

Least Squares for a linear model

$$y = a_1 x_1 + \dots + a_k x_k + \varepsilon \quad \text{where } \varepsilon \sim N(0, \sigma)$$

in other words

$$\hat{y} = a_1 x_1 + \dots + a_k x_k$$

Problem: given a dataset $D = \{(\vec{x}_n, y_n)\}$
 where $x_n = (x_{n1}, \dots, x_{nk}) \in \mathbb{R}^k$
 find $a_1, \dots, a_k \in \mathbb{R}$ which
 minimize the squared error

$$SE = \sum_n (y_n - (a_1 x_{n1} + \dots + a_k x_{nk}))^2$$

SE is a quadratic function of a_1, \dots, a_k
 and the coefficient in front of the quadratic term
 is positive hence the minimum is in the point
 (a_1, \dots, a_k) where

$$\forall_i \frac{\partial}{\partial a_i} SE(a_1, \dots, a_k) = 0$$

For simplicity of notation consider the 2D case
 where we need to find (a_1, a_2)

We need to find a_1, a_2 where

$$\begin{cases} \frac{\partial}{\partial a_1} \sum_n (y_n - (a_1 x_{n1} + a_2 x_{n2}))^2 = 0 \\ \frac{\partial}{\partial a_2} \sum_n (y_n - (a_1 x_{n1} + a_2 x_{n2}))^2 = 0 \end{cases}$$

We have

$$\begin{aligned} \frac{\partial}{\partial a_1} \sum_n (y_n - (a_1 x_{n1} + a_2 x_{n2}))^2 &= \sum_n \frac{\partial}{\partial a_1} (y_n - (a_1 x_{n1} + a_2 x_{n2}))^2 \\ &= \sum_n 2 (y_n - (a_1 x_{n1} + a_2 x_{n2})) (-x_{n1}) \\ &= -2 \sum_n x_{n1} (y_n - (a_1 x_{n1} + a_2 x_{n2})) \end{aligned}$$

$$\frac{\partial}{\partial a_2} \sum_n (y_n - (a_1 x_{n1} + a_2 x_{n2}))^2 = -2 \sum_n x_{n2} (y_n - (a_1 x_{n1} + a_2 x_{n2}))$$

Hence we have to solve the following system of linear equations:

$$\begin{cases} \sum_n x_{n1} (y_n - (a_1 x_{n1} + a_2 x_{n2})) = 0 \\ \sum_n x_{n2} (y_n - (a_1 x_{n1} + a_2 x_{n2})) = 0 \end{cases}$$

\Leftrightarrow

$$\begin{cases} \sum_n x_{n1} y_n = \sum_n x_{n1}^2 a_1 + \sum_n x_{n1} x_{n2} a_2 \\ \sum_n x_{n2} y_n = \sum_n x_{n1} x_{n2} a_1 + \sum_n x_{n2}^2 a_2 \end{cases}$$

\Leftrightarrow

$$\begin{cases} \sum_n x_{n1} y_n = \alpha_1 \sum_n x_{n1}^2 + \alpha_2 \sum_n x_{n1} x_{n2} \\ \sum_n x_{n2} y_n = \alpha_1 \sum_n x_{n1} x_{n2} + \alpha_2 \sum_n x_{n2}^2 \end{cases} \quad (*)$$

\Updownarrow

$$\begin{cases} \sum_n x_{n1} y_n \sum_n x_{n2}^2 = \alpha_1 \sum_n x_{n1}^2 \sum_n x_{n2}^2 + \alpha_2 \sum_n x_{n1} x_{n2} \sum_n x_{n2}^2 \\ \sum_n x_{n2} y_n \sum_n x_{n1} x_{n2} = \alpha_1 \sum_n x_{n1} x_{n2} \sum_n x_{n1} x_{n2} + \alpha_2 \sum_n x_{n2}^2 \sum_n x_{n1} x_{n2} \end{cases}$$

\Downarrow

$$\alpha_1 = \frac{\sum_n x_{n2} y_n \sum_n x_{n1} x_{n2} - \sum_n x_{n1} y_n \sum_n x_{n2}^2}{\left(\sum_n x_{n1} x_{n2}\right)^2 - \sum_n x_{n1}^2 \sum_n x_{n2}^2}$$

Denoting $X_1 = (x_{11}, x_{21}, \dots, x_{n1})$ - the first coordinate of every datapoint
 $X_2 = (x_{12}, x_{22}, \dots, x_{n2})$ - the second coordinate of every datapoint

we can write the formula for α_1 in a concise way

$$\alpha_1 = \frac{(x_2 \cdot y)(x_1 \cdot x_2) - (x_1 \cdot y) \|x_2\|^2}{x_1 \cdot x_2 - \|x_1\|^2 \|x_2\|^2}$$

and for α_2 :

$$\alpha_2 = \frac{(x_1 \cdot y)(x_1 \cdot x_2) - (x_2 \cdot y) \|x_1\|^2}{x_1 \cdot x_2 - \|x_1\|^2 \|x_2\|^2}$$

The system of equations (*) can be written in a more generic way

$$\begin{cases} \sum_n x_{n1} y_n = \alpha_1 \sum_n x_{n1}^2 + \alpha_2 \sum_n x_{n1} x_{n2} \\ \sum_n x_{n2} y_n = \alpha_1 \sum_n x_{n1} x_{n2} + \alpha_2 \sum_n x_{n2}^2 \end{cases}$$



$$X^T y = X^T X A$$

where

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad A = \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix}$$

Then

$$X^T y = X^T X A \quad | \cdot (X^T X)^{-1}$$



$$A = (X^T X)^{-1} X^T y$$

This form holds also for k -dimensional x vectors

$X^T X$ is a $k \times k$ dimensional matrix
 \Rightarrow if k is small, inverting $X^T X$ is not costly

ALS - Alternating Least Squares

Recall that the matrix factorization problem is given by

$$\min_{p_u, q_i \in \mathbb{R}^d} \sum_{(u,i) \in K} (r_{u,i} - q_i^T p_u)^2$$

If we fix the item representations q_i , then the problem becomes

54 ...

$$\min_{p_u \in \mathbb{R}^d} \sum_i (r_{u,i} - (q_{i,1} p_{u,1} + \dots + q_{i,d} p_{u,d}))^2$$

where this expression has to be minimized for every user over all possible values of $p_u = (p_{u,1}, p_{u,2}, \dots, p_{u,d})$

This is the Linear Least Squares problem!

Therefore

$$p_u = (X^T X)^{-1} X^T y$$

where

$$X = \begin{bmatrix} q_{1,1} & q_{1,2} & \dots & q_{1,d} \\ q_{2,1} & q_{2,2} & \dots & q_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ q_{m,1} & q_{m,2} & \dots & q_{m,d} \end{bmatrix} \quad y = \begin{bmatrix} r_{u,1} \\ r_{u,2} \\ \vdots \\ r_{u,m} \end{bmatrix}$$

ALS

1. Initialize all user and item representation vectors p_u and q_i with random values
2. Iterate until convergence
(i.e. changing of representations less than ϵ)
 - 2.a. Set all item representations q_i and solve the Linear Least Squares problem for user representations p_u
 - 2.b. Set all user representations p_u and solve the Linear Least Squares problem for item representations q_i

Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation (MLE)

Consider again the following linear model

$$y = a_1 x_1 + \dots + a_k x_k + \varepsilon \quad \text{where } \varepsilon \sim N(0, \sigma)$$

Assuming this model is true, the likelihood of observing a datapoint $(x_1, x_2, \dots, x_k, y)$ in the data is equal to

$$L(\varepsilon) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{\varepsilon}{\sigma}\right)^2} = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y - (a_1 x_1 + a_2 x_2 + \dots + a_k x_k)}{\sigma}\right)^2}$$

The idea behind MLE is that for a given set of observed datapoints $\{(\vec{x}_n, y_n)\} = \{(x_{n1}, x_{n2}, \dots, x_{nk}, y_n)\}$ we want to find such model parameters a_1, a_2, \dots, a_k that the likelihood of observing such dataset is maximal, i.e. we want to solve

$$\max_{a_1, \dots, a_k} \prod_n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y - (a_1 x_1 + a_2 x_2 + \dots + a_k x_k)}{\sigma}\right)^2}$$

This expression can be further simplified since:

$$\begin{aligned} \arg \max_{a_1, \dots, a_k} \prod_n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y - (a_1 x_1 + a_2 x_2 + \dots + a_k x_k)}{\sigma}\right)^2} \\ &= \arg \max_{a_1, \dots, a_k} \prod_n e^{-\frac{1}{2} \left(\frac{y - (a_1 x_1 + a_2 x_2 + \dots + a_k x_k)}{\sigma}\right)^2} \\ &= \arg \max_{a_1, \dots, a_k} e^{-\frac{1}{2} \sum_n \left(\frac{y - (a_1 x_1 + a_2 x_2 + \dots + a_k x_k)}{\sigma}\right)^2} \quad \left(\text{since } e^a e^b = e^{a+b} \right) \\ &= \arg \min_{a_1, \dots, a_k} \sum_n \left(\frac{y - (a_1 x_1 + a_2 x_2 + \dots + a_k x_k)}{\sigma}\right)^2 \quad \left(\text{since } e^{-x} \text{ is decreasing} \right) \\ &= \arg \min_{a_1, \dots, a_k} \sum_n (y - (a_1 x_1 + a_2 x_2 + \dots + a_k x_k))^2 \end{aligned}$$

$$= \underset{a_1, \dots, a_n}{\operatorname{argmin}} \sum_n (y - (a_1 x_1 + a_2 x_2 + \dots + a_n x_n))^2$$

But this is exactly Least Squares !

MLE and Least Squares are equivalent
if the noise in the data is normal (Gaussian)

Note

MLE is a powerful and general method
which can be used with any probability
distribution, for instance Bernoulli, Binomial,
Poisson, Exponential, Gamma, Beta