

Лекция 1

Введение в статистику и ML

Курс: Введение в DS на УБ и МиРА (весна, 2022)

Преподаватель: Владимир Омелюсик

4 апреля 2022 г.

- Всевышкинский экзамен по анализу данных.
- Три уровня: начальный, **базовый**, продвинутый.
- На базовом уровне нужно знать основы статистики (проверка гипотез, линейные модели), машинного обучения (постановка задач, простые модели) и анализа данных (визуализация, обработка).

- 10 лекций (онлайн), 14 семинаров (онлайн / очно).
- На лекциях – теория, на семинарах – практика на компьютерах. Программировать будем только на Python.
- Лекции – в Тимс, с записью. Записи будут в Тимс и на YT (ссылку сообщим). Конспекты лекций будут (где – чуть позже). Лекции где-то парные, где-то одинарные.
- Семинары – на усмотрение семинаристов.
- Семинары рассинхронизированы, и мы будем это учитывать.

$$\text{Итог} = 0.5 \times \text{ДЗ} + 0.2 \times \text{Квизы} + 0.3 \times \text{Экзамен}$$

- 5 домашних заданий на компьютере, оценка за ДЗ – это средняя оценка по ним.
- Квизы проводятся на лекциях в Google-формах, предупреждаем заранее (первый квиз будет на лекции 11 апреля).
- Письменный экзамен в конце курса, не блокирующий. Автоматов не предусмотрено.
- Любая форма контроля на усмотрение преподавателя может быть представлена к устной защите.
- Округление арифметическое, округляется только итог.

- [Страничка на Wiki](#) (здесь агрегируются вся информация и материалы).
- [GitHub с материалами](#) (сюда ведут ссылки со странички на Wiki).
- [Телеграм-канал](#) (тут выкладываются важные объявления).
- Телеграм-чаты групп (ссылки выложим в канал).
- [Анонимная Google-форма](#) (сюда можно писать текущий фидбэк).

И ещё две ссылки, чтобы всё было в одном месте:

- [Ссылка на Тимс.](#)
- [Ссылка на плейлист с записями лекций.](#)

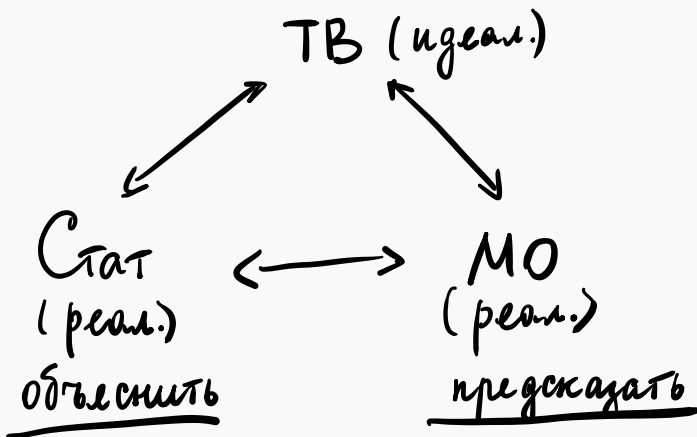
План лекций

1. (4 апреля) Введение в статистику и ML. Выборка.
2. (4 апреля) Теория вероятностей. Распределения. Описательные статистики.
3. (11 апреля) Ковариация и корреляция. Кейсы по визуализации.
4. (18 апреля) Тестирование гипотез (Z- и t-тесты для одной и двух выборок). p-value.
5. (25 апреля) Непараметрика. χ^2 -критерий согласия.
6. (25 апреля) Линейная регрессия (оценка и метрики).
7. (16 мая) Линейная регрессия (интерпретация и гипотезы).
8. (16 мая) Задачи машинного обучения. Кросс-валидация. Градиентный спуск.
9. (30 мая) Задача классификации. Метод k ближайших соседей.
10. (6 июня) Логистическая регрессия.

- Выглядит технично, но это так кажется.
- Будем максимально использовать математику 9-10 класса.
- Много практических кейсов, чтобы было интереснее.
- Running gag про то, что машинное обучение – это `.fit()` и `.predict()`.
- [Understanding Machine Learning through Memes](#).

Теория вероятностей, статистика и машинное обучение

- Теория вероятностей: идеальный мир, кубики и урны.
- Статистика: цель – построить модель, которая объясняет мир вокруг нас.
- Машинное обучение: цель – построить модель, которая хорошо предсказывает какую-то величину.



Теория вероятностей (кубики и урны)

Играл. кубик, 6 граней, все грань равнов-ог

$$P(\text{грань } \cdot) = \frac{1}{6}$$

$$P(\text{грань } \cdot \cdot) = \frac{1}{6}$$

\vdots

$$P(\text{грань } \ddots) = \frac{1}{6}$$

? А если 30 раз, то
сколько в среднем

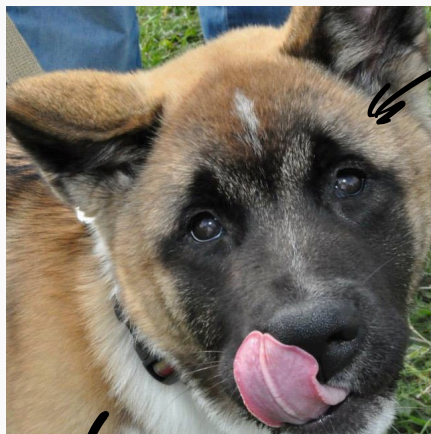
? А если 2 раза подряд,
то \ddots , \ddots

? ...

Статистика: пример с ростом собак

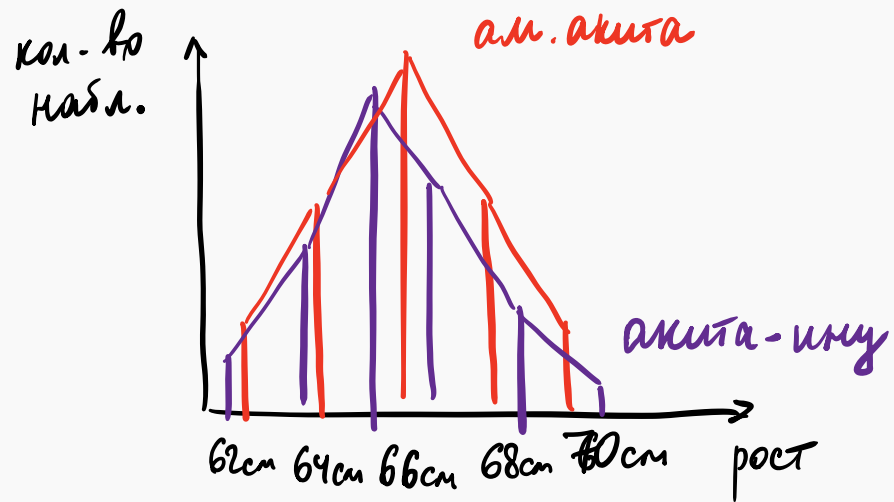
- Выборка из 1000 независимых измерений роста собак породы американская акита. Средний рост – 66 см
- Выборка из 1000 независимых измерений роста собак породы акита-ину. Средний рост – 65 см

акита-
ину

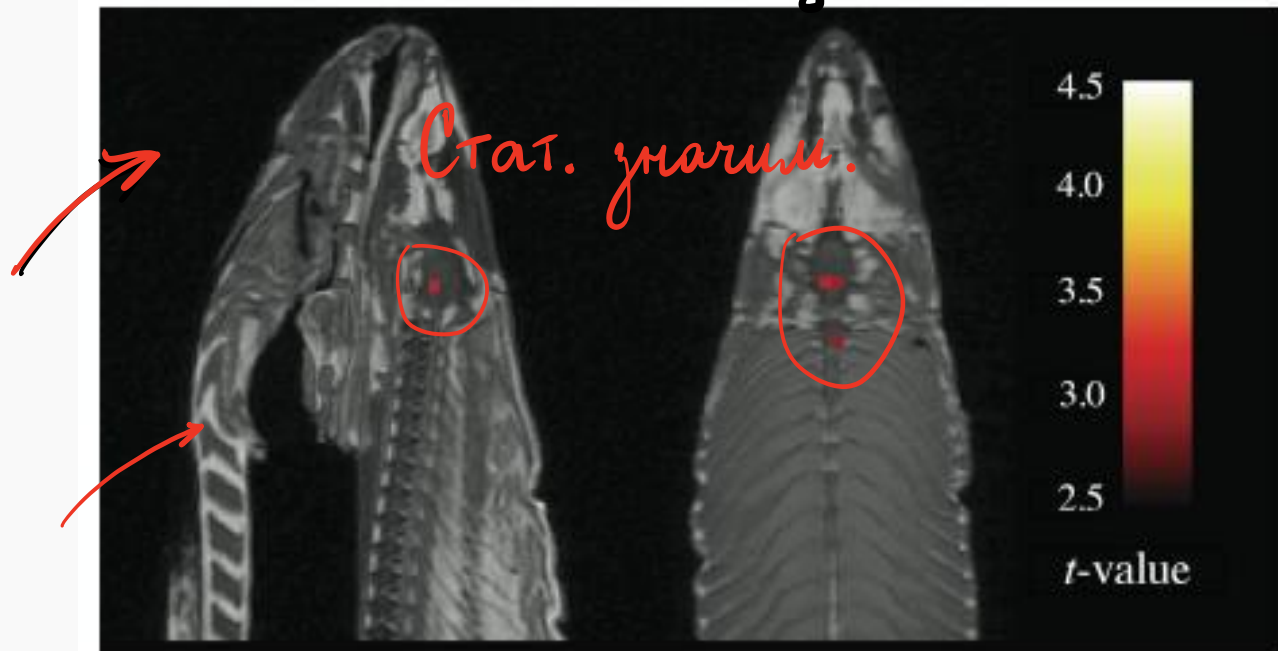


амер.
акита

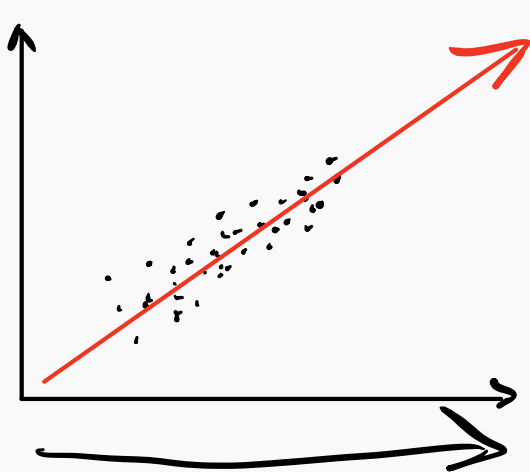
Статистика: пример с ростом собак



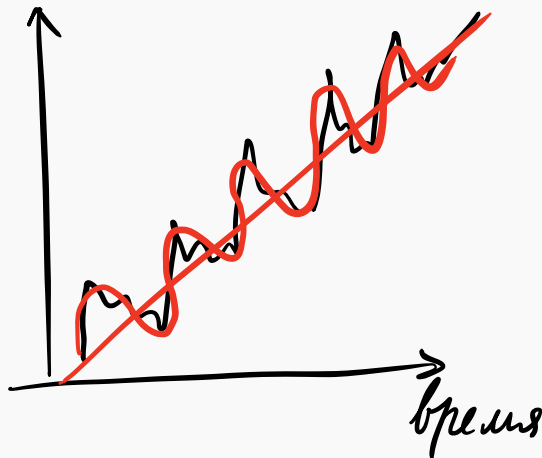
- Почитать Dead Salmon Study (Bennet et al.)



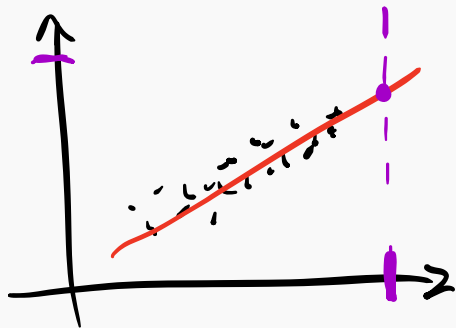
Статистика: пример с регрессией и временными рядами



ВВП



Машинное обучение: пример со стоимостью квартиры

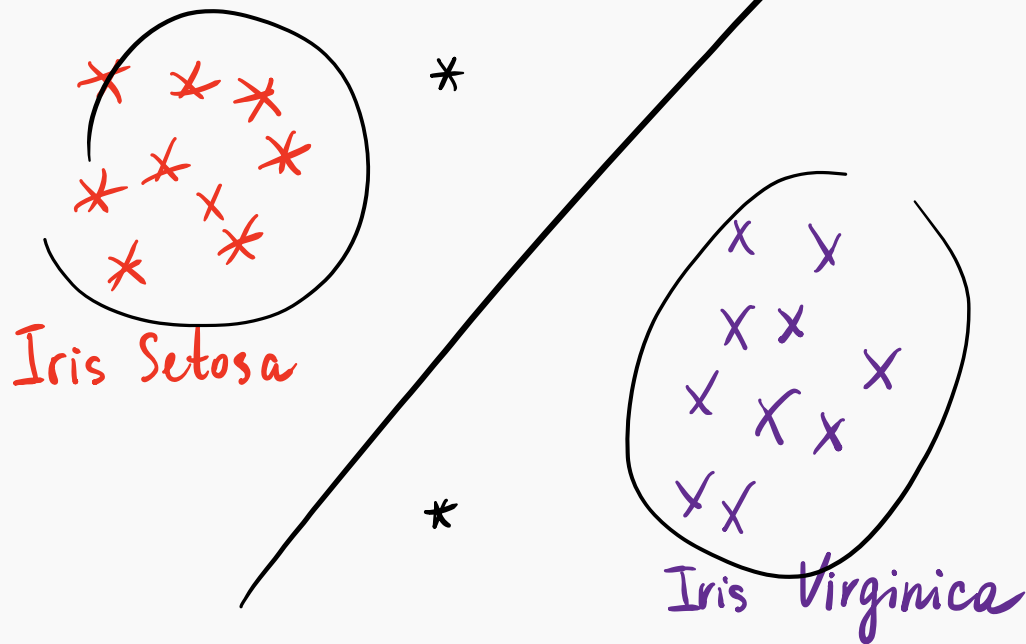


Набор данных:

стоим. кв. (у.е.)	S, m^2	расст. до метро
10	15	.
20	45	.
70	100	.
15	.	.
...

1000 набл.

Машинное обучение: пример с классификацией цветков ириса



ТВ
(инструм.)
? ↗ ↘
СТАТ \Leftrightarrow МО
(объяс.) (предск.)

Генеральная совокупность

Совокупность всех мыслимых объектов, относительно которых предстоит делать выводы в исследуемой задаче.

Выборка

Часть генеральной совокупности, используемая в исследовании.

$$X = [3, 4, 5, 1, 7, 8, \dots, 9, 10]$$

$$Y = [1000, 1500, 700, 900] \text{ "вектор"}$$

$$Z =$$

\equiv
"списки
чисел"

$$A =$$

$$B =$$

Глубоков. калымаро - Ген. Сов.

L, ① К. в Фих. ок.

② К. с губ. вулк.

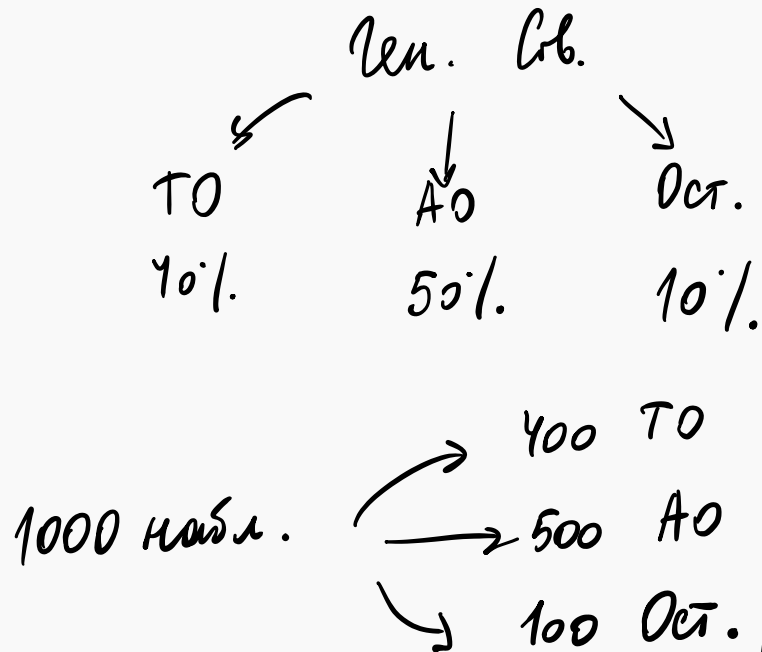
Все ли выборки хороши?

Репрезентативность выборки

Соответствие характеристик выборки характеристикам генеральной совокупности.

- * Опред. факторы / хар-и ИС
- * Сегм. ИС
- * Служ.

Случайная выборка



Репр.
=>
случ. набл. в
соотв. с
хар. ген.
стохастич.

Типы данных

	\downarrow X	\downarrow Y	\downarrow Z
→	40	41	0
→	50	12	0
→	100	8	1

матрица "наблюдения - признаки"

$$X = [40, 50, 100]$$

	дискр.	текущ. пер.	бинар.	ном.	ранж.
	X	Y	Z	W	P
→	40	1.01	1	"red"	max.
→	50	2.12	0	"green"	max.
→	100	3.14	1	"red"	всех.

дискретн.

числовые

бинар.

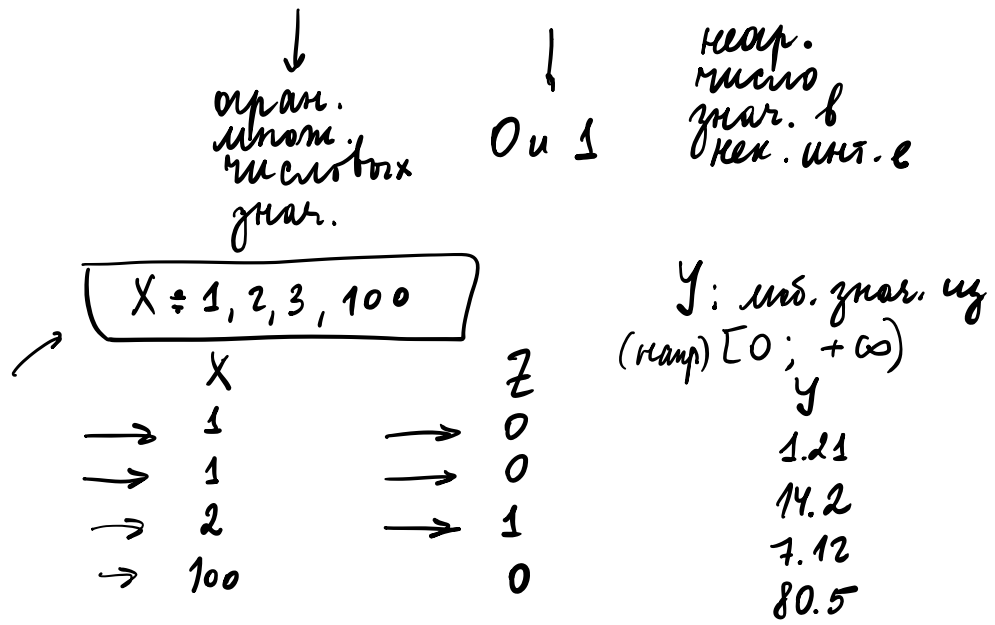
текуществ.

категор.

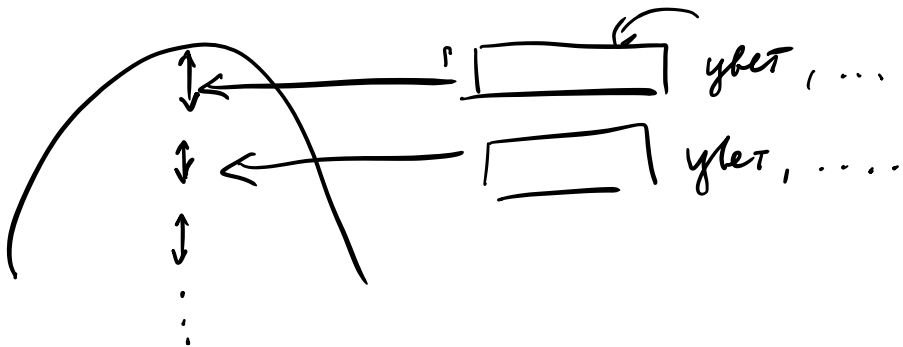
номин.

(нет порядка)

ранж. (порядок)



Полны: цвета, материалы, назв. городов



Рам: степени отр., физик. зв.

