

# Linear regression

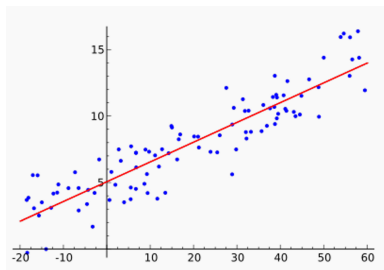
Alexander Sirotkin

HSE University  
November 14, 2023

# Linear regression

- Linear model: consider a linear function

$$y(x, w) = w_0 + \sum_{j=1}^p x_j w_j = x^\top w, \quad x = (1, x_1, \dots, x_p).$$



- How can we find optimal parameters  $\hat{w}$  by training data of the form  $(x_i, y_i)_{i=1}^N$ ?

# Linear regression

- How can we find optimal parameters  $\hat{w}$  by training data of the form  $(x_i, y_i)_{i=1}^N$ ?
- Least squares estimation: we will minimize

$$\text{RSS}(w) = \sum_{i=1}^N (y_i - x_i^\top w)^2.$$

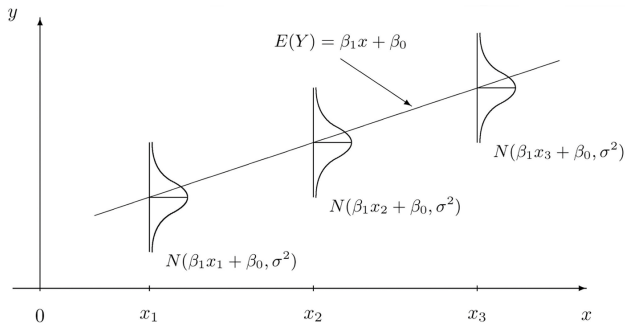
- There is even an exact solution, but that's not important right now.

# Linear regression

- What is important: suppose that noise (error in the data) has a normal distribution, i.e., observed variable  $t$  is

$$t = y(x, w) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2), \text{ то есть}$$

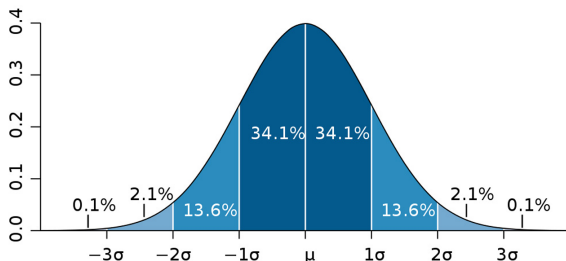
$$p(t \mid x, w, \sigma^2) = \mathcal{N}(t \mid y(x, w), \sigma^2).$$



# Linear regression

- Aside – normal distribution:

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$



- Why is it so important?

# Linear regression

- Consider a dataset  $X = \{x_1, \dots, x_N\}$  with correct answers  $t = \{t_1, \dots, t_N\}$ .
- We assume that the data points are independent identically distributed:

$$p(t | X, w, \sigma^2) = \prod_{n=1}^N \mathcal{N}(t_n | w^\top \phi(x_n), \sigma^2).$$

- We take the logarithm (we omit  $X$  below for brevity):

$$\ln p(t | w, \sigma^2) = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2.$$

# Linear regression

- And we see that to maximize the likelihood w.r.t.  $w$  we need to minimize mean squared error!

$$\nabla_w \ln p(t | w, \sigma^2) = \frac{1}{\sigma^2} \sum_{n=1}^N \left( t_n - w^\top \phi(x_n) \right) \phi(x_n).$$

- We can also get a posterior distribution, introducing prior distributions (also normal).
- And then the predictive distribution

$$p(y | x, D) = \int_w p(y | x, w) p(w | D) dw$$

...but that's beside the point right now.

- Main conclusion: in many regression problems it makes sense to minimize the sum of squared deviations, this corresponds to normally distributed noise.

# Bayesian regularization

- And now let us look at regression from the pure Bayesian perspective.
- Recall that in Bayesian inference, we
  - 1 find the posterior distribution on the hypothesis/params:

$$p(\theta \mid D) \propto p(D|\theta)p(\theta)$$

(and/or find the maximal a posteriori hypothesis  
 $\arg \max_{\theta} p(\theta \mid D)$ );

- 2 find the predictive distribution:

$$p(x \mid D) \propto \int_{\theta \in \Theta} p(x \mid \theta) p(D|\theta) p(\theta) d\theta.$$



# Bayesian regularization

- We have not yet had any priors in our study of linear regression.
- Let us introduce a prior; e.g., the normal distribution:

$$p(w) = \mathcal{N}(w \mid \mu_0, \Sigma_0).$$

- Consider a dataset  $X = \{x_1, \dots, x_N\}$  with values  $t = \{t_1, \dots, t_N\}$ ; we again assume that they are independent and identically distributed:

$$p(t \mid X, w, \sigma^2) = \prod_{n=1}^N \mathcal{N}(t_n \mid w^\top \phi(x_n), \sigma^2).$$

# Bayesian regularization

- Then the problem is to compute

$$\begin{aligned} p(w \mid t) &\propto p(t \mid X, w, \sigma^2) p(w) \\ &= \mathcal{N}(w \mid \mu_0, \Sigma_0) \prod_{n=1}^N \mathcal{N}(t_n \mid w^\top \phi(x_n), \sigma^2). \end{aligned}$$

- Let us compute!

# Bayesian regularization

- We get

$$p(\mathbf{w} \mid \mathbf{t}) = \mathcal{N}(\mathbf{w} \mid \boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N),$$

$$\boldsymbol{\mu}_N = \boldsymbol{\Sigma}_N \left( \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 + \frac{1}{\sigma^2} \boldsymbol{\Phi}^\top \mathbf{t} \right),$$

$$\boldsymbol{\Sigma}_N = \left( \boldsymbol{\Sigma}_0^{-1} + \frac{1}{\sigma^2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \right)^{-1}.$$

- And now the log likelihood.

# Bayesian regularization

- If we take the prior distribution around zero:

$$p(w) = \mathcal{N}(w \mid 0, \frac{1}{\alpha}I),$$

we get the log likelihood as

$$\ln p(w \mid t) = -\frac{1}{2\sigma^2} \sum_{n=1}^N \left( t_n - w^\top \phi(x_n) \right)^2 - \frac{\alpha}{2} w^\top w + \text{const},$$

i.e., precisely ridge regression!

# Generalization

- A slight generalization – a more general prior distribution:

$$p(w \mid \alpha) = \left[ \frac{q}{2} \left( \frac{\alpha}{2} \right)^{1/q} \frac{1}{\Gamma(1/q)} \right]^M e^{-\frac{\alpha}{2} \sum_{j=1}^M |w_j|^q}.$$

Try to compute the log likelihood.

# Regularization again

- We know that least squares do not always work well. Two reasons:
  - ① bad predictive power – often better to regularize, trading bias for variance;
  - ② hard to interpret – we often want to understand what is going on, and if we have lots of different nonzero numbers, it's hard.
- Hence, we'd like to get more nonzero components in the vector  $w$ .

# Subset selection

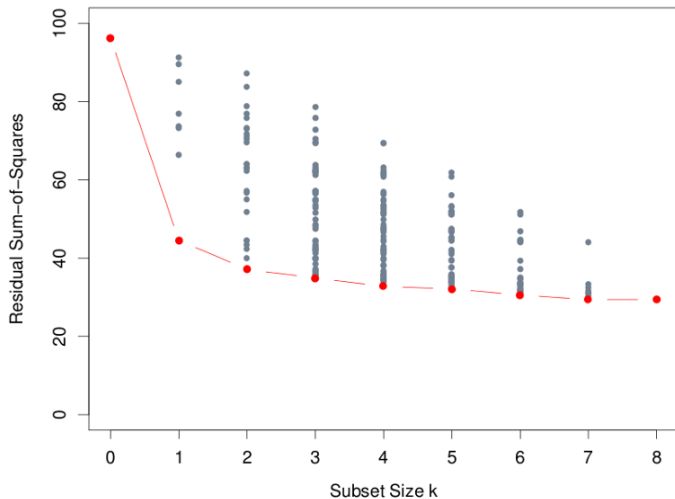
- What if we do it directly? Simply presume most coefficients are zero and find the nonzero ones.
- This is called *subset selection*.
- Best subset selection: choose the subset of  $k$  input variables that gives the best results

## Subset selection

- Naturally, this does not work computationally: there are lots of subsets.
- Forward-stepwise selection: start from the intercept, then add one best predictor per step.
- Backward-stepwise selection: start from full regression and remove the predictor that influences the error the least.



# Subset selection



# Lasso

- Let us now consider lasso regression:

$$L(w) = \frac{1}{2} \sum_{i=1}^N (f(x_i, w) - y_i)^2 + \lambda \sum_{j=0}^p |w_j|.$$

- The main difference is that the form of the constraints is now such that it is much more probable to get strictly zero  $w_j$ .
- Btw, what do I mean by “form of the constraints”?

# Lasso

- We can rewrite the regression with regularizer in a different way:

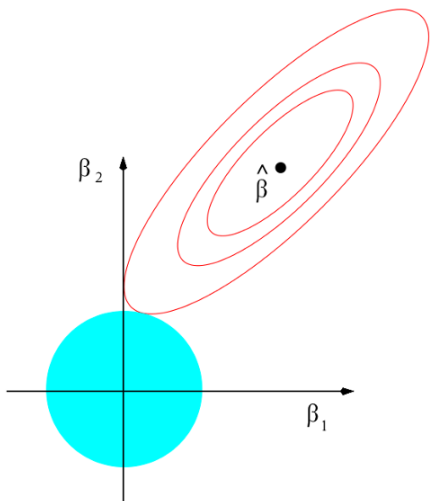
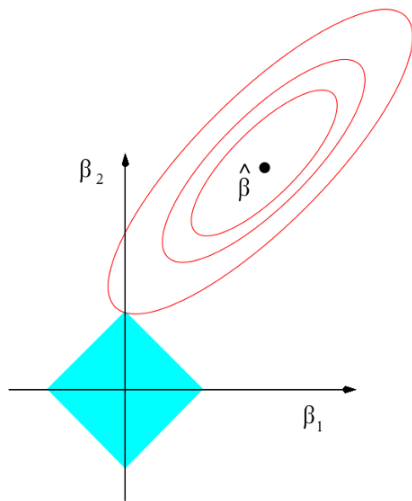
$$w^* = \arg \min_w \left\{ \frac{1}{2} \sum_{i=1}^N (f(x_i, w) - y_i)^2 + \lambda \sum_{j=0}^p |w_j| \right\},$$

is equivalent to

$$w^* = \arg \min_w \left\{ \frac{1}{2} \sum_{i=1}^N (f(x_i, w) - y_i)^2 \right\} \text{ for } \sum_{j=0}^p |w_j| \leq t.$$

**Prove it.**

# Ridge and lasso



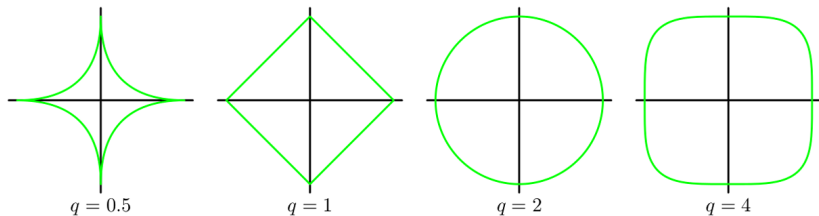
# Generalization

- We can generalize ridge and lasso regression to

$$L(w) = \frac{1}{2} \sum_{i=1}^N (f(x_i, w) - y_i)^2 + \lambda \sum_{j=0}^p (|w_j|)^q.$$

**Which prior distribution on  $w$  does this correspond to?**

# Different $q$



Thank you!

**Thank you for your attention!**