



ІІТМО

Basic ML
Не магия наука

Что такое Машинное обучение?

Машинное обучение — это наука, изучающая алгоритмы, автоматически улучшающиеся благодаря опыту.



Какие задачи решает? Те, которые “легки” для человека, но крайне трудно запрограммировать:

- Перевод текста
- Постановка диагноза
- Ранжирование документов по поиску
- Классификация картинок

Данные, предсказания модели.



Данные бывают разные, но
в общем и целом для
моделей это:

- примеры
- предсказания

Набор примеров – датасет.

Часть датасета – выборка.

Примеры → предсказания – модель.

Какие задачи мы решаем?

Базовые задачи



$Y = R$ // $Y = R^M$ – Регрессия – пытаемся предсказать число.

$Y = 0,1$ // $Y = [1, K]$ – Классификация – предсказываем класс или классы

Задачи со звездочкой



$Y = [0,1]^K$ – Многоклассовая классификация с пересечением

Y – конечное упорядоченное множество – Ранжирование

Задачи сложнее и без учителя

ІІТМО

Сегментация.
Перевод текста.



Создание новых объектов из ничего -
диффузионки.

Кластеризация, “угадай следующее
слово” (в сыром тексте) – без учителя.

А если все соединить?

Ответ:



При обучении модели в разработке используют разные метрики:

- бизнес метрики
- онлайн метрики
- ассесоры
- оффлайн метрики

Мы не бизнес. Нам нужно понять, как мы хорошо обучили модель на текущих датасетах:

- Доли правильных ответов
- Разница между предсказанным и истинным значением
- Для ранжировки - доля пар документов, которые упорядочены неправильно

Все это можно назвать “функциями потерь”.



Начальные метрики

TP - сделали верное +предсказание
FP - сделали НЕверное +предсказание

TN - сделали верное -предсказание
FN - сделали НЕверное -предсказание

| | | Actual class | |
|-----------------|---|-----------------------|-----------------------|
| | | + | - |
| Predicted class | + | TP True Positives | FP Flase Positives |
| | - | FN Flase Negatives | TN True Negatives |

Серьезные метрики

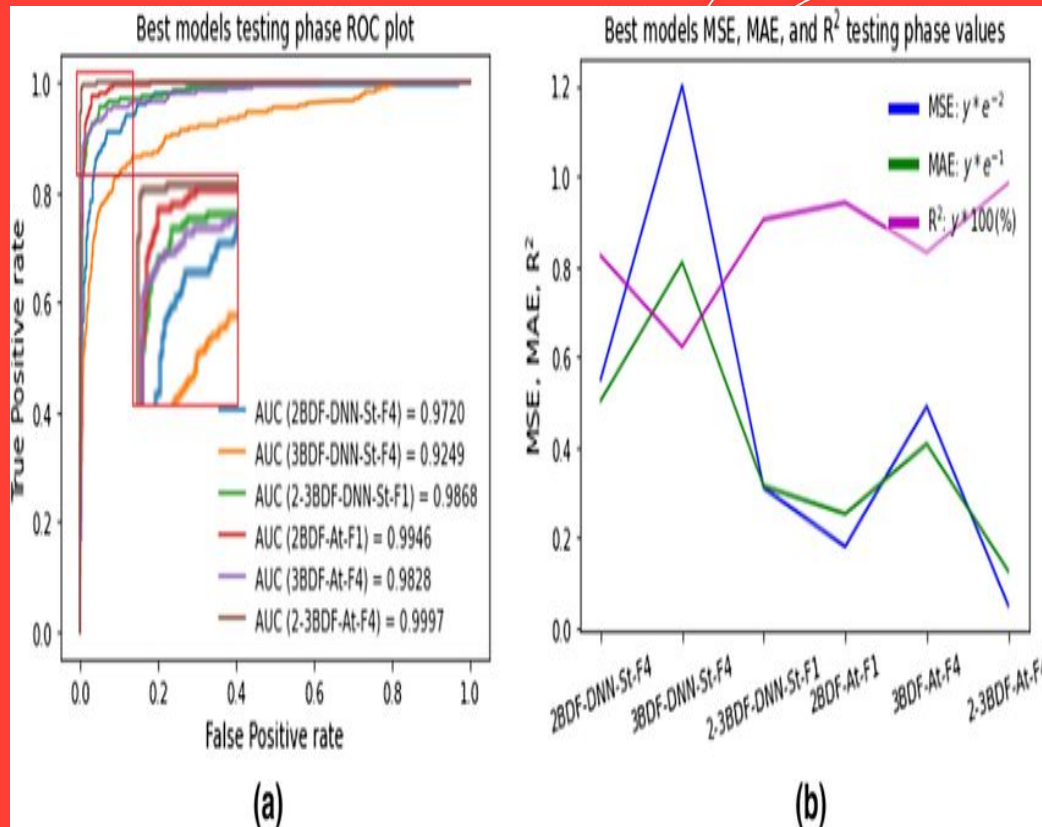
$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

$$\text{Precision} = \frac{TP}{TP+FP}$$

$$\text{Recall} = \frac{TP}{TP+FN}$$

$$F1 = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

“Крутые” метрики



$$MAE = (1/n) * \sum |y[i] - y'[i]|$$

$$MSE = (1/n) * \sum (y[i] - y'[i])^2$$

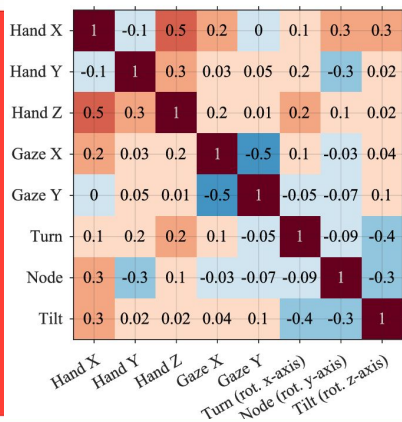
$$RMSE = \sqrt{MSE}$$

$$MAPE = (100\% / n) * \sum | (y[i] - y'[i]) / y[i] |$$

$$LOG LOSS = - (1/n) * \sum [y[i] * \log(p[i]) + (1 - y[i]) * \log(1 - p[i])]$$

Преобразование данных

| PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|-------------|----------|--------|---|--------|------|-------|-------|------------------|---------|-------|----------|
| 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |



1. Машина не любит тип string.
2. Машина не любит разброс в датасете.
3. Машина не любит СИЛЬНО коррелирующие данные
4. Машина не любит слабо коррелирующие данные
5. Машина не любит выбросы
6. Машина не любит пропуски
7. Машина любит, когда данные находятся в интервале $[-1;1]$
8. Машина любит, когда данных МНОГО (очень. НУ ОЧЕНЬ МНОГО)
9. Машина любит, когда вы сделали датасет информативным, каждый признак несет в себе суть для таргета.

Так я модель сделаю уже?

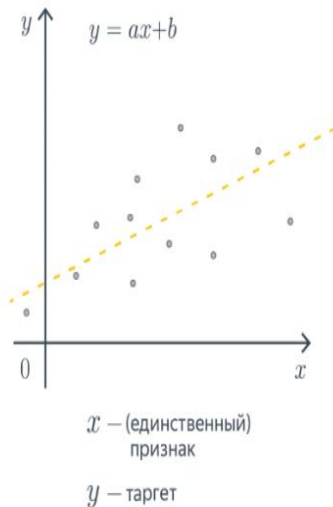
ІІТМО

У линейных моделей достаточно понятная и “прямая” формула:
 $y = w_1 \cdot x_1 + \dots w_i \cdot x_i + w_0 == \langle x, w \rangle + w_0$

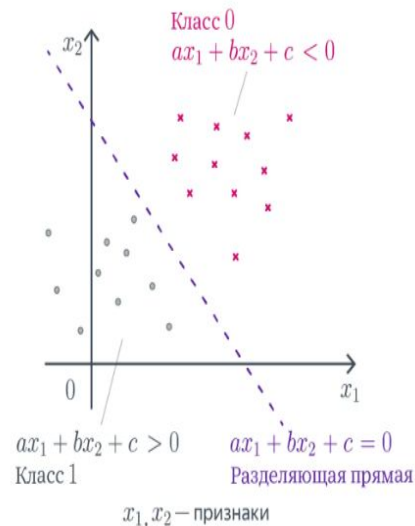
В задачах регрессии нам важно расположить “прямую” настолько близко к точкам, насколько возможно.

В задачах классификации нам надо разделить два класса, допуская выбросы.

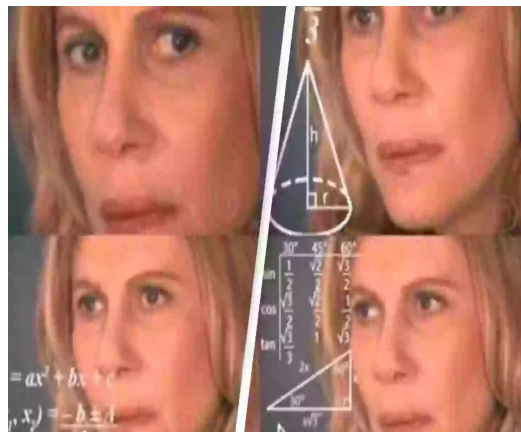
Регрессия



Классификация



Oh, no, it's a math!



Как нам подкрутить веса, чтобы все завелось?

$Xw = w_1 * x_1 + \dots w_i * x_i$, w - веса, x - признаки

Пусть $y = y_{pred} + y_{ort}$, т.е.

$y_{ort} = y - Xw \perp x_1, \dots, x_i \rightarrow$ в матрички $\rightarrow X^T * (y - Xw) = 0 \rightarrow$

$$\rightarrow w = (X^T * X)^{-1} X^T * y$$

Проблемы:

1. Обращение матрицы $1e20 \times 1e20$
2. $X^T X$ в 99.9% обратима, но чем больше признаков, тем чаще мы встречаем линейные зависимости. Малые “возмущения” таргета вызовут большие изменения w .

Немного схитрим

ІТМО

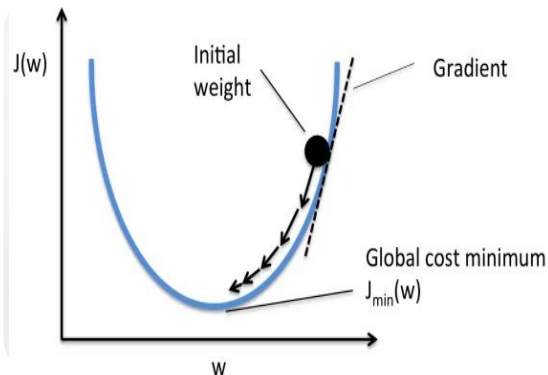


Посчитаем, как сильно мы ошибаемся и попробуем это уменьшать?

$y_{\text{pred}} = Xw$
 $\text{err} = y_{\text{pred}} - y_{\text{true}}$
 $\text{changes} = (2 * X^T * \text{err}) / N$
 $w \leftarrow w + \alpha * \text{changes}$

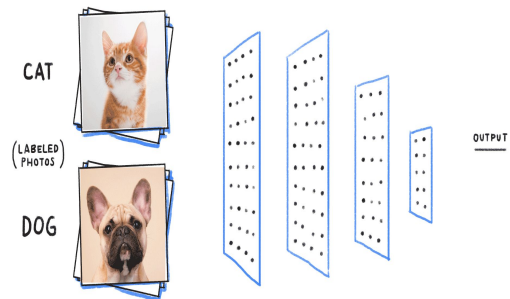
Поздравляю! А теперь повторим раз 500 на одном датасете!

Что за колдунство?



На самом деле, мы нашли минимум

Да, псевдокод из прошлого слайда, по сути – решение всех проблем линейной регрессии *базово*.



Ладно, а с классами что?

Тут немного другой псевдокод... Нам надо ошибку завернуть в сигмоиду.

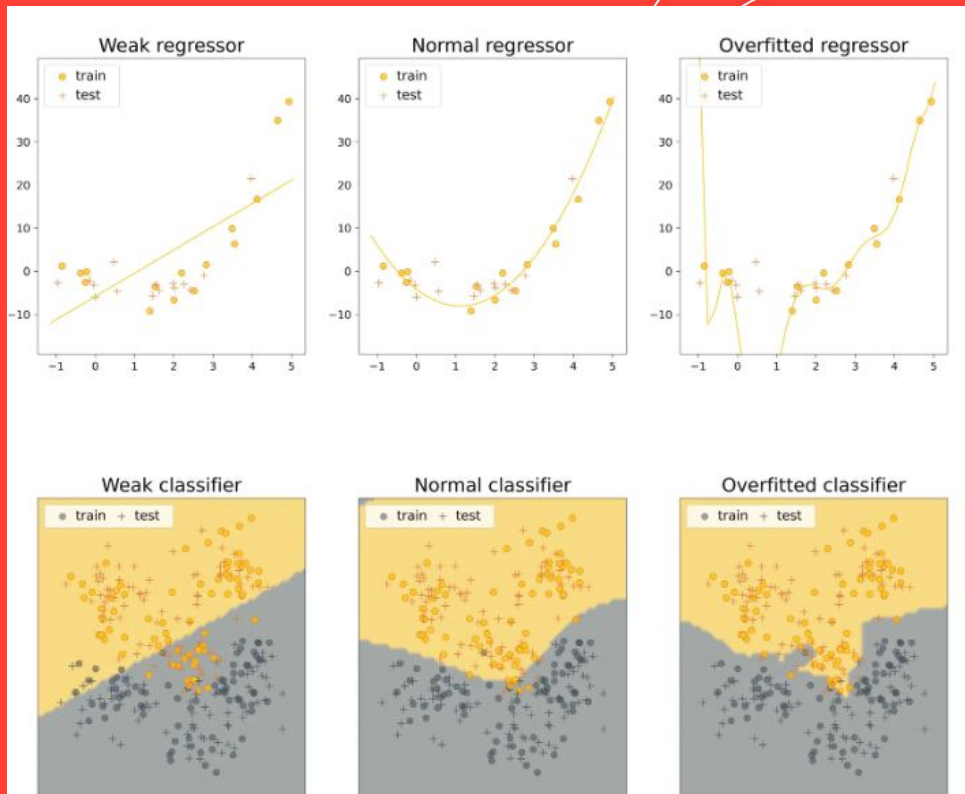
$$-\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

Binary Cross-Entropy / Log Loss

А “изменения” той страшной штукой?

Если взять производную от LogLoss’а, то получится, на деле, тот же самый mse :)

Первые проблемы



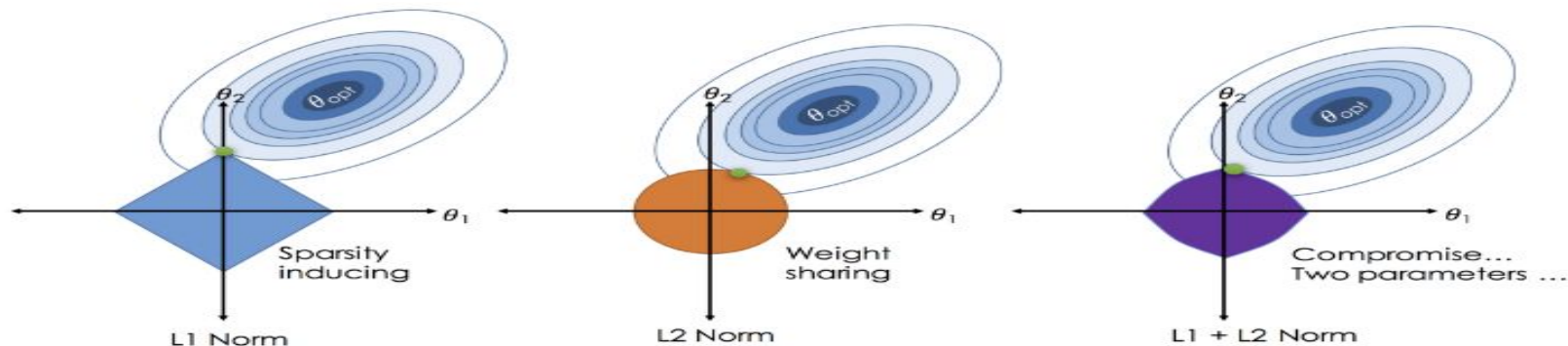
Недообучение – модель не нашла закономерности и решает задачу на “авось”.

Переобучение - модель просто запомнила весь датасет.



Для решения таких проблем используется деление на подвыборки (train, test, validate) и регуляризация

Регуляризация



L1 или lasso

Штрафует за сумму абсолютных весов.

Функция потерь:

$$\text{changes} = \text{our_err} + \lambda * \sum |y[i]|$$

L2 или Ridge

Штрафует за сумму абсолютных весов.

Функция потерь:

$$\text{changes} = \text{our_err} + \lambda * \sum (y[i])^2$$

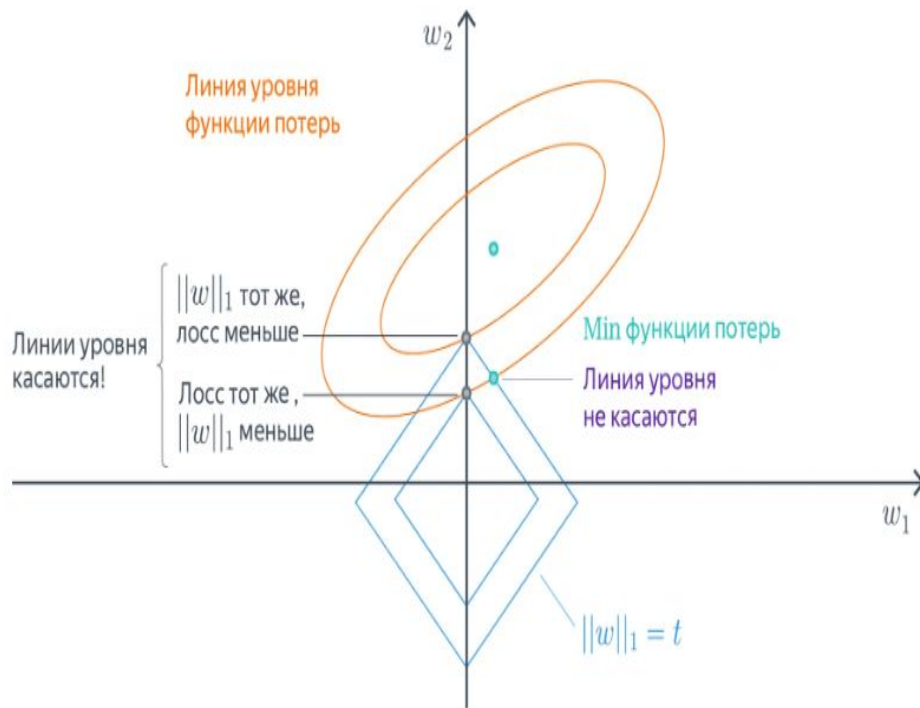
Combine them! or ElasticNet

Пытаемся найти компромисс между L1 и L2

Функция потерь:

$$\text{changes} = \text{our_err} + \lambda_1 * \sum |y[i]| + \lambda_2 * \sum (y[i])^2$$

Немного подробнее про L1

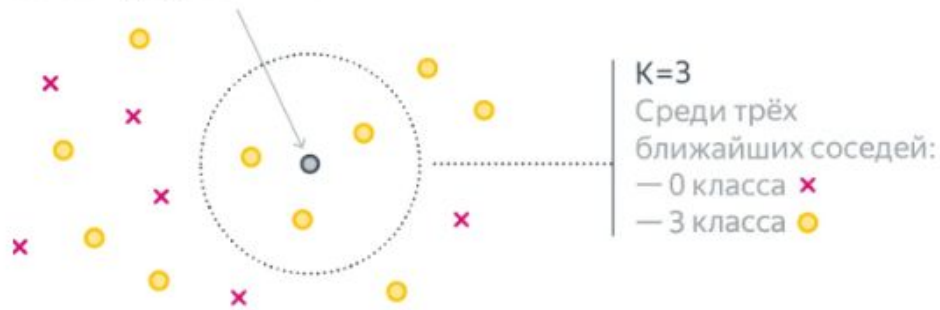


При применении l1, веса по некоторым признакам зануляются → они почти не влияют на наш таргет → мы автоматом “лечимся” от мультиколлинеарности.



Если твой сосед - я, то ты == я. KNN ІІТМО

Хотим предсказать класс



У нас есть выборка из N записей. Запомнили ее - X



Прилетает какой-то объект u . Ищем K соседей (по сути, наименьшие K расстояний от u до каких-то объектов из X)

Какие они в *большинстве* своем? Значит и он такой же.

Глубже в соседей.

Как математически найти метку класса?



1. Пройдись по всем классам (y), которые лежат в выборке.
2. Посчитай сколько среди k ближайших соседей объекта u имеют метку этого класса.
3. Присвой объекту u тот класс, каких соседей большинство

$$P(u \sim y) = \frac{\sum_{i=1}^k \mathbb{I}(y_u^{(i)} = y)}{k}$$

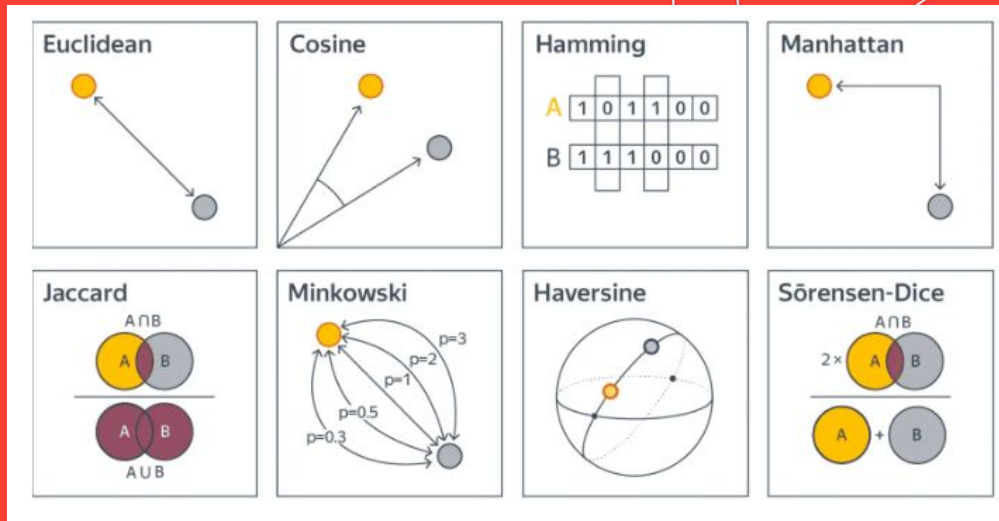
А оценить вероятность?



1. Выполняем шаги 1-2 из прошлого пункта
2. Каждый полученный результат делим на наше k .

$$a(u) = \operatorname{argmax}_{y \in Y} \sum_{i=1}^k \mathbb{I}(y_u^{(i)} = y)$$

Метрики в KNN



База это Евклид. Он будет работать отлично



$$\rho(x,y) = \sqrt{(\sum_i (x_i - y_i)^2)}$$

Манхеттенское:

$$\rho(x,y) = \sum_i |x_i - y_i|$$

Минковский

$$\rho(x,y) = (\sum_i |x_i - y_i|^p)^{1/p}$$

Косинусное

$$\rho(x,y) = 1 - \cos(\theta) = 1 - (xy) / (||x|| * ||y||)$$

Делаем KNN круче!

$$a(u) = \underset{y \in Y}{\operatorname{argmax}} \sum_{i=1}^k w_i \mathbb{I}(y_u^{(i)} = y)$$

$$w_i = \frac{k+1-i}{k} \quad w_i = q^i, 0 < q < 1$$

$$a(u) = \underset{y \in Y}{\operatorname{argmax}} \sum_{i=1}^k K\left(\frac{p(u, x_u^i)}{h}\right) \mathbb{I}(y_u^{(i)} = y)$$

1. Добавляем веса. Вес тем больше, чем ближе объект к цели.

2. Чаще всего у нас затухающие веса берутся линейно или экспоненциально.

А можем еще круче? Да.
На ядрах.

Какие бывают ядра?

- $K(x) = \frac{1}{2} \mathbb{I}[|x| \leq 1]$ — прямоугольное ядро;
- $K(x) = (1 - |x|) \mathbb{I}[|x| \leq 1]$ — треугольное ядро (непрерывное);
- $K(x) = \frac{3}{4} (1 - x^2) \mathbb{I}[|x| \leq 1]$ — ядро Епанечникова (гладкое везде, кроме -1 и 1);
- $K(x) = \frac{15}{16} (1 - x^2)^2 \mathbb{I}[|x| \leq 1]$ — биквадратное ядро (гладкое везде);
- $K(x) = \frac{1}{\sqrt{2\pi}} e^{-2x^2}$ — гауссовское ядро (бесконечно гладкое везде).

От ядра зависит гладкость аппроксимации.



Чем ниже, тем более гладкое.

На практике, нужно либо прямоугольное (быстро и просто), либо гауссовское (супер гладко)

h , ширина окна, как и learning rate подбирается от задачи к задаче

Не хочу делить котов. Хочу квартиру! ИТМО

Логичные варианты



$$a(u) = 1/k * \sum_{i=1}^k y_u^{(i)}$$

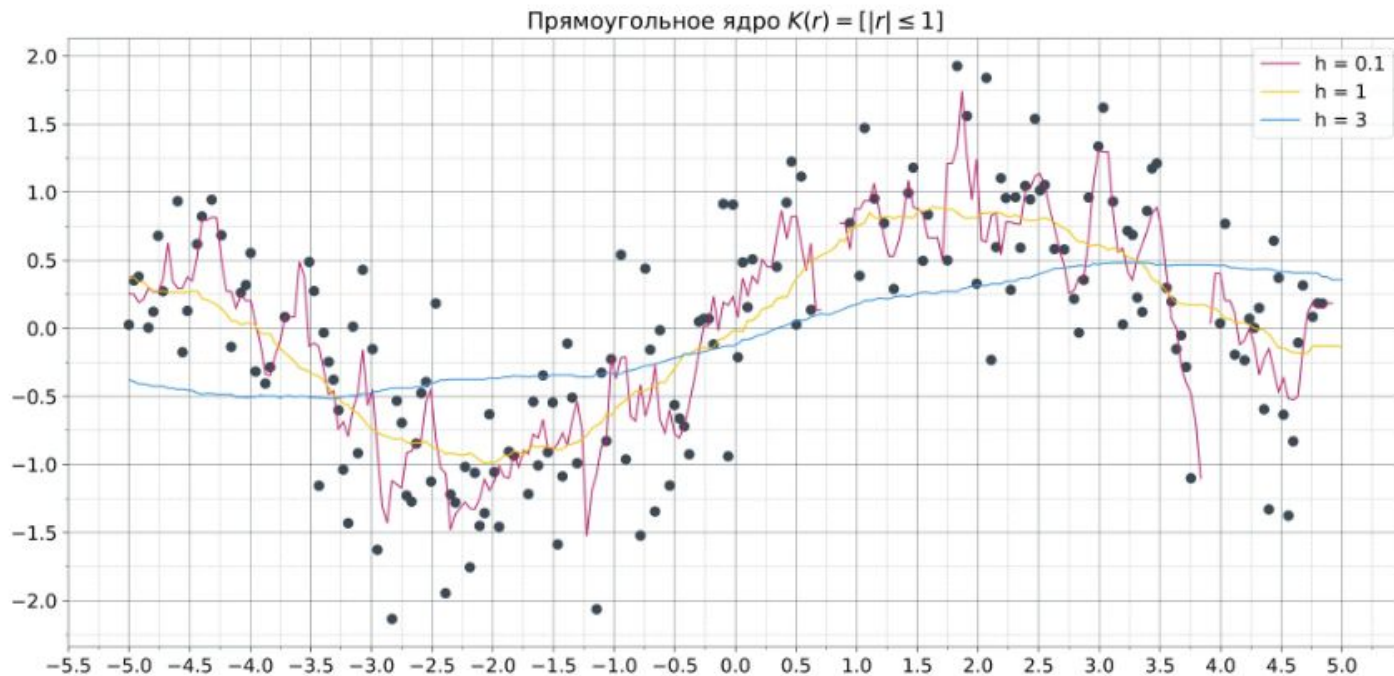
$$a(u) = (\sum_{i=1}^k K(\frac{p(u, x_u^i)}{h}) y_u^{(i)}) / (\sum_{i=1}^k K(\frac{p(u, x_u^i)}{h}))$$

Вариант спустя пару логичных замечаний

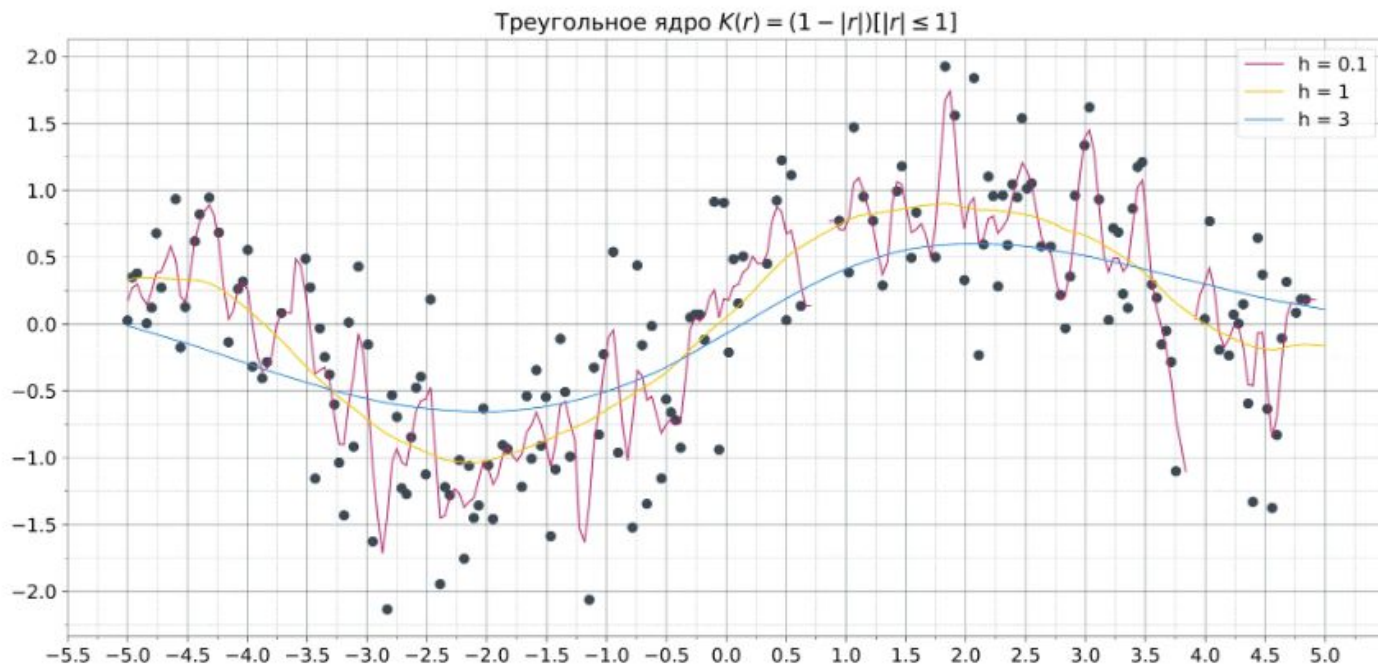


$$a(u) = \operatorname{argmin}_{y \in R} \sum_{i=1}^k K(\frac{p(u, x_u^i)}{h}) * (y - y_u^{(i)})^2$$

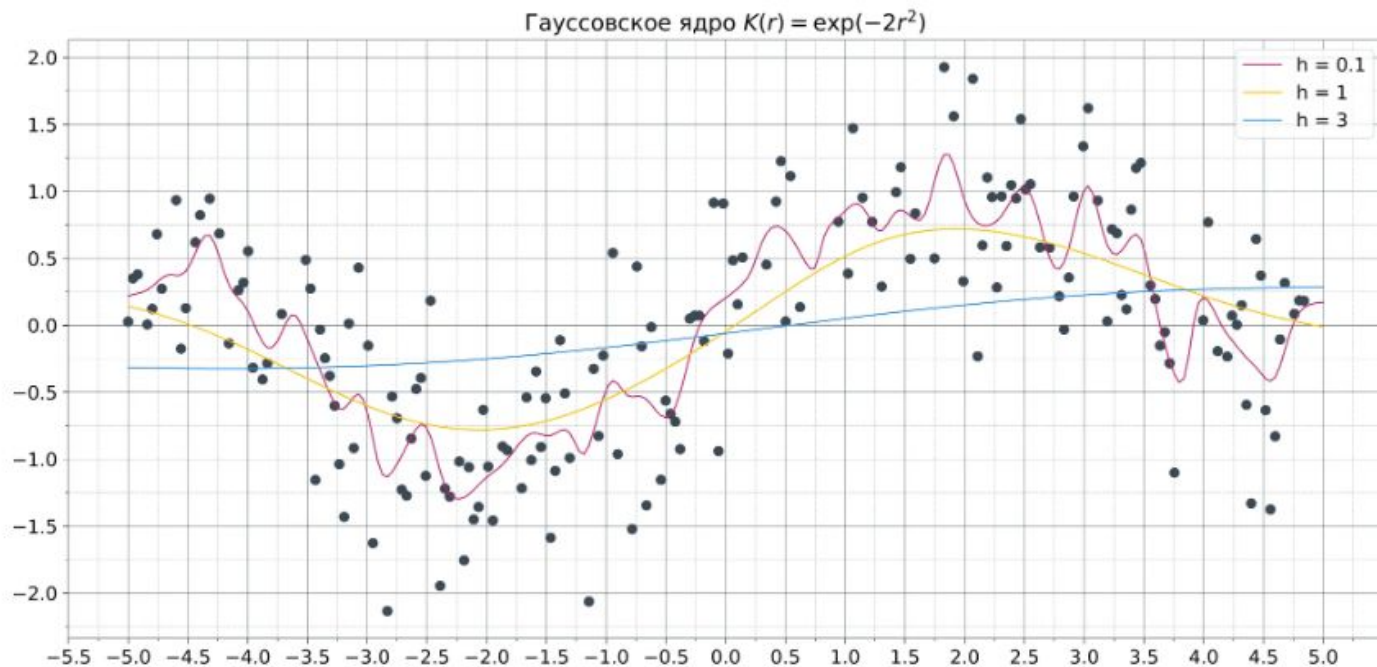
Разная ширина окна



Разная ширина окна



Разная ширина окна



Круче KNN (с оговоркой)

ИТМО

Деревья

Не ML:

- K-d деревья
- Random projection trees

ML:

- “Леса”

Хэши

- LSH
- HNSW

**Спасибо
за внимание!**

it'sMO *re than a*
UNIVERSITY