc21(hidden)
40         z_scale = torch.exp(self.fc22(hidden))
41         return z_loc, z_scale
42
43 # define the PyTorch module that parameterizes the
44 # observation likelihood p(x|z)
45 class Decoder(nn.Module):
46     def __init__(self, z_dim, hidden_dim):
47         super().__init__()
48         # setup the two linear transformations used
49         self.fc31 = nn.Linear(z_dim, hidden_dim)
50         self.fc32 = nn.Linear(hidden_dim, 784)
51         # setup the non-linearities
52         self.softplus = nn.Softplus()
53
54     def forward(self, z):
55         # define the forward computation on the latent z
56         # first compute the hidden units
57         hidden = self.softplus(self.fc31(z))
58         # return the parameter for the output Bernoulli
59         # each is of size batch_size x 784
60         loc_log = torch.sigmoid(self.fc32(hidden))
61         return loc_log
62
63 # define a PyTorch module for the VAE
64 class VAE(nn.Module):
65     def __init__(self, z_dim=50, hidden_dim=400, use_cuda=False):
66         # by default our latent space is 50-dimensional
67         # and we use 400 hidden units
68         super().__init__()
69         # create the encoder and decoder networks
70         self.encoder = Encoder(z_dim, hidden_dim)
71         self.decoder = Decoder(z_dim, hidden_dim)
72
73     if use_cuda:
74         # calling cuda[] here will put all the parameters of
75         # the encoder and decoder networks into gpu memory
76         self.encoder.cuda()
77         self.decoder.cuda()
78         self.use_cuda = use_cuda
79         self.z_dim = z_dim
80
81     def model(self, x):
82         # register PyTorch module "decoder" with Pyro
83         pyro.module("decoder", self.decoder)
84         with pyro.plate("data", x.shape[0]):
85             # setup hyperparameters for prior p(z)
86             z_loc = torch.randn(x.shape[0], self.z_dim, dtype=x.dtype, dev
87             z_scale = torch.ones(x.shape[0], self.z_dim, dtype=x.dtype, de
88             # sample from prior (z will be sampled by guide when compo
92         # decode the latent code z
93         loc_log = self.decoder.forward(z)
94         # score against actual images
95         pyro.sample("obs", dist.Bernoulli(loc_log).to_event(1), obs=x)
96         # return the loc so we can visualize it later
97         return loc_log
98
99     # define the guide (i.e. variational distribution) q(z|x)
100     def guide(self, x):
101         # register PyTorch module "encoder" with Pyro
102         pyro.module("encoder", self.encoder)
103         with pyro.plate("data", x.shape[0]):
104             # use the encoder to get the parameters used to define q(z|x)
105             z_loc, z_scale = self.encoder.forward(x)
106             # sample the latent code z
107             pyro.sample("latent", dist.Normal(z_loc, z_scale).to_event(1))
108
109     # define a helper function for reconstructing images
110     def reconstruct_img(self, x):
111         # encode image x
112         z_loc, z_scale = self.encoder(x)
113         # sample in latent space
114         z = dist.Normal(z_loc, z_scale).sample()
115         # decode the image (note we don't sample in image space)
116         loc_log = self.decoder(z)
117         return loc_log
118
119 def main(args):
120     # clear param store
121     pyro.clear_param_store()
122
123     # setup MNIST data loaders
124     train_loader, test_loader
125     train_loader, test_loader = setup_data_loaders(MNIST, use_cuda=args.c
126
127     # setup the VAE
128     vae = VAE(use_cuda=args.cuda)
129
130     # setup the optimizer
131     adam_args = {"lr": args.learning_rate}
132     optimizer = Adam(adam_args)
133
134     # setup the inference algorithm
135     elbo = JitTrace_ELBO() if args.jit else Trace_ELBO()
136     svi = SVI(vae.model, vae.guide, optimizer, loss=elbo)
137
138     # setup visdom for visualization
139     if args.visdom_flag:
140         vis = visdom.Visdom()
141
142     train_elbo = []
143     test_elbo = []
144     # training loop
145     for epoch in range(args.num_epochs):
146         # initialize loss accumu
147         epoch_loss = 0.
148         # do a training epoch over
149         # by the data loader
150         for x, _ in train_loader:
151             # if on GPU put mini
152             if args.cuda:
153                 x = x.cuda()
154             # do ELBO gradient ar
155             epoch_loss += svi.step
156
157     # report training diagnos
158     normalizer_train = len(tr
159     total_epoch_loss_train = e
160     train_elbo.append(total_elo
161     print("epoch %d loss %f" %
162
163     if args.test_flag:
164         # initialize loss accu
165         test_loss = 0.
166         # compute the loss ove
167         for i, (x, _) in enumerate
168             # if on GPU put ac
169             if args.cuda:
170                 x = x.cuda()
171             # compute ELBO est
172             test_loss += svi.e
173
174     # pick three rando
175     # visualize how w
176     if i == 0:
177         if args.visdom:
178             plot_vae_s
179             recs_initi
180             for index
181                 test_l
182                 recs_l
183                 vis.is
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500

```





# Какие вообще инструменты могут потребоваться



PyTorch



Yandex  
CatBoost



plotly



Вывод:  
и математика  
и программирование



# Анализ данных это

процесс поиска закономерностей  
в данных при помощи

- ▶ средств визуализации данных,
- ▶ математических методов,
- ▶ программных алгоритмов.



Искусственный интеллект

## Отличительная особенность:

нет четко зафиксированного ответа на каждый входящий объект.

## Что можно почитать:

Анализ данных — основы и терминология

<https://habr.com/ru/post/352812/>

Всё, что вам нужно знать об ИИ — за несколько минут

<https://habr.com/ru/post/416889/>



# Сравним задачи

## Алгоритмы и структуры данных

*Задача:* дан массив  $x$ , нужно его отсортировать.

Ровно один правильный ответ, можно получить с помощью четких алгоритмов.

## Комбинаторика

*Задача:* Сколько имеется способов раздать 11 разных цветков, трём девушкам: какой-то – 5, а остальным – по 3 цветка? [ОКТЧ 2019]

Ровно один правильный ответ.

## Анализ данных

*Задача:* Имеются данные  $(x_1, y_1), \dots, (x_n, y_n)$ .

Восстановите по ним функцию  $f : x \mapsto y$ .

*Особенности:* нет четкого ответа, требуется только приближение, но есть критерии качества.



## Пример — распознавание рукописных цифр

Вход: 

Ожидается на выходе: 5

Но как четко алгоритмически определить границу между 6 и 8?



— 2 или 9?



— 4 или 7?



## Актуальность в научной среде

### Число статей по запросам в Google Scholar с 2020:

- |   |   |
|---|---|
| ▶ <i>statistics</i> $\approx$ 1 240 000           | ▶ <i>статистика</i> $\approx$ 14 600 статей       |
| ▶ <i>machine learning</i> $\approx$ 1 040 000     | ▶ <i>машинное обучение</i> $\approx$ 15 200       |
| ▶ <i>computer vision</i> $\approx$ 537 000        | ▶ <i>компьютерное зрение</i> $\approx$ 12 100     |
| ▶ <i>deep learning</i> $\approx$ 236 000          | ▶ <i>глубокое обучение</i> $\approx$ 15 900       |
| ▶ <i>generative models</i> $\approx$ 143 000      | ▶ <i>генеративные модели</i> $\approx$ 15 800     |
| ▶ <i>language models</i> $\approx$ 137 000        | ▶ <i>языковые модели</i> $\approx$ 15 800         |
| ▶ <i>neural network</i> $\approx$ 108 000         | ▶ <i>нейронные сети</i> $\approx$ 15 600          |
| ▶ <i>artificial intelligence</i> $\approx$ 91 900 | ▶ <i>искусственный интеллект</i> $\approx$ 15 900 |



# Обзор задач анализа данных



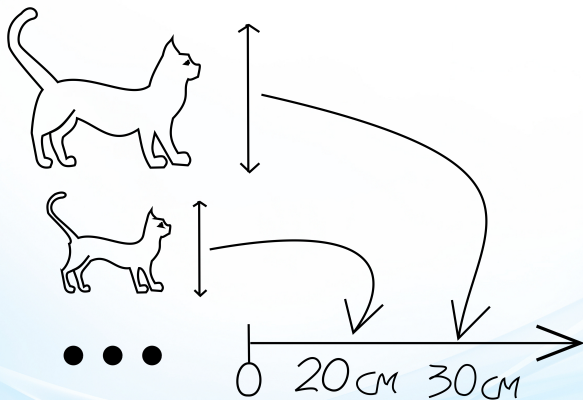
# Мурмурландия



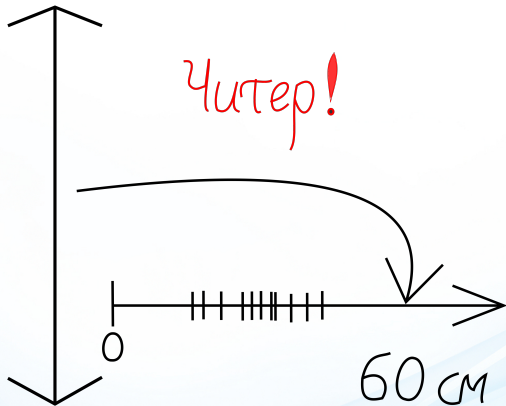




Каков средний рост котиков?



**Точечное оценивание**

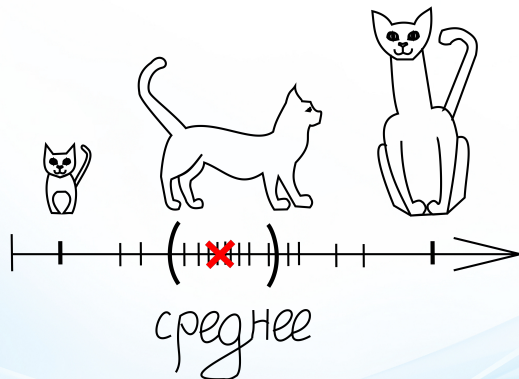


Выбросы



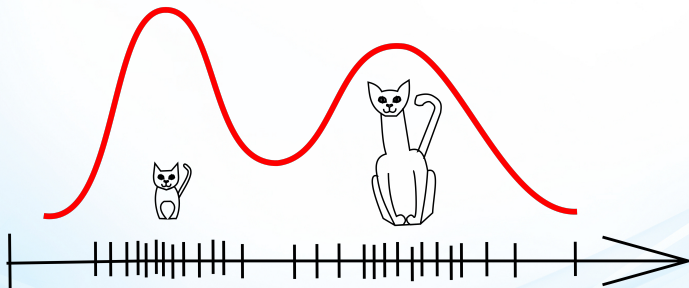


Среднее определяется неточно



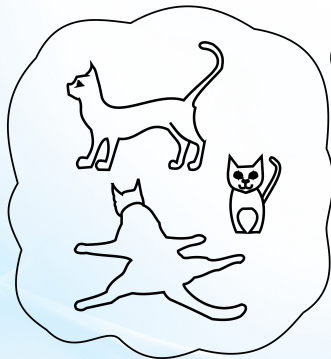
Интервальное оценивание

# Характер распределения



Непараметрическое оценивание

низкие



высокие

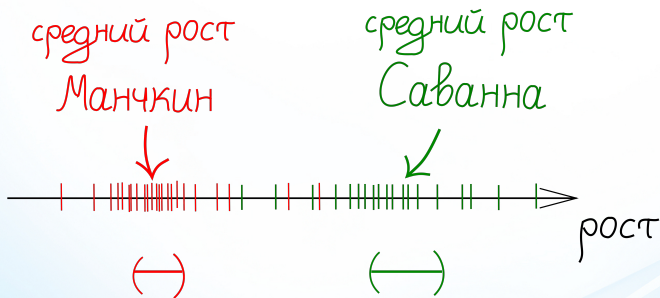




Отличается ли их средний рост?



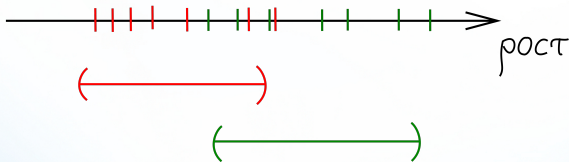
# Собираем данные



отличается

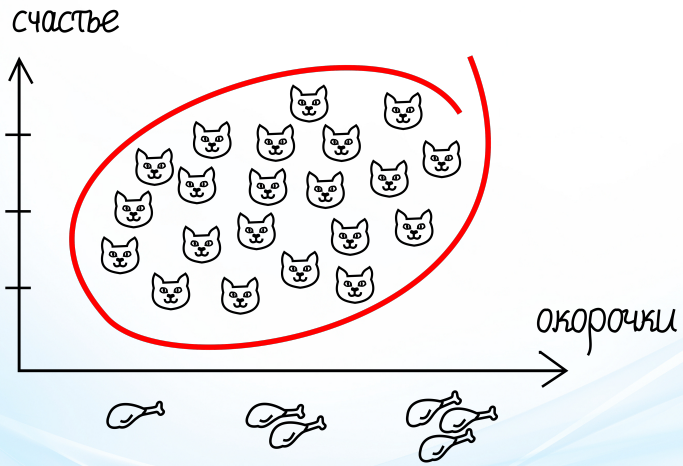


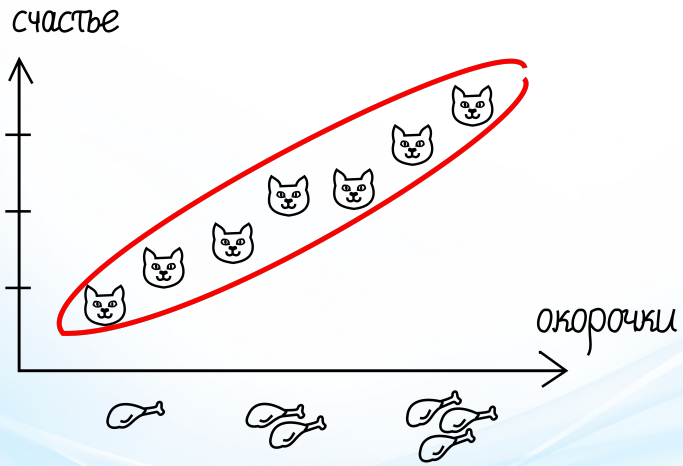
Если данных мало

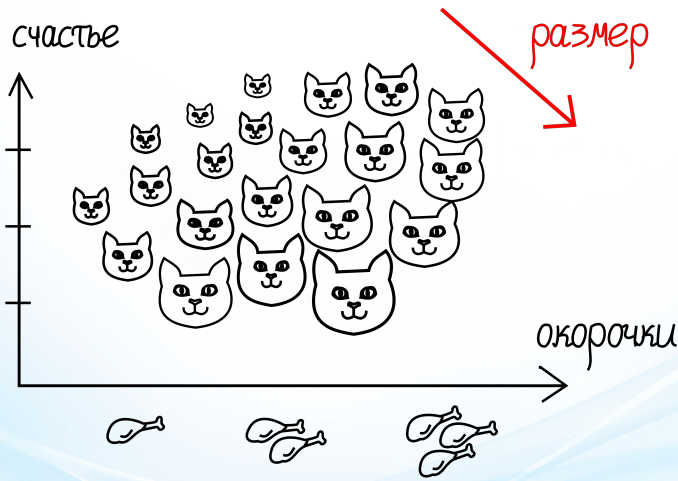


непонятно

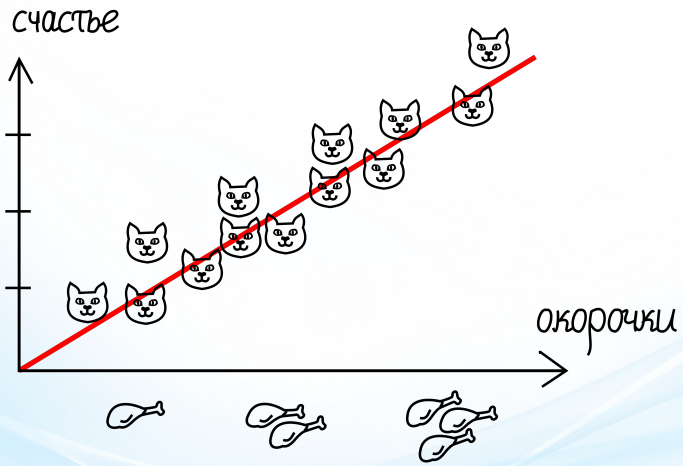
Статистические гипотезы, АВ-тесты







Корреляционный анализ





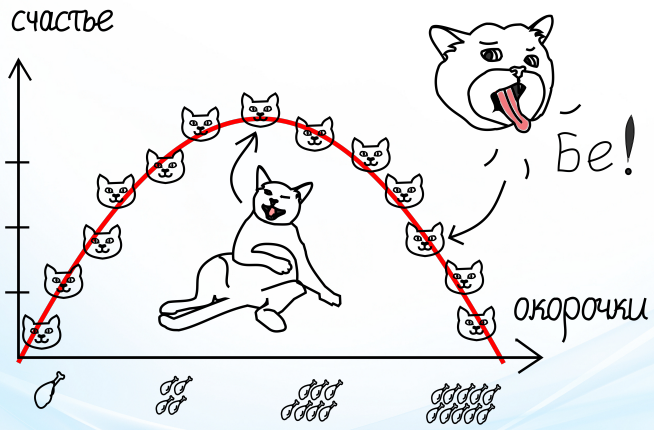
# Формула счастья



$$= \theta_0 + \theta_1 \times \text{кол-во} \begin{array}{c} \text{курицы} \\ \text{и} \\ \text{кости} \end{array} + \text{погрешности}$$



# Больше окорочков





# Формула счастья



$$= \theta_0 + \theta_1 \times \text{кол-во} - \theta_2 \times (\text{кол-во})^2 + \text{погрешности}$$





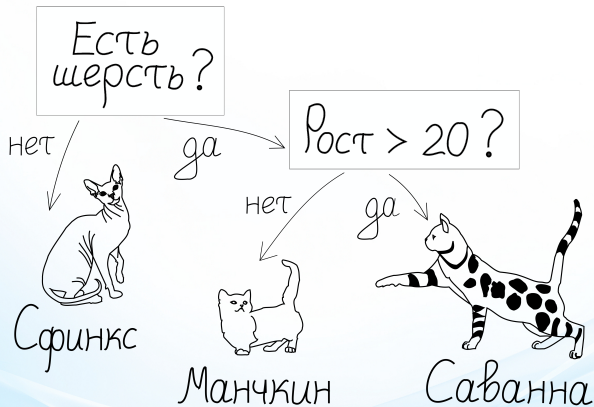
## Другие факторы

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

Регрессионный анализ







# Классификация котиков



**Классификация**

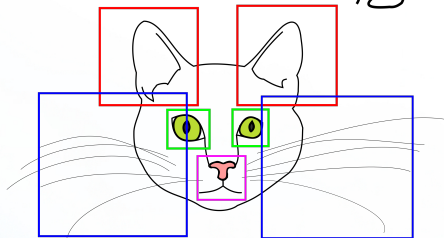


# Собираем данные

котик	порода	рост	шерсть
	Саванна	50 см	да
	Сфинкс	30 см	нет
	Манчкин	15 см	да
	Саванна	40 см	да



# Распознавание мордочек



Нейронные сети



# Художник курса



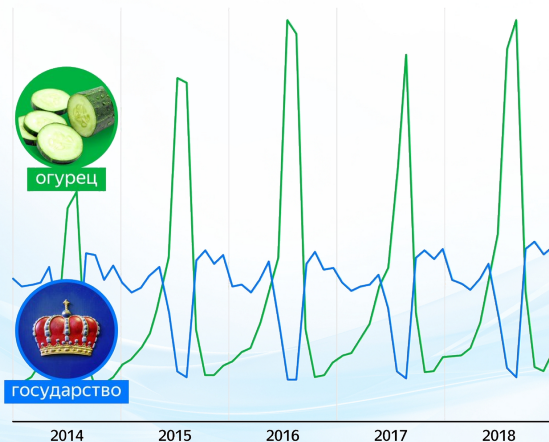
*Евгений*



Книга с похожим содержанием



Когда в Поиске растёт интерес к **огурцам**, снижается доля запросов со словом **государство**





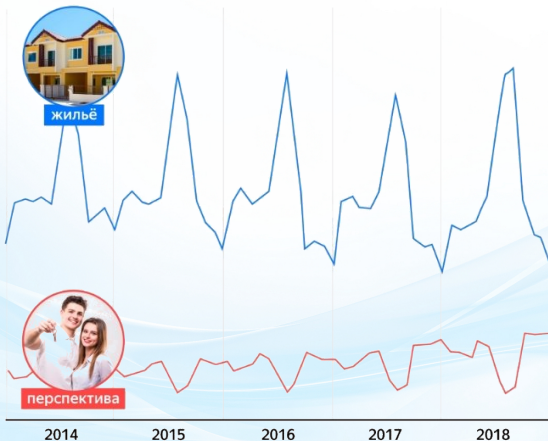
Когда в Поиске растёт интерес к **тату**,  
снижается доля запросов со словом **СМЫСЛ**







Когда в Поиске растёт интерес к **жилью**, снижается доля запросов со словом **перспектива**





# Президентские выборы в США в 1948 г.

Гарри Труман (демократы) vs. Томас Дьюи (республиканцы)

В ночь на оглашение результатов газета Chicago Tribune опубликовала заголовок: **DEWEY DEFEATS TRUMAN**



После закрытия участков газета провела опрос, обзвонив большое число избирателей, все предвещало оглушительную победу Дьюи.



## Президентские выборы в США в 1948 г.

Смеющийся Труман, победитель выборов 1948 года.



Что же пошло не так?

В 1948 году телефон был доступен только людям определенного достатка и редко встречался у людей с небольшим заработком.

Выборка не учитывала достаточно широкий пласт избирателей Трумана, т.к. как правило демократы имеют большую долю голосов среди бедного населения, которым телефон в свою очередь был недоступен.



Попробуем решить задачу



А ты кто?

Перед нами домашнее животное. Кто это — собака или кот?

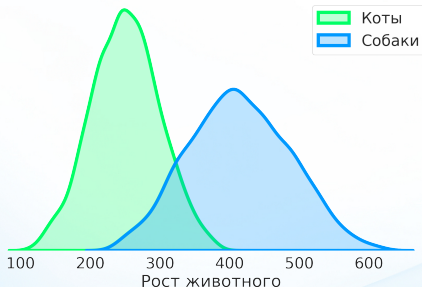




## Классификация: собака vs кот

Попробуем сначала извлечь какой-то *признак*.

Построим вероятностные плотности для каждого класса.



При каких-то значениях роста мы уже можем с большой уверенностью сказать ответ.

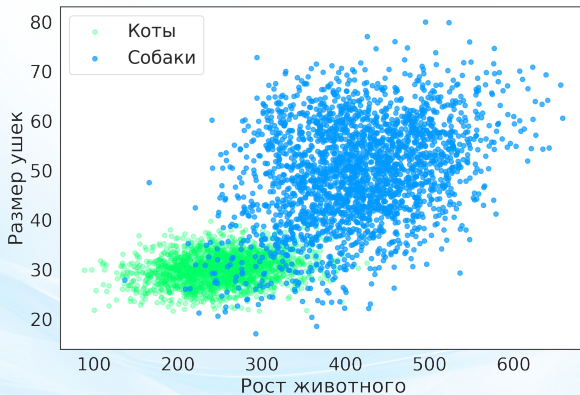
Но есть большое пересечение, это не очень здорово.



## Классификация: собака vs кот

Извлечем еще один признак — размер ушек.

Теперь классы лучше разделяются.

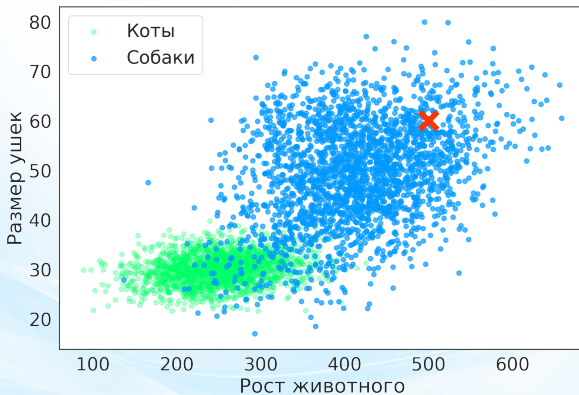




# Классификация: собака vs кот

Попробуем классифицировать новое животное.

Кто отмечен красным?



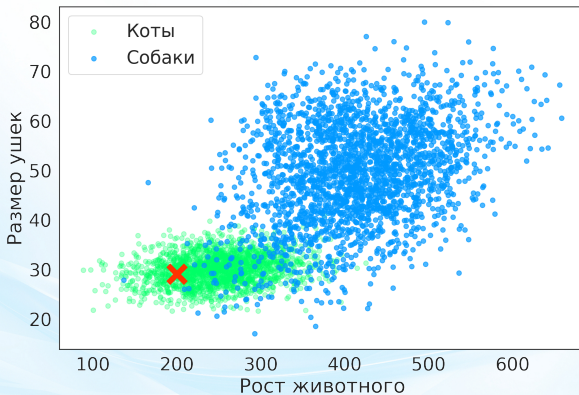




# Классификация: собака vs кот

Попробуем классифицировать новое животное.

Кто отмечен красным?

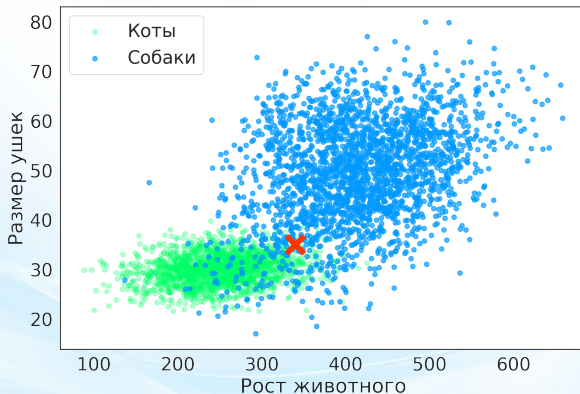




## Классификация: собака vs кот

Попробуем классифицировать новое животное.

Кто отмечен красным?



На основе чего вы сделали все выводы?



## Метод ближайших соседей (kNN)

### Дано:

$X_1, \dots, X_n$  — набор размеченных объектов.

$Y_1, \dots, Y_n$  — соответствующие метки класса.

### Задача:

Пусть  $x$  — исследуемый объект. Какого он класса?

### Решение:

Будем смотреть на свойства  $k$  ближайших соседей.

$x_{(1)}, \dots, x_{(k)}$  —  $k$  его соседей в порядке удаления от  $x$ .

$Y_1, \dots, Y_k$  — соответствующие им классы.

Ответ — наиболее часто встречаемый класс среди  $x_{(1)}, \dots, x_{(k)}$ .

### Свойства:

1.  $k$  — гиперпараметр модели;
2. Не редко на практике показывает хорошие результаты.

### 3. Дорогое применение:

для каждого  $x$  результат вычисляется за  $O(n \ln n)$ .



## Взвешенный метод ближайших соседей

Пусть  $x$  — исследуемый объект.

$x_{(1)}, \dots, x_{(k)}$  —  $k$  его соседей в порядке удаления от  $x$ .

$Y_1, \dots, Y_k$  — соответствующий класс.

$w_1, \dots, w_k$  — вклад  $k$ -го соседа, определяемый пользователем.

Способы определения веса:

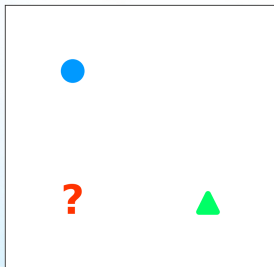
- ▶  $w_j = 1 - j/k$  — зависящий от номера соседа;
- ▶  $w_j = \|x - x_{(j)}\|^{-1}$  — зависящий от расстояния до соседа.

$$\hat{y}(x) = \arg \max_y \sum_{j=1}^k w_j I\{Y_j = k\} \text{ — классификация}$$



# Особенности

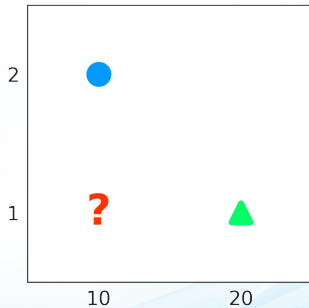
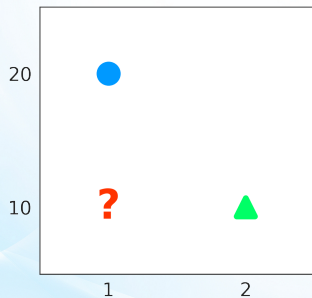
Классифицируйте объект "?".





# Особенности

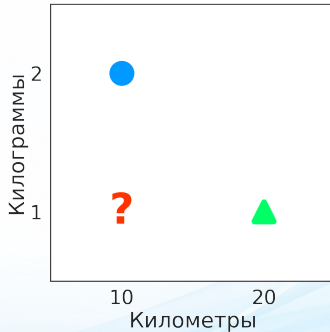
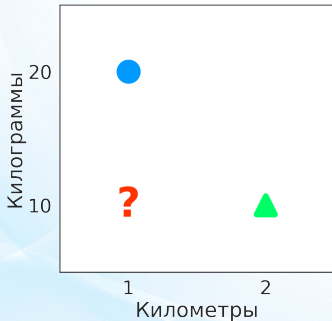
Классифицируйте объект "?".





# Особенности

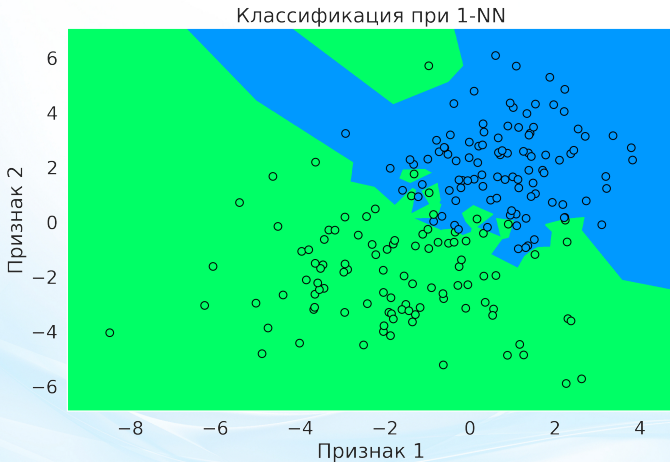
Классифицируйте объект "?".



**Вывод:** результат сильно зависит от используемой метрики между точками в пространстве. Не складывайте *кг* с *км*!



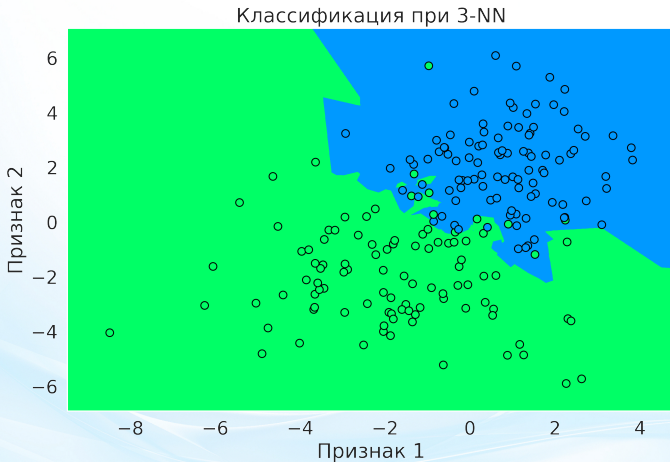
# Что происходит при разных $k$ ?





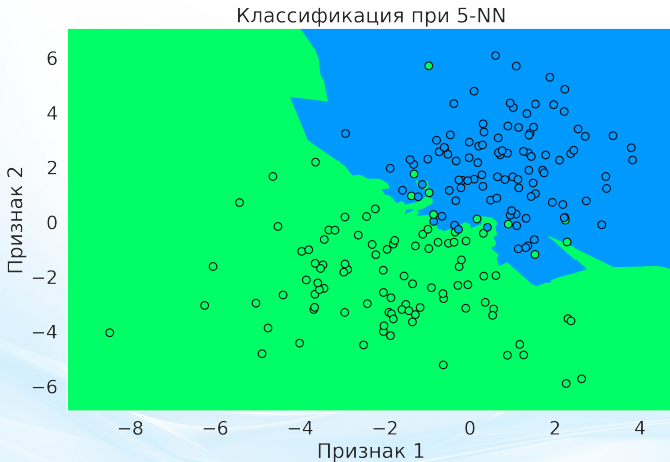


# Что происходит при разных $k$ ?



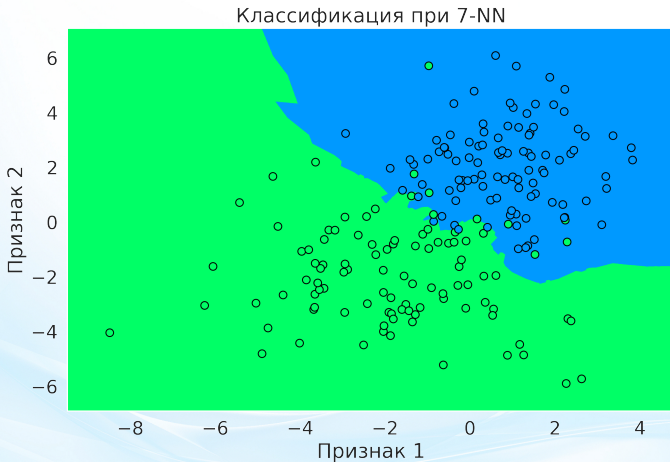


# Что происходит при разных $k$ ?



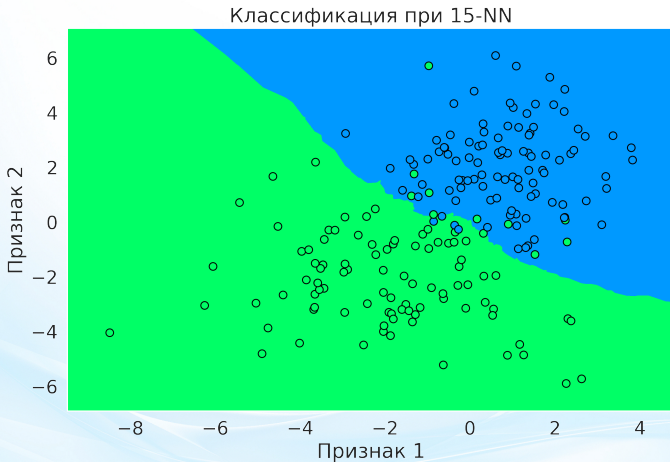


# Что происходит при разных $k$ ?



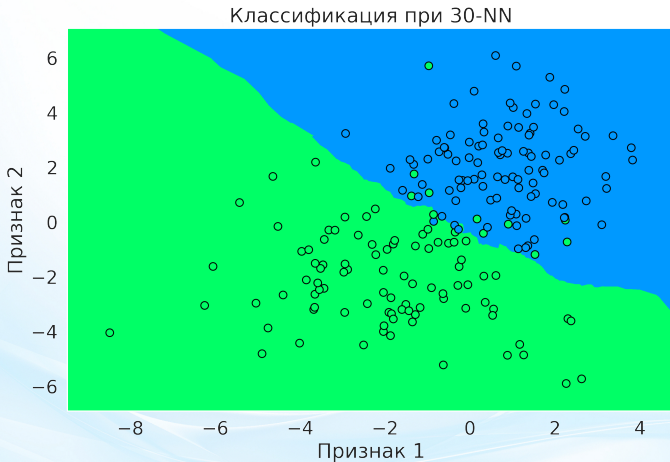


# Что происходит при разных $k$ ?





# Что происходит при разных $k$ ?





## Как оценить качество классификации?

Пусть  $\hat{y}(x)$  — оценка класса для объекта  $x$ .

Можем посчитать **точность** — доля правильно угаданных классов

$$A = \frac{1}{n} \sum_{i=1}^n I\{Y_i = \hat{y}(x_i)\}$$

Оценка качества называется **метрикой** (не путать с метр. пр-вами).

**Какое число соседей оптимизирует эту метрику?**

Ответ:  $k = 1$ , т.к. при вычислении  $\hat{y}(x_i)$  берем сам  $Y_i$ .

Поэтому данные делят случайно на **две непересекающиеся части**:

1. на одной определяют правило классификации,
2. на другой — считают оценку качества классификации.

**Точность 90% это много или мало?**

Кажется, круто. А если в данных 85% котов? Тогда отвечая всегда "кот" сможем добиться точности 85%, и 90% уже не так круто...



## А что если по картинке?

Хорошо, но что если объект — изображение кота или собаки?

Изображение  $100 \times 100$  состоит из  $10^4$  пикселей,

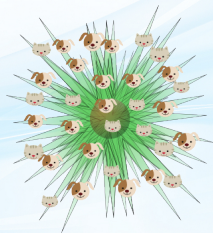
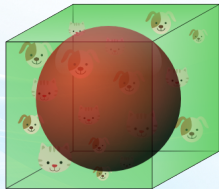
в каждом по 3 числа. Какой размерности получается объект?

Ответ:  $100 \times 100 \times 3 = 30\,000$  чисел в одной картинке.

### Проблема:

в пр-ве больших размерностей расстояния неинформативны.

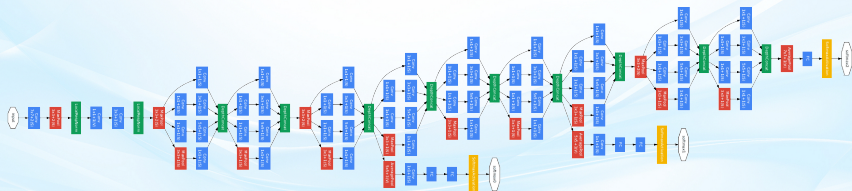
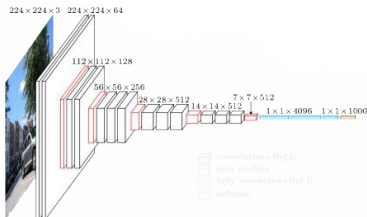
Например, среди фиксированного количества случайных точек в единичном кубе в пространстве большой размерности почти все точки будут лежать около границы куба.





# А что в сложных случаях?

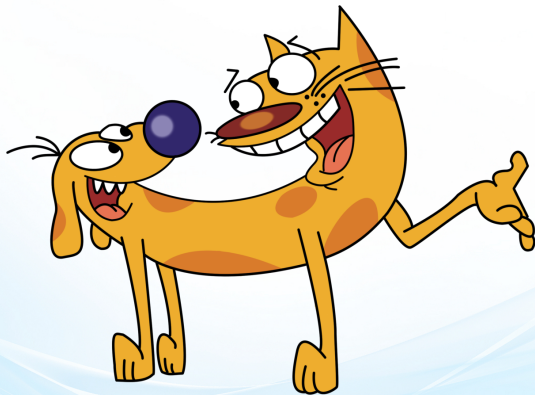
Нейросети! Но об этом позже :)







А потом приходит кто?





# Примеры прикладных задач

Рекомендательная система

Продуктовая аналитика

Распознавание лиц

Генерация изображений

Синтез речи

Где еще?



# Рекомендательная система


## Постановка задачи:



Спроектировать рекомендательную систему для муз. сайта.  
Система должна рекомендовать музыку каждому пользователю  
в соответствии с его предпочтениями.

Данные:

- ▶ оценки пользователей различным трекам;
- ▶ прослушивание треков пользователями.



Яндекс Музыка  Главное Подкасты Жанры Радио

 **ИСПОЛНИТЕЛЬ**  
**David Guetta**  
Нравится слушателям: Avicii, Rihanna, Calvin Harris

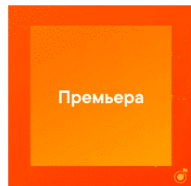

[▶ Слушать](#) 2 577 575  Радио 

**ГЛАВНОЕ** [ТРЕКИ](#) [АЛЬБОМЫ](#) [КЛИПЫ](#) [ПОХОЖИЕ](#) [ИНФО](#)

**Популярные треки** Смотреть всё >

Трек	Длительность
 Memories (2021 Remix; feat. Kid Cudi) — David Guetta, K... <span>♥</span>	2:42
 Let's Love — David Guetta, Sia <span>♥</span>	3:20

**Недавний релиз**



Премьера

Только новинки, подобранные  
по вашим предпочтениям

Обновлён 7 февраля



# Рекомендательная система

Пусть  $U$  — множество треков,  $I$  — множество пользователей.

$R = (r_{ui})_{u \in U, i \in I}$  — матрица лайков/прослуш. пользователями треков.

**Идея:**  $T$  — небольшое множество интересов.

$P = (p_{tu})_{t \in T, u \in U}$  — матрица интересов пользователей.

$Q = (q_{ti})_{t \in T, i \in I}$  — матрица соответствия треков интересам.

Если величины  $p_{tu}$  и  $q_{ti}$  неотрицательны, то их можно нормировать по темам, получив вероятности:

- ▶  $p_{tu} / \sum_{t \in T} p_{tu} = P(\text{пользователю } u \text{ интересно } t),$
- ▶  $q_{ti} / \sum_{t \in T} q_{ti} = P(\text{трек } i \text{ соответствует } t).$



# Рекомендательная система

## Неотрицательные матричные разложения

Решаем поочередно для каждой  $t$ .

$$\begin{cases} \|R_t - p_t^T q_t\|^2 \rightarrow \min_{p_t} \\ p_t \geq 0 \end{cases} \implies p_t = \left( \frac{q_t R_t^T}{q_t q_t^T} \right)^+$$

$$\begin{cases} \|R_t - p_t^T q_t\|^2 \rightarrow \min_{q_t} \\ q_t \geq 0 \end{cases} \implies q_t = \left( \frac{p_t R_t^T}{p_t p_t^T} \right)^+$$

## Что нужно помнить при реализации

- ▶ Пользователей и треков миллионы;
- ▶ Матрица  $R$  сильно разрежена — лайков очень малая доля.



# Примеры прикладных задач

Рекомендательная система

Продуктовая аналитика

Распознавание лиц

Генерация изображений

Синтез речи

Где еще?

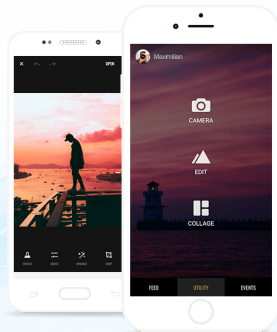


# Продуктовая аналитика

Представьте, что вы продуктовый аналитик в фоторедакторе.  
Вы всегда хотите, чтобы продукт приносил больше денег.

**Реальный вопрос:**

пользователи какой страны платят больше?





# Продуктовая аналитика

## Постановка задачи:

Оказывает ли страна пользователя влияние на его LifeTime Value?

**LifeTime Value** — доход, который принес пользователь за все время жизни в продукте.

Данные:

- ▶ страна пользователя — фактор;
- ▶ LifeTime Value каждого пользователя.





# Продуктовая аналитика

## **Решение задачи:**

Однофакторный дисперсионный анализ и Post-Hoc анализ.

Данные методы помогают понять какие значения фактора влияют на исследуемую величину и как именно.

Т.е., можно узнать, пользователи каких стран платят больше.

Как же выглядят эти методы?



# Однофакторный дисперсионный анализ

## Независимые выборки

1	2	...	k
$X_{11}$	$X_{12}$	...	$X_{1k}$
$X_{21}$	$X_{22}$	...	$X_{2k}$
...	...	...	...
$X_{n_1 1}$	$X_{n_2 2}$	...	$X_{n_k k}$

$$X_{ij} = \overbrace{\mu + \beta_j}^{\mu_j} + \varepsilon_{ij},$$

$i = 1, \dots, n_j$  — номер наблюдения в выборке

$j = 1, \dots, k$  — номер выборки

$\mu$  — неизвестное общее среднее

$\beta_j$  — неизвестный эффект воздействия фактора для  $j$ -й выборки

$\varepsilon_{ij}$  — случайная ошибка

### Предположение:

$\varepsilon_{ij}$  независимы и имеют одинаковое непрерывное распределение.

$H_0: \mu_1 = \dots = \mu_k$  vs.  $H_1: \exists j_1, j_2$  т.ч.  $\mu_{j_1} \neq \mu_{j_2}$

## Оценка контраста

Пусть  $H_{rs}$  отвергается  $\implies$  оцениваем контраст  $\Delta_{rs} = \mu_r - \mu_s$ .

$V_{rs} = \text{med}\{X_{ir} - X_{is}, i = 1..n_r, j = 1..n_s\}$  — первичная оценка

$W_r = \frac{1}{N} \sum_{s=1}^k n_s V_{rs}$ , где  $V_{rr} = 0$

$\hat{\Delta} = W_r - W_s$  — уточненная оценка контраста

### Свойства:

1. Первичные оценки могут быть несогласованными:  $V_{12} \neq V_{13} + V_{32}$
2. Уточненные оценки согласованы и состоятельны.
3. Уточненные оценки зависят от всех выборок.

Для применения данных методов нужно знать теорию вероятностей и математическую статистику.



## Результат

Теперь мы знаем, в какой стране завести маркетинговую кампанию!

Мы принесли компании деньги, мы молодцы!





# Примеры прикладных задач

Рекомендательная система

Продуктовая аналитика

Распознавание лиц

Генерация изображений

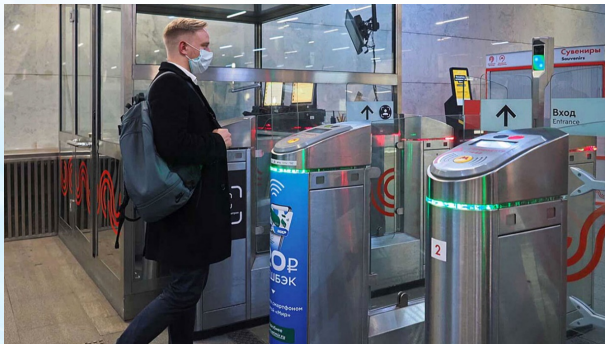
Синтез речи

Где еще?



# Распознавание лиц

Фото человека → Модель детекции лиц → Координаты лица →  
Обработка фото → Сиамская сеть → Идентификатор человека





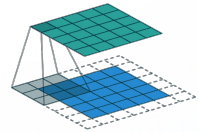
# Распознавание лиц

Пример модели детекции лиц: TinaFace.

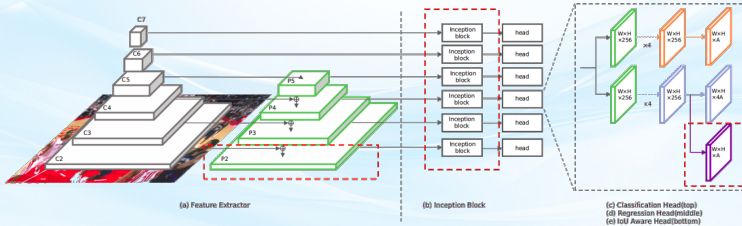
- ▶ Основу модели составляют сверточные слои.

Формула свертки для изображения  $X$  и фильтра с весами  $W$  и сдвигом  $b$ :

$$\sum_{i=1}^M \sum_{j=1}^M w_{ij} \cdot x_{m+i-1, n+j-1} + b$$



- ▶ Архитектура модели содержит множество блоков из сверточных слоев, функций активаций и т.д.





## Распознавание лиц

- ▶ **Сиамская сеть** для двух объектов  $X_1$  и  $X_2$  определяет, принадлежат ли они одному классу, оценивая близость между ними.
- ▶ Архитектура модели представляет собой сверточную сеть.
- ▶ Для оптимизации параметров модели минимизируется **contrastive loss**.

Пусть  $Y_1$  и  $Y_2$  — классы объектов  $X_1$  и  $X_2$  соотв.,  
 $d$  — функция расстояния,

$L_{sim}$  и  $L_{dissim}$  — функции штрафующие за близость объектов одного класса и дальность объектов разных классов соотв.

Тогда лосс равен:

$$L(X_1, X_2, Y_1, Y_2) = I\{Y_1 = Y_2\}L_{sim}(d(f(X_1), f(X_2))) \\ + I\{Y_1 \neq Y_2\}L_{dissim}(d(f(X_1), f(X_2)))$$



# Примеры прикладных задач

Рекомендательная система

Продуктовая аналитика

Распознавание лиц

**Генерация изображений**

Синтез речи

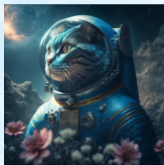
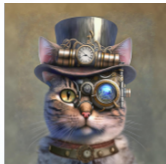
Где еще?





# Генерация изображений

**Задача** — научиться генерировать разнообразные правдоподобные изображения, например котов.





# Генерация изображений

Для этого построим **диффузионную модель**.

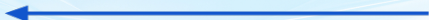
Она моделирует 2 процесса:

- ▶ Прямой процесс — постепенно добавляем шум ко входу.
- ▶ Обратный процесс — модель постепенно восстанавливает данные из шума.

Прямой диффузионный процесс



Обратный диффузионный процесс





# Генерация изображений

## Прямой диффузионный процесс

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \varepsilon, \text{ где } \varepsilon \sim \mathcal{N}(0, I)$$

$$X_t | X_{t-1} \sim \mathcal{N}(\sqrt{1 - \beta_t} X_{t-1}, \beta_t I)$$

Прямой диффузионный процесс



$x_0$

$x_1$

$x_2$

$x_3$

$x_4$

...

$x_T$



# Генерация изображений


## Обратный диффузионный процесс

Обучается нейросетевая модель таким образом, чтобы минимизировать  $ELBO$ :

$$ELBO = E_q \log \frac{p_\theta(X_0, \dots, X_T)}{q(X_1, \dots, X_T | X_0)}$$

$$\simeq const - \sum_{t=2}^T \frac{\tilde{\alpha}_{t-1} \beta_t^2}{2\tilde{\beta}_t(1 - \alpha_t)^2} E_q \|X_0 - x_\theta(X_t, t)\|^2$$



  
 Обратный диффузионный процесс



# Примеры прикладных задач

Рекомендательная система

Продуктовая аналитика

Распознавание лиц

Генерация изображений

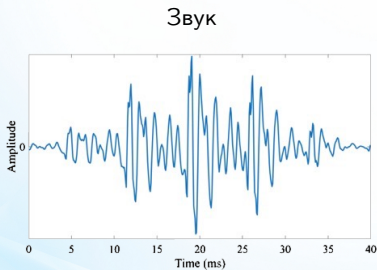
**Синтез речи**

Где еще?



# Синтез речи

## Что такое звук?



**Звук** — композиция волн с разными амплитудами и частотой.

**Волна** — периодич. ф-ия, имеющая амплитуду, период и частоту.



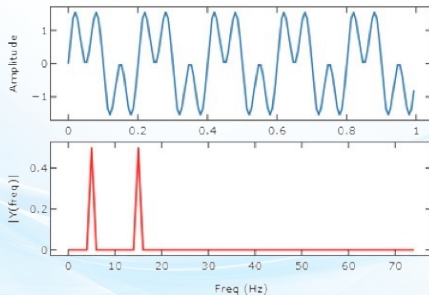
# Синтез речи

## Как работать со звуком?

Посмотрим на распределение частот волн в звуке.

Для этого применяется дискретное преобразование Фурье:

$$X_k = \sum_{n=0}^{N-1} x_n \cdot e^{-i\frac{2\pi}{N}kn} = \sum_{n=0}^{N-1} x_n \left[ \cos\left(\frac{2\pi}{N}kn\right) - i \sin\left(\frac{2\pi}{N}kn\right) \right]$$



Пример для звука из двух волн



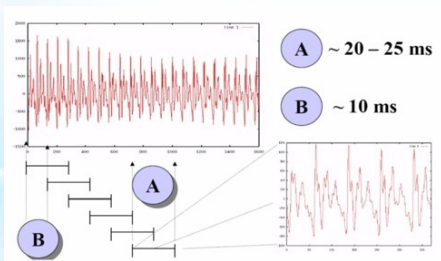
# Синтез речи

## Как работать со звуком?

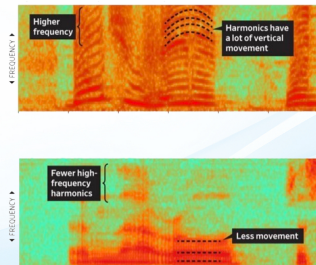
У результата преобразования Фурье ко всем данным нет времени.

Посчитаем распределение на окнах из звука.

Окна



Спектрограмма



- ▶ Берем маленькие окна (20-25 мс) от исходных данных.
- ▶ К каждому окну применяем преобразования Фурье.
- ▶ Стакаем распределения частот вместе, получаем спектрограмму.





# Синтез речи

## Постановка задачи синтеза речи

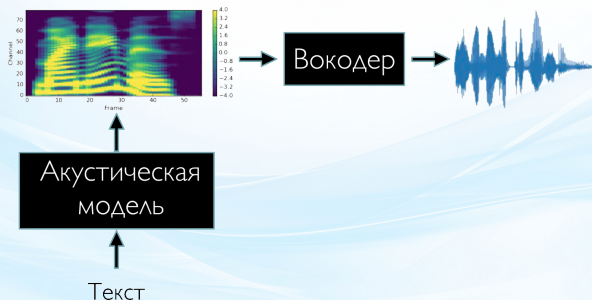
Для введенного получить соответствующее аудио.

## Другие задачи, связанные со звуком:

- ▶ Определить, была ли речь на записи.
- ▶ Споттер — распознавание определенной фразы или определенного события.  
*Например, "Ok, Google", "Алиса" или определение выстрела.*
- ▶ Идентификация говорящего или его признаков.  
*Например, мужчина/женщина, ребенок/взрослый.*
- ▶ Распознавание речи  
*По аудио вернуть полную его расшифровку в виду текста.*

# Синтез речи

1. Акустическая модель по тексту предсказывает мел-спектограмму.  
Акустическая модель — это какая-то нейросеть.
2. Вокодер по мел-спектограмме предсказывает аудио.  
Вокодер может быть как алгоритмом, так и нейросетью.  
Некоторые известные вокодеры используют байесовские методы.





# Примеры прикладных задач

Рекомендательная система

Продуктовая аналитика

Распознавание лиц

Генерация изображений

Синтез речи

Где еще?



- ▶ Беспилотные автомобили
- ▶ Ранжирование результатов в поисковой системе
- ▶ Поиск по картинкам и видео на основе содержания
- ▶ Распознавание спама/фрода
- ▶ Прогноз погоды
- ▶ Прокладывание маршрутов в навигаторах
- ▶ AI-камеры в телефонах
- ▶ Генерация картин подобно художникам
- ▶ Оплата проезда в метро взглядом
- ▶ Решение о выдаче кредита
- ▶ Персонализированная реклама
- ▶ Определение причин оттока клиентов
- ▶ Прогнозирование спроса на товар
- ▶ Определение месторасположения новой торговой точки
- ▶ Расстановка полок в магазинах
- ▶ Автоопределение свежести товаров в магазинах
- ▶ Распознавание патологий на медицинских снимках
- ▶ Персонализированная медицина
- ▶ Прогнозирование поломок оборудования
- ▶ Автоматическая система оптимального управления оборудованием

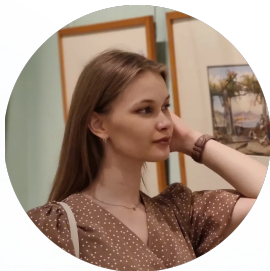


# Анонс следующей лекции

См. ноутбук



## Слово студентам DS-потока



Анна Ефимова

3 курс, DS-поток

tg: @anna\_rrchy



**ВСЁ!**