

курс «Прикладные задачи анализа данных»

# **Функции ошибки / функционалы качества**

## **Часть 1: Задачи регрессии**

**Александр Дьяконов**

**02 октября 2020 года**



## **План на несколько лекций**

**задача регрессии**

**задача бинарной классификации**

- **чёткая классификация**
- **скоринговые функции**

**задача классификации с несколькими классами**

**задачи ранжирования**

**задачи кластеризации**

## Задача – ДНК

**Дано**

**Найти**

**Критерий**

**Построить алгоритм легко!**

**Чтобы улучшить... надо уметь оценивать.**

**Метрики**

- **функции ошибки**
- **функционалы качества**

## Функции ошибки / функционалы качества

Пожалуй, **самое главное** при решении задачи...  
иногда важнее данных!

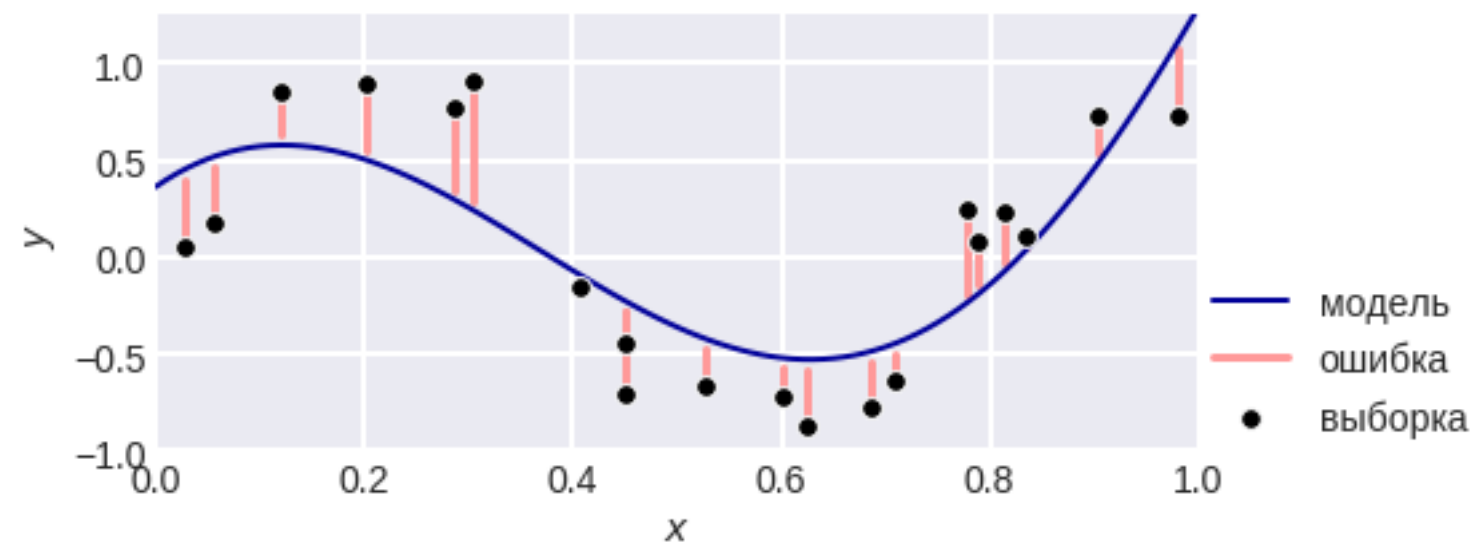
а что такое решение!

**В анализе данных:**

- формализация ответа (формат)
- как ответ оценивается (критерий качества)

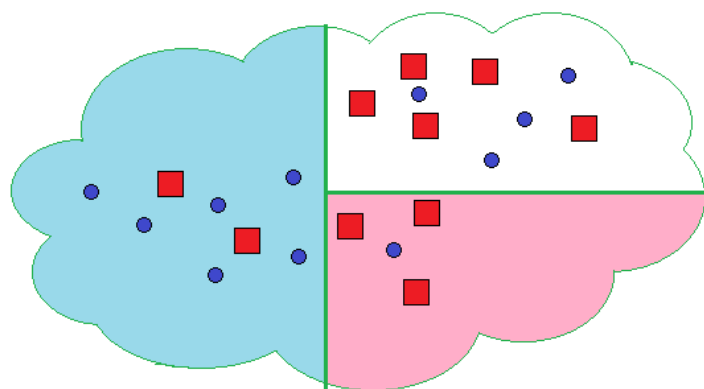
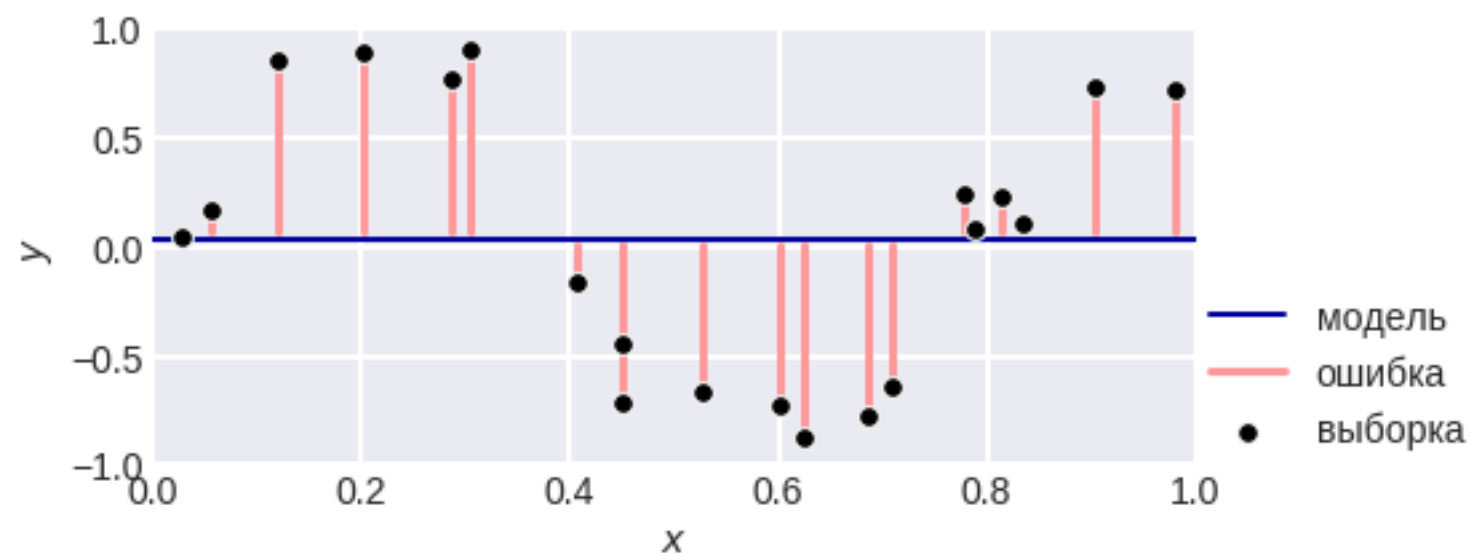
**Случай из практики: задача про траектории зрачка**  
(задача с 3 классами, а не с двумя)

## Задача регрессии



## Задача регрессии

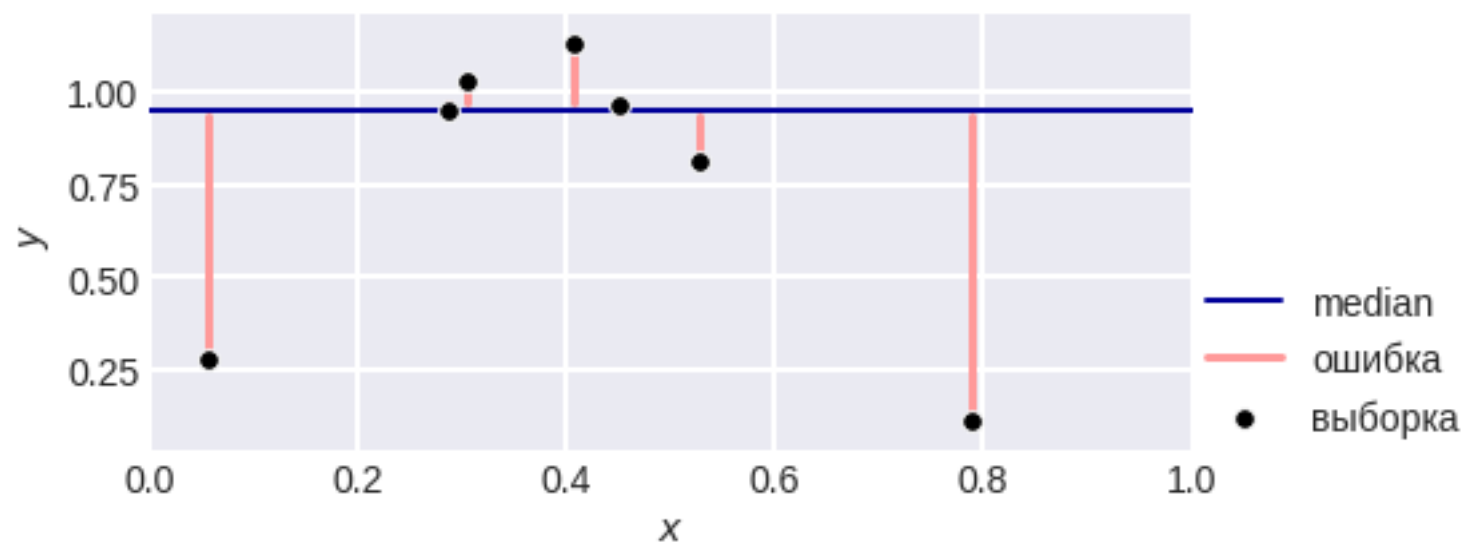
**Будем дальше пытаться всё решать в классе констант**



**1. Простейшее решение**

**2. Примерно это и происходит в листьях решающих деревьев**

**3. Раскрывает природу функционалов**

**Средний модуль отклонения – Mean Absolute Error (MAE), Mean Absolute Deviation (MAD)**

$$\text{MAE} = \frac{1}{m} \sum_{i=1}^m |a_i - y_i|$$

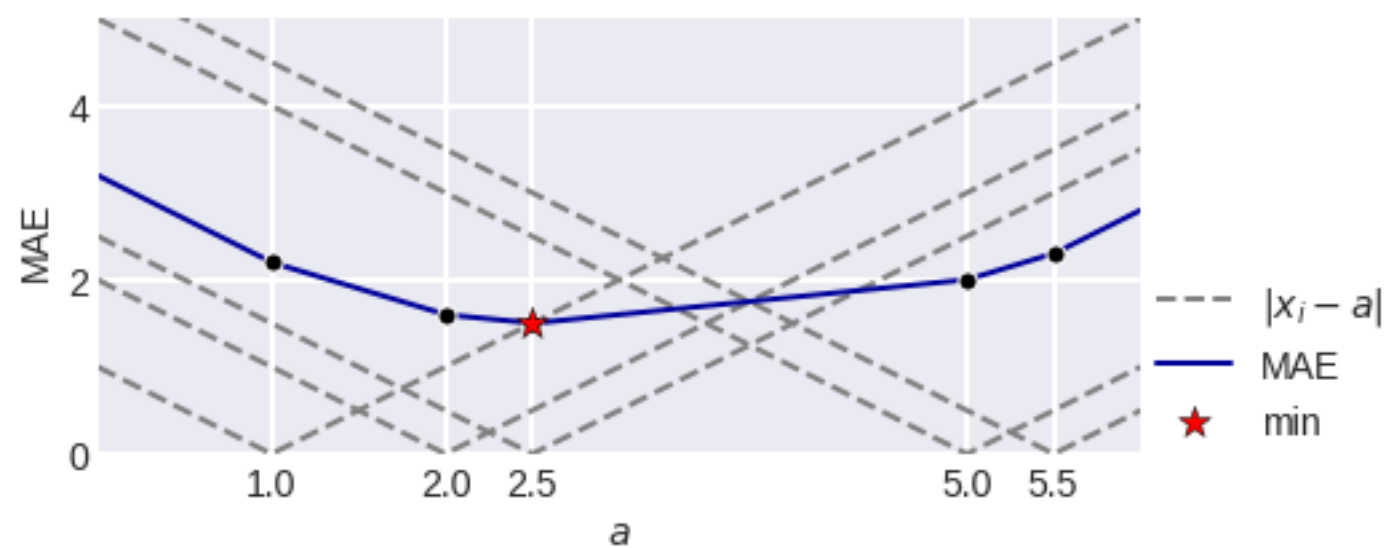
**Напоминание:**

$$\frac{1}{m} \sum_{i=1}^m |a - y_i| \rightarrow \min$$

$$a = \text{median}(\{y_i\}_{i=1}^m)$$

**Это открывает смысл решений!**

## Средний модуль отклонения





## **Средний модуль отклонения**

### **Способы использования тайных знаний:**

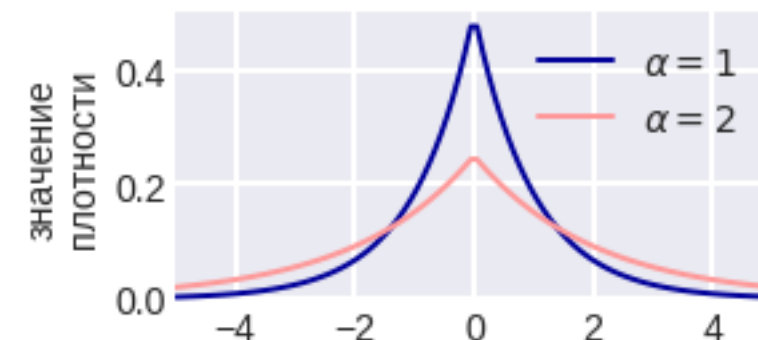
- **медиана, вместо усреднения, в ансамбле**
- **округление ответа (если целевой вектор целочисленный)**

## Откуда берётся MAE

$$y = a_w(x) + \varepsilon$$

$w$  – параметры алгоритма  $a_w(x)$

$$\varepsilon \sim \text{laplace}(0, \alpha)$$



**Для оценки параметров выписываем правдоподобие модели**

$$p(y | x, w) = \frac{\alpha}{2} \exp[-\alpha |y - a_w(x)|]$$

**Метод максимального правдоподобия:**

$$\begin{aligned} \log L(w) &= \log \prod_{i=1}^m p(y_i | x_i, w) = \\ &= \sum_{i=1}^m \left[ \log \frac{\alpha}{2} - \alpha |y_i - a_w(x_i)| \right] \rightarrow \max \end{aligned}$$

## Откуда берётся MAE

Получаем

$$\alpha \sum_{i=1}^m |y_i - a_w(x_i)| \rightarrow \min$$

**т.е. задачу минимизации MAE!**

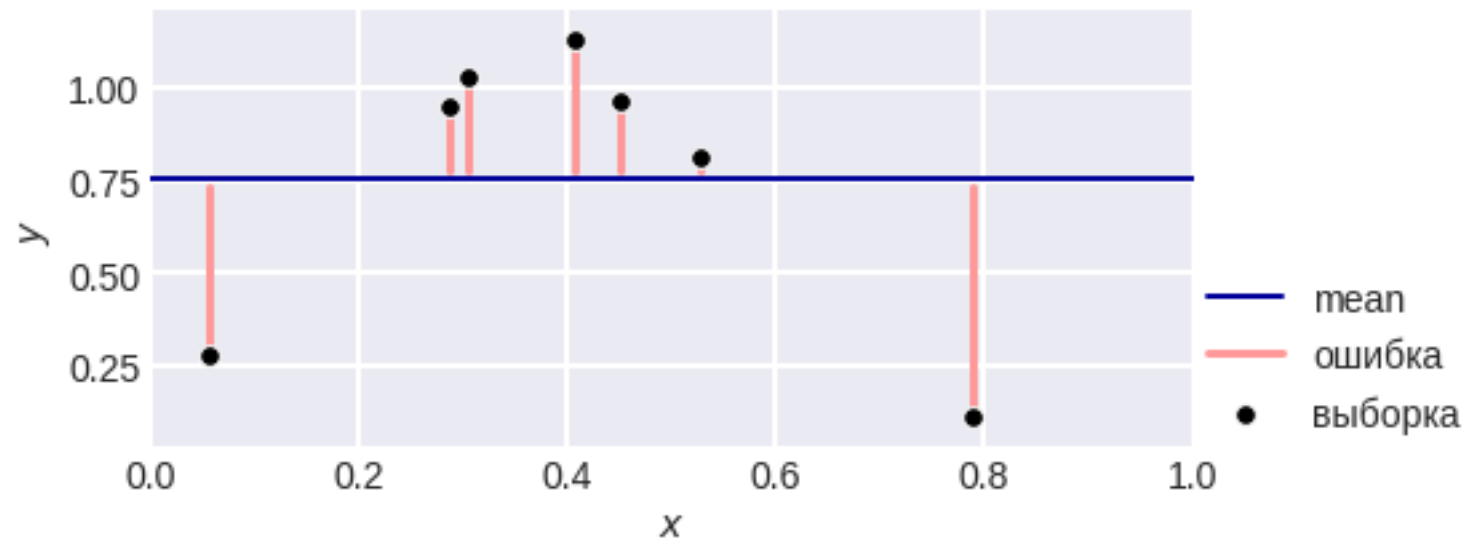
- не зависит от природы модели
- зависит от распределения ошибок  
(почему Residual Plots)

**Максимизация правдоподобия эквивалентна минимизации MAE!**

**Чему соответствует минимизация весового MAE?**

## Средний квадрат отклонения ~ Mean Squared Error (MSE)

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m |a_i - y_i|^2$$



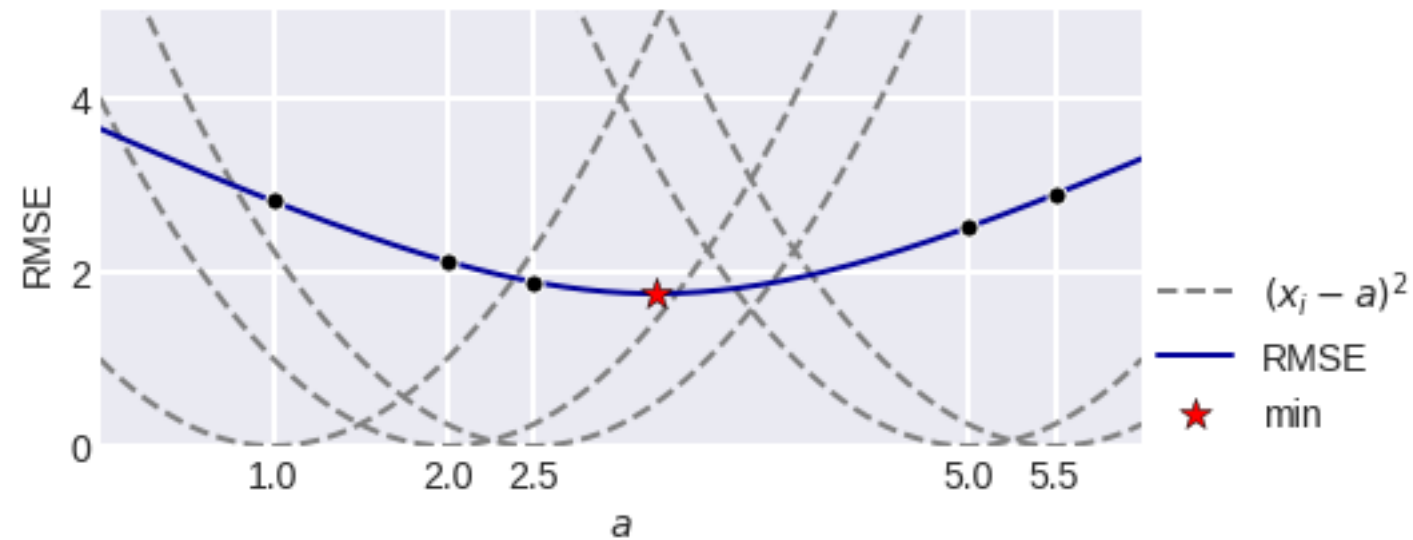
$$\frac{1}{m} \sum_{i=1}^m |a - y_i|^2 \rightarrow \min$$

$$a = \frac{1}{m} \sum_{i=1}^m y_i$$

**Root Mean Squared Error (RMSE) / Root Mean Square Deviation (RMSD)**

$$\text{RMSE} = \sqrt{\frac{1}{m} \sum_{i=1}^m |a_i - y_i|^2}$$

## Средний квадрат отклонения ~ Mean Squared Error (MSE)



### Способы использования тайных знаний

- ничего не делать (в RF, GBM и т.д. всё равно усредняют)
- метод НСКО – классическая регрессия!

## Нормированная версия: коэффициент детерминации $R^2$ (Coefficient of Determination)

$$R^2 = 1 - \frac{\sum_{i=1}^m |a_i - y_i|^2}{\sum_{i=1}^m |\bar{y} - y_i|^2}$$

$$\bar{y} = \frac{1}{q} \sum_{i=1}^m y_i$$

**В общем случае (в статистике) коэффициент детерминации:**

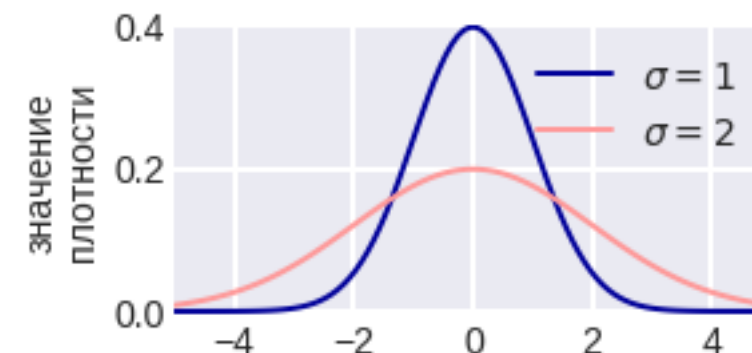
$$R^2 = 1 - \frac{\mathbf{D}(y | x)}{\mathbf{D}(y)}$$

## Откуда берётся (R)MSE

$$y = a_w(x) + \varepsilon$$

$w$  – параметры алгоритма  $a_w(x)$

$$\varepsilon \sim \text{norm}(0, \sigma^2)$$



Для оценки параметров выписываем правдоподобие модели

$$p(y | x, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y - a_w(x))^2}{2\sigma^2}\right]$$

**Метод максимального правдоподобия:**

$$\begin{aligned} \log L(w) &= \log \prod_{i=1}^m p(y_i | x_i, w) = \\ &= \sum_{i=1}^m \left[ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(y_i - a_w(x_i))^2}{2\sigma^2} \right] \rightarrow \max \end{aligned}$$

## Откуда берётся (R)MSE

Получаем

$$\frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - a_w(x_i))^2 \rightarrow \min$$

**т.е. задачу минимизации MSE!**

- не зависит от природы модели
- зависит от распределения ошибок  
(почему Residual Plots)

**Максимизация правдоподобия эквивалентна минимизации среднеквадратичной ошибки!**

**Д3 Каким ещё распределениям какие ошибки соответствуют?**



## Откуда берётся (R)MSE: ещё одно «оправдание»

Пусть функция ошибки  $l(y, a) = g(y - a)$

**Что логично?**

**1.**  $g(0) = 0$

**2.**  $|z_1| \leq |z_2| \Rightarrow g(z_1) \leq g(z_2)$

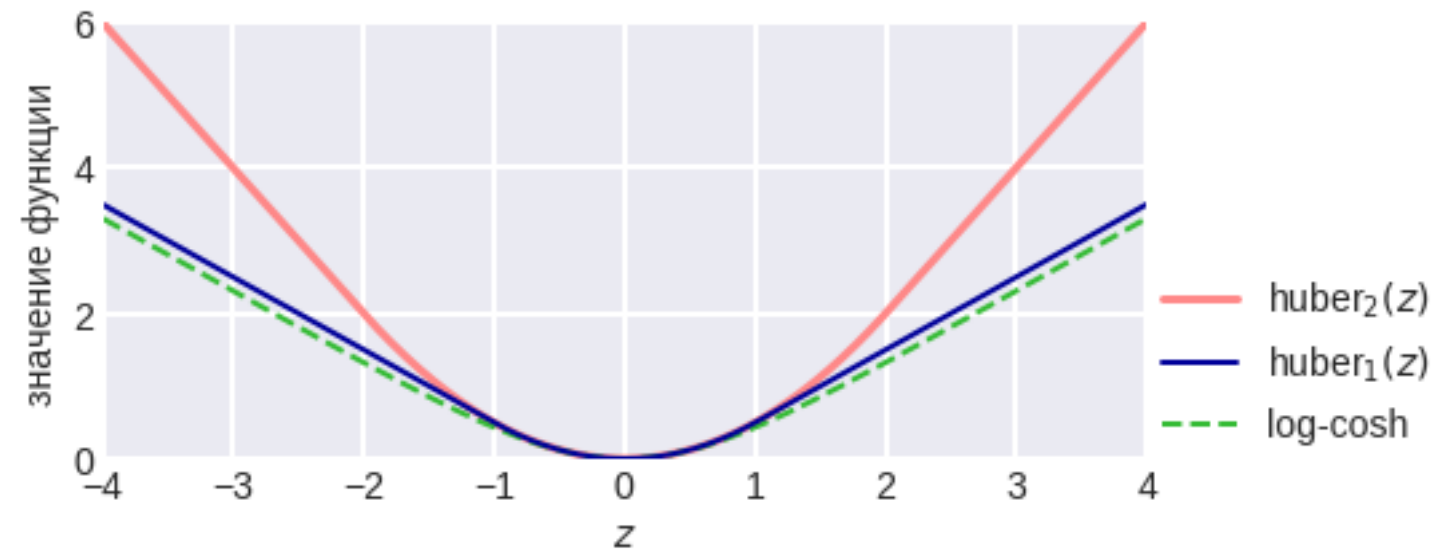
**3. достаточно гладкая...**

$$g(z) = g(0) + g'(0)z + \frac{g''(0)}{2}z^2 + o(z^2)$$

**но тогда**

$$l(y, a) = g(y - a) \approx \underbrace{g(0)}_{=0(1)} + \underbrace{g'(0)(y - a)}_{=0(2)} + \frac{g''(0)}{2}(y - a)^2 = \underbrace{C}_{>0}(y - a)^2$$

## Функция Хьюбера и logcosh

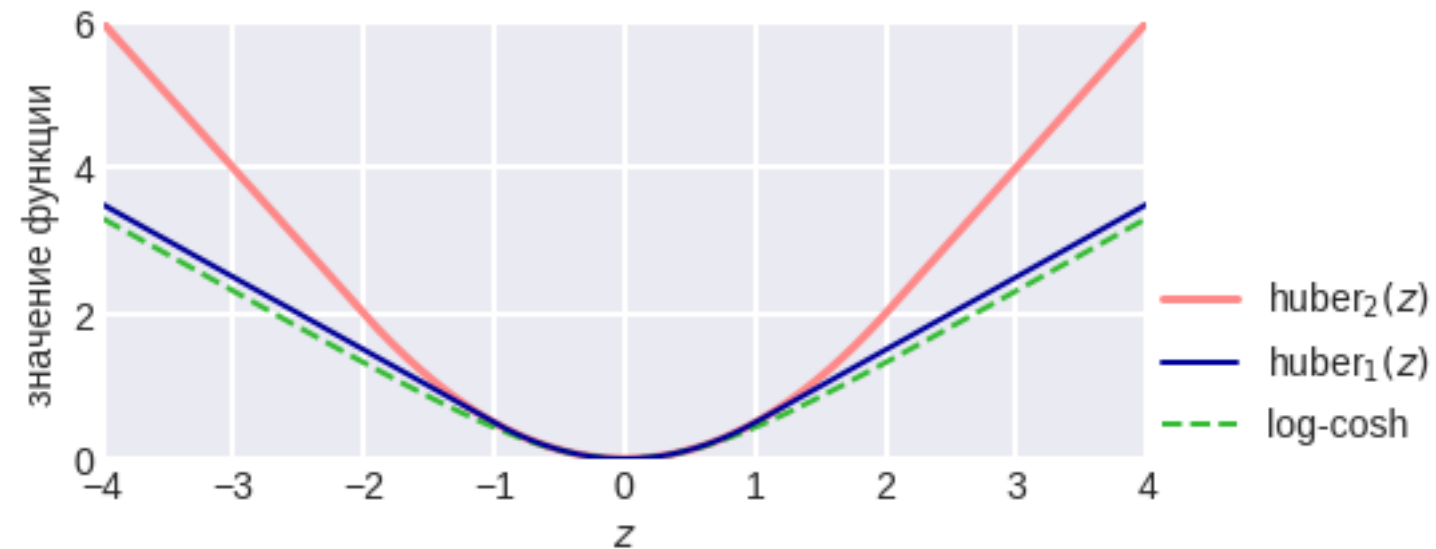


$$\text{huber}(z) = \begin{cases} \frac{1}{2}z^2, & |z| \leq \delta, \\ \delta \left( |z| - \frac{1}{2}\delta \right), & |z| > \delta. \end{cases}$$

**Как только что вывели:**

**когда отклонение мало – ошибка квадратичная**  
**когда велико (в т.ч. выбросы) – линейная**

## Функция Хьюбера и logcosh

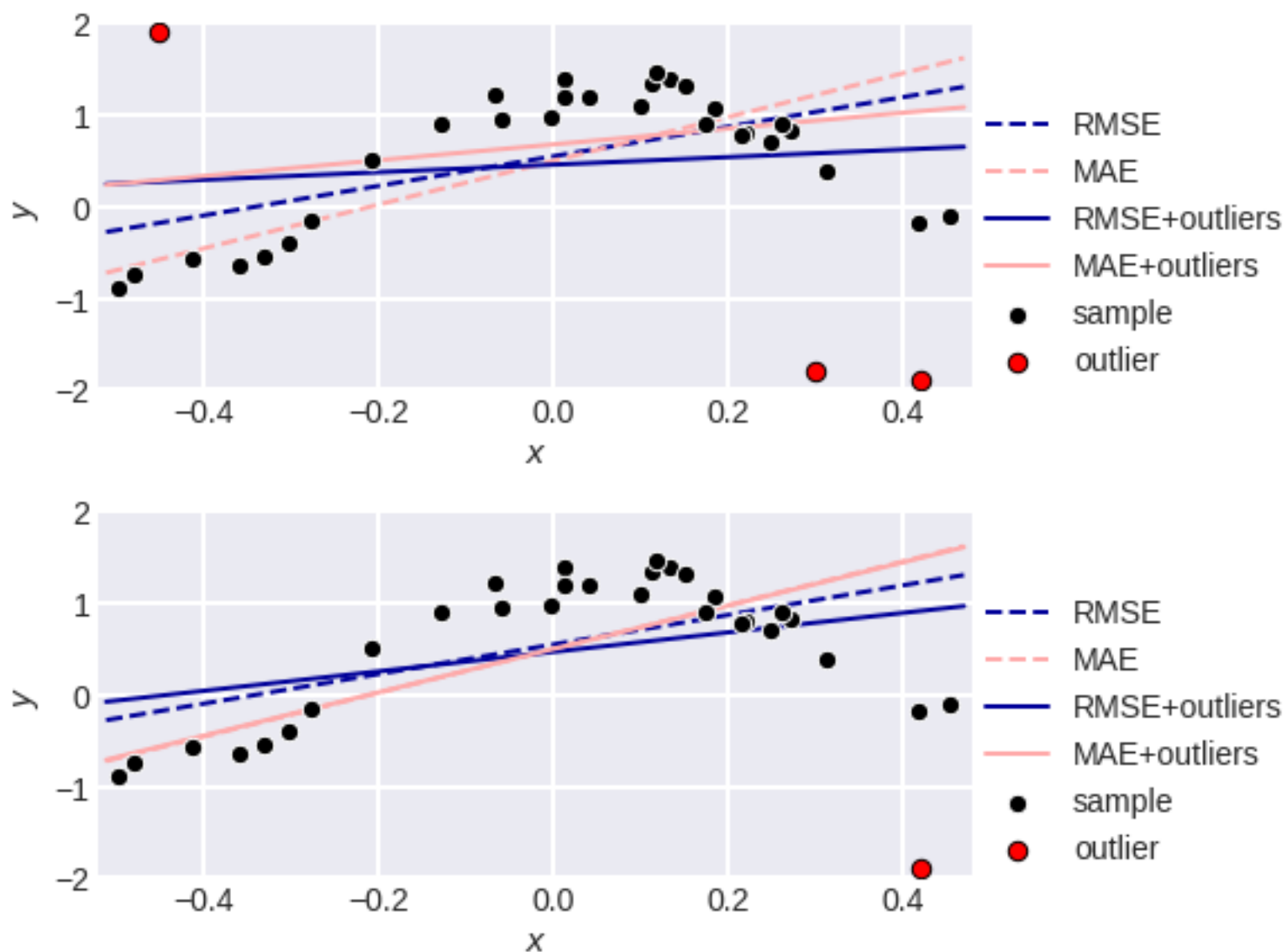


$$\text{logcosh} = \log\left(\frac{\exp(z) + \exp(-z)}{2}\right)$$

**непараметрическая,  
но используется редко.**

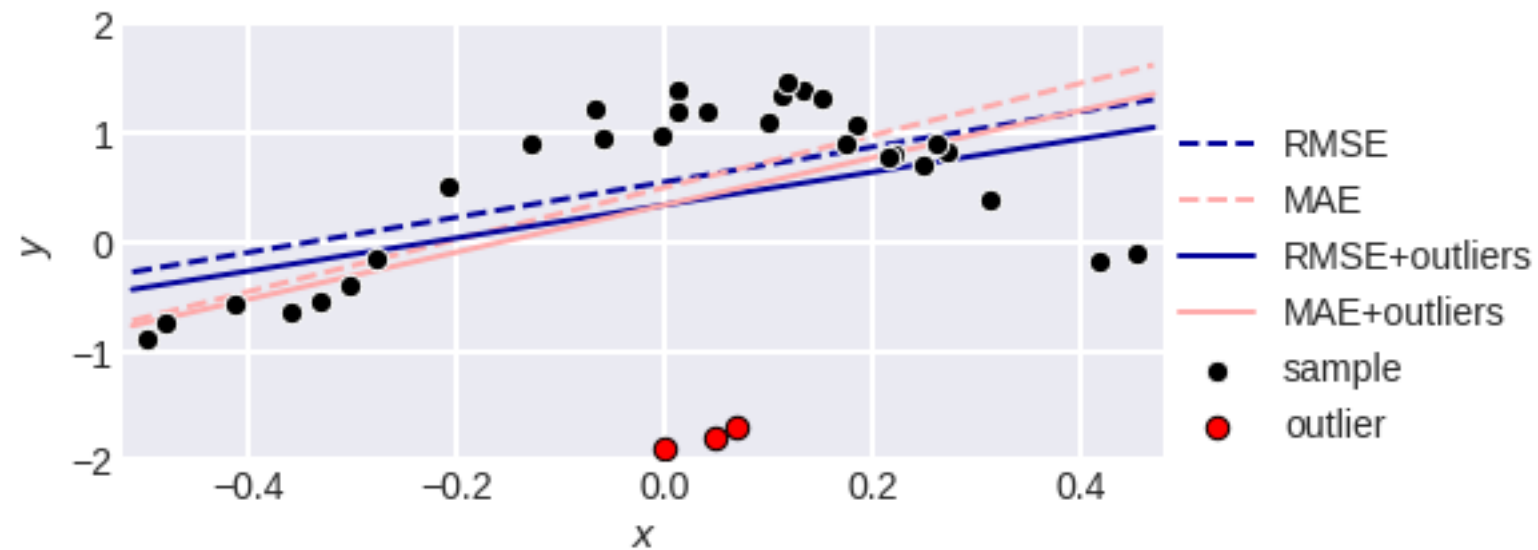
## Различия MSE и MAE

## Устойчивость к выбросам...



## Различия MSE и MAE

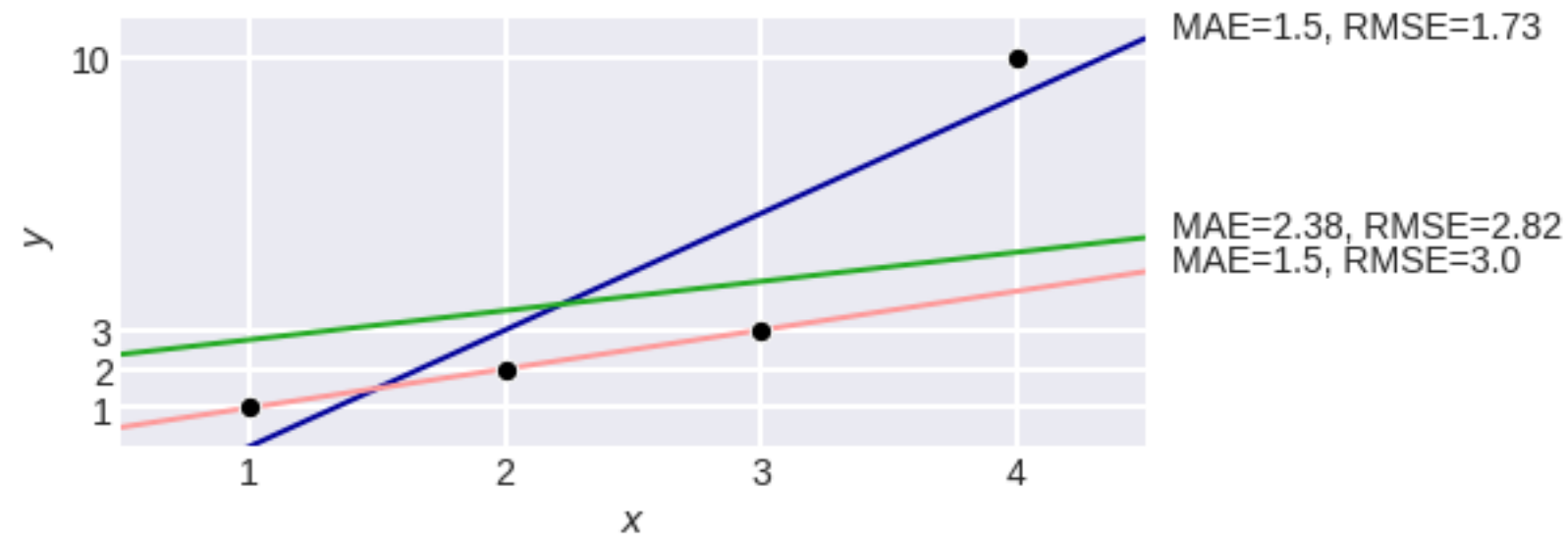
### Устойчивость к выбросам...



**Считается, что MAE устойчивее к выбросам...**

**Д3 Честный эксперимент: зависимость результата от функций ошибок / выбросов**

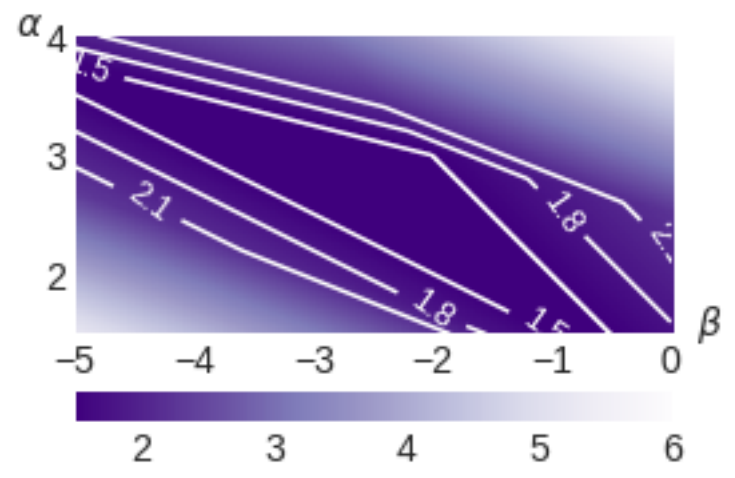
Различия MSE и MAE



посмотрим на неконстантное решение:

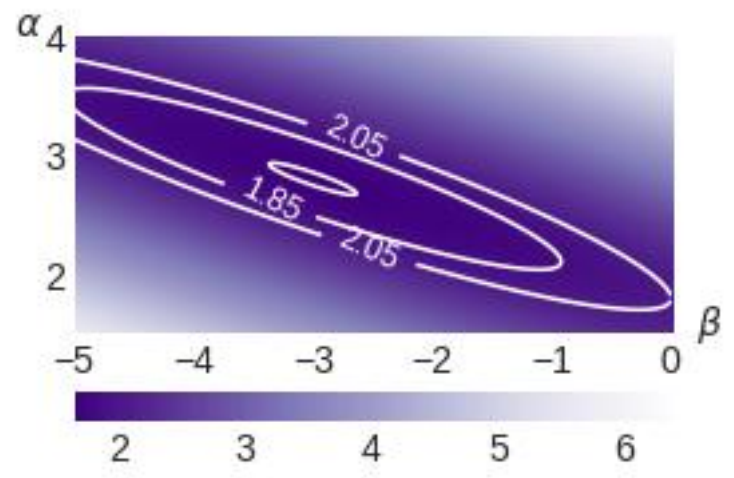
$$\sum_{i=1}^m |y_i - a(x_i)|^p \rightarrow \min,$$
$$a(x) = \alpha x + \beta$$

MAE

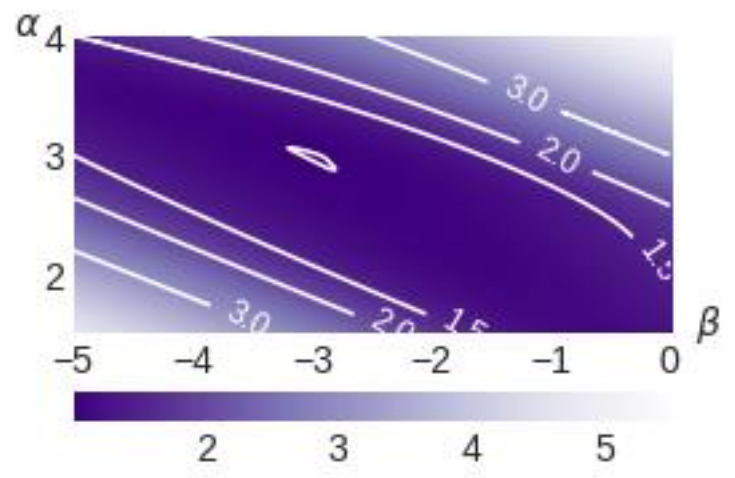


внутри «треугольника»  
одинаковый MAE=1.5

RMSE



Huber ( $\delta = 1$ )



## Различия MSE и MAE

**внутри «треугольника» одинаковый  $MAE=1.5$**

**можно привести примеры, когда MAE меняется слабо,  
а RMSE значительно**

**Д3 Хороший нетривиальный пример / может ли быть наоборот?**

## Обобщения

$$\sqrt[p]{\frac{1}{m} \sum_{i=1}^m w_i |\varphi(a_i) - \varphi(y_i)|^p}$$

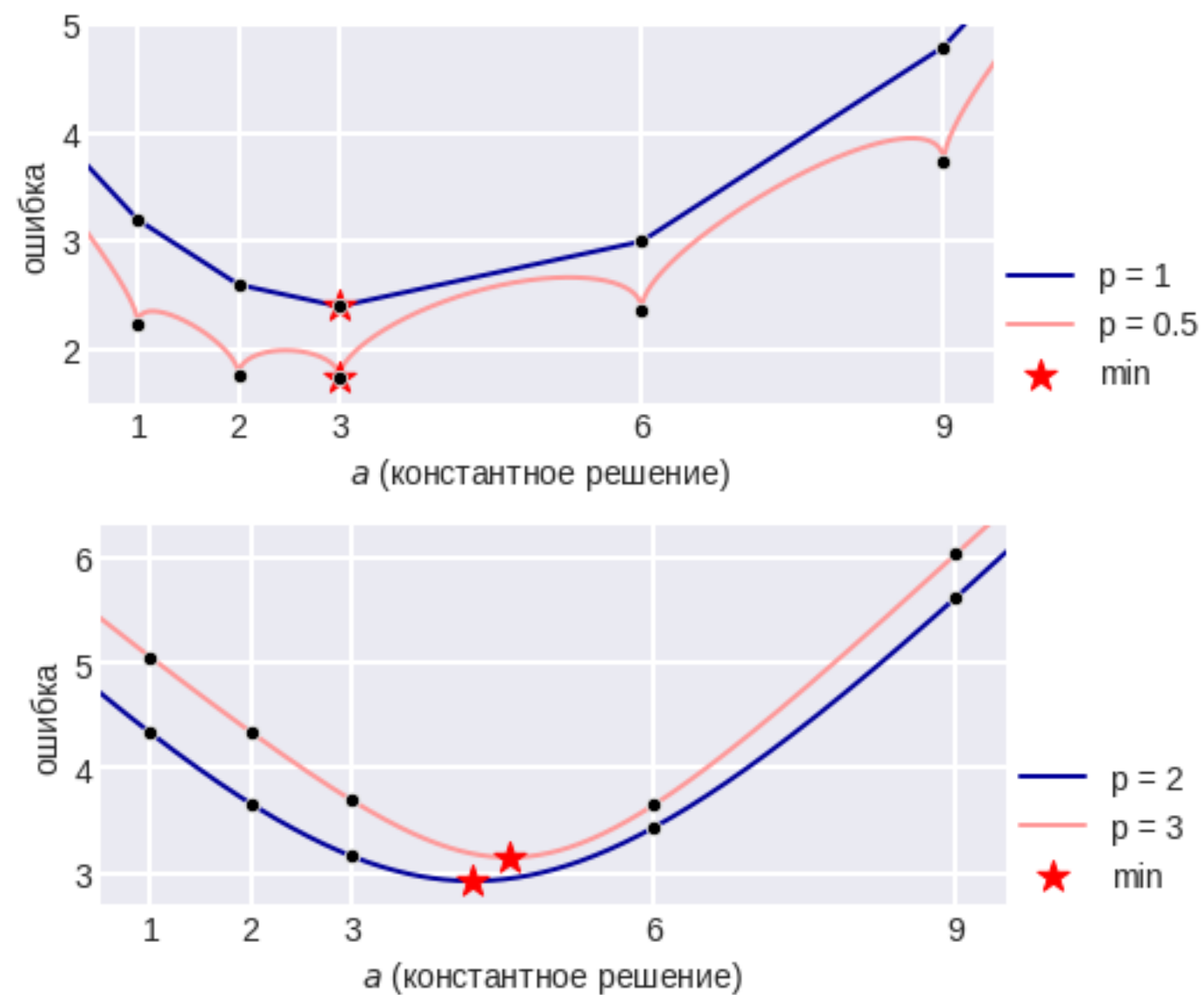
## Рецепты

1. Преобразование целевого вектора  $\varphi(y)$
2. Веса ~ вероятности появления объектов в сэмплировании  
Некоторые модели поддерживают веса объектов
3. В случае нетривиальных  $p$  – прямая настройка

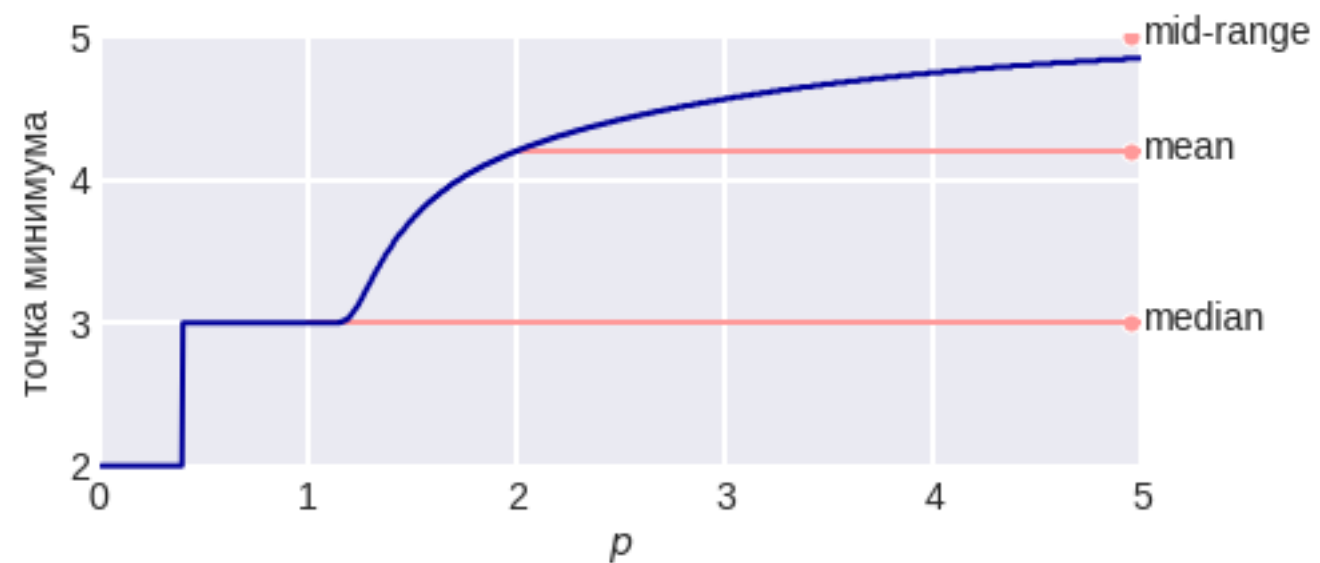
**Дальше к этому вернёмся...**



Про нетривиальные  $p$



## Как точка минимума зависит от степени



**Разные постановки задачи регрессии**

$$\|Xw - y\|_p \rightarrow \min$$

$p < 1$	<b>Это не норма, NP-сложная задача</b>
$p = 1$	<b>Линейное программирование</b>
$1 < p < 2$	<b>Нет стандартных методов</b>
$p = 2$	<b>Аналитическое решение (линейная алгебра)</b>
$p > 2$	<b>Градиентные методы</b>
$p = \infty$	<b>Линейное программирование</b>

**Symmetric mean absolute percentage error (SMAPE or sMAPE)**

$$\text{SMAPE} = \frac{2}{m} \sum_{i=1}^m \frac{|y_i - a_i|}{y_i + a_i} = 100\% \cdot \frac{1}{m} \sum_{i=1}^m \frac{|y_i - a_i|}{(y_i + a_i) / 2}$$

**Когда надо интерпретировать погрешность как проценты  
– плохо, если есть нули (и отрицательные значения)**

**1 – 2**  
**SMAPE = 67%**

**100 – 101**  
**SMAPE = 1%**

**0 – 1**  
**SMAPE = 200%**

Начальники не знают, что такое проценты...

Применение SMAPE – прогноз временных рядов

## Mean Absolute Percent Error (MAPE)

$$\text{MAPE} = \frac{1}{m} \sum_{i=1}^m \frac{|y_i - a_i|}{|y_i|}$$

**Чем MAPE явно лучше SMAPE на практике?**

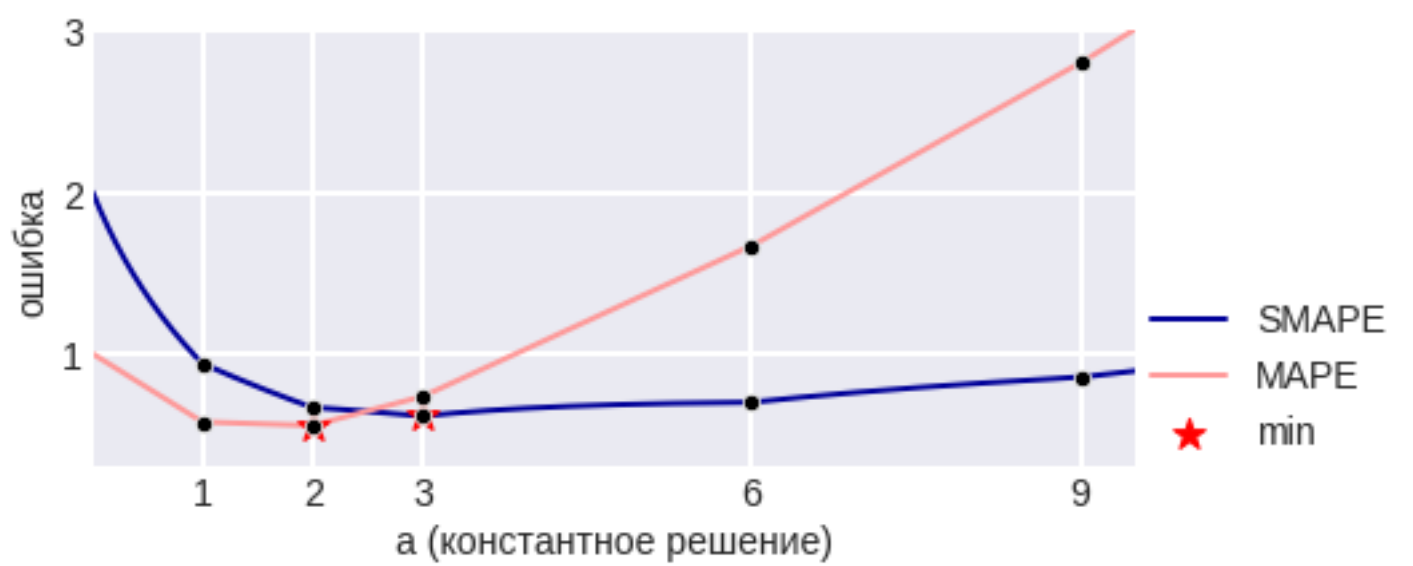
**Mean Absolute Percent Error (MAPE)**

$$\text{MAPE} = \frac{1}{m} \sum_{i=1}^m w_i |y_i - a_i|$$

$$w_i = \frac{1}{|y_i|}$$

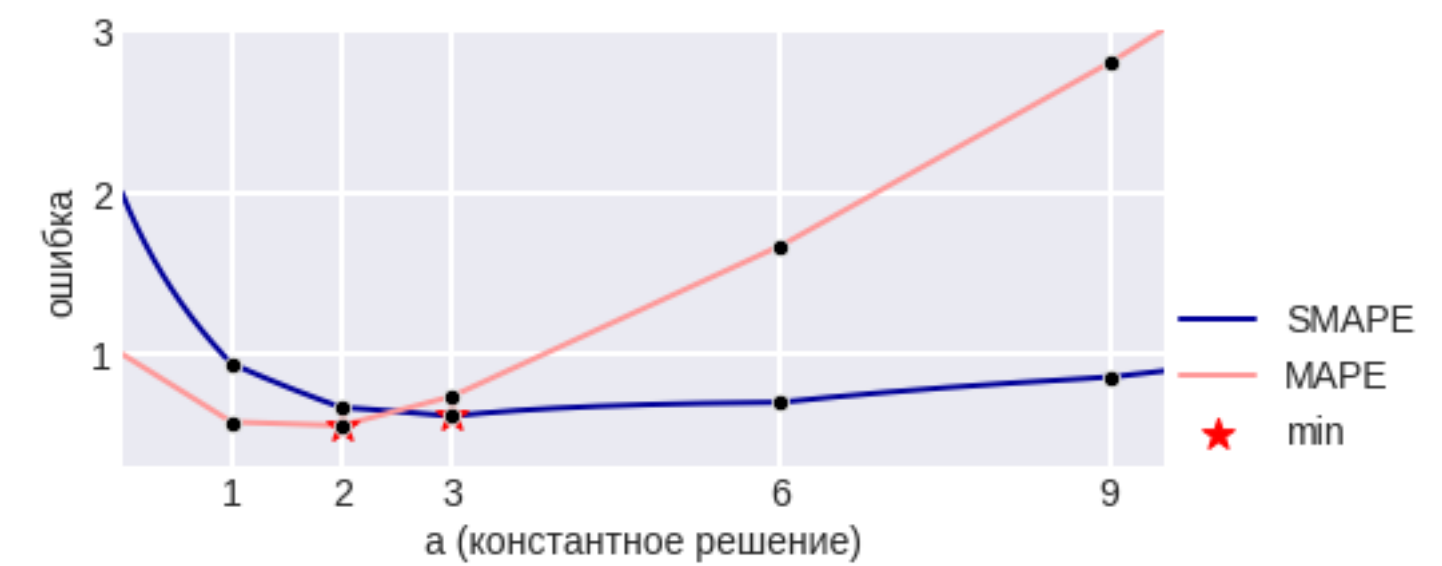
**Просто весовой MAE!**  
как оптимизировать? дальше...

MAPE и SMAPE

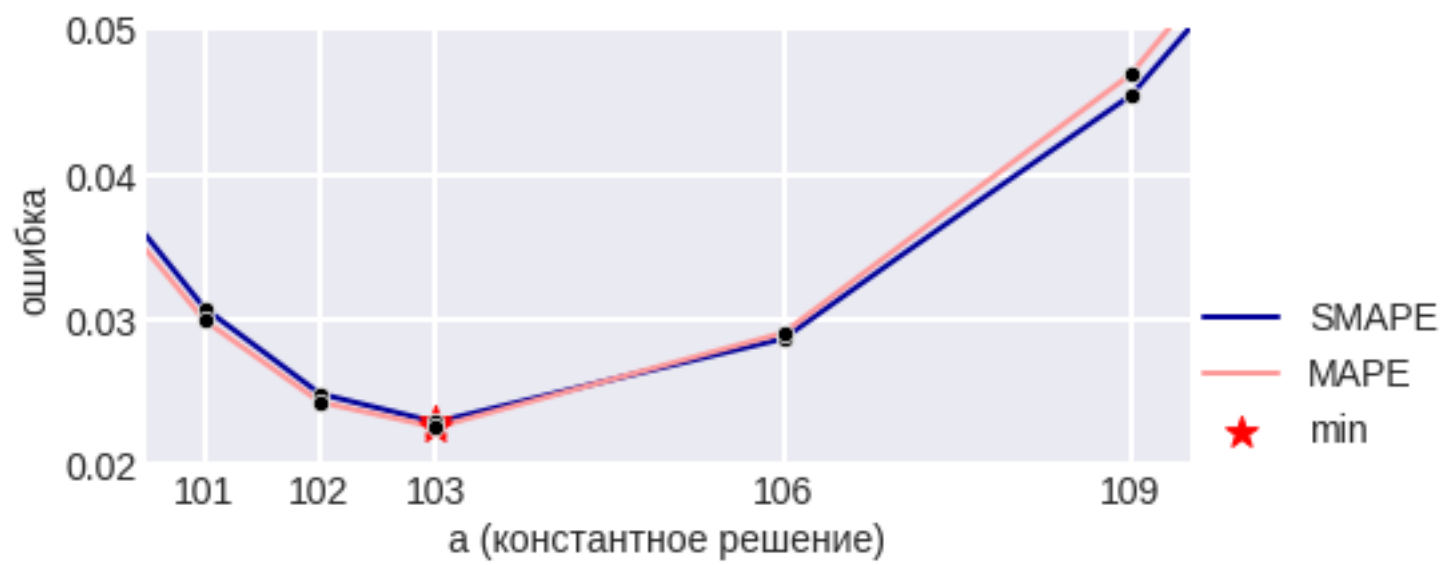


Что настораживает в этом графике?

MAPE и SMAPE

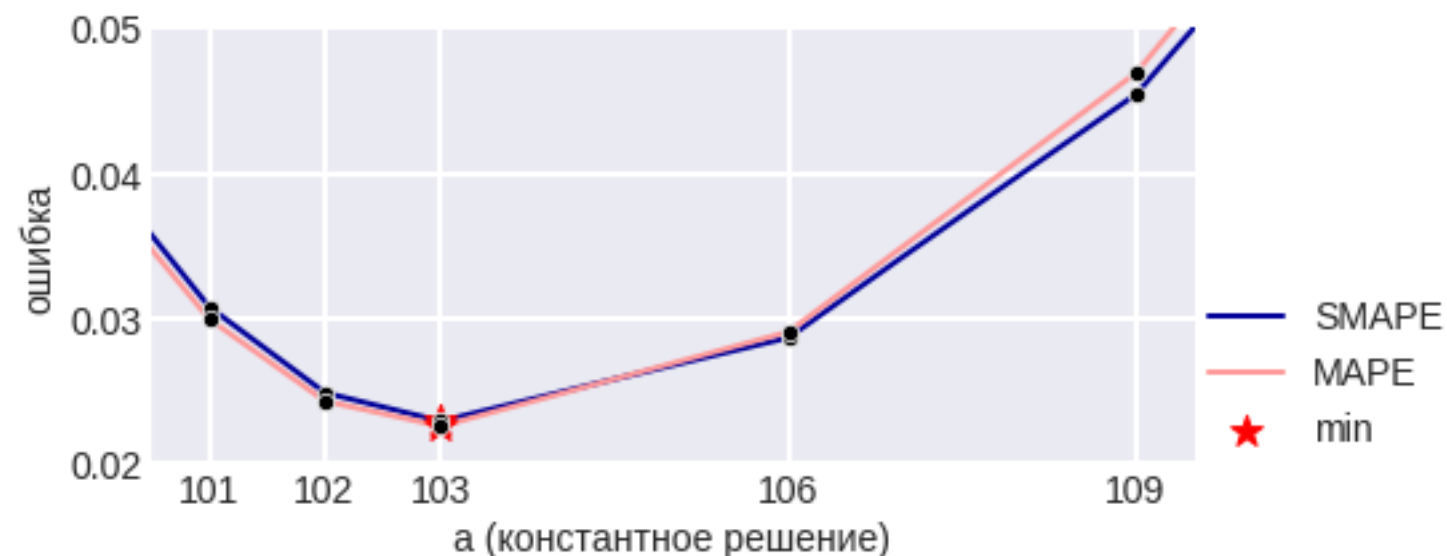


Масштаб! Типичная ошибка (и во многих курсах)





## MAPE и SMAPE



**Например, MAPE – весовой MAE,  
но на практике веса не сильно отличаются!**

**Поэтому решение около медианы**

**Д3 Предложить минимизацию для MAPE и SMAPE (обосновать в экспериментах)**

**PMAD**

**Другой способ нормировки ошибки...**

$$\text{PMAD} = \frac{\sum_{i=1}^m |y_i - a_i|}{\sum_{i=1}^m |y_i|}$$

**эквивалентен MAE**

**ДЗ Как на типичных и специальных выборках соотносятся решения задач минимизации перечисленных функций ошибки?**

## Меры на сравнении с бенчмарком

**Классная идея:**

**сделать простой алгоритм и смотреть ошибку относительно него**

**Mean Relative Absolute Error  
(MRAE)**

$$\text{MRAE} = \frac{1}{m} \sum_{i=1}^m \frac{|y_i - a_i|}{|y_i - a'_i|}$$

**REL\_MAE**

$$\text{REL\_MAE} = \frac{\sum_{i=1}^m |y_i - a_i|}{\sum_{i=1}^m |y_i - a'_i|}$$

**Percent Better**

$$\text{PB(MAE)} = \frac{1}{m} \sum_{i=1}^m I[|y_i - a_i| < |y_i - a'_i|]$$

## Меры на сравнении с бенчмарком

**Как выбрать бенчмарк в задачах прогнозирования?**

## Нормированные ошибки

Не зависят от шкалы...

### Mean Absolute Scaled Error

$$\text{MASE} = \frac{1}{\frac{m}{m-1} \sum_{i=2}^m |y_i - y_{i-1}|} \sum_{t=1}^m |y_t - a_t|$$

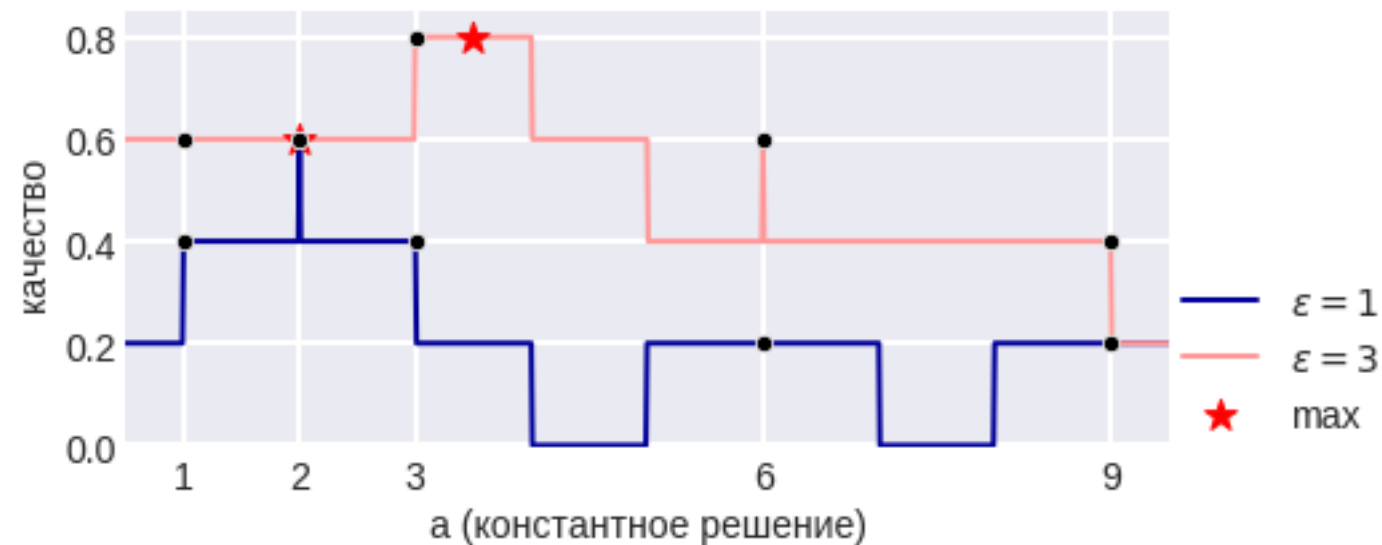
**Какие ещё бывают функционалы в регрессии?**

**С точностью до порога****функция ошибки**

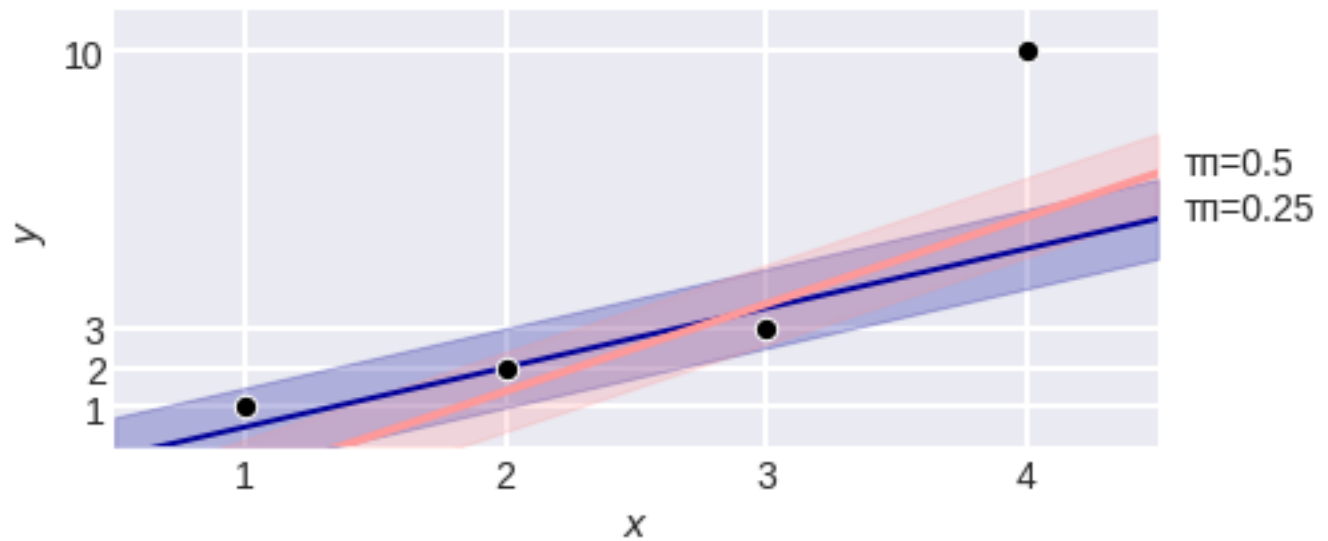
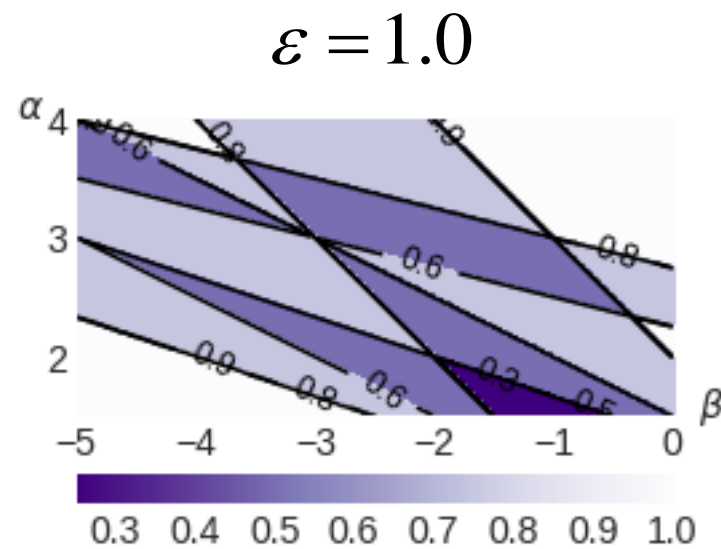
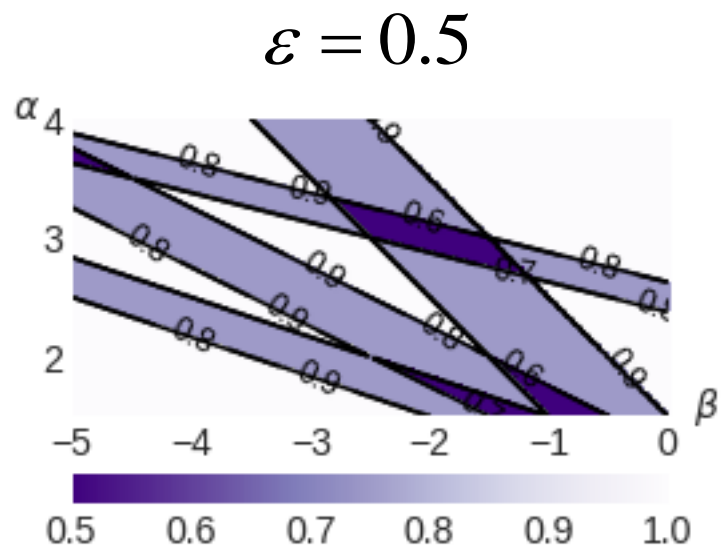
$$eB = \frac{1}{m} \sum_{i=1}^m I[|y_i - a_i| > \varepsilon]$$

**функционал качества**

$$eB = \frac{1}{m} \sum_{i=1}^m I[|y_i - a_i| < \varepsilon]$$

**был в задаче Dunnhumby****Оптимальное решение – мода парzenовской плотности**

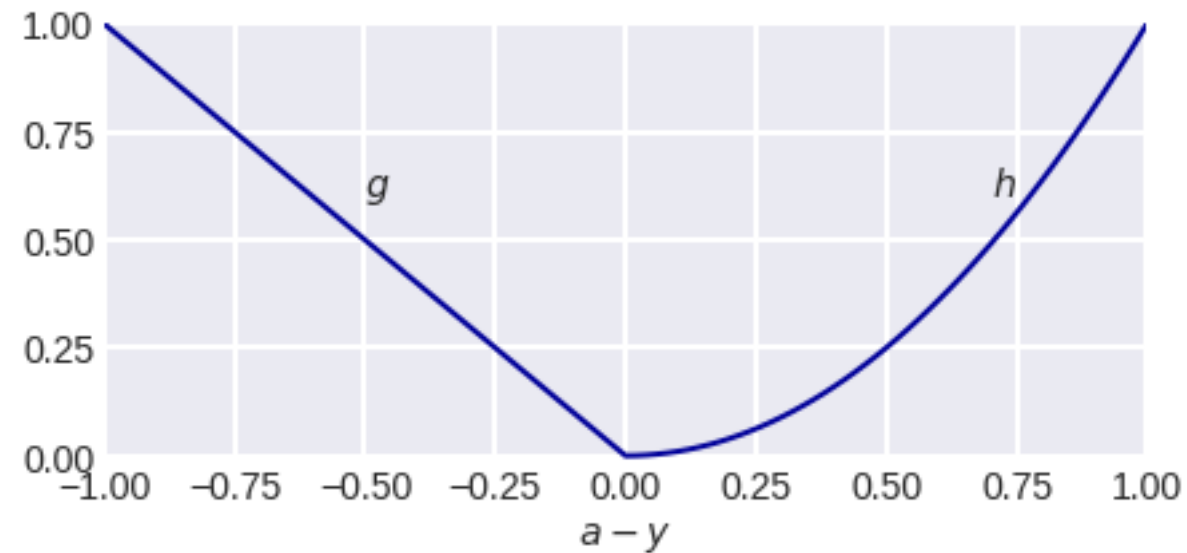
С точностью до порога



**Д3 Реализуйте многомерную линейную регрессию, оптимизирующую  $\epsilon_V$ .**

## Несимметричные функции потерь

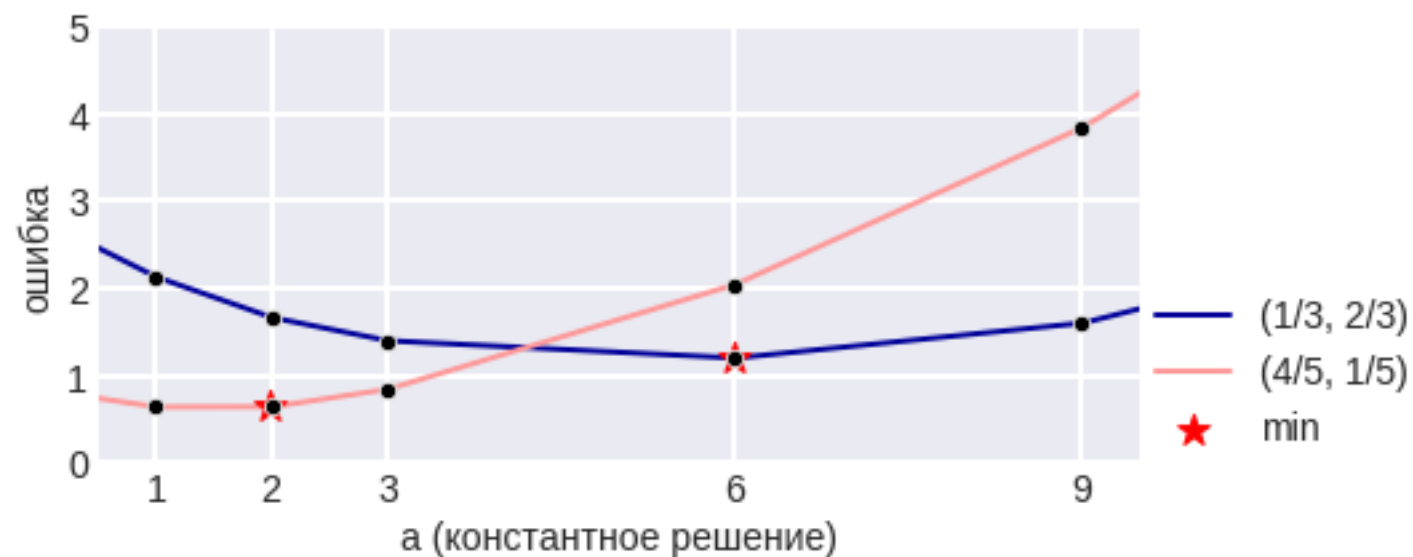
$$\frac{1}{m} \sum_{i=1}^m \begin{cases} g(|y_i - a_i|), & y_i < a_i, \\ h(|y_i - a_i|), & y_i \geq a_i, \end{cases}$$



**Зачем нужны такие функции?**



## Несимметричные функции потерь



$$\frac{1}{m} \sum_{i=1}^m \begin{cases} k_1 |y_i - a_i|, & y_i < a_i, \\ k_2 |y_i - a_i|, & y_i \geq a_i, \end{cases}$$

**Д3 Реализуйте многомерную линейную регрессию, оптимизирующую такую функцию.**

## Совет

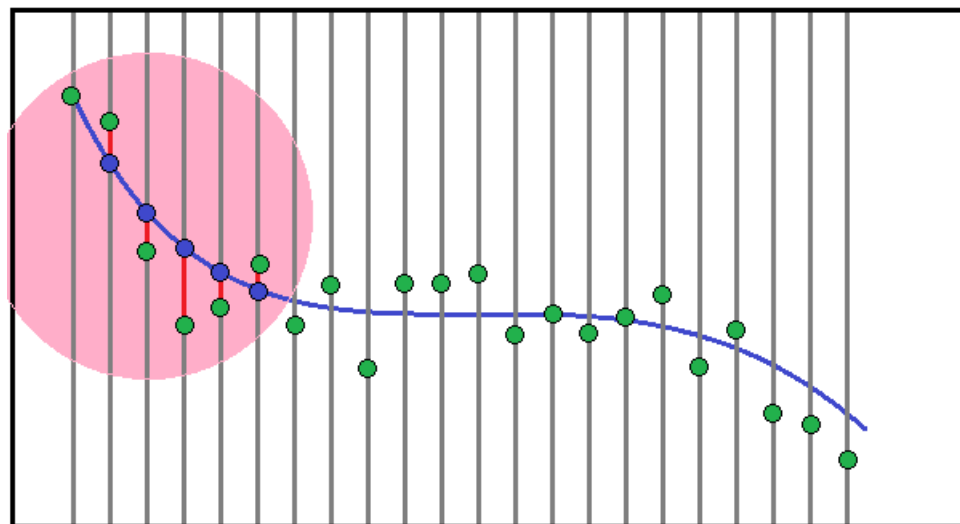
**Функции ошибок иногда и классные признаки...**

**Пример:** в Casualty придумываем бенчмарки  
(восстановление одной переменной по другой),  
признаки – их относительные ошибки,  
т.к. абсолютные брать нельзя

**Почему?**

**Совет**

**Аналогично во многих задачах с сигналами...**



**Признак – не только коэффициенты в приближении,  
но и ошибка приближения!**

**~ отклонение от типичного поведения**

## Монотонное изменение функции ошибки

**Формально задачи эквивалентные:**

$$\text{MSE} \rightarrow \min$$

$$\frac{1}{m} \sum_{i=1}^m |a - y_i|^2 \rightarrow \min$$

$$\text{RMSE} \rightarrow \min$$

$$\sqrt{\frac{1}{m} \sum_{i=1}^m |a_i - y_i|^2} \rightarrow \min$$

**Решения на практике могут отличаться...**

**В методе градиентного спуска разные производные**

$$\frac{\partial \text{MSE}}{\partial a} = \frac{2}{m} \sum_{i=1}^m (a - y_i)$$

$$\frac{\partial \text{RMSE}}{\partial a} = \frac{1}{m \text{RMSE}} \sum_{i=1}^m (a_i - y_i)$$

**Д3 На что это влияет на практике? что лучше минимизировать?**

**Рассмотреть ещё подобные случаи в ML!**

## Метрики в регрессии: минутка кода

```
from sklearn.metrics import r2_score
from sklearn.metrics import mean_absolute_error
from sklearn.metrics import mean_squared_error
from sklearn.metrics import mean_squared_log_error
from sklearn.metrics import median_absolute_error
from sklearn.metrics import explained_variance_score

# R^2
print (r2_score(y, a),
       1 - np.mean((y - a) ** 2) / np.mean((y - np.mean(y)) ** 2))

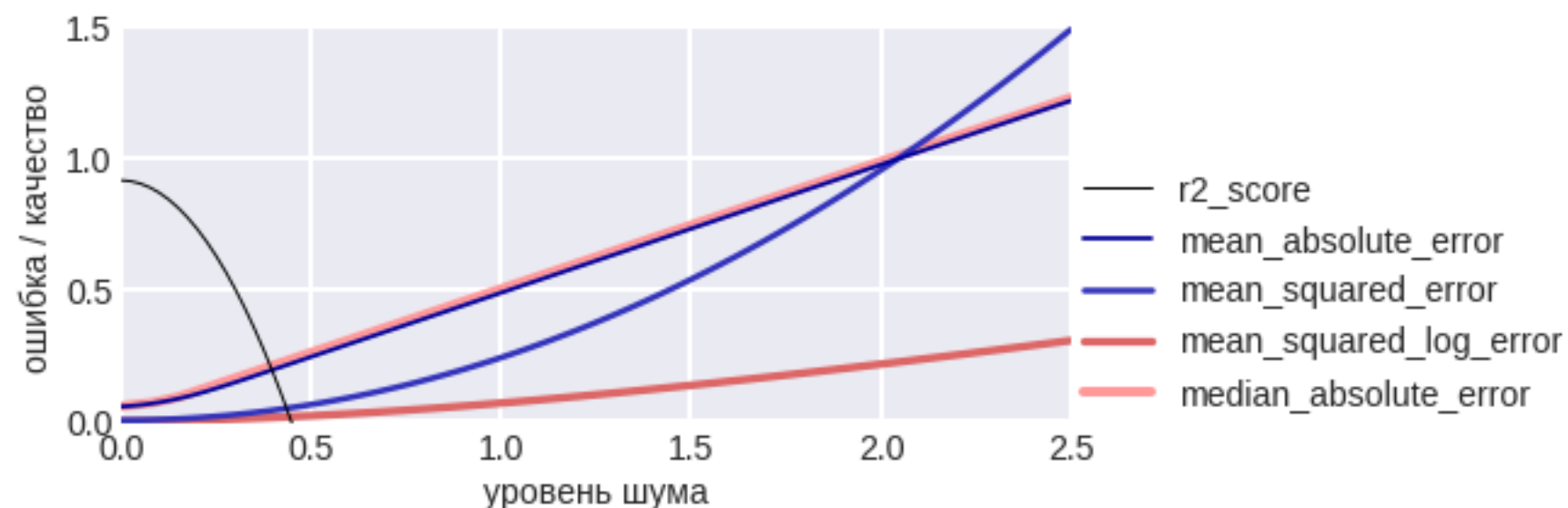
# MAE
print (mean_absolute_error(y, a),
       np.mean(np.abs(y - a)))

# MSE
print (mean_squared_error(y, a),
       np.mean((y - a) ** 2))

# MSLp1E
print (mean_squared_log_error(y, a),
       np.mean((np.log1p(y) - np.log1p(a)) ** 2))

# MedAE
print (median_absolute_error(y, a),
       np.median(np.abs(y - a)))
```

## Сравнение метрик в одном эксперименте



**Д3 Что за эксперимент? Почему ошибки ведут себя так?  
(попробовать восстаносить)**

**Д3 Как число фолдов влияет на CV-оценку ошибки?**

**Д3 Как шум влияет на выбор оптимального решения?**

## Итоги

**Функции ошибки имеют вероятностное обоснование**  
(через правдоподобие)

**средний модуль отклонения MAE (MAD)**

**средний квадрат отклонения MSE**

**+ RMSE, коэффициент детерминации  $R^2$ , функция Хьюбера, Logcosh**  
**Можно невероятно обосновать для малых отклонений**

**Иногда попадают обобщения MAE и RMSE**

## Итоги

**Процентные функции ошибок**  
(SMAPE, MAPE, PMAD)

**Основанные на сравнении с бенчмарком**  
(MRAE, REL\_MAE, PB)

**Нормированные ошибки**  
(MASE)

**Несимметричные ошибки**

**Ошибки с точностью до порога**

**Есть нетрадиционные применения функций ошибок**  
для генерации признаков



## Литература

**Стрижов В.В. Функция ошибки в задачах восстановления регрессии //**  
**Заводская лаборатория, 2013, 79(5): 65-73.**

<http://strijov.com/papers/Strijov2012ErrorFn.pdf>

**«How to Win a Data Science Competition: Learn from Top Kagglers»**

<https://ru.coursera.org/learn/competitive-data-science>