

# Чудесный мир регрессии

## Машинное обучение, 2017

Спасибо К. В. Воронцову, МФТИ, Data Factory Яндекса, O.D.S. и кофеину.

Малютин Е. А.

## Сегодня в программе:

- Регрессия с точки зрения алгебры
- Регрессия с точки зрения логики
- Метод градиентного спуска
- Трюки и финты
- Проблемы регрессии
- Практика

# Регрессия. Начало



Рис.: Мемасик для вдохновения

# Регрессия. Начало

## Back to the roots

Пусть  $\exists \{x_1, x_2 \dots x_l\} \in X$  – множество объектов;  $\exists \{y_i\}_{i=1}^l \in Y$  – множество допустимых ответов. Пары  $(x_i, y_i)$  – называются прецедентами, а совокупность пар  $X^l = (x_i, y_i)_{i=1}^l$  – обучающая выборка, а так же существует зависимость (алгоритм):  $y^* : X \rightarrow Y$  – его-то и необходимо восстановить.

## Регрессия

- $a(x) = x_0 + \sum_{i=1}^l \omega_i x_i \rightarrow a(x) = \sum_{i=0}^l \omega_i x_i \rightarrow \vec{y} = X \vec{\omega} + \epsilon$ 
  - $\vec{y} \in R^n$  – объясняемая (или целевая) переменная;  $\epsilon$  – случайная переменная.
  - $\omega$  – вектор параметров модели (в машинном обучении эти параметры часто называют весами);
  - $X$  – матрица наблюдений и признаков размерности  $n$  строк на  $m + 1$  столбцов (включая фиктивную единичную колонку слева) с полным рангом по столбцам:  $\text{rank}(X) = m + 1$ ;

## Условности

Можно выписать модель регрессии явным образом для отдельного объекта:

$$y_i = \sum_{j=0}^n w_j X_{ij} + \epsilon$$

И, соответственно, мы приходим к некоторым условиям:

- матожидание:  $\forall i : E[\epsilon_i] = 0$ ;
- гомоскедастичность:  $\forall : Var(\epsilon_i) = \sigma^2 < \infty$
- некоррелированы:  $\forall i \neq j : Cov(\epsilon_i, \epsilon_j) = 0$

# Регрессия. Начало

## Веса

Оценка  $\bar{\omega}_i$  весов  $\omega_i$  называется **линейной**, если:

$$\bar{\omega}_i = \omega_{1i}y_1 + \omega_{2i}y_2 + \dots + \omega_{ni}y_n;$$

Несмещённая оценка:  $E[\bar{\omega}_i] = \omega_i$ ;

МНК:  $\mathcal{L}(X, y, \omega) = \frac{1}{2n} \sum_{i=1}^n (y_i - \omega^T x_i)^2$ ;

Найдем минимум данной ошибки:

$$\frac{d\mathcal{L}}{d\omega} = \frac{1}{2n}(-2X^T y + 2X^T X \omega)$$

И, соответственно:

$$\frac{d\mathcal{L}}{d\omega} = 0 \Leftrightarrow \omega = (X^T X)^{-1} X^T y$$

# Регрессия. Начало

## Метод максимального правдоподобия

Правдоподобие:  $L = \prod p(y_i|x_i, \alpha)$

**Log-likelihood:**  $W(\alpha) = \log(L) = \sum \ln(P(y|X, \alpha))$

Положим ошибки нормально-распределёнными:  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

И тогда для нашей модели:  $p(y_i|x_i, w) = \mathcal{N}(\sum_j w_j X_{ij}, \sigma^2)$

$$\log(p(X, y|w)) = \log(\prod \mathcal{N}(\sum_j w_j X_{ij}, \sigma^2)) = -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum (y_i - \omega^T x_i)^2$$

И, максимизируя это вот всё:

$$w = \arg \max_w p(X, y|w) = \arg \max_w -\mathcal{L}(X, y, w)$$

**Вывод:** минимизация МНК эквивалентна максимизации МП.

## Разложение ошибки на смещение и разброс

- истинное значение целевой переменной складывается из некоторой детерминированной функции  $f(\vec{x})$  и случайной ошибки  $y = f(\vec{x}) + \epsilon$ ;
- ошибка распределена нормально с центром в нуле и некоторым разбросом:  $\epsilon \sim \mathcal{N}(0, \sigma^2)$
- мы пытаемся приблизить детерминированную, но неизвестную функцию  $f(\vec{x})$  линейной функцией от регрессоров  $f(\vec{x})$ , которая, в свою очередь, является точечной оценкой функции  $f$  в пространстве функций.



# Регрессия. Начало

## Разложение ошибки

Рассмотрим ошибку в точке  $\vec{x}$ :

$$E[\vec{x}] = E[(x - \hat{f}(\vec{x}))] = E[y^2] + E[\hat{f}^2] - 2E[y\hat{f}];$$

По отдельности:

$$E[y^2] = \text{Var}(y) + E[y]^2 = \sigma^2 + f^2; \quad E[\hat{f}^2] = \text{Var}(\hat{f}) + E[f]^2$$

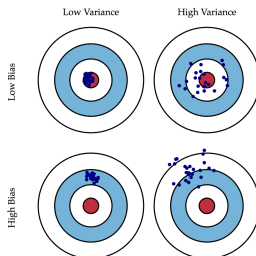
$$E[y\hat{f}] = E[(f + \epsilon)\hat{f}] = f E[\hat{f}]$$

$$\begin{aligned} \text{Err}(\vec{x}) &= E[(y - \hat{f}(\vec{x}))^2] = \sigma^2 + f^2 + \text{Var}(\hat{f}) + E[\hat{f}^2] - 2fE[f] = \\ &= (f - E[\hat{f}])^2 + \text{Var}(\hat{f}) + \sigma^2 = \text{Bias}(\hat{f}^2) + \text{Var}(\hat{f}) + \sigma^2 \end{aligned}$$

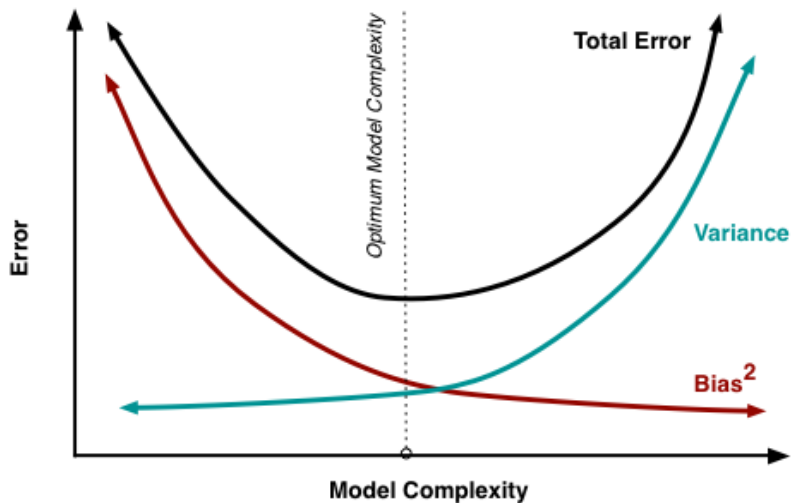
# Регрессия. Начало

## Ошибки:

- квадрат смещения:  $Bias(\hat{f})$  – средняя ошибка по всевозможным наборам данных;
- дисперсия:  $Var(\hat{f})$  – на сколько ошибка будет отличаться, если обучать модель на разных наборах данных;
- неустранимой ошибки:  $\sigma^2$



# Регрессия. Начало



## Регуляризация

- При мультиколлинеарности в данных получаем  $(X^T X)^{-1}$  – экстремально большие значения собственных чисел  $(\frac{1}{\lambda_i})$

## Регуляризация

- При мультиколлинеарности в данных получаем  $(X^T X)^{-1}$  – экстремально большие значения собственных чисел  $(\frac{1}{\lambda_i})$
- Решение? Регуляризация по Тихонову:

$$\mathcal{L}(X, \vec{y}, \vec{w}) = \frac{1}{2n} \|\vec{y} - X\vec{w}\|^2 + \|\mathcal{G}\vec{w}\|^2$$

## Регуляризация

- При мультиколлинеарности в данных получаем  $(X^T X)^{-1}$  – экстремально большие значения собственных чисел  $(\frac{1}{\lambda_i})$
- Решение? Регуляризация по Тихонову:

$$\mathcal{L}(X, \vec{y}, \vec{w}) = \frac{1}{2n} \|\vec{y} - X\vec{w}\|^2 + \|\mathcal{G}\vec{w}\|^2$$

- Часто используется в таком виде:  $\mathcal{G} = \frac{\lambda}{2} E$

## Регуляризация

- При мультиколлинеарности в данных получаем  $(X^T X)^{-1}$  – экстремально большие значения собственных чисел  $(\frac{1}{\lambda_i})$
- Решение? Регуляризация по Тихонову:

$$\mathcal{L}(X, \vec{y}, \vec{w}) = \frac{1}{2n} \|\vec{y} - X\vec{w}\|^2 + \|\mathcal{G}\vec{w}\|^2$$

- Часто используется в таком виде:  $\mathcal{G} = \frac{\lambda}{2} E$
- Точное решение:

$$\vec{w} = (X^T X + \lambda E)^{-1} X^T \vec{y}$$

# Регрессия. Начало

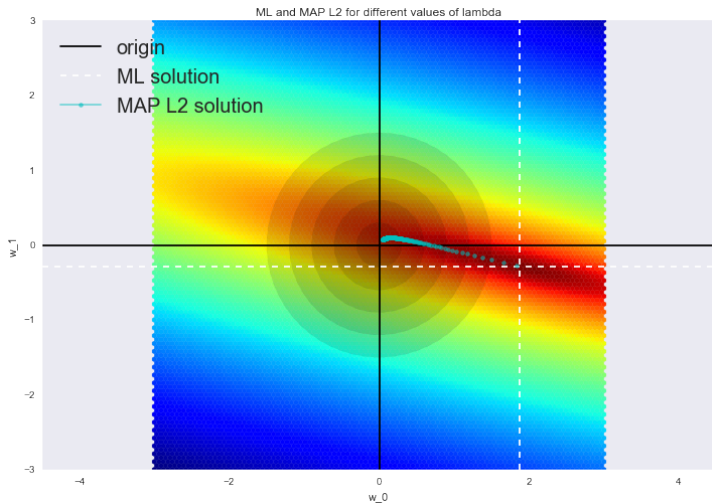
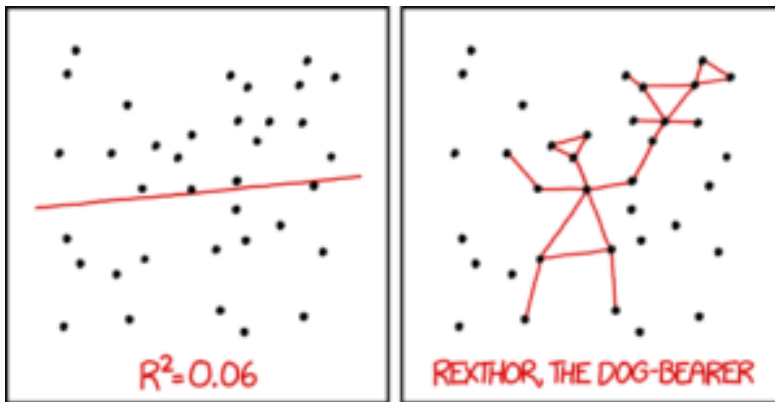


Рис.: Особенности регуляризации



# Мемасик для вдохновения



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER  
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE  
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

# Логистическая регрессия

бинарный классификатор на основе регрессии:  $a(\vec{x}) = \text{sign}(\vec{w}^T x)$

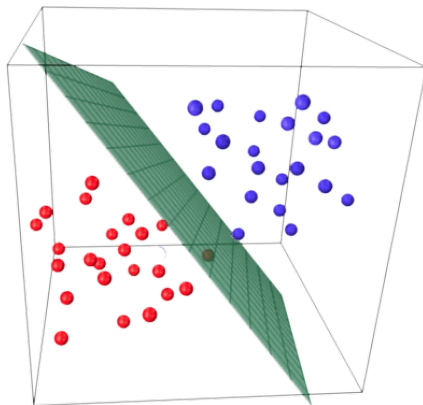


Рис.: Линейно-разделимая выборка

# Логистическая регрессия

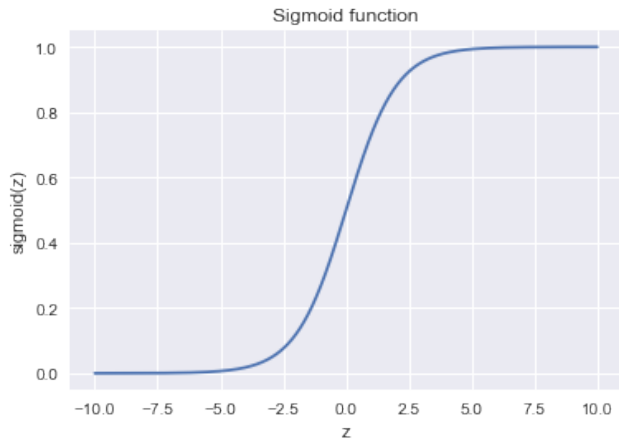
$$p_+ = P(y_i = 1 | \vec{x}_i, \vec{\omega})$$

Клиент	Вероятность невозврата
Mike	0.78
Jack	0.45
Larry	0.13
Kate	0.06
William	0.03
Jessica	0.02

Рис.: Пример бинаризации

# Логистическая регрессия

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



# Логистическая регрессия

$\exists X, P(X); OR(X) \equiv \frac{P(X)}{1-P(X)}$ ; (отношение вероятностей)  $P(X) \in [0, 1]; OR(X) \in R$

## Вычисляем лог. регрессию

- Вычислить значение  $w_0 + w_1x_1 + w_2x_2 + \dots = \vec{w}^T \vec{x}$ . (уравнение  $\vec{w}^T \vec{x} = 0$  задает гиперплоскость, разделяющую примеры на 2 класса);
- Вычислить логарифм отношения шансов:  $\log(OR_+) = \vec{w}^T \vec{x}$ .
- Вычисляем вероятность:  $p_+ = \frac{OR_+}{1+OR_+} = \frac{\exp^{\vec{w}^T \vec{x}}}{1+\exp^{\vec{w}^T \vec{x}}} = \frac{1}{1+\exp^{-\vec{w}^T \vec{x}}} = \sigma(\vec{w}^T \vec{x})$

# Логистическая регрессия

## Начало

- $p_+ = P(y_i = 1 | \vec{x}_i, \vec{\omega}) = \sigma(\vec{\omega}^T, \vec{x})$

# Логистическая регрессия

## Начало

- $p_+ = P(y_i = 1 | \vec{x}_i, \vec{w}) = \sigma(\vec{w}^T, \vec{x})$
- $p_- = P(y_i = -1 | \vec{x}_i, \vec{w}) = 1 - \sigma(\vec{w}^T, \vec{x}) = \sigma(-\vec{w}^T, \vec{x})$

# Логистическая регрессия

## Начало

- $p_+ = P(y_i = 1 | \vec{x}_i, \vec{w}) = \sigma(\vec{w}^T, \vec{x})$
- $p_- = P(y_i = -1 | \vec{x}_i, \vec{w}) = 1 - \sigma(\vec{w}^T, \vec{x}) = \sigma(-\vec{w}^T, \vec{x})$
- следим за руками...



# Логистическая регрессия

## Начало

- $p_+ = P(y_i = 1 | \vec{x}_i, \vec{\omega}) = \sigma(\vec{\omega}^T, \vec{x})$
- $p_- = P(y_i = -1 | \vec{x}_i, \vec{\omega}) = 1 - \sigma(\vec{\omega}^T, \vec{x}) = \sigma(-\vec{\omega}^T, \vec{x})$
- следим за руками...
- $P(y = y_i | \vec{x}_i, \vec{\omega}) = \sigma(y_i \vec{\omega}^T, \vec{x})$

# Логистическая регрессия

## Начало

- $p_+ = P(y_i = 1 | \vec{x}_i, \vec{\omega}) = \sigma(\vec{\omega}^T, \vec{x})$
- $p_- = P(y_i = -1 | \vec{x}_i, \vec{\omega}) = 1 - \sigma(\vec{\omega}^T, \vec{x}) = \sigma(-\vec{\omega}^T, \vec{x})$
- следим за руками...
- $P(y = y_i | \vec{x}_i, \vec{\omega}) = \sigma(y_i \vec{\omega}^T, \vec{x})$
- $M(\vec{x}_i) = y_i \vec{\omega} \vec{x}$  – отступ

# Логистическая регрессия

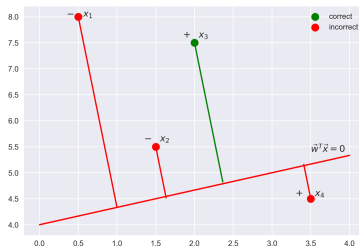


Рис.: Отступы на объектах

## Отступы

- $M \gg 0$  – метка поставлена правильно;  $M \ll 0$  – выброс, шум
- $M \simeq 0$  – ну не знаю...

# Логистическая регрессия

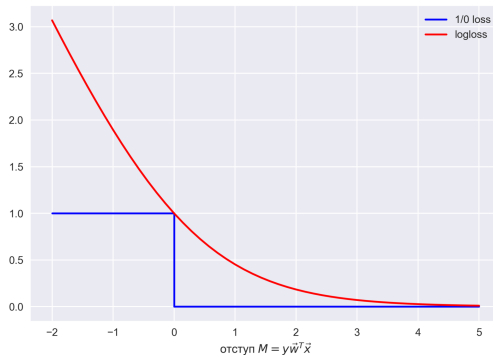
Угадайте что?

$$P(\vec{y} \mid X, \vec{w}) = \prod_{i=1}^{\ell} P(y = y_i \mid \vec{x}_i, \vec{w}),$$

$$\begin{aligned} \log P(\vec{y} \mid X, \vec{w}) &= \log \prod_{i=1}^{\ell} P(y = y_i \mid \vec{x}_i, \vec{w}) \\ &= \log \prod_{i=1}^{\ell} \sigma(y_i \vec{w}^T \vec{x}_i) \\ &= \sum_{i=1}^{\ell} \log \sigma(y_i \vec{w}^T \vec{x}_i) \\ &= \sum_{i=1}^{\ell} \log \frac{1}{1 + \exp^{-y_i \vec{w}^T \vec{x}_i}} \\ &= - \sum_{i=1}^{\ell} \log(1 + \exp^{-y_i \vec{w}^T \vec{x}_i}) \end{aligned}$$

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

# Логистическая регрессия



# Градиентный спуск

Обобщенный алгоритм градиентного спуска:

- Минимизируем эмпирический риск:  $Q(w) = \sum_{i=1}^I \mathcal{L}_i(w) \rightarrow \min_w$

# Градиентный спуск

## Обобщенный алгоритм градиентного спуска:

- Минимизируем эмпирический риск:  $Q(w) = \sum_{i=1}^I \mathcal{L}_i(w) \rightarrow \min_w$
- Начальное приближение:  $w^0 = \text{start}$

# Градиентный спуск

## Обобщенный алгоритм градиентного спуска:

- Минимизируем эмпирический риск:  $Q(w) = \sum_{i=1}^l \mathcal{L}_i(w) \rightarrow \min_w$
- Начальное приближение:  $w^0 = start$
- $t = 1 \dots n$  итеративно пересчитываем:

$$w^{(t+1)} = w^{(t)} - h * \nabla Q(w^{(t)}), \quad \nabla Q(w) = \left( \frac{\partial Q(w)}{\partial w_i} \right)_{i=0}^n$$



# Градиентный спуск

## Обобщенный алгоритм градиентного спуска:

- Минимизируем эмпирический риск:  $Q(w) = \sum_{i=1}^l \mathcal{L}_i(w) \rightarrow \min_w$
- Начальное приближение:  $w^0 = start$
- $t = 1 \dots n$  итеративно пересчитываем:

$$w^{(t+1)} = w^{(t)} - h * \nabla Q(w^{(t)}), \quad \nabla Q(w) = \left( \frac{\partial Q(w)}{\partial w_i} \right)_{j=0}^n$$

- Всё вместе:

$$w^{(t+1)} = w^{(t)} - h * \sum_{i=1}^l \nabla \mathcal{L}_i(w^{(t)})$$

где  $h$  – это шаг градиента(learning rate).

# Градиентный спуск

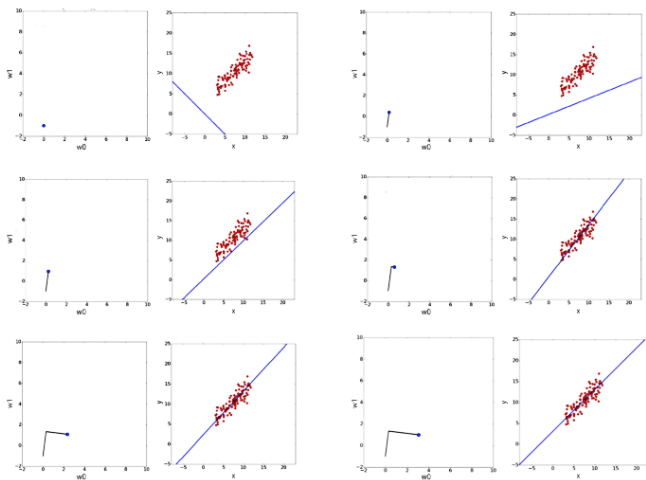


Рис.: Пример Gradient Descent на одномерной регрессии

# Градиентный спуск

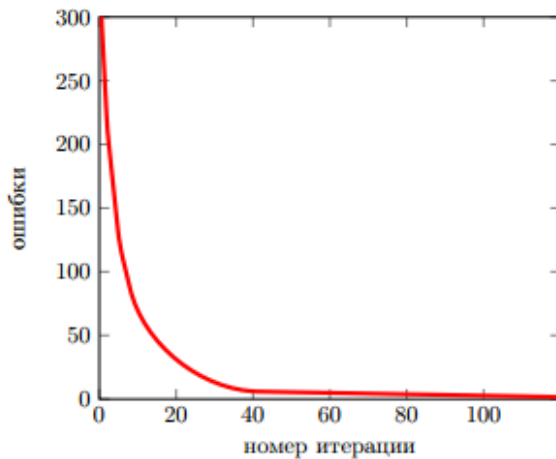


Рис.: Ошибка Gradient Descent

# Градиентный спуск

## Инициализация весов:

- как выбрать инициализацию?

# Градиентный спуск

## Инициализация весов:

- как выбрать инициализацию?
- Случайно:  $w_j = \text{random}(-\frac{1}{2n}, \frac{1}{2n})$

# Градиентный спуск

## Инициализация весов:

- как выбрать инициализацию?
- Случайно:  $w_j = \text{random}(-\frac{1}{2n}, \frac{1}{2n})$
- Эвристика:

$$w_j = \frac{\langle y, f_j \rangle}{\langle f_j, f_j \rangle}$$

# Градиентный спуск

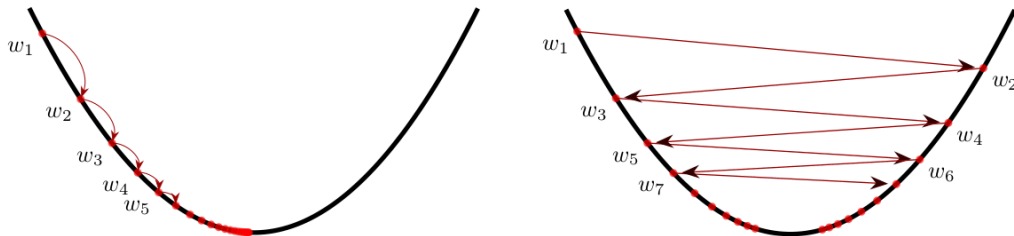


Рис.: Случай маленького и большого шага

$$h_t = \frac{k}{t} - \text{адаптивный шаг}$$

# Градиентный спуск

## Пример с регрессией

- Задача оптимизации:

$$Q(\omega, X) = \frac{1}{I} \|X\omega - y\|^2 \rightarrow \min_{\omega}$$

- Градиентный спуск:

$$\nabla_{\omega} Q(\omega, X) = \frac{2}{I} X^T (X\omega - y)$$

- $\|X\omega - y\|^2$  – вектор ошибок



problems?

- j-я компонента для градиента:

$$\frac{\partial Q}{\partial \omega_j} = \frac{2}{l} \sum_{i=1}^l x_i^j (\langle w_i, x_i \rangle - y_i)^2$$

problems?

- j-я компонента для градиента:

$$\frac{\partial Q}{\partial \omega_j} = \frac{2}{l} \sum_{i=1}^l x_i^j (\langle w_i, x_i \rangle - y_i)^2$$

- Выход – стохастический градиент:  $w^0 = 0$

problems?

- j-я компонента для градиента:

$$\frac{\partial Q}{\partial \omega_j} = \frac{2}{l} \sum_{i=1}^l x_i^j (\langle w_i, x_i \rangle - y_i)^2$$

- Выход – стохастический градиент:  $w^0 = 0$
- Шаг:  $\omega^t = \omega^t - h_t \nabla Q(w^{t-1}, \{x_i\})$

problems?

- j-я компонента для градиента:

$$\frac{\partial Q}{\partial \omega_j} = \frac{2}{l} \sum_{i=1}^l x_i^j (\langle w_i, x_i \rangle - y_i)^2$$

- Выход – стохастический градиент:  $w^0 = 0$
- Шаг:  $\omega^t = \omega^t - h_t \nabla Q(w^{t-1}, \{x_i\})$
- Остановка:  $\|w^t - w^{t-1}\| < \epsilon$

# Градиентный спуск

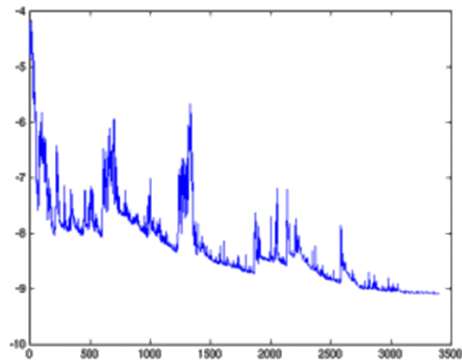
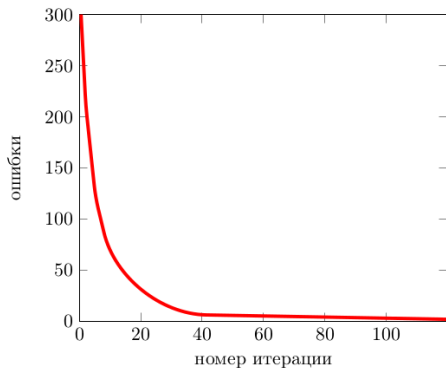


Рис.: Сходимость классического и стохастического градиентного спуска

# Градиентный спуск

## Преимущества и недостатки

### Преимущества:

- Легко обобщается
- Поточковый
- Работает с большими данными
- Можно заканчивать, даже не предъявив всю выборку

### Недостатки:

# Градиентный спуск

## Преимущества и недостатки

### Преимущества:

- Легко обобщается
- Поточковый
- Работает с большими данными
- Можно заканчивать, даже не предъявив всю выборку

### Недостатки:

- Многоэкстремальный
- Переобучение
- Если функция потерь имеет горизонтальные асимптоты, то процесс может попасть в состояние «паралича».

## Крутые эвристики

- Нормализация данных:

$$x_j = \frac{x_j - x_{min}}{x_{max} - x_{min}} \quad \text{либо} \quad x_j = \frac{x_j - x_{med}}{x_{var}}$$



# Градиентный спуск

## Крутые эвристики

- Нормализация данных:

$$x_j = \frac{x_j - x_{min}}{x_{max} - x_{min}} \quad \text{либо} \quad x_j = \frac{x_j - x_{med}}{x_{var}}$$

- Порядок предъявления объектов:

## Крутые эвристики

- Нормализация данных:

$$x_j = \frac{x_j - x_{min}}{x_{max} - x_{min}} \quad \text{либо} \quad x_j = \frac{x_j - x_{med}}{x_{var}}$$

- Порядок предъявления объектов:

- перетасовка объектов (shuffling)

## Крутые эвристики

- Нормализация данных:

$$x_j = \frac{x_j - x_{min}}{x_{max} - x_{min}} \quad \text{либо} \quad x_j = \frac{x_j - x_{med}}{x_{var}}$$

- Порядок предъявления объектов:

- перетасовка объектов (shuffling)
- предъявлять те объекты, на которых была допущена ошибка

## Крутые эвристики

### ■ Нормализация данных:

$$x_j = \frac{x_j - x_{min}}{x_{max} - x_{min}} \quad \text{либо} \quad x_j = \frac{x_j - x_{med}}{x_{var}}$$

### ■ Порядок предъявления объектов:

- перетасовка объектов (shuffling)
- предъявлять те объекты, на которых была допущена ошибка
- сравнить величину ошибки на предъявленном объекте с некоторым порогом

# Градиентный спуск

## Квадратичная регуляризация:

- Штраф за норму:  $Q_\tau(w) = Q(w) + \tau/2 \|w\|^2$
- $w = w(1 - h * \tau) + h * Q'(w)$

## Ещё эвристики

- выбор величины шага
  - От  $t$ :  $h = 1/t$
  - Решить оптимизацию:  $Q(w - hQ'(w)) \rightarrow \min_w$
- выбивание из локальных минимумов
- ранний останов.

To be continued...