



Введение в АД





О курсе



ThetaHat



Проводимые нами учебные курсы

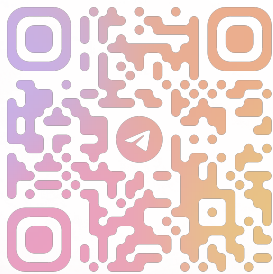
- ▶ Введение в анализ данных;
- ▶ DS-поток;
- ▶ Phystech@DataScience (в т.ч. статистика на ЛФИ);
- ▶ Статистика и машинное обучение на ФБМФ БТ;
- ▶ Мат. статистика на ФБМФ и ФЭФМ (семинары);
- ▶ Прикладная статистика на кафедрах в магистратуре;
- ▶ Машинное обучение в ШАД (частично);
- ▶ АБ-тестирование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



Организационная информация

Telegram-бот

@thetahat_ds25_bot



Код регистрации **Loss=0.193**

Сайт команды thetahat.ru

Почта thetahat@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
08.02	О <i>Введение</i>	 Никита Волков	01.03
15.02	<i>Занятия нет, делаем домашки</i>		
22.02	О <i>AI-инструменты</i>	 Валерия Я. Елизавета Д.	01.03
01.03	О <i>Линейные модели</i>	 Никита Волков	15.03
8.03	<i>Занятия нет, делаем домашки</i>		
15.03	О <i>Введение в нейросети</i>	 Елизавета Дахова	22.03
22.03	Ф <i>Компьютерное зрение и генеративные модели</i>	 Елизавета Дахова	29.03



План занятий (О — обязательная, Ф — факультативная)

Дата	Тема	Лектор	Дедлайн по ДЗ
29.03	Ф <i>Обработка текстов</i>	 Денис Титов	05.04
05.04	Ф <i>Кластеризация и пониж. размерности</i>	 Никита Волков	19.04
12.04	О <i>Гостевая лекция от Яндекса</i>		
19.04	О <i>Теория вероятностей на практике</i>	 Дарья Плотникова	26.04
26.04	Ф <i>Статистика</i>	 Дарья Плотникова	03.05
03.05	Ф <i>Аналитика</i>	 Валерия Якупова	10.05



ChatGPT

Gemini



deepseek



**GIGA
CHAT**



Claude



ChatGPT

Gemini

YaGPT



deepseek

Сопровождение бесплатно



Claude



SIGMA
CHAT

Не можешь предотвратить — возглавь!



Использование ИИ-инструментов в курсе

При выполнении *технической* работы в домашних заданиях **рекомендуется** использовать ИИ-инструменты!

Как? Узнаете на второй лекции!

Ограничения

1. ИИ может ошибаться. Его ошибка в вашей работе — ваша ошибка. Вам необходимо понимать и перепроверять ответы ИИ.
Аргументы "мне так сказал ИИ" не принимаются.
2. Всю содержательную работу по задаче вам необходимо делать самостоятельно.
3. Злоупотребление ИИ приравнивается к списыванию.
4. Ваша цель — обучиться.
Используйте ИИ для выполнения этой цели.



Система оценивания обязательной части

Официальное название: Введение в анализ данных

Кафедра дискретной математики

Активности:

- ▶ Л — доля выполнения легких заданий (их немного);
- ▶ С — доля выполнения сложных заданий (их много);
- ▶ ЛС — доля выполнения всех домашних заданий;
- ▶ В — доля правильных ответов на вопросы в боте на занятии;
- ▶ Т — доля выполнения тестов (на "удовл", в мае);

Списывания:

- ▶ Штраф -2 балла за каждый случай всем участникам;
- ▶ Объяснение "мы просто общались" не прокатит;
- ▶ Злоупотребление ИИ приравнивается к списыванию.



Система оценивания обязательной части

Правила:

Ставится максимальная оценка X , для которой $A \& (B \mid C) = \text{True}$.

Оценка X	Условие А	Условие В	Условие С
3	$T > 34\%$		
4	$T > 67\%$		
5	$L > 25\%$		
6	$L > 50\%$		
7	$L > 75\%$		
8	$L > 25\%$ и $C > 25\%$	$ЛС > 50\%$ и $В > 50\%$	$ЛС > 75\%$
9	$L > 25\%$ и $C > 25\%$	$ЛС > 65\%$ и $В > 65\%$	$ЛС > 85\%$
10	$L > 25\%$ и $C > 25\%$	$ЛС > 80\%$ и $В > 80\%$	$ЛС > 95\%$



Система оценивания факультативной части

Официальное название: Введение в анализ данных: доп. главы

Кафедра дискретной математики

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 — курс обязателен.
- ▶ ПМФ, ИВТ, гр. 351, 352 — 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.
- ▶ Применяем разработанный нами ML-пайплайн, который 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.



ThetaGrader

ThetaGrader — система автоматической проверки домашних заданий с помощью технологий искусственного интеллекта, разрабатываемая командой ThetaHat.

Она будет помогать проверять ваши домашки!



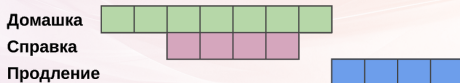
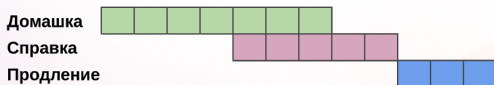
Переносы дедлайнов по уважительным причинам

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



Анна Бурханова

tg: @Anches_here

- ▶ 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
5	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	
6	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
Кафедра АД	Курс ШАД	Методы прикладной статистики
7	Байесовский подход в анализе данных	Курс по выбору x 2
	<i>Практика</i>	
8	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
5	Математическая статистика	Математическая статистика
	Машинное обучение	Вычислительная математика
	<i>Практика</i>	Курс по выбору
	Основы прикладной статистики	
6	Дискр. случ. процессы и временные ряды	Случайные процессы
	Глубокое обучение и его приложения	Квантовая механика, ч. 2
	Прикладная статистика и анализ данных	Теория и реализация языков программирования
	<i>Практика</i>	Курс по выбору x 2
Кафедра АД	Курс ШАД	Методы прикладной статистики
7	Байесовский подход в анализе данных	Машинное обучение
	<i>Практика</i>	Курс по выбору
8	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. Лекторы, применяющие анализ данных на практике;
4. Разбор статей на тему анализа данных;
5. Разработка DS-проектов.



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.Классика и ПМФ.КТ.

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, создавать систем хранения и обработки больших данных и др.. Независима по отношению к обучению в МФТИ.
Подробнее: 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. Вступительный экзамен (математика, алгоритмы)
2. Техническое собеседование (математика, алгоритмы)
3. Мотивационное собеседование

Техническое собеседование засчитывается автоматом, если

1. Экзамен написан достаточно хорошо
2. Оценки ОТЛ за обе части курса Введение в АД
3. 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 *экзамен и техническое собеседование засчитываются*

Кафедра по своему усмотрению может назначить собеседование.



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 Бен-Джурена

Выявление зависимостей с $\alpha = \frac{\alpha}{m}$

Тестовые. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

④. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑤. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑥. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑦. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑧. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑨. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑩. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑪. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑫. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑬. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑭. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑮. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑯. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑰. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑱. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑲. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

⑳. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉑. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉒. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉓. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉔. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉕. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉖. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉗. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉘. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉙. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉚. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉛. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉜. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉝. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉞. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㉟. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊱. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊲. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊳. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊴. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊵. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊶. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊷. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊸. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊹. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

㊺. Если $H_0: P(X_j) \sim U(0,1)$, тогда $FWER(\beta, \alpha) \leq \alpha$

Статистика, прикладной поток 11. 7-е учеб. Практическая значимость. Множественная проверка

Критерии (напоминание)

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

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

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

$Pr(t)$ $Pr(t)$ $Pr(t)$

c_α c_α c_α c_α

Статистика, прикладной поток 9. Бутстрап. Задание оценки близости. Проверка статистического гипотезы

Метод бутстрапа

Этап 2.

Процедуру генерации выборки повторить B раз: $X_b^* = (X_{b1}^*, \dots, X_{bn}^*)$, где $1 \leq b \leq B$.

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

Этап 3.

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

Например, бутстрапная оценка дисперсии имеет вид

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

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

Тогда $\hat{F}_n(x) - F(x) \leq \hat{F}_n(U_{(1)} - \alpha) - F(U_{(1)} - \alpha) = \hat{F}_n(U_{(1)} - \alpha) - F(U_{(1)} - \alpha) + F(U_{(1)} - \alpha) - F(U_{(1)} - \alpha) \leq \hat{F}_n(U_{(1)} - \alpha) - F(U_{(1)} - \alpha) + \frac{1}{N} \leq \frac{1}{N}$

Аналогично $\hat{F}_n(x) - F(x) \geq \hat{F}_n(U_{(1)} - \alpha) - F(U_{(1)} - \alpha) - \frac{1}{N}$

Итого $|\hat{F}_n(x) - F(x)| \leq \frac{1}{N}$

или $P_{boot} \leq P \leq P_{boot}$

$\alpha \rightarrow 0$

Аналогично

Итого x — эквивалентные

$|\hat{F}_n(x) - F(x)| \leq \frac{1}{N}$

$N \rightarrow \infty$



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 Δ 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  $\leftarrow$   $\{([x, y], 0) \mid y \in X\}$ 
  for all  $y \in$  knn do
     $d_{x,y} \leftarrow$  max $\{0, \text{dist}(x, y) - \rho\}/\sigma$ 
    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  $\sum_{i=1}^n \exp(-(\text{knn-dists}_i - \rho)/\sigma) = \log_2(n)$ 
  return  $\sigma$ 
  
```

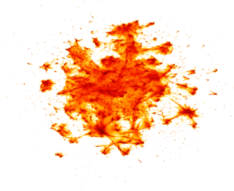


Figure 8: Visualization of the full 3 million word vectors from the GoogleNews dataset as embedded by UMAP.

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 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:

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

We want to differentiate and optimize the lower bound $\mathcal{L}(\theta, \phi; \mathbf{x}$ parameters ϕ and generative parameters θ . 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 a 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 (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

$$\mathbb{E}_{p(\zeta)} [f(h_{\phi}(\zeta))] = \mathbb{E}_{q_{\phi}(\mathbf{z})} [f(\mathbf{z})] \quad (20)$$

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

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

(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 ϕ , therefore the estimator can be differentiated w.r.t. ϕ .

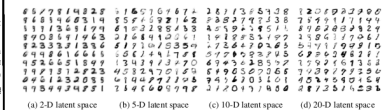


Figure 5: Random samples from learned generative models of MNIST for different dimensionalities of latent space.

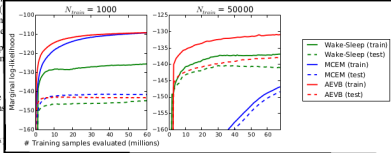


Figure 6: Marginal log-likelihood vs. # Training samples evaluated (millions) for $N_{\text{train}} = 1000$ and $N_{\text{train}} = 50000$. The plot compares four methods: Wake-Sleep (train) (green solid), Wake-Sleep (test) (green dashed), MCEM (train) (blue solid), and AEVB (train) (red solid). AEVB (train) consistently achieves the highest marginal log-likelihood, followed by MCEM (train). Wake-Sleep methods are significantly lower.



Посмотрим на примере 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 torchvision
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 * 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 * 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=torch.float)
87             z_scale = torch.ones(x.shape[0], self.z_dim, dtype=torch.float)
88             # sample from prior (z will be sampled by guide when compo
91
92     def 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 code z", 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 accumulator
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 diagnostic
158     normalizer_train = len(tr
159     total_epoch_loss_train = e
160     train_elbo.append(total_los
161     print("epoch %d" % epoch)
162
163 if __name__ == '__main__':
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.s
173
174     # pick three rank
175     # visualize how
176     if i == 0:
177         if args.visdom
178             plot_vae.i
179             recd_inde
180             for index
181                 test_
182                 rec_1
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
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```



Какие вообще инструменты могут потребоваться



PyTorch



Yandex
CatBoost



Numba



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 считать:

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

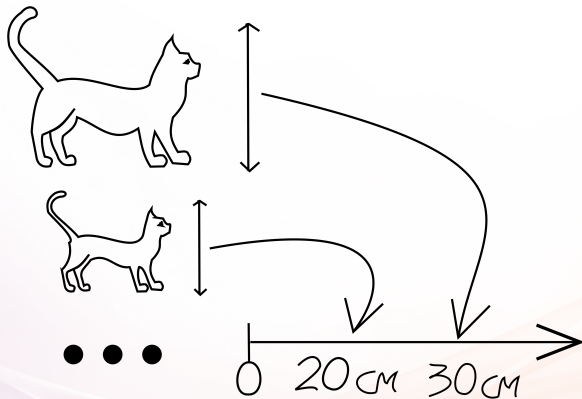


Мурмурландия

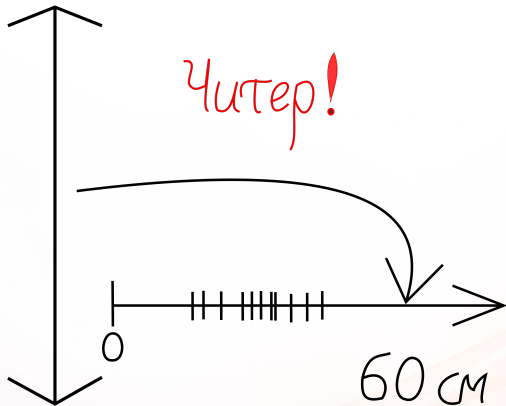




Каков средний рост котиков?



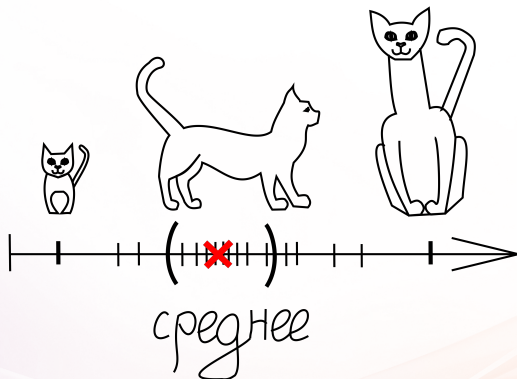
Точечное оценивание



Выбросы



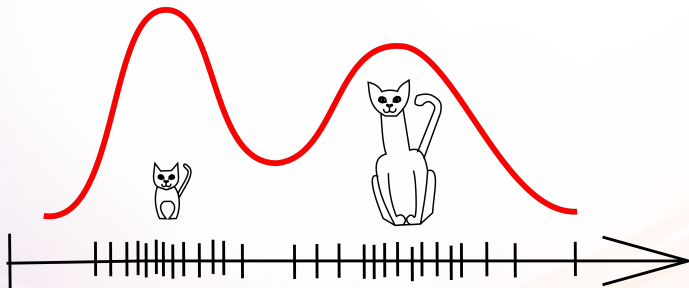
Среднее определяется неточно



Интервальное оценивание



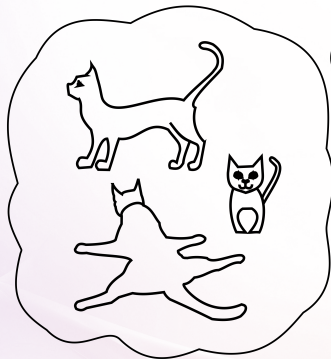
Характер распределения



Непараметрическое оценивание



низкие



высокие

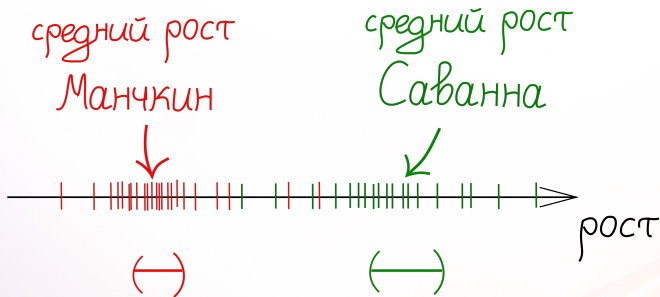




Отличается ли их средний рост?



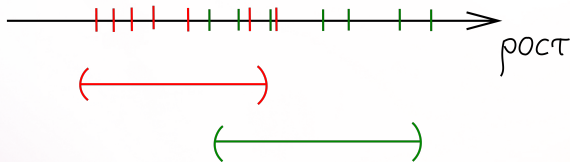
Собираем данные



отличается

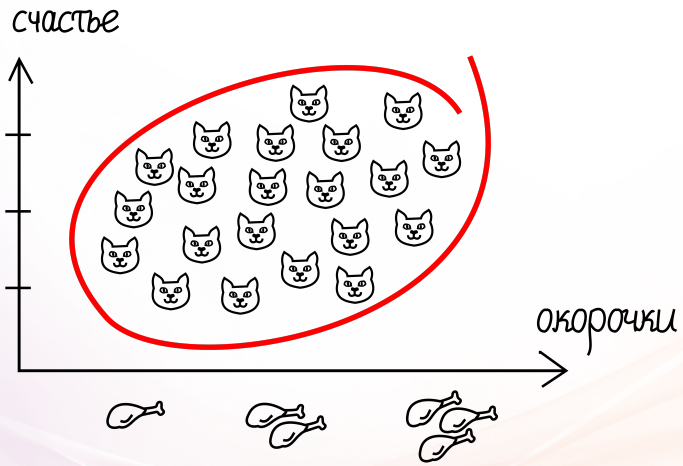


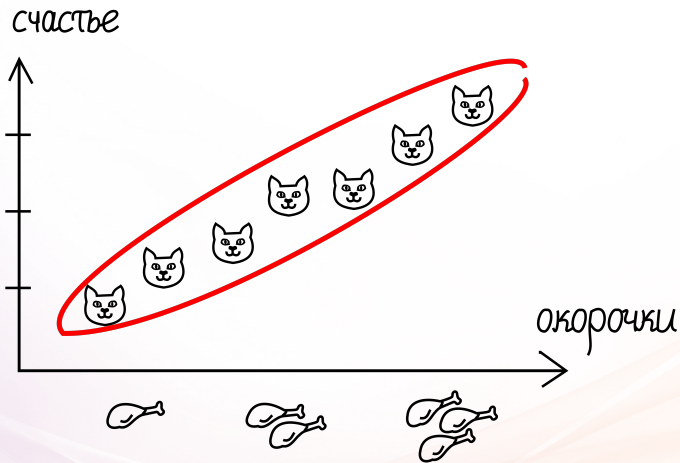
Если данных мало

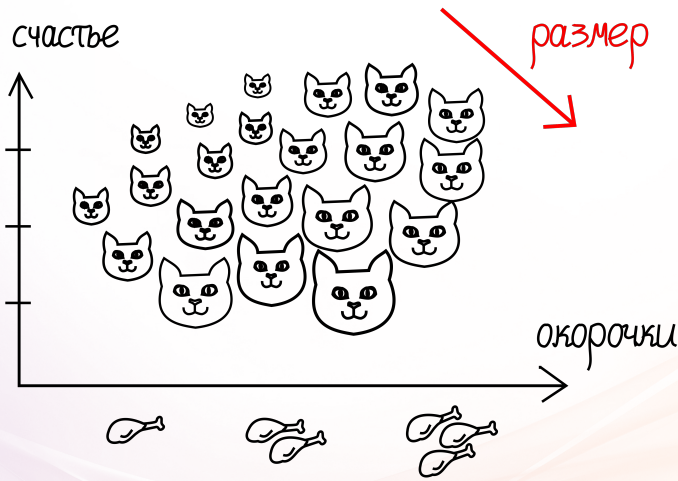


непонятно

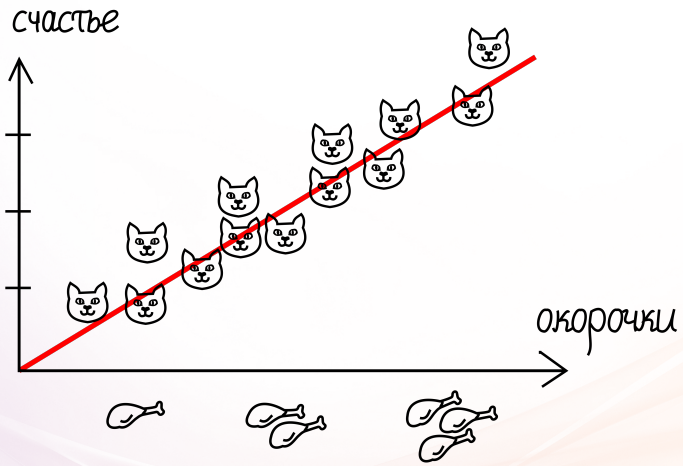
Статистические гипотезы, АВ-тесты







Корреляционный анализ





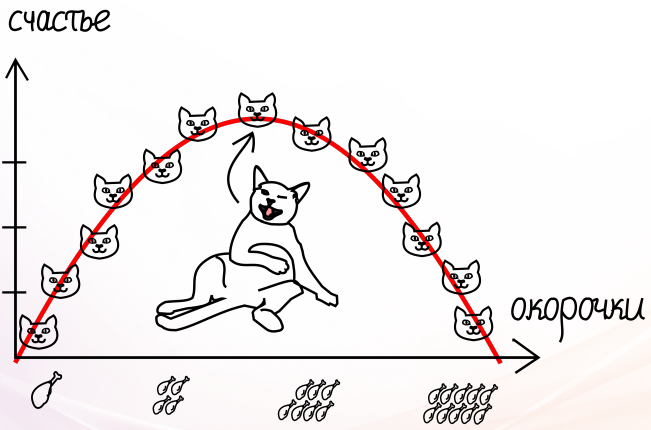
Формула счастья



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

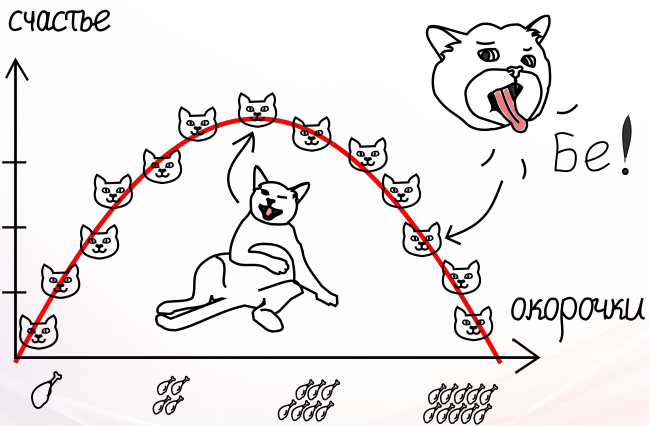


Больше окорочков





Больше окорочков





Формула счастья



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



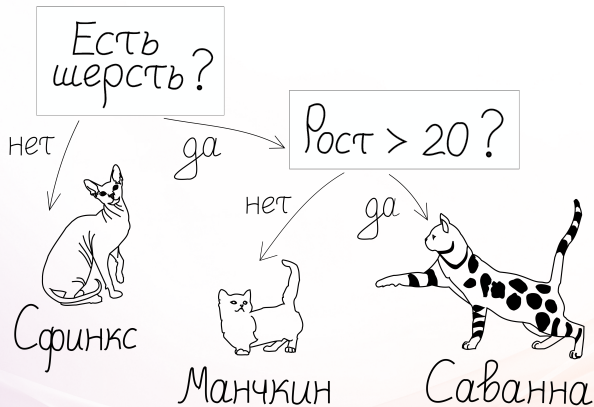
Другие факторы

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

Регрессионный анализ







Классификация котиков



Классификация

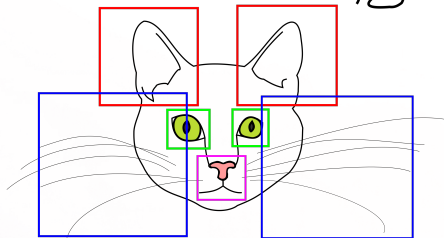


Собираем данные

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



Распознавание мордочек



Нейронные сети



Книга с похожим содержанием





Попробуем решить задачу



А ты кто?

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

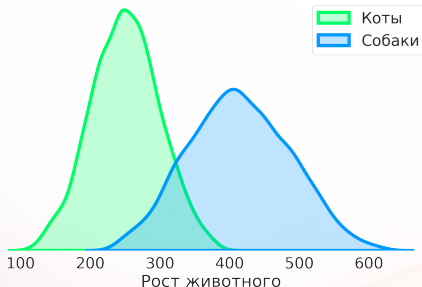




Классификация: собака vs кот

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.

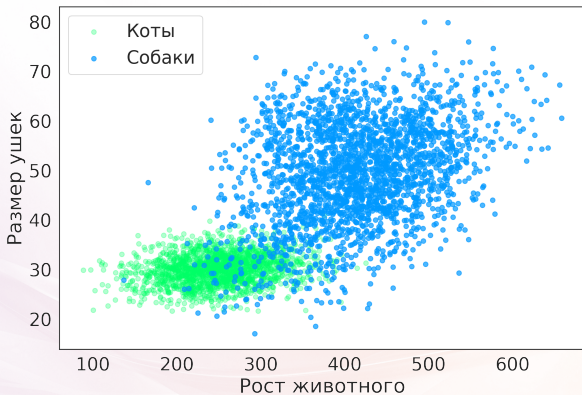
Но есть большое пересечение, это не очень здорово.



Классификация: собака vs кот

Извлечем еще один признак — размер ушек.

Теперь классы лучше разделяются.

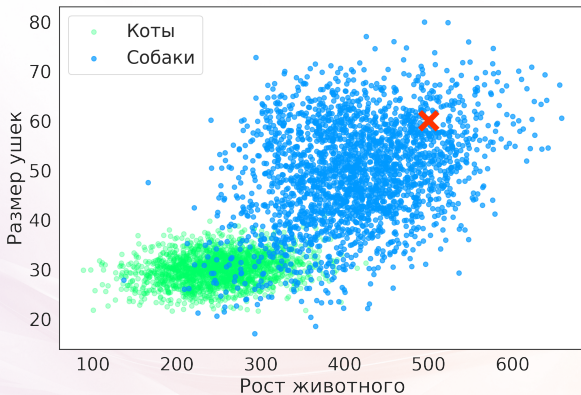




Классификация: собака vs кот

Попробуем классифицировать новое животное.

Кто отмечен красным?

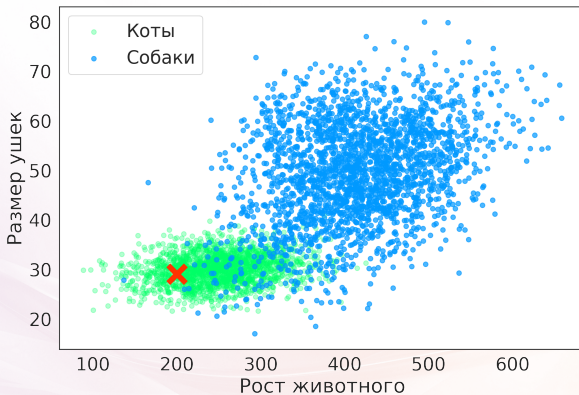




Классификация: собака vs кот

Попробуем классифицировать новое животное.

Кто отмечен красным?

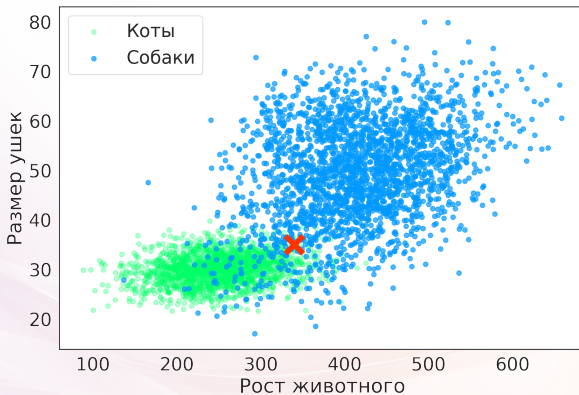




Классификация: собака vs кот

Попробуем классифицировать новое животное.

Кто отмечен красным?

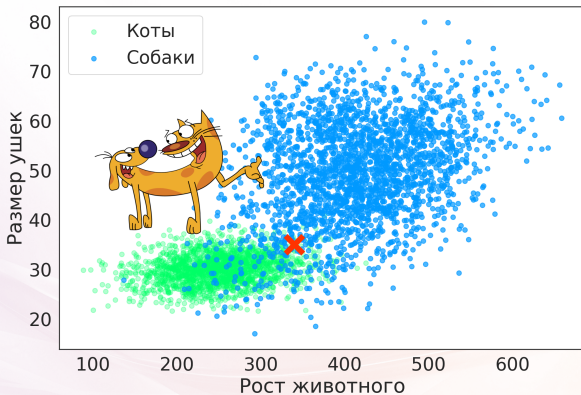




Классификация: собака vs кот

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:

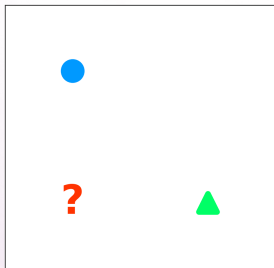
- ▶ $w_j = 1 - j/k$ — зависящий от номера соседа;
- ▶ $w_j = \|x - x_{(j)}\|^{-1}$ — зависящий от расстояния до соседа.

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



Особенности

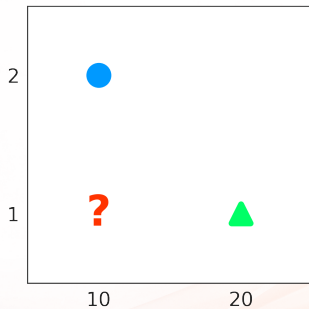
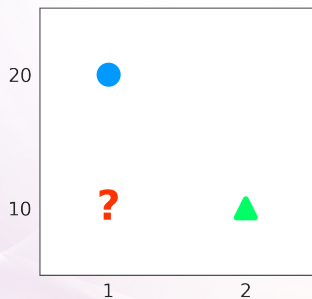
Классифицируйте объект "?".





Особенности

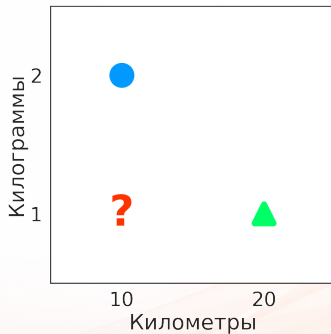
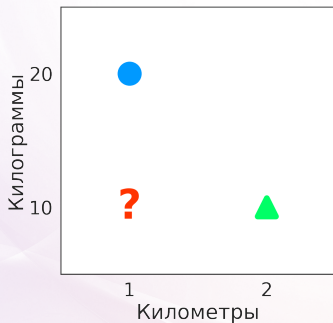
Классифицируйте объект "?".





Особенности

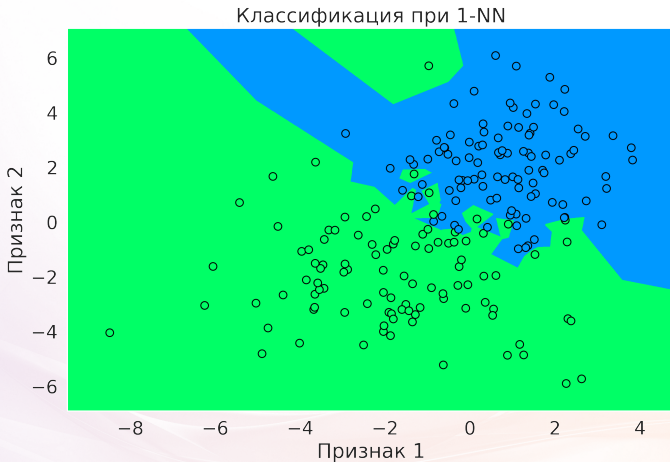
Классифицируйте объект "?".



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

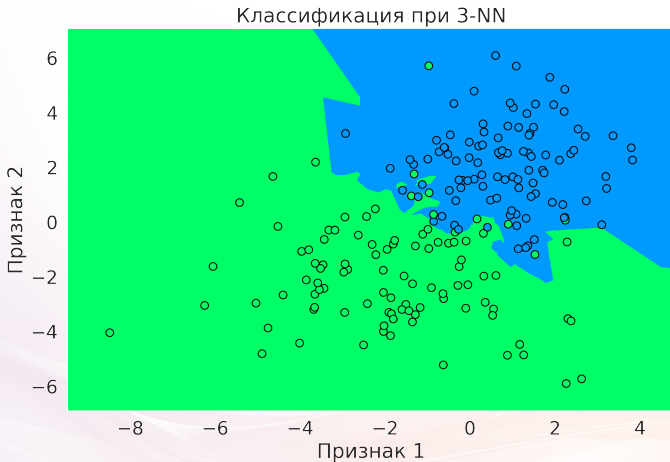


Что происходит при разных k ?



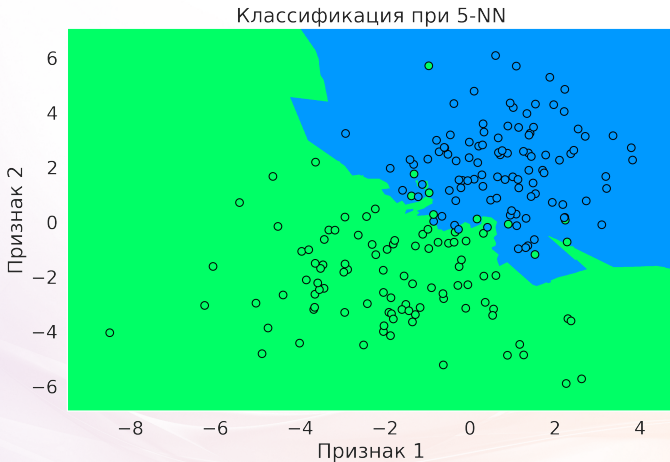


Что происходит при разных k ?



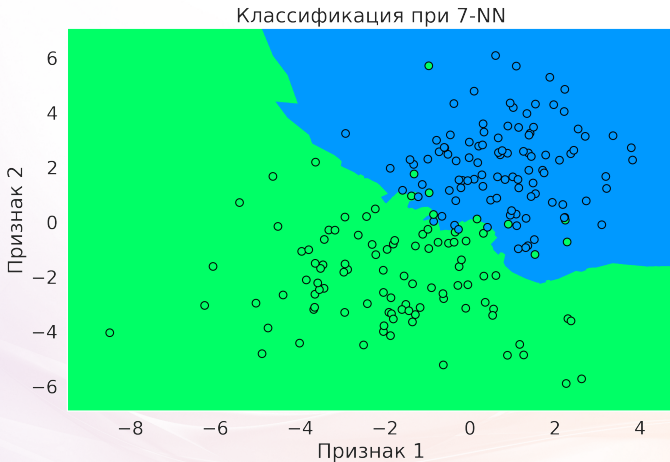


Что происходит при разных k ?



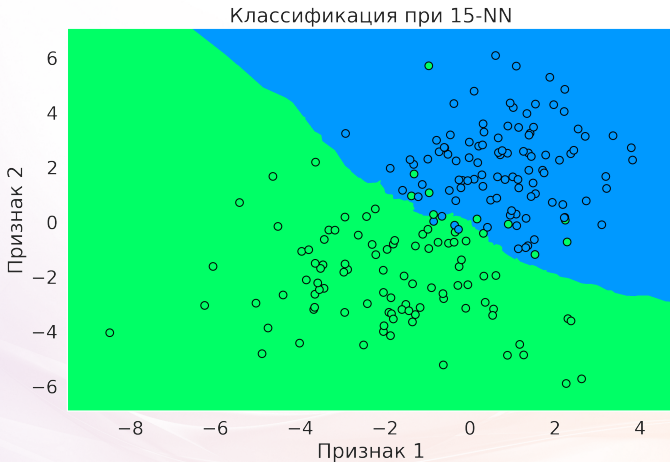


Что происходит при разных k ?



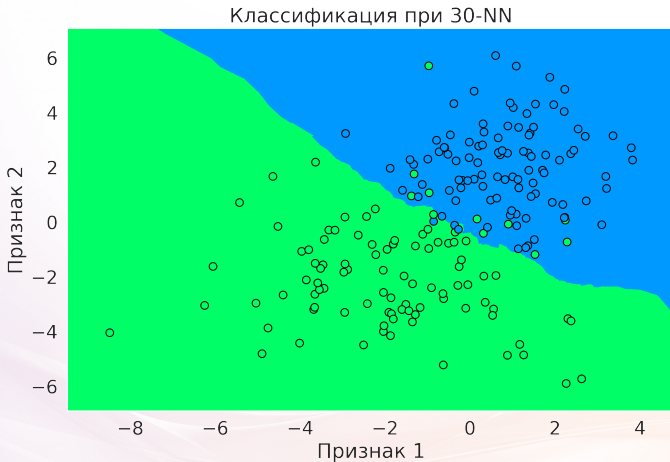


Что происходит при разных k ?





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 100×100 состоит из 10^4 пикселей,

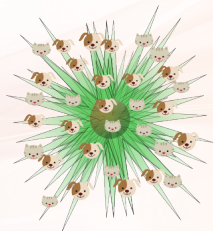
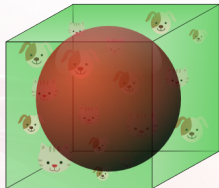
в каждом по 3 числа. Какой размерности получается объект?

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

Проблема:

в пр-ве больших размерностей расстояния неинформативны.

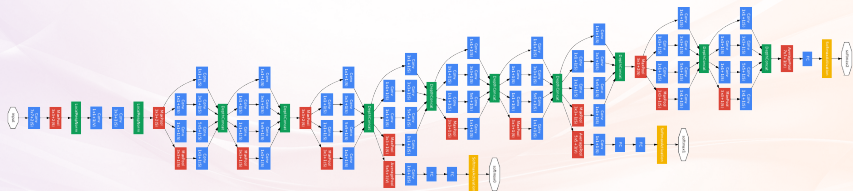
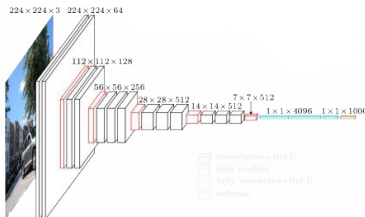
Например, среди фиксированного количества случайных точек в единичном кубе в пространстве большой размерности почти все точки будут лежать около границы куба.





А что в сложных случаях?

Нейросети! Но об этом позже :)





Примеры прикладных задач

АБ-тестирование

Распознавание лиц

Генерация изображений

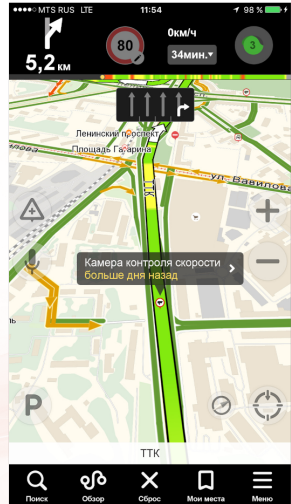
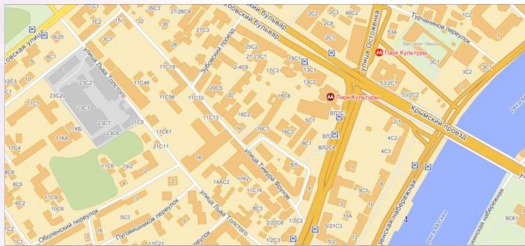
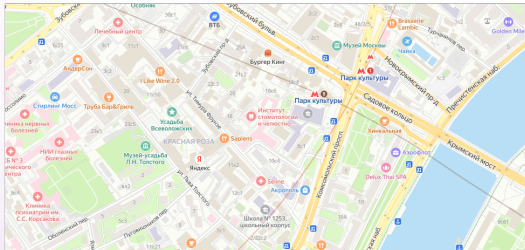
Синтез речи

Обучение с подкреплением

Где еще?



Какой стиль карт удобнее пользователям?





Какая концепция магазина приносит больше выручки?





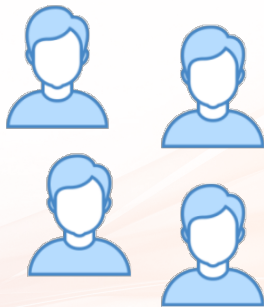
Контрольная группа Группа А

Пользователи видят
прежнюю версию сервиса



Тестовая группа Группа В

Пользователи видят
новую версию сервиса





Классический способ проверки

Пусть X_1, \dots, X_n и Y_1, \dots, Y_m — значения целевой метрики (выручка, клики, рейтинг и т.д.) для контрольной и тестовой групп.

Постановка задачи:

$H_0: EX_1 = EY_1$ vs. $H_1: EX_1 \{<, \neq, >\} EY_1$

Справедлива сходимость

$$T(X, Y) = \frac{\bar{X} - \bar{Y}}{\sqrt{S_X^2/n + S_Y^2/m}} \xrightarrow{d_0} \mathcal{N}(0, 1).$$

Статистический критерий $S = \{|T(X, Y)| > z_{1-\alpha/2}\}$

Доверительный интервал для $EX_1 - EY_1$ уровня доверия $1 - \alpha$

$$\left(\bar{X} - \bar{Y} \pm z_{1-\alpha/2} \sqrt{S_X^2/n + S_Y^2/m} \right).$$



При наличии дополнительных данных

Введем обозначения

$$Y = \begin{pmatrix} Y_{11} \\ \dots \\ Y_{n1} \\ Y_{12} \\ \dots \\ Y_{m2} \end{pmatrix}, \quad X = \begin{pmatrix} 1 & 0 \\ \dots & \dots \\ 1 & 0 \\ 1 & 1 \\ \dots & \dots \\ 1 & 1 \end{pmatrix}, \quad \theta = \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \dots \\ \varepsilon_{n+m} \end{pmatrix}.$$

Линейная регрессия предполагает зависимость $Y = X\theta + \varepsilon$.

Тогда $Y_{i1} = \theta_0 + \varepsilon_i$ и $Y_{i2} = \theta_0 + \theta_1 + \varepsilon_{n+i}$, следовательно

- ▶ θ_0 — среднее в группе А,
- ▶ θ_1 — эффект от эксперимента.

Вывод: для проверки АВ-теста нужно построить интервал для θ_1 и проверить гипотезу $H_0: \theta_1 = 0$ критерием Стьюдента.

Важно использовать оценку дисперсии, устойчивую к гетероскедастичности, т.к. группы могут иметь разную дисперсию.



Примеры прикладных задач

АБ-тестирование

Распознавание лиц

Генерация изображений

Синтез речи

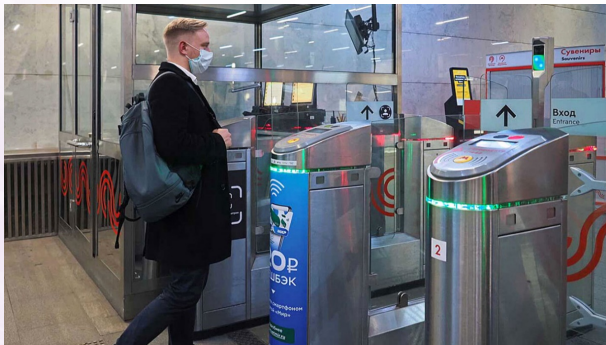
Обучение с подкреплением

Где еще?



Распознавание лиц

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





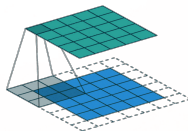
Распознавание лиц

Пример модели детекции лиц: TinaFace.

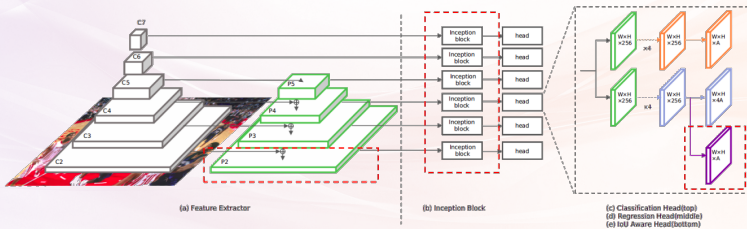
- ▶ Основу модели составляют сверточные слои.

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:

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



Примеры прикладных задач

АБ-тестирование

Распознавание лиц

Генерация изображений

Синтез речи

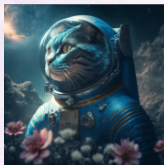
Обучение с подкреплением

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 $ELBO$:

$$ELBO = E_q \log \frac{p_\theta(X_0, \dots, X_T)}{q(X_1, \dots, X_T | X_0)}$$
$$\simeq const - \sum_{t=2}^T \frac{\tilde{\alpha}_{t-1} \beta_t^2}{2\tilde{\beta}_t(1 - \alpha_t)^2} E_q \|X_0 - x_\theta(X_t, t)\|^2$$

 x_0 x_1 x_2 x_3 x_4

...

 x_T 

Обратный диффузионный процесс



Примеры прикладных задач

АБ-тестирование

Распознавание лиц

Генерация изображений

Синтез речи

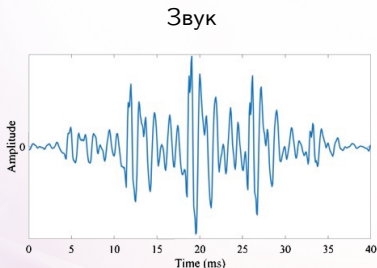
Обучение с подкреплением

Где еще?



Синтез речи

Что такое звук?



Звук — композиция волн с разными амплитудами и частотой.

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



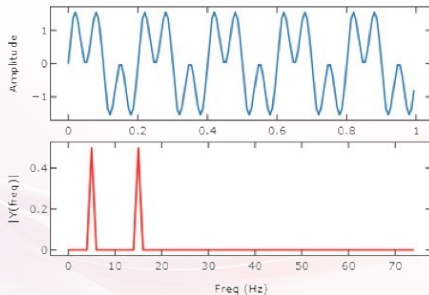
Синтез речи

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



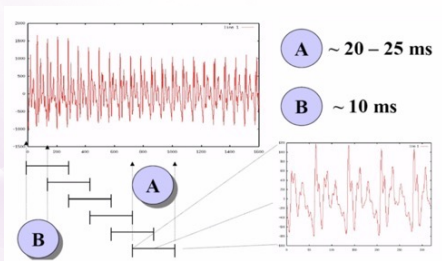
Синтез речи

Как работать со звуком?

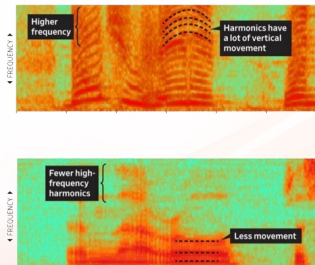
У результата преобразования Фурье ко всем данным нет времени.

Посчитаем распределение на окнах из звука.

Окна



Спектрограмма

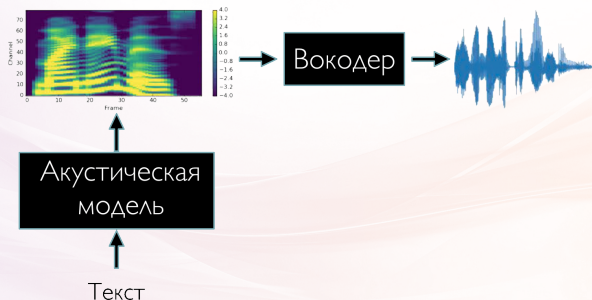


- ▶ Берем маленькие окна (20-25 мс) от исходных данных.
- ▶ К каждому окну применяем преобразования Фурье.
- ▶ Стакаем распределения частот вместе, получаем спектрограмму.



Синтез речи

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

АБ-тестирование

Распознавание лиц

Генерация изображений

Синтез речи

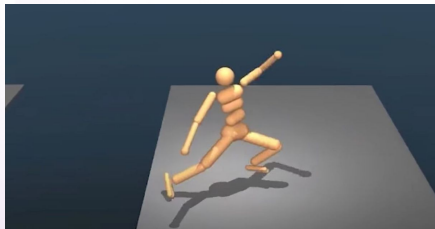
Обучение с подкреплением

Где еще?



Учимся ходить

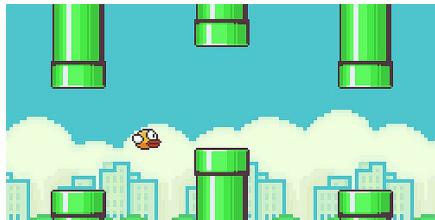
Хотим научить компьютерную симуляцию человека ходить.



Как человек учится ходить?
Делаем много попыток,
постоянно улучшая свои навыки.

Играем в игры

Хотим научить компьютер играть
в игру Flappy Bird.



Как бы мы сами учились играть?
Играем много раз, в каждый из
которых улучшаем свои навыки игры.

Похожим образом устроен **Reinforcement learning**.



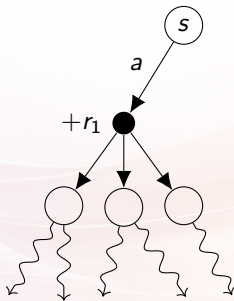
Важные определения RL

Стратегия π — способ выбора **действия** a в **среде** (стиль игры).

Награда $R(s, a)$ — величина поощрения в действии a и состоянии s .

Q-функция — **ожидаемый** итоговый выигрыш в будущем при условии **состоянии** s и совершения **действия** a при следовании стратегии π .

$$Q_{\pi}(s, a) = E_{\pi, P, R}(R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots \mid S_t = s, A_t = a)$$





Свойства Q-функции

Теорема (об оптимальности). Существует π^* такая, что она даёт выигрыш не хуже любой другой стратегии π :

$$Q_{\pi^*}(s, a) = \max_{\pi} Q_{\pi}(s, a) = Q_{opt}(s, a)$$

Теорема (уравнение Беллмана). В произвольной среде и для произвольной стратегии π выполнено соотношение

$$Q_{\pi}(s, a) = ER(s, a) + \gamma \sum_{s' \in \mathcal{S}, a' \in \mathcal{A}} P(s'|s, a) \pi(a'|s') Q_{\pi}(s', a').$$

Уравнение показывает, как Q-функция обновляется: текущая ценность действия складывается из непосредственной награды и лучших возможных будущих действий.



Q-learning

Если знаем оптимальную Q-функцию, можем определить **лучшую стратегию действий**. Воспользуемся итерационной процедурой:

$$Q(s_t, a_t) := Q(s_t, a_t) + \alpha \left(r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right)$$

Этот процесс называется **Q-learning**:

- ▶ Агент **совершает действие** и получает **награду**.
- ▶ **Обновляет** Q-значение, учитывая **лучшие будущие действия**.
- ▶ Со временем стратегия становится всё более оптимальной.

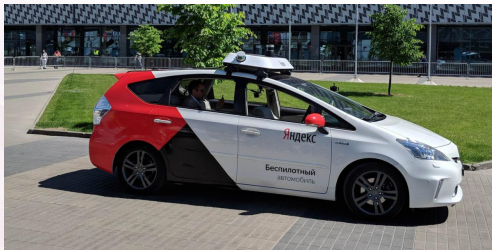
Пример: Учимся играть в новую игру. Сначала ты пробуем случайные действия, но постепенно запоминаем, какие из них приводят к победе.



Q-learning

Теорема. В случае конечного MDP в алгоритме Q-learning $Q(s, a)$ сходится к $Q_{opt}(s, a)$ с вероятностью 1, если

1. $\forall s \in \mathcal{S}, a \in \mathcal{A} \mu(a|s) \neq 0$, μ — стратегия для сессий в обучении
2. $\sum_{k=1}^{\infty} \alpha_k = \infty$, где α_k — степень **exploration** среды на шаге k
3. $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$.





Примеры прикладных задач

АБ-тестирование

Распознавание лиц

Генерация изображений

Синтез речи

Обучение с подкреплением

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



Слово студентам DS-потока



Вероника Прохорова

3 курс, DS-поток

tg: @roorpoop2



ВСЁ!