

# Лекция 4.

# Линейные модели

Введение в машинное обучение

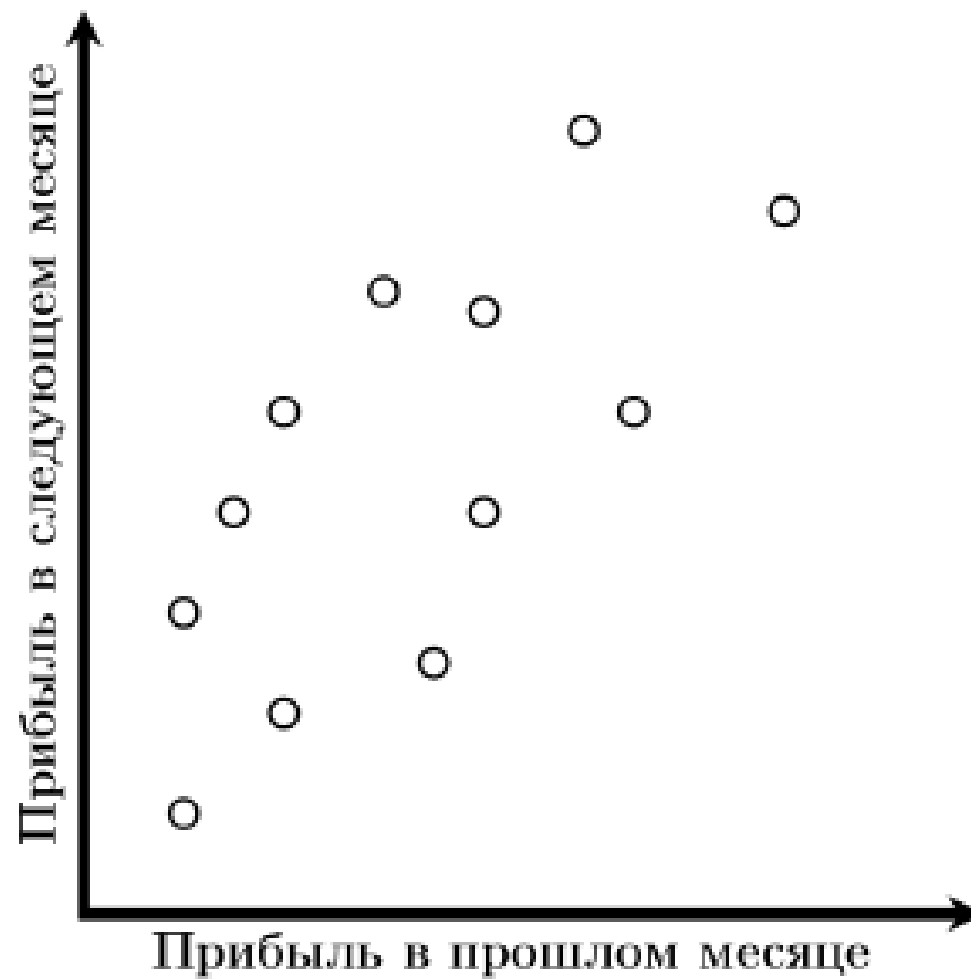
# Обозначения

- $X$  — пространство объектов
- $Y$  — пространство ответов
- $x = (x^1, \dots, x^d)$  — признаковое описание объекта
- $X = (x_i, y_i)_{i=1}^l$  — обучающая выборка
- $a(x)$  — алгоритм, модель
- $Q(a, X)$  — функционал ошибки алгоритма  $a$  на выборке  $X$
- Обучение:  $a(x) = \operatorname{argmin}_{a \in A} Q(a, X)$
- $Y = \mathbb{R}$

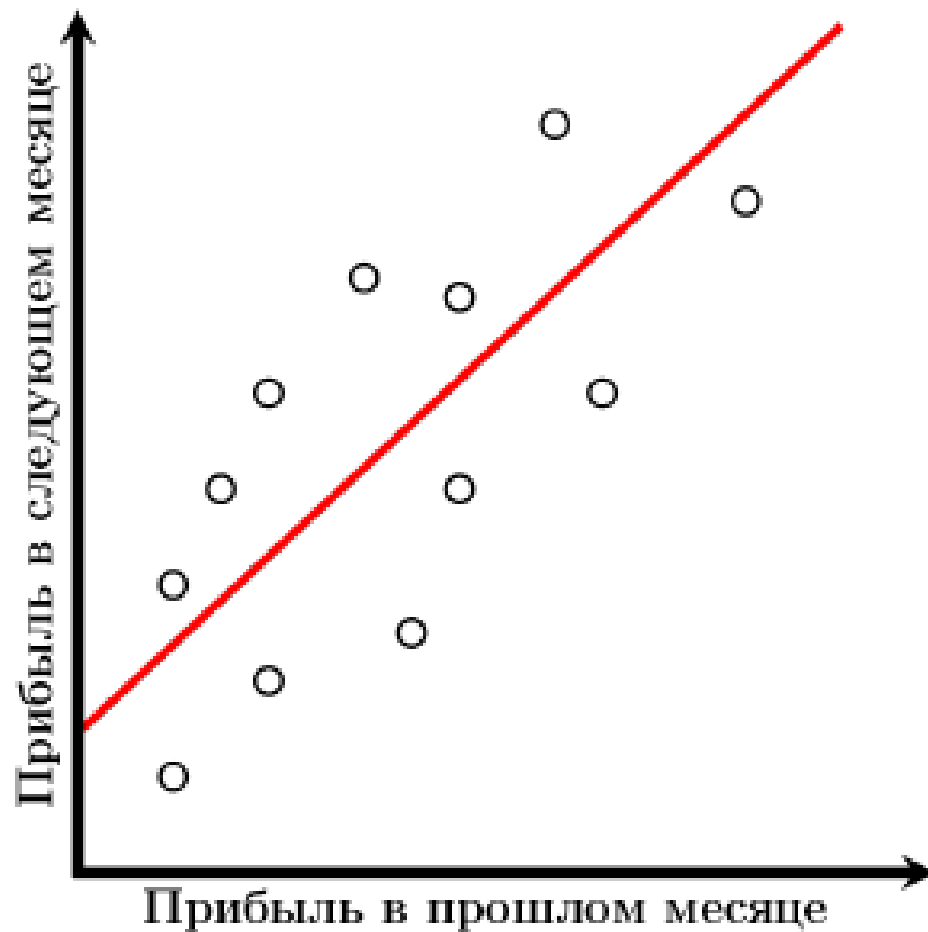
# Напоминание

- **Функционал ошибки  $Q$ :** способ измерения того, хорошо или плохо работает алгоритм на конкретной выборке
- **Семейство алгоритмов  $A$  :** как выглядит множество алгоритмов, из которых выбирается лучший
- **Метод обучения :** как именно выбирается лучший алгоритм из семейства алгоритмов.

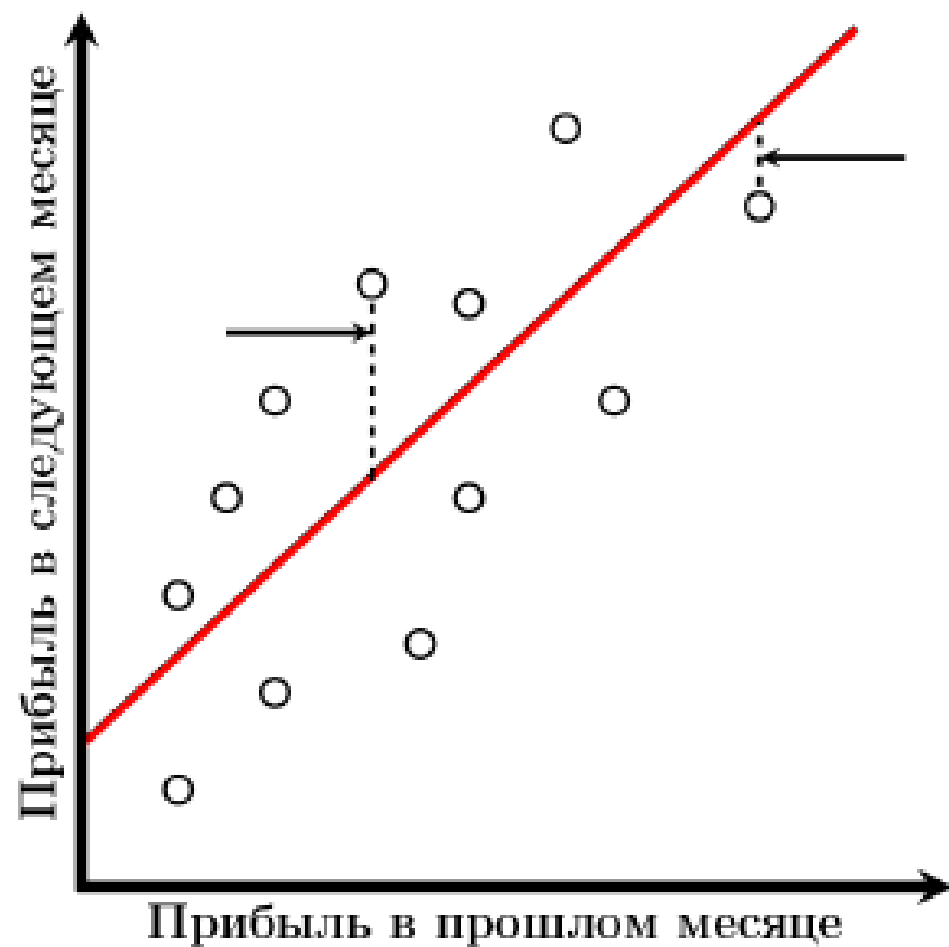
# Задача регрессии



# Задача регрессии



# Задача регрессии



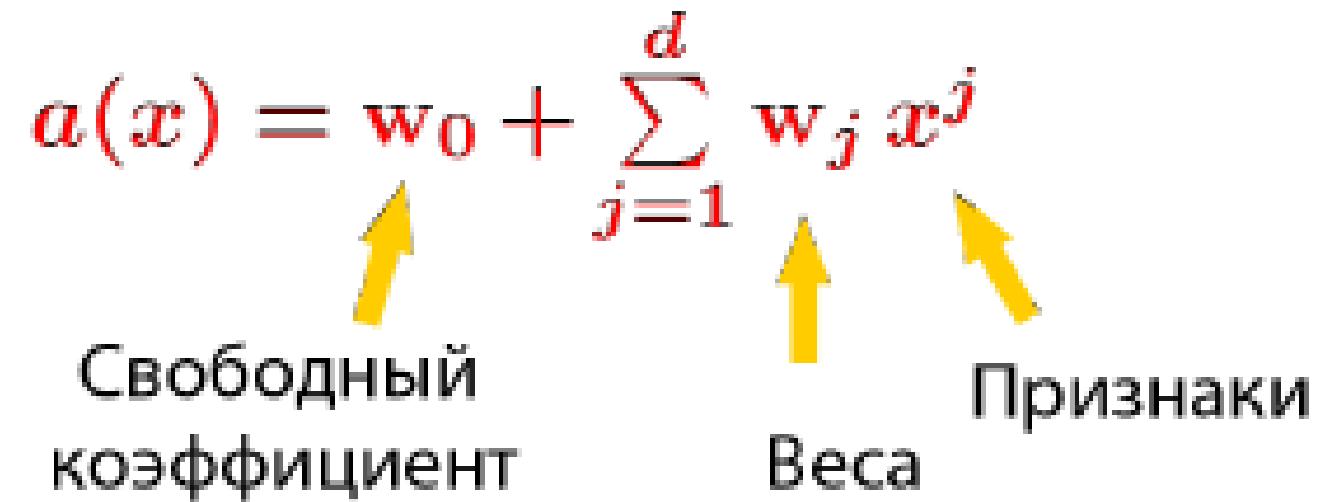
# Семейство алгоритмов

$$a(x) = w_0 + \sum_{j=1}^d w_j x^j$$

Свободный коэффициент

Весы

Признаки



# Семейство алгоритмов

$$a(x) = w_0 + \sum_{j=1}^d w_j x_j$$

$$a(x) = \sum_{i=0}^m w_i x_i = \langle w, x \rangle$$

$$y = XW + \varepsilon$$



# Линейная регрессия

$$y = \sum_{j=0}^m X_j W_j + \varepsilon$$
$$y_i = \sum_{j=0}^m w_j X_{ij} + \varepsilon_i$$

- где  $y \in R^n$  — целевая переменная
- $w \in R^{m+1}$  — вектор параметров модели(весы)
- $X$  — матрица наблюдений и признаков размерности  $n$  строк на  $m + 1$  столбцов с полным рангом по столбцам:  $rank(X) = m + 1$
- $\varepsilon$  — случайная переменная, соответствующая случайной, непрогнозируемой ошибке модели

# Линейная регрессия

На модель накладываются ограничения:

1. Матожидание случайных ошибок равно нулю:  $\forall i: E[\varepsilon_i] = 0$
2. Дисперсия случайных ошибок одинакова и конечна:  
 $\forall i: Var(\varepsilon_i) = \sigma^2 < \infty$
3. Случайные ошибки не скоррелированы:  $\forall i \neq j: Cov(\varepsilon_i, \varepsilon_j) = 0$

# Линейная оценка

- Оценка весов называется линейной если :
- $\widehat{w}_i = w_{1i}y_1 + w_{2i}y_2 + \dots + w_{1n}y_n$ ,
- Где  
 $\forall k w_{ki}$  зависит только от наблюдаемых данных  $X$  и почти наверняка нелинейно
- Так как решением задачи поиска оптимальной весов будет именно линейная оценка, то и модель называется линейной регрессией
- Один из способов вычислить значения параметров является метод наименьших квадратов (МНК), который минимизирует среднеквадратичную ошибку между реальным значением зависимой переменной и прогнозом, выданной моделью

# Метод наименьших квадратов

$$\begin{aligned}\mathcal{L}(\mathbf{X}, \mathbf{y}, \mathbf{w}) &= \frac{1}{2n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \\ &= \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \\ &= \frac{1}{2n} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})\end{aligned}$$

# Шпаргалка по матричным производным

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{X}^T \mathbf{A} = \mathbf{A}$$

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{X}^T \mathbf{A} \mathbf{X} = (\mathbf{A} + \mathbf{A}^T) \mathbf{X}$$

$$\frac{\partial}{\partial \mathbf{A}} \mathbf{X}^T \mathbf{A} \mathbf{y} = \mathbf{X}^T \mathbf{y}$$

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{A}^{-1} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \mathbf{X}} \mathbf{A}^{-1}$$

# МНК. Дифференцирование

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \frac{1}{2n} (\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}) \\ &= \frac{1}{2n} (-2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{w})\end{aligned}$$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 &\Leftrightarrow \frac{1}{2n} (-2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{w}) = 0 \\ &\Leftrightarrow -\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \mathbf{w} = 0 \\ &\Leftrightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \\ &\Leftrightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

# Аналитическое решение

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

- Нужно обращать матрицу  $d \times d$  – сложность  $d^3$
- Могут возникнуть численные проблемы

# Обучение линейной регрессии

$$Q(\mathbf{w}, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2 \rightarrow \min_{\mathbf{w}}$$

- $d$  неизвестных
- Есть константный признак
- Выпуклая функция



# Матричная запись

- Матрица объекты-признаки

$$X = \begin{pmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{\ell 1} & \cdots & x_{\ell d} \end{pmatrix} \text{ Объект}$$

$$X = \begin{pmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{\ell 1} & \cdots & x_{\ell d} \end{pmatrix}$$

Признак

Вектор ответов

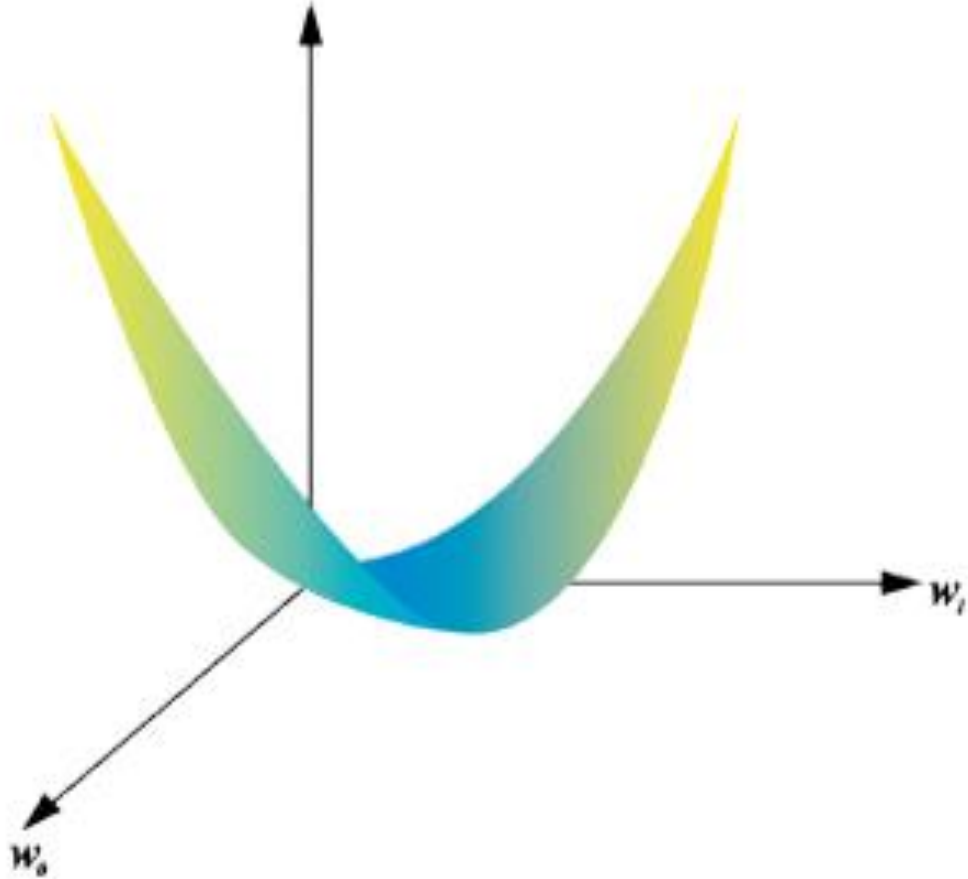
$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_\ell \end{pmatrix}$$

# Матричная запись

$$Q(\mathbf{w}, X) = \frac{1}{\ell} \|X \mathbf{w} - \mathbf{y}\|^2 \rightarrow \min_{\mathbf{w}}$$

# Градиентный спуск

- Функция ошибки гладкая и выпуклая



# Градиентный спуск: алгоритм

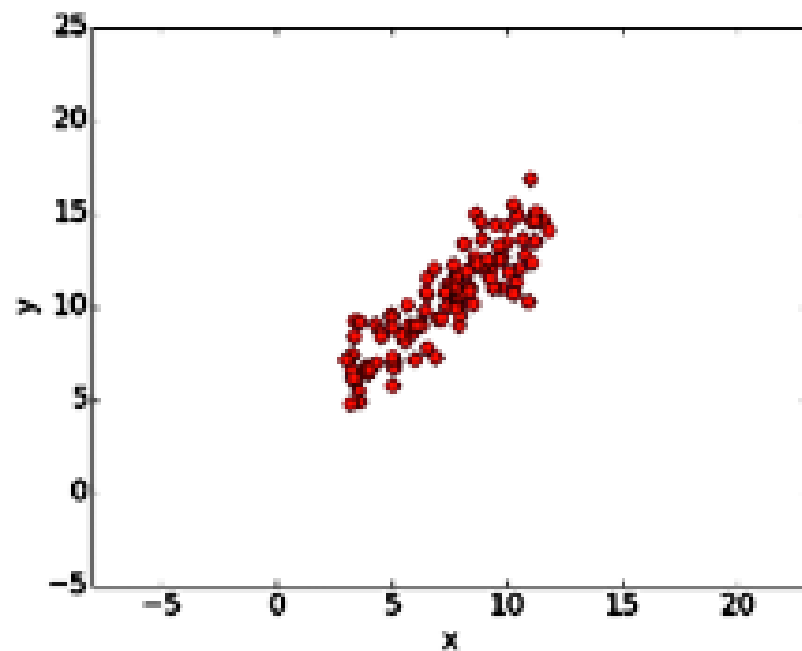
- Инициализация:  $w^0 = 0$
- Цикл по  $t = 1, 2, 3, \dots$ :
- $w^t = w^{t-1} - \eta_t \nabla Q(w^{t-1}, X)$
- Если  $\|w^t - w^{t-1}\| < \varepsilon$ , то завершить

# Парная регрессия

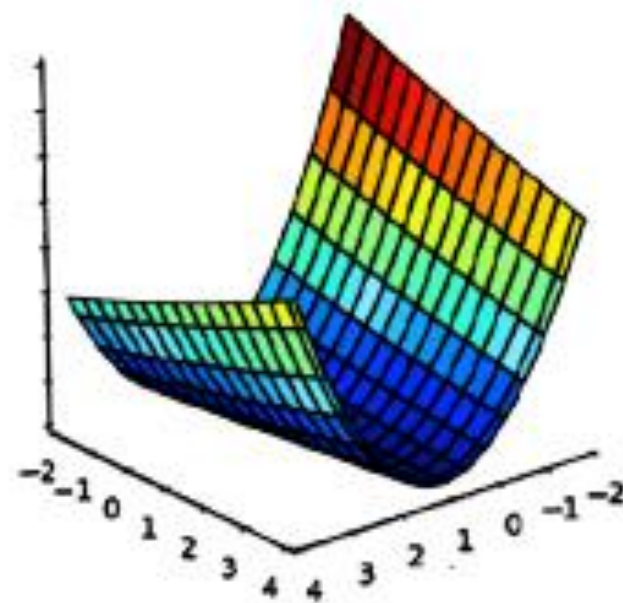
- Простейший случай: один признак
- Модель:  $a(x) = w_1x + w_0$
- Два параметра:  $w_1$  и  $w_0$
- Функционал:

$$Q(w_0, w_1, X) = \frac{1}{l} \sum_{i=1}^l (w_1x_i + w_0 - y_i)^2$$

# Парная регрессия



Выборка



Функционал качества

# Градиентный спуск

- Инициализация:  $w^0 = 0$
- Цикл по  $t = 1, 2, 3, \dots$ :
- $w^t = w^{t-1} - \eta_t \nabla Q(w^{t-1}, X)$
- Если  $\|w^t - w^{t-1}\| < \varepsilon$ , то завершить

# Градиент для парной регрессии

$$Q(w_0, w_1, X) = \frac{1}{l} \sum_{i=1}^l (w_1 x_i + w_0 - y_i)^2$$

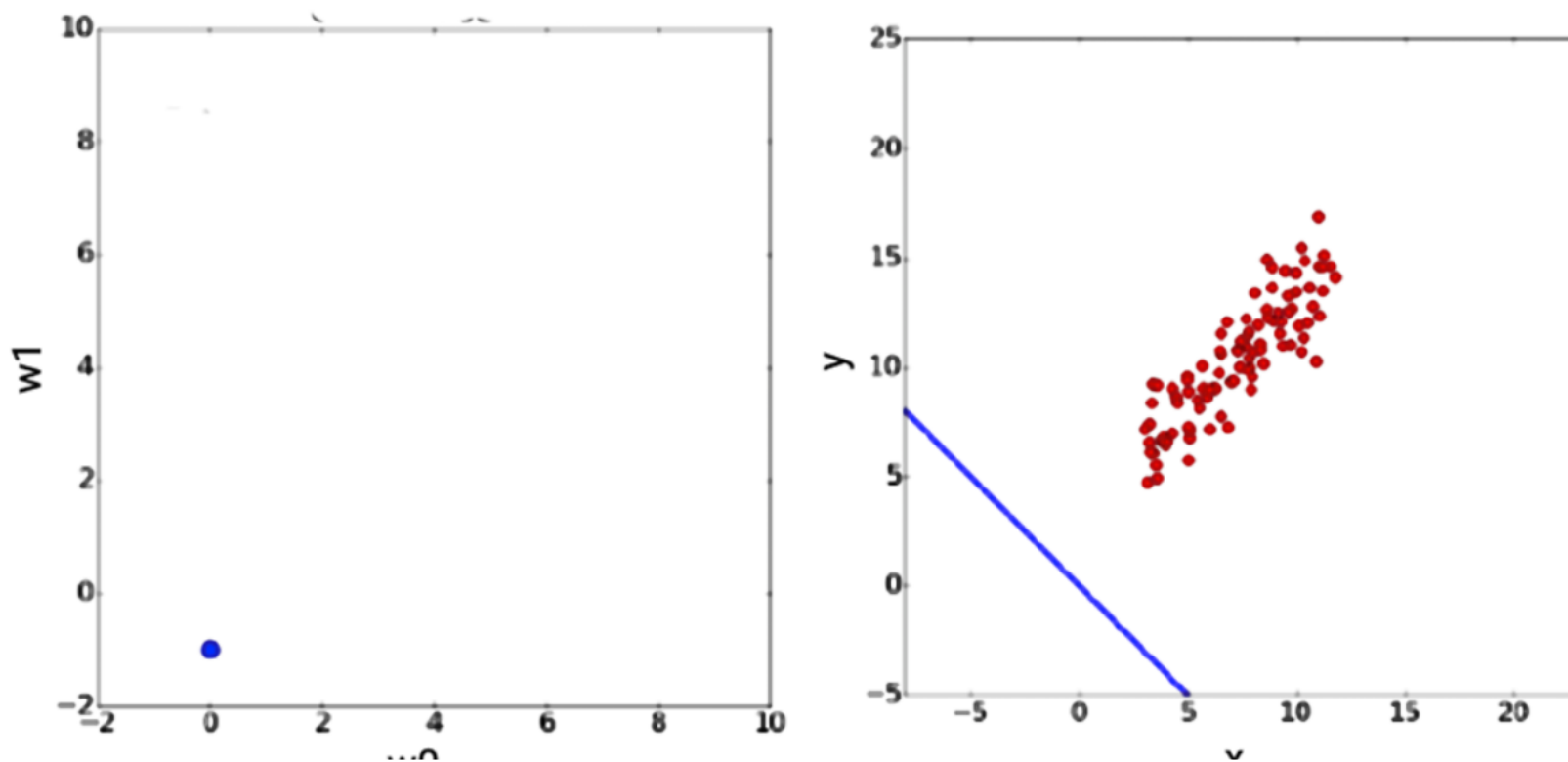
- Частные производные:

$$\frac{\partial Q}{\partial w_1} = \frac{2}{l} \sum_{i=1}^l (w_1 x_i + w_0 - y_i) x_i$$

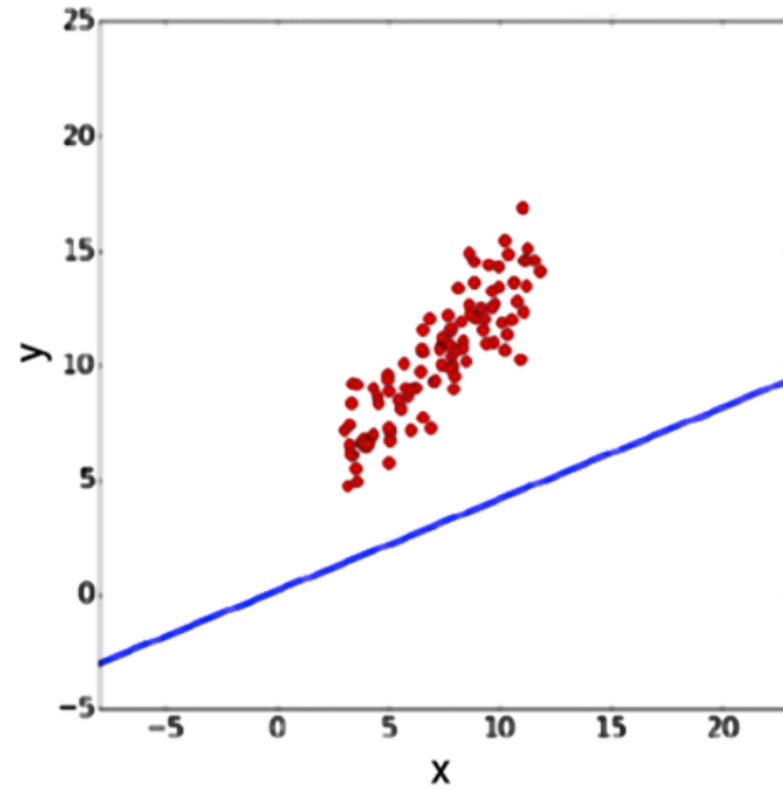
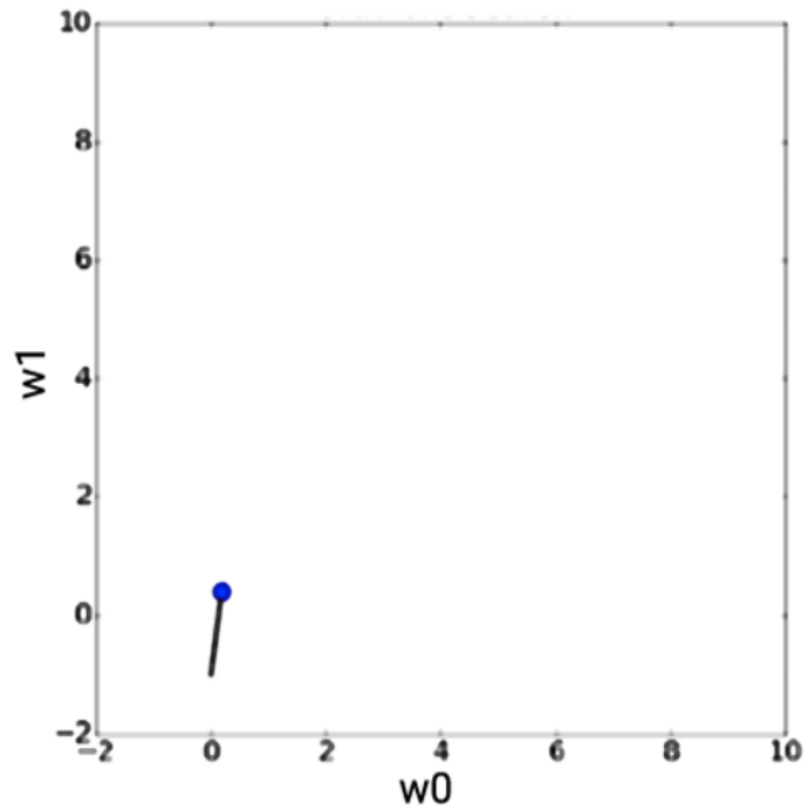
$$\frac{\partial Q}{\partial w_0} = \frac{2}{l} \sum_{i=1}^l (w_1 x_i + w_0 - y_i)$$



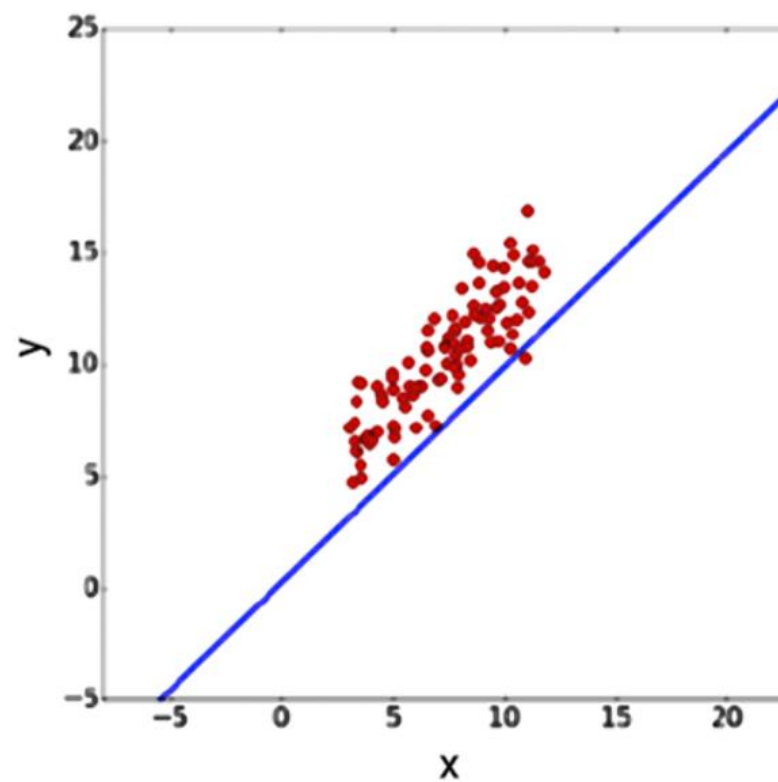
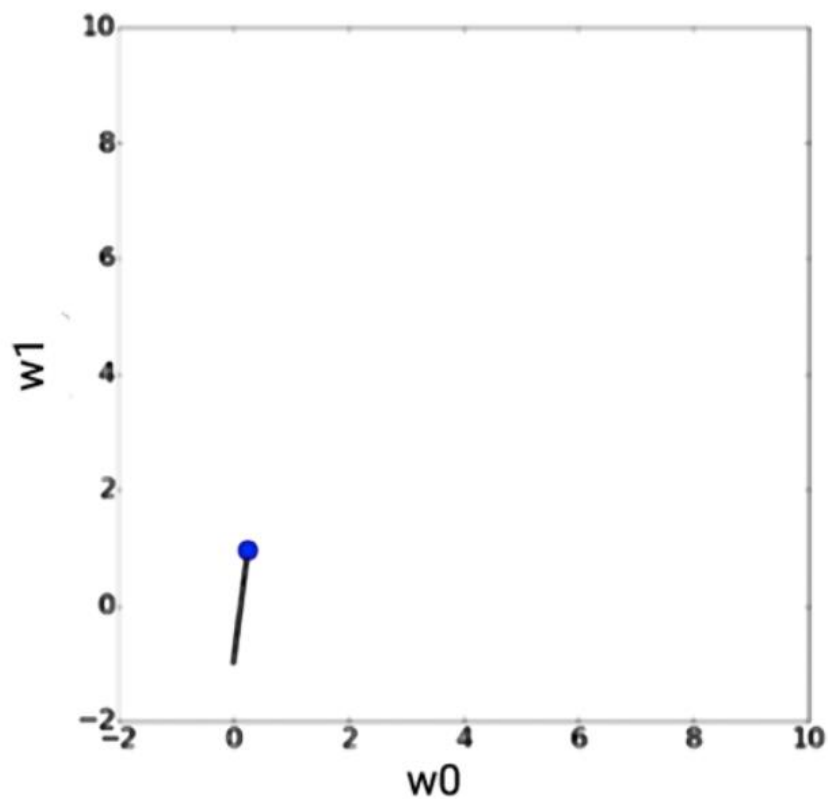
# Парная регрессия – 1 итерация



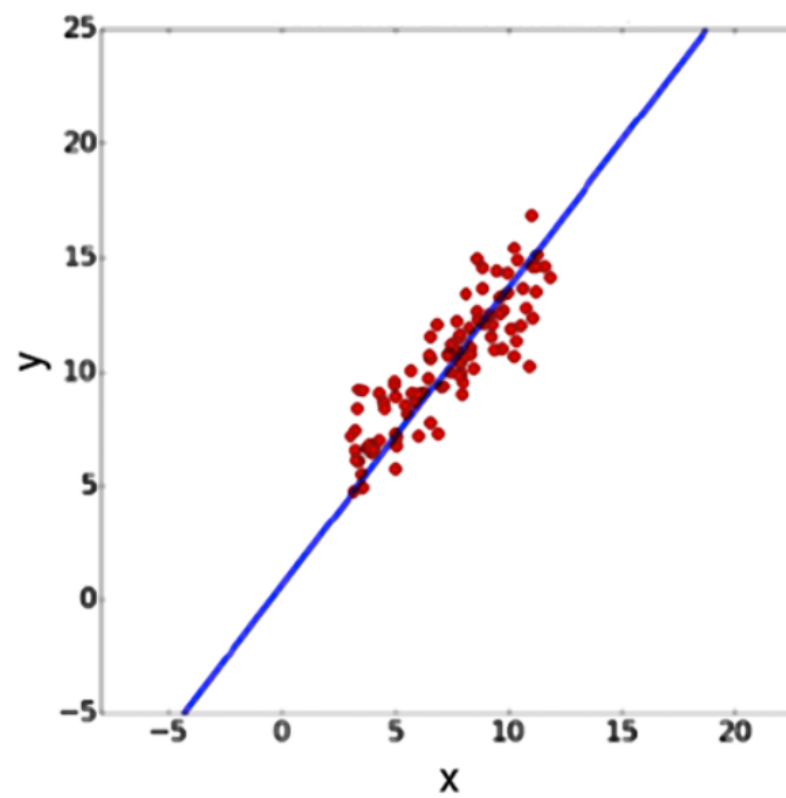
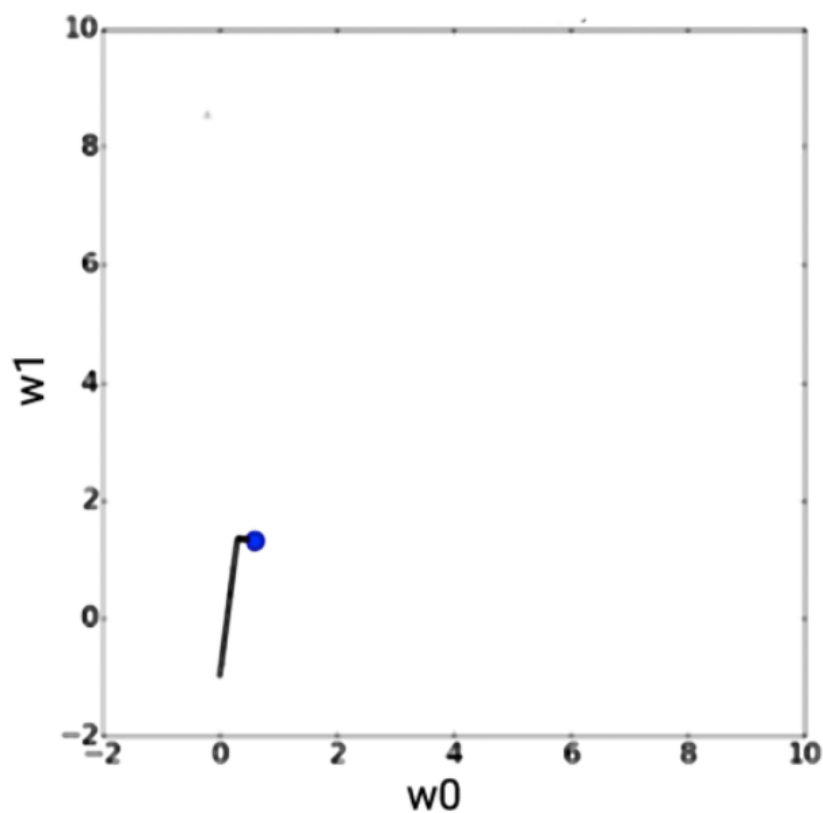
# Парная регрессия – 2 итерация



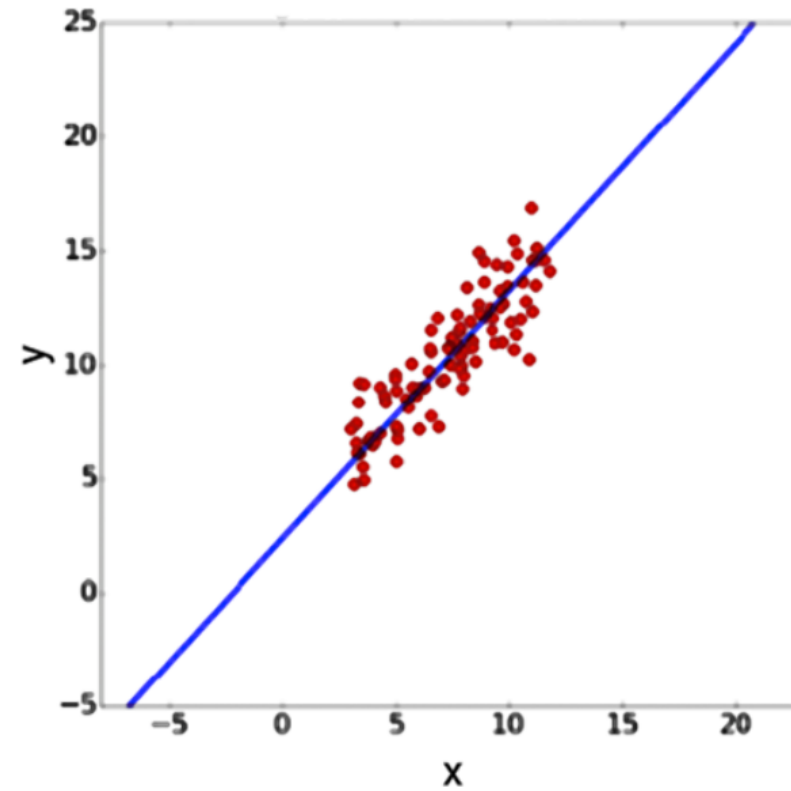
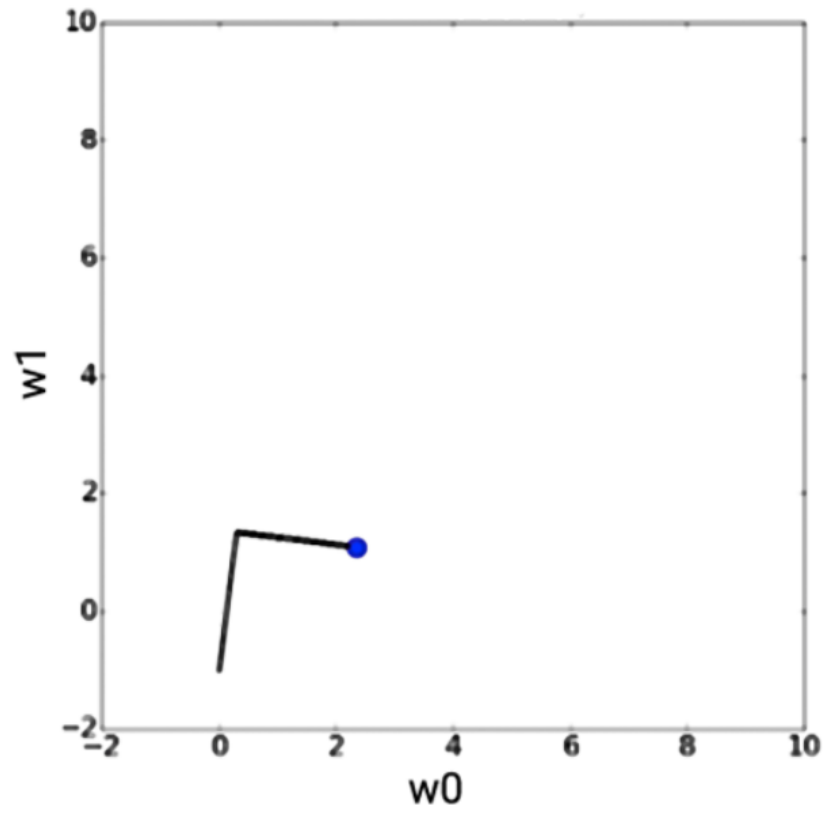
# Парная регрессия – 3 итерация



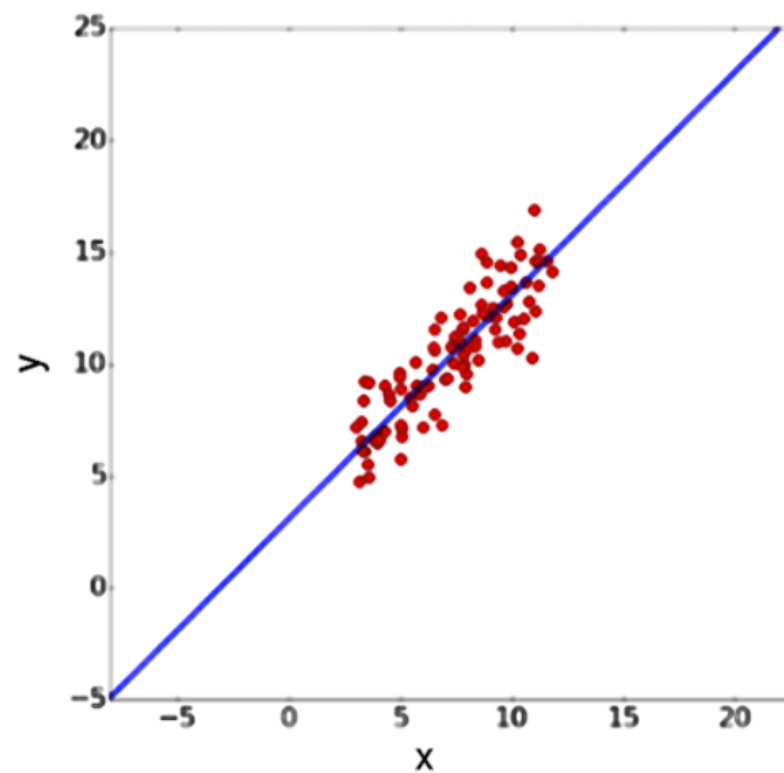
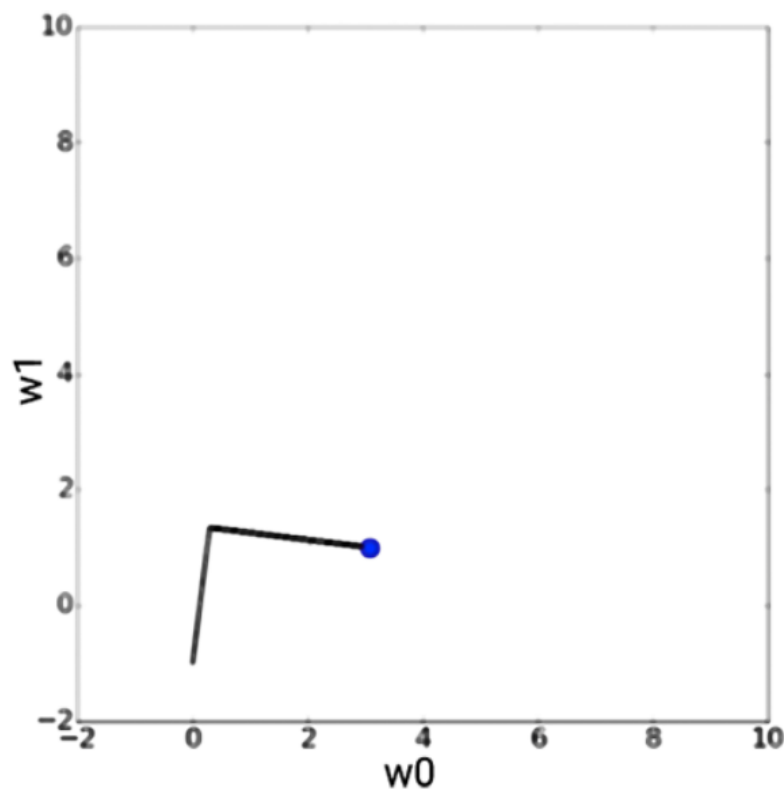
# Парная регрессия – 4 – я итерация



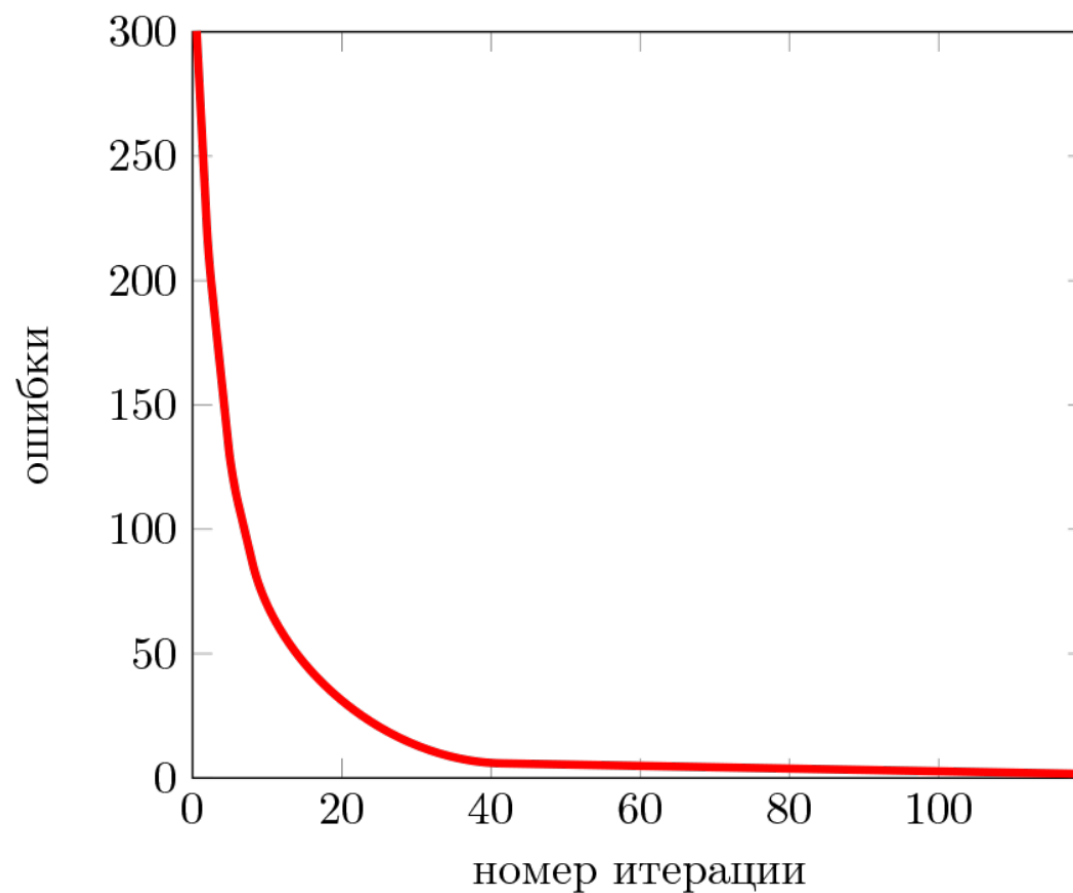
# Парная регрессия – 5-я итерация



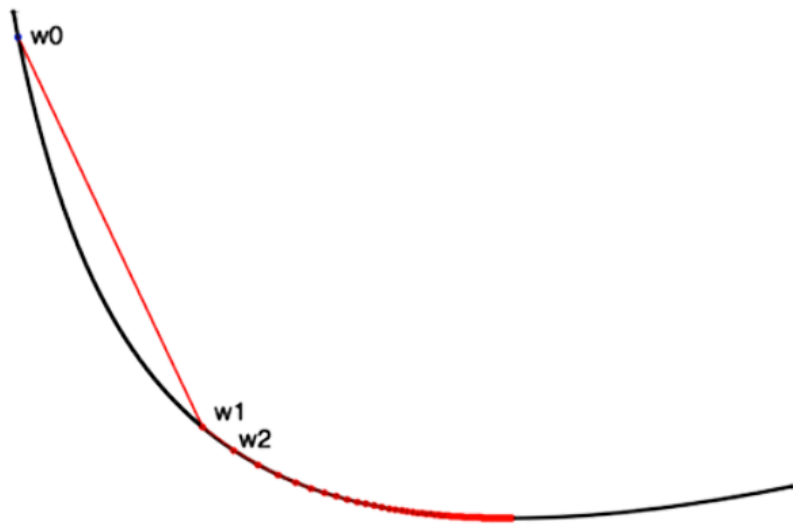
# Парная регрессия – 6-я итерация



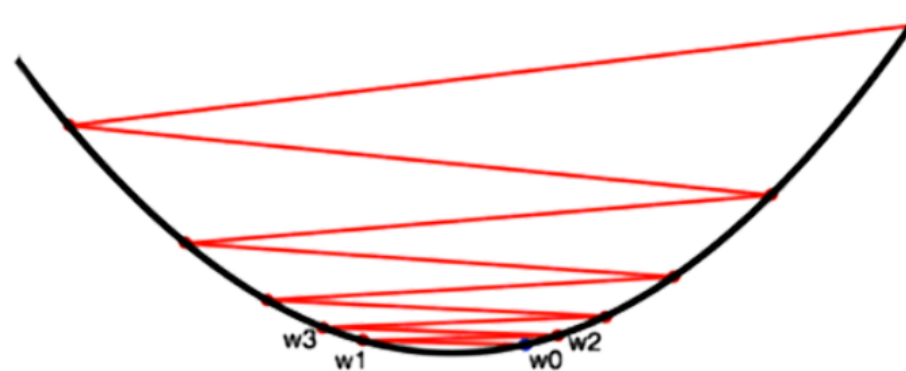
# Функционал качества



# Размер шага



Маленький шаг

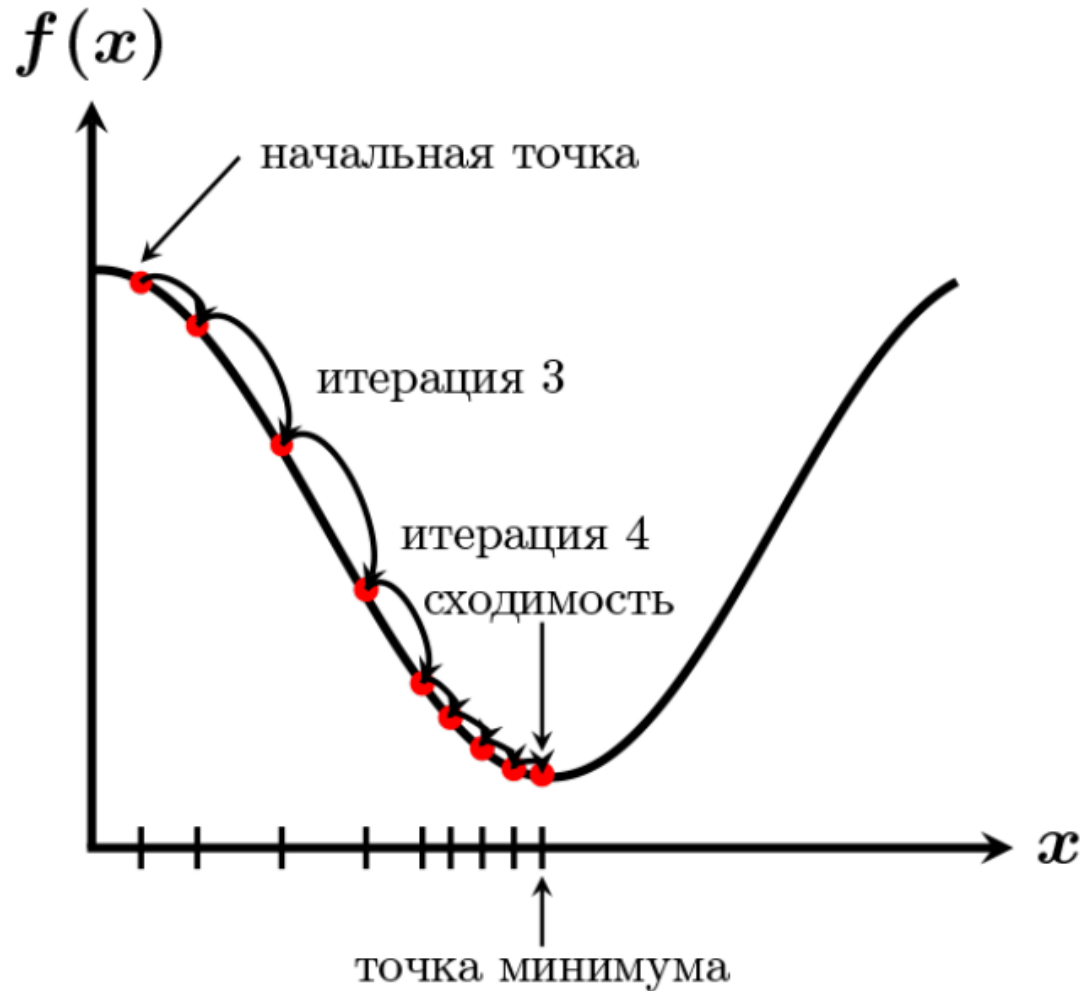


Большой шаг



# Размер шага

- Размер шага, гиперпараметр, нужно подбирать
- Обычно пользуются эвристиками
- Чем ближе к минимуму, тем меньше надо шагать



# Многомерная линейная регрессия

$$Q(\mathbf{w}, X) = \frac{1}{\ell} \|X \mathbf{w} - \mathbf{y}\|^2 \rightarrow \min_{\mathbf{w}}$$

- Градиент:

$$\nabla_{\mathbf{w}} Q(\mathbf{w}, X) = \frac{2}{\ell} X^T (X \mathbf{w} - \mathbf{y})$$

# Стохастический градиентный спуск

- Инициализация:  $w^0 = 0$
- Цикл по  $t = 1, 2, 3, \dots$ :
- $w^t = w^{t-1} - \eta_t \nabla Q(w^{t-1}, X)$
- Если  $\|w^t - w^{t-1}\| < \varepsilon$ , то завершить

# Градиент функционала

$$\nabla_{\mathbf{w}} Q(\mathbf{w}, X) = \frac{2}{\ell} X^T (X \mathbf{w} - \mathbf{y})$$

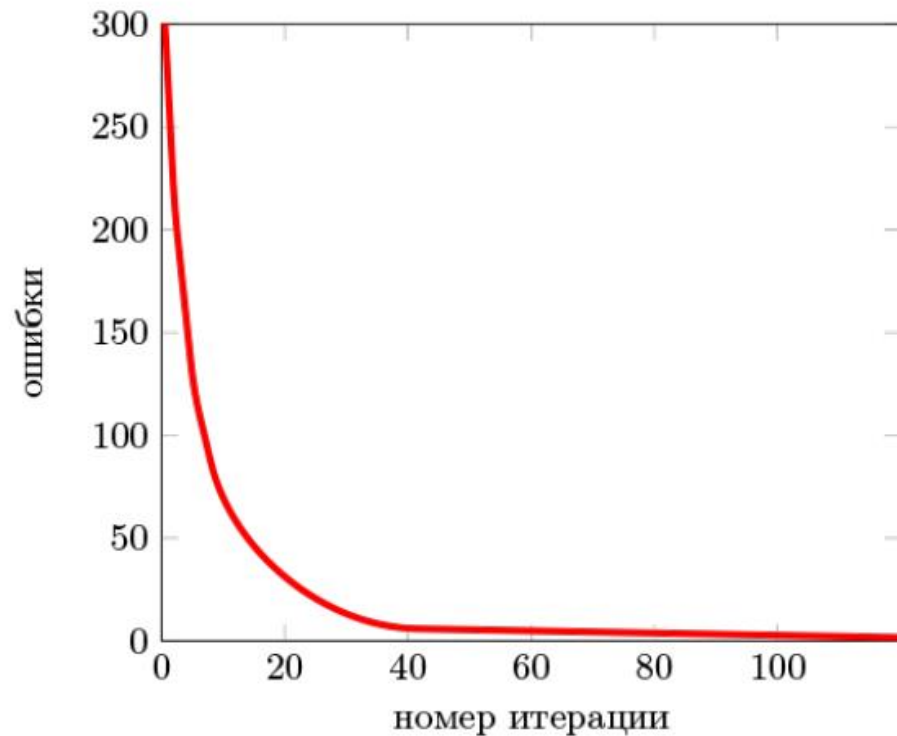
$$\frac{\partial Q}{\partial w_j} = \frac{2}{\ell} \sum_{i=1}^{\ell} x_i^j (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)$$

Градиентный спуск требует вычисления полного градиента!

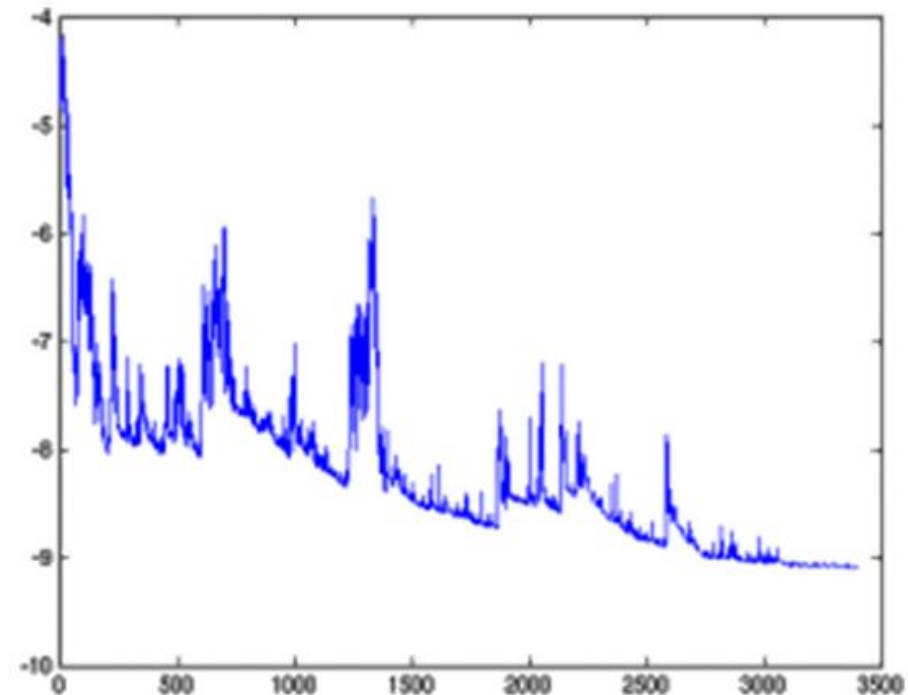
# Стохастический градиентный спуск

- Инициализация:  $w^0 = 0$
- Цикл по  $t = 1, 2, 3, \dots$ :
- Выбрать случайный объект  $x_i$  из  $X$
- $w^t = w^{t-1} - \eta_t \nabla Q(w^{t-1}, \{x_i\})$
- Если  $\|w^t - w^{t-1}\| < \varepsilon$ , то завершить

# Стохастический градиентный спуск



Градиентный спуск



Стохастический  
градиентный спуск

## Преимущества SGD

- Быстрее выполняется один шаг
- Не требует хранения выборки в памяти
- Подходит для онлайн обучения

# Метод максимального правдоподобия

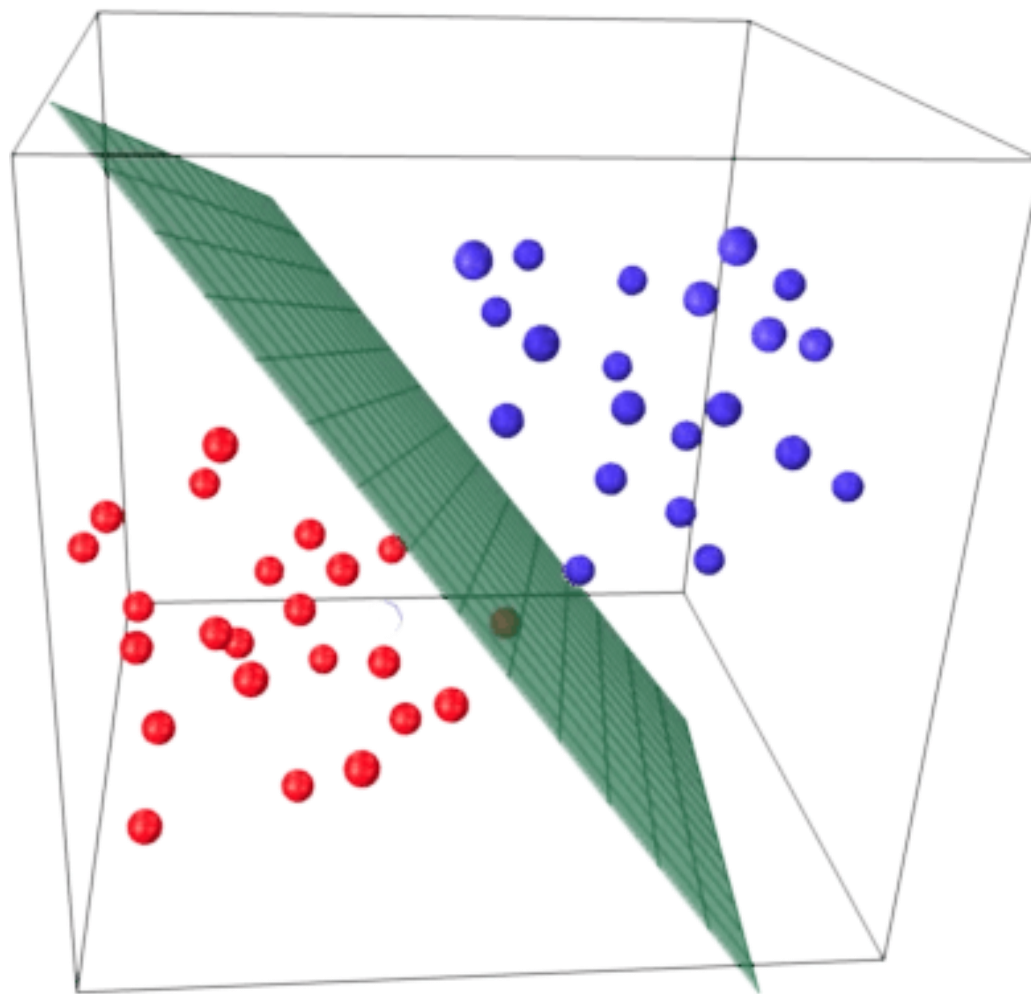
- Интуитивная оценка вероятности  $117/400 = 29\%$
- Такая оценка является оценкой максимального правдоподобия
- Разберемся откуда берется данная оценка, вспомни  
Распределение Бернулли : случайная величина  $X$  имеет  
распределение Бернулли, если она принимает всего два значения  
(1 и 0 с вероятностями  $\theta$  и  $1 - \theta$ ) и имеет следующую функцию  
распределения вероятности :

$$p(\theta, x) = \theta^x (1 - \theta)^{(1-x)}, x \in \{0, 1\}$$



Ноутбук Jupyter

# Линейный классификатор



# Линейный классификатор

$$a(x) = \text{sign} \left( w_0 + \sum_{j=1}^d w_j x^j \right)$$

Свободный коэффициент

Весы

Признаки

# Линейный классификатор

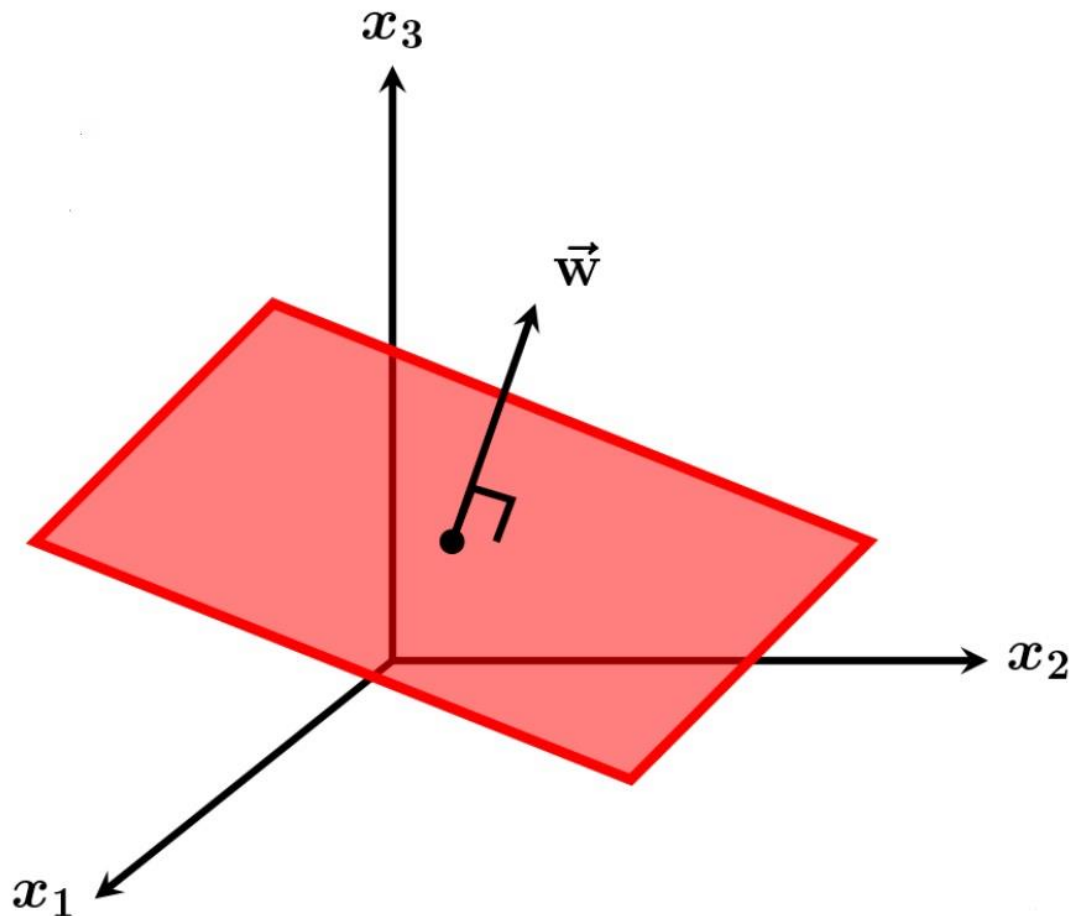
Добавим единичный признак:

$$a(x) = \text{sign} \sum_{j=1}^{d+1} w_j x^j = \text{sign} \langle w, x \rangle$$

# Уравнение гиперплоскости

$$\langle \mathbf{w}, \mathbf{x} \rangle = 0$$

Уравнение  
гиперплоскости



# Расстояние от гиперплоскости

- Чтобы получить расстояние до гиперплоскости, нужно решить почти классическую задачу из курса линейной алгебры: найти расстояние от точки с радиус вектором  $x_A$  до плоскости, которая задается уравнением  $w^T x = 0$

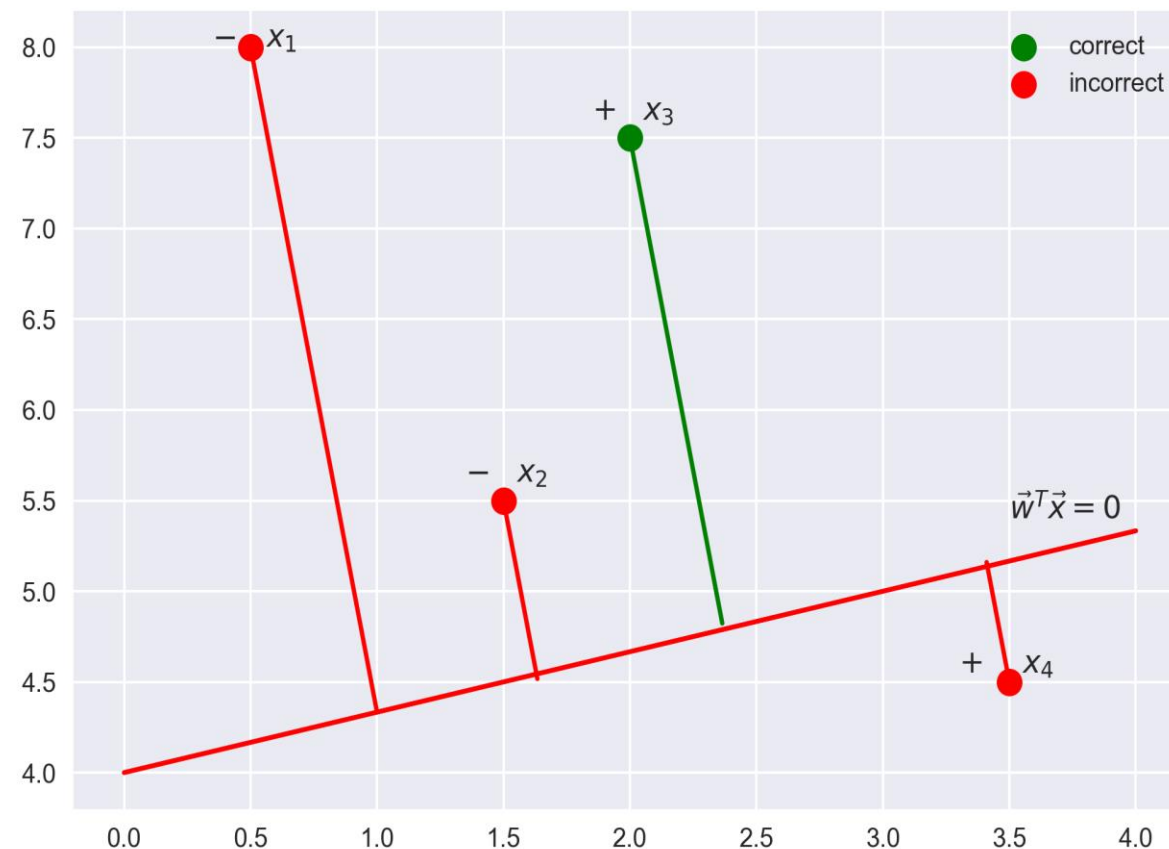
$$\rho(x_A, w^T x = 0) = \frac{w^T x_A}{\|w\|}$$

Когда получим ответ, то поймем, что чем больше по модулю выражение  $w^T x_i$ , тем дальше точка  $x_i$  находится от плоскости  $w^T x = 0$

# Отступ

$M(x_i) = y_i w^T x_i$  — уверенность модели в классификации объекта  $x_i$

- если отступ большой (по модулю) и положительный, это значит, что метка класса поставлена правильно, На рисунке —  $x_3$ .
- если отступ большой (по модулю) и отрицательный, значит метка класса поставлена неправильно (скорее всего такой объект — аномалия, например, его метка в обучающей выборке поставлена неправильно). На рисунке —  $x_1$ .
- если отступ малый (по модулю), то объект находится близко к разделяющей гиперплоскости, а знак отступа определяет, правильно ли объект классифицирован. На рисунке —  $x_2$  и  $x_4$ .



# Логистическая регрессия

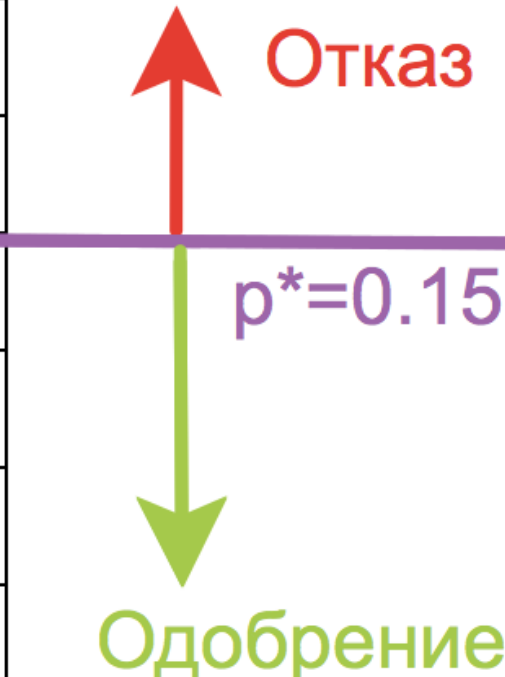
Логистическая регрессия является частным случаем линейного классификатора, но она обладает способностью оценивать вероятность принадлежности объекта  $x_i$  к классу '+'

$$p_+ = P(y_i = 1 \mid \mathbf{x}_i, \mathbf{w})$$



# Бизнес предпосылки

| Клиент  | Вероятность<br>невозврата |
|---------|---------------------------|
| Mike    | 0.78                      |
| Jack    | 0.45                      |
| Larry   | 0.13                      |
| Kate    | 0.06                      |
| William | 0.03                      |
| Jessica | 0.02                      |



Отказ

$p^*=0.15$

Одобрение

# Логистическая регрессия

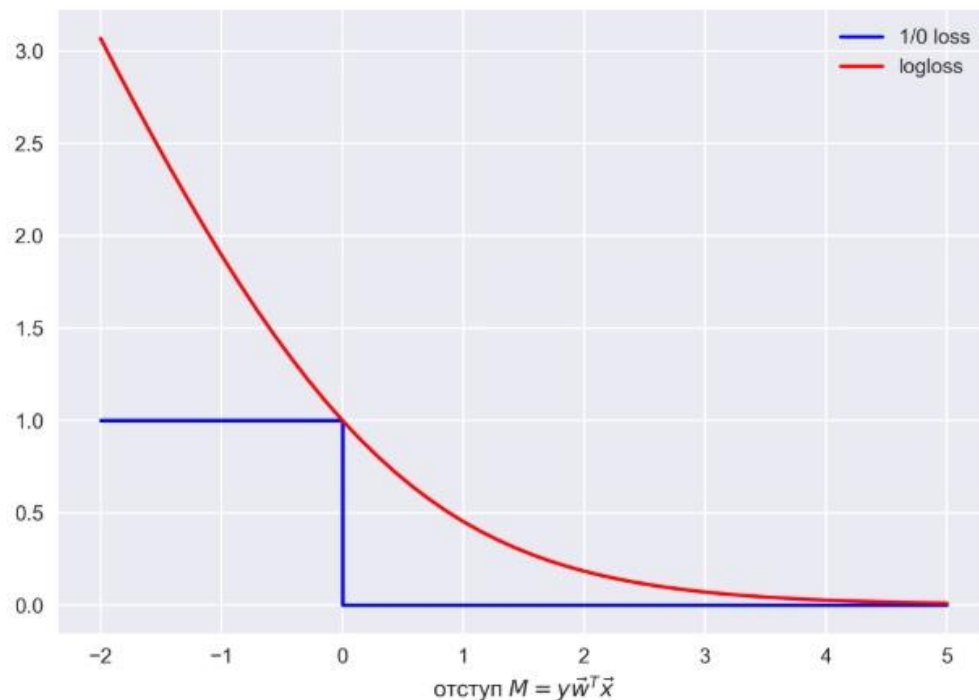
- Мы хотим прогнозировать вероятность  $p_+ \in [0, 1]$
- Умеем строить линейный прогноз с помощью МНК:  $a(x) = w^T x \in R$
- Каким образом преобразовать полученное значение в вероятность, пределы которых –  $[0, 1]$
- Для этого нужна функция  $f : \mathbb{R} \rightarrow [0, 1]$
- В логистической регрессии берется функция сигмоид
- $\sigma(z) = \frac{1}{1+e^{-z}}$

# Логистическая регрессия

$$\mathcal{L}_{log}(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \sum_{i=1}^{\ell} \log(1 + \exp^{-y_i \mathbf{w}^T \mathbf{x}_i}).$$

# Верхняя оценка

- Посмотрим на новую функцию как функцию от отступа :  $L(M) = \log(1 + e^{-M})$
- Нарисуем график 1/0 функций потерь (zero one loss), которая просто штрафует модель на 1 за ошибку на каждом объекте (отступ отрицательный) :  $L_{1/0}(M) = [M < 0]$



# Верхняя оценка

- Картинка отражает общую идею, что в задаче классификации, не умея напрямую минимизировать число ошибок (по крайней мере, градиентными методами это не сделать – производная  $1/0$  функций потерь в нуле обращается в бесконечность), мы минимизируем некоторую ее верхнюю оценку.

$$\mathcal{L}_{1/0}(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \sum_{i=1}^{\ell} [M(\mathbf{x}_i) < 0] \leq \sum_{i=1}^{\ell} \log(1 + \exp^{-y_i \mathbf{w}^T \mathbf{x}_i}) = \mathcal{L}_{\log}(\mathbf{X}, \mathbf{y}, \mathbf{w}),$$

# Преимущества

- Линейные модели быстро работают
- Мало параметров
- Хорошо применимы для больших данных

# Недостатки: пример XOR проблема

XOR

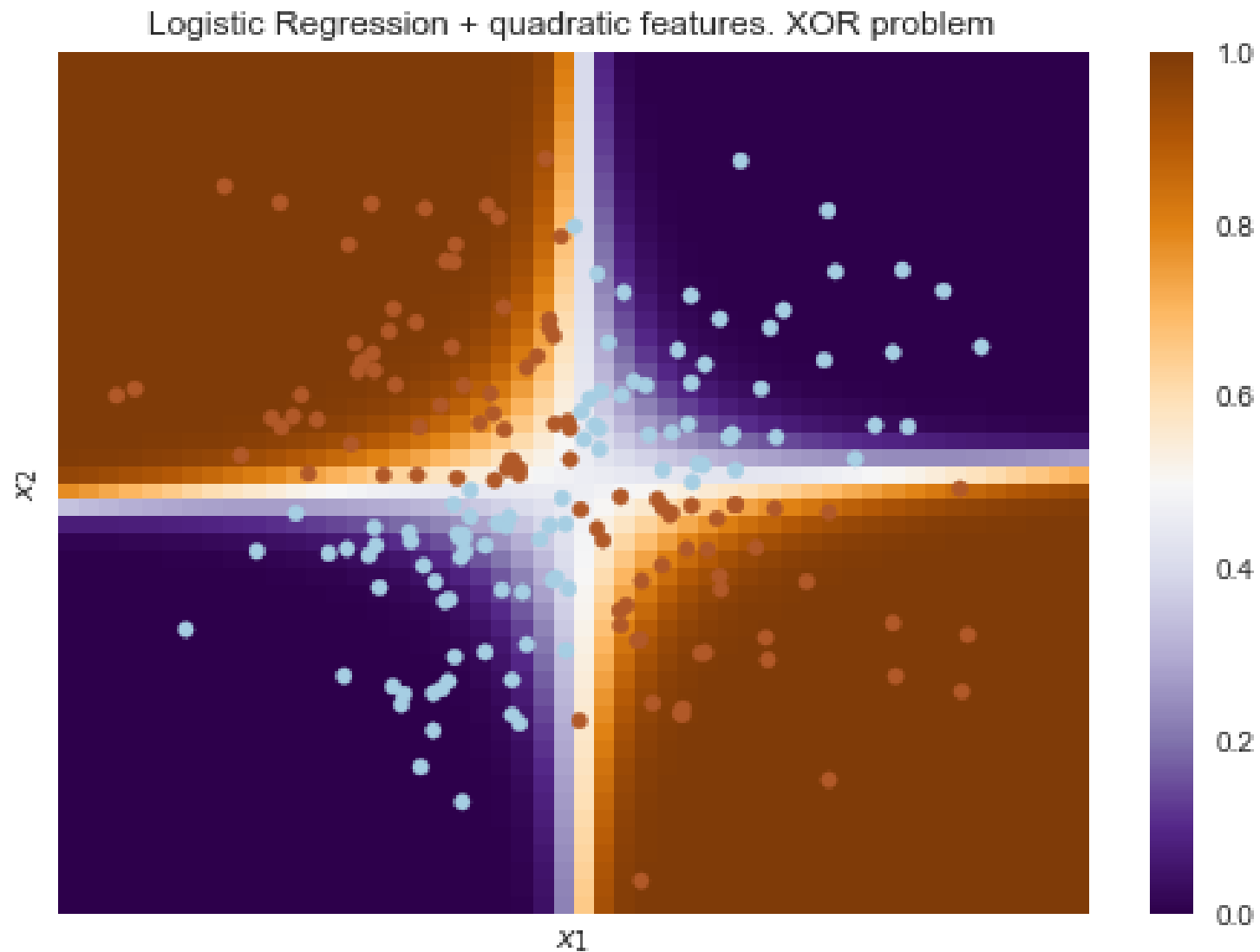
Input 1

Input 2

|   |   |   |
|---|---|---|
|   | 0 | 1 |
| 0 | 0 | 1 |
| 1 | 1 | 0 |



# Полиномиальные признаки





# Недостатки

- Линейные модели требуют тщательной обработки данных, для достижения хороших результатов
- Линейные модели чувствительны к выбросам и масштабированию
- В идеале, нужно, чтобы выборка соответствовала условиям теоремы Гаусса – Маркова, но это практически никогда не случается.

Спасибо за внимание !!!