

# Linear regression

Victor Kitov

[v.v.kitov@yandex.ru](mailto:v.v.kitov@yandex.ru)

Yandex School of Data Analysis



# Table of Contents

- 1 Linear regression
- 2 Nonlinear transformations
- 3 Regularization & restrictions.
- 4 Different loss-functions
- 5 Weighted account for observations

# Linear regression

- Linear model  $f(x, \beta) = \langle x, \beta \rangle = \sum_{i=1}^D \beta_i x^i$
- Define  $X \in \mathbb{R}^{N \times D}$ ,  $\{X\}_{ij}$  defines the  $j$ -th feature of  $i$ -th object,  $Y \in \mathbb{R}^n$ ,  $\{Y\}_i$  - target value for  $i$ -th object.
- Ordinary least squares (OLS) method:

$$\sum_{n=1}^N (f(x_n, \beta) - y_n)^2 = \sum_{n=1}^N \left( \sum_{d=1}^D \beta_d x_n^d - y_n \right)^2 \rightarrow \min_{\beta}$$

# Solution

Stationarity condition:

$$2 \sum_{n=1}^N x_n \left( \sum_{d=1}^D \beta_d x_n^d - y_n \right) = 0$$

In matrix form:

$$2X^T(X\beta - Y) = 0$$

so

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

This is the global minimum, because the optimized criteria is convex.

- Geometric interpretation of linear regression, estimated with OLS.

# Linearly dependent features

- Solution  $\hat{\beta} = (X^T X)^{-1} X^T Y$  exists when  $X^T X$  is non-degenerate
- Using property
$$\text{rank}(X) = \text{rank}(X^T) = \text{rank}(X^T X) = \text{rank}(X X^T)$$
  - problem occurs when one of the features is a linear combination of the other
    - example: constant unity feature  $c$  and one-hot-encoding  $e_1, e_2, \dots, e_K$ , because  $\sum_k e_k \equiv c$
    - interpretation: non-identifiability of  $\hat{\beta}$
  - solved using:
    - feature selection
    - extraction (e.g. PCA)
    - regularization.

# Analysis of linear regression

## Advantages:

- single optimum, which is global (for the non-singular matrix)
- analytical solution
- interpretability algorithm and solution

## Drawbacks:

- too simple model assumptions (may not be satisfied)
- $X^T X$  should be non-degenerate (and well-conditioned)

# Table of Contents

- 1 Linear regression
- 2 Nonlinear transformations
- 3 Regularization & restrictions.
- 4 Different loss-functions
- 5 Weighted account for observations

## Generalization by nonlinear transformations

Nonlinearity by  $x$  in linear regression may be achieved by applying non-linear transformations to the features:

$$x \rightarrow [\phi_0(x), \phi_1(x), \phi_2(x), \dots, \phi_M(x)]$$

$$f(x) = \langle \phi(x), \beta \rangle = \sum_{m=0}^M \beta_m \phi_m(x)$$

The model remains to be linear in  $w$ , so all advantages of linear regression remain.



# Typical transformations

$\phi_k(x)$	comments
$\exp \left\{ -\frac{\ x-\mu\ ^2}{s^2} \right\}$	closeness to point $\mu$ in feature space
$x^i x^j$	interaction of features
$\ln x_k$	the alignment of the distribution with heavy tails
$F^{-1}(x_k)$	conversion of atypical continuous distribution to uniform <sup>1</sup>

---

<sup>1</sup>why?

# Table of Contents

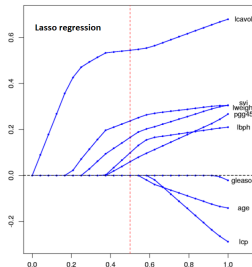
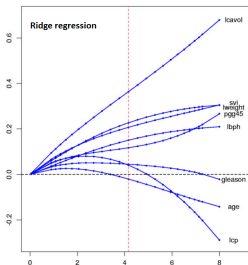
- 1 Linear regression
- 2 Nonlinear transformations
- 3 Regularization & restrictions.**
- 4 Different loss-functions
- 5 Weighted account for observations

# Regularization

- Variants of target criteria  $Q(\beta)$  with regularization<sup>2</sup>:

$$\begin{aligned} \sum_{n=1}^N (x_n^T \beta - y_n)^2 + \lambda \|\beta\|_1 & \quad \text{Lasso} \\ \sum_{n=1}^N (x_n^T \beta - y_n)^2 + \lambda \|\beta\|_2^2 & \quad \text{Ridge} \\ \sum_{n=1}^N (x_n^T \beta - y_n)^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2 & \quad \text{Elastic net} \end{aligned}$$

- Dependency of  $\beta$  from  $\frac{1}{\lambda}$ :



<sup>2</sup>Derive solution for ridge regression. Will it be uniquely defined for correlated features?

## Linear monotonic regression

- We can impose restrictions on coefficients such as non-negativity:

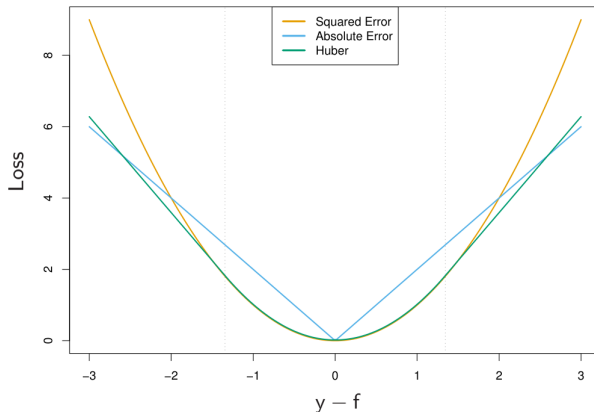
$$\begin{cases} Q(\beta) = \|X\beta - Y\|^2 \rightarrow \min_{\beta} \\ \beta_i \geq 0, \quad i = 1, 2, \dots, D \end{cases}$$

- Example: averaging of forecasts of different prediction algorithms
- $\beta_i = 0$  means, that  $i$ -th component does not improve accuracy of forecasting.

# Table of Contents

- 1 Linear regression
- 2 Nonlinear transformations
- 3 Regularization & restrictions.
- 4 Different loss-functions**
- 5 Weighted account for observations

# Non-quadratic loss functions<sup>34</sup>



<sup>3</sup>What is the value of constant prediction, minimizing sum of squared errors?

<sup>4</sup>What is the value of constant prediction, minimizing sum of absolute errors?

## Conditional non-constant optimization

- For  $x, y \sim P(x, y)$  and prediction being made for fixed  $x$ :

$$\arg \min_{f(x)} \mathbb{E} \left\{ (f(x) - y)^2 \mid x \right\} = \mathbb{E}[y|x]$$

$$\arg \min_{f(x)} \mathbb{E} \{ |f(x) - y| \mid x \} = \text{median}[y|x]$$

# Minimization of expected squared error

- Let  $x, y \sim P(x, y)$  and  $\mathbb{E}[y|x]$  exist. Then

$$\arg \min_{f(x)} \mathbb{E} \left\{ (f(x) - y)^2 \middle| x \right\} = \mathbb{E}[y|x]$$

$$\begin{aligned} \mathbb{E} \{ (f(x) - y)^2 | x \} &= \mathbb{E} \{ (f(x) - \mathbb{E}[y|x] + \mathbb{E}[y|x] - y)^2 | x \} \\ &= \mathbb{E} \{ (f(x) - \mathbb{E}[y|x])^2 | x \} + \mathbb{E} \{ (\mathbb{E}[y|x] - y)^2 | x \} \\ &\quad + 2\mathbb{E} \{ (f(x) - \mathbb{E}[y|x]) (\mathbb{E}[y|x] - y) | x \} = \\ &= (f(x) - \mathbb{E}[y|x])^2 + \mathbb{E} \{ (\mathbb{E}[y|x] - y)^2 | x \} \end{aligned} \quad (1)$$



## Minimization of expected squared error

We used

$$\begin{aligned}\mathbb{E} \{ (f(x) - \mathbb{E}[y|x]) (\mathbb{E}[y|x] - y) | x \} = \\ (f(x) - \mathbb{E}[y|x]) \mathbb{E} \{ \mathbb{E}[y|x] - y | x \} \equiv 0\end{aligned}$$

Minimum of (1) is achieved at  $f(x) = \mathbb{E}[y|x]$ .

# Table of Contents

- 1 Linear regression
- 2 Nonlinear transformations
- 3 Regularization & restrictions.
- 4 Different loss-functions
- 5 Weighted account for observations**

## Weighted account for observations<sup>5</sup>

- Weighted account for observations

$$\sum_{n=1}^N w_n (x_n^T \beta - y_n)^2$$

- Weights may be:
  - increased for incorrectly predicted objects
    - algorithm becomes more oriented on error correction
  - decreased for incorrectly predicted objects
    - they may be considered outliers that break our model

---

<sup>5</sup>Derive solution for weighted regression.

# Robust regression

- Initialize  $w_1 = \dots = w_N = 1/N$
- Repeat:
  - estimate regression  $\hat{y}(x)$  using observations  $(x_i, y_i)$  with weights  $w_i$ .
  - for each  $i = 1, 2, \dots, N$ :
    - re-estimate  $\varepsilon_i = \hat{y}(x_i) - y_i$
    - recalculate  $w_i = K(|\varepsilon_i|)$
  - normalize weights  $w_i = \frac{w_i}{\sum_{n=1}^N w_n}$

**Comments:**  $K(\cdot)$  is some *decreasing* function, repetition may be

- predefined number of times
- until convergence of model parameters.

## Robust classification

- Initialize  $w_1 = \dots = w_N = 1/N$
- Repeat:
  - estimate classifier discriminant functions  $\{g_y(\cdot)\}_{y=1,\dots,C}$  using observations  $(x_i, y_i)$  with weights  $w_i$ .
  - for each  $i = 1, 2, \dots, N$ :
    - re-estimate  $M_i = g_{y_i}(x_i) - \max_{y \neq y_i} g_y(x_i)$
    - recalculate  $w_i = K(M_i)$
  - normalize weights  $w_i = \frac{w_i}{\sum_{n=1}^N w_n}$

**Comments:**  $K(\cdot)$  is some *increasing* function, repetition may be

- predefined number of times
- until convergence of model parameters.