

## Least squares and Maximum Likelihood Estimation (MLE)

Friday, 30 April 2021 07:42

### Least Squares for a linear model

$$y = \alpha_1 x_1 + \dots + \alpha_k x_k + \varepsilon \quad \text{where } \varepsilon \sim N(0, \sigma^2)$$

in other words

$$\hat{y} = \alpha_1 x_1 + \dots + \alpha_k x_k$$

Problem: given a dataset  $D = \{(\vec{x}_n, y_n)\}$   
 where  $\vec{x}_n = (x_{n1}, \dots, x_{nk}) \in \mathbb{R}^k$   
 find  $\alpha_1, \dots, \alpha_k \in \mathbb{R}$  which  
 minimize the squared error

$$SE = \sum_n (y_n - (\alpha_1 x_{n1} + \dots + \alpha_k x_{nk}))^2$$

$SE$  is a quadratic function of  $\alpha_1, \dots, \alpha_k$   
 and the coefficient in front of the quadratic term  
 is positive hence the minimum is in the point  
 $(\alpha_1, \dots, \alpha_k)$  where

$$\nabla_{\alpha_i} SE(\alpha_1, \dots, \alpha_k) = 0$$

For simplicity of notation consider the 2D case  
 where we need to find  $(\alpha_1, \alpha_2)$

We need to find  $\alpha_1, \alpha_2$  where

$$\left\{ \begin{array}{l} \frac{\partial}{\partial \alpha_1} \sum_n (y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2}))^2 = 0 \\ \frac{\partial}{\partial \alpha_2} \sum_n (y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2}))^2 = 0 \end{array} \right.$$

we have

$$\begin{aligned} \frac{\partial}{\partial \alpha_1} \sum_n (y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2}))^2 &= \sum_n \frac{\partial}{\partial \alpha_1} (y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2}))^2 \\ &= \sum_n 2(y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2}))(-x_{n1}) \\ &= -2 \sum_n x_{n1} (y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2})) \end{aligned}$$

$$\frac{\partial}{\partial \alpha_2} \sum_n (y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2}))^2 = -2 \sum_n x_{n2} (y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2}))$$

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

$$\left\{ \begin{array}{l} \sum_n x_{n1} (y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2})) = 0 \\ \sum_n x_{n2} (y_n - (\alpha_1 x_{n1} + \alpha_2 x_{n2})) = 0 \end{array} \right.$$

↑

$$\left\{ \begin{array}{l} \sum_n x_{n1} y_n = \sum_n x_{n1}^2 \alpha_1 + \sum_n x_{n1} x_{n2} \alpha_2 \\ \sum_n x_{n2} y_n = \sum_n x_{n1} x_{n2} \alpha_1 + \sum_n x_{n2}^2 \alpha_2 \end{array} \right.$$

↑

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



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



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

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  $\alpha_1$  in a concise way

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

and for  $\alpha_2$ :

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

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

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad A = \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix}$$

Then

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

↑

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

Thus 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

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

If we fix the item representations  $q_i$ , then the problem becomes independent for every user and takes the following form:

$$\min_{\beta_u \in \mathbb{R}^d} \sum_{i \in K(u)} (r_{ui} - (q_{i1}\beta_{u1} + \dots + q_{id}\beta_{ud}))^2$$

where this expression has to be minimized for every user over all possible values of  $\beta_u = (\beta_{u1}, \beta_{u2}, \dots, \beta_{ud})$

This is the Linear Least Squares problem!

Therefore

$$\beta_u = (X^\top X)^{-1} X^\top y$$

where

$$X = \begin{bmatrix} q_{11} & q_{12} & \dots & q_{1d} \\ q_{21} & q_{22} & \dots & q_{2d} \\ \vdots & \vdots & & \vdots \\ q_{m1} & q_{m2} & \dots & q_{md} \end{bmatrix} \quad y = \begin{bmatrix} r_{u1} \\ r_{u2} \\ \vdots \\ r_{um} \end{bmatrix}$$

ALS

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

+ Scalability

Because the Linear Least Squares problem is independent for every user/item, the algorithm can be highly parallelized

## Maximum Likelihood Estimation (MLE)

Consider again the following linear model

$$y = \alpha_1 x_1 + \dots + \alpha_k x_k + \varepsilon \quad \text{where } \varepsilon \sim N(0, \sigma^2)$$

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

$$L(\varepsilon) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{\varepsilon}{\sigma}\right)^2} = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y - (\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k)}{\sigma}\right)^2}$$

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  $\alpha_1, \alpha_2, \dots, \alpha_k$  that the likelihood of observing such dataset is maximal, i.e. we want to solve

$$\max_{\alpha_1, \dots, \alpha_k} \prod_n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y - (\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k)}{\sigma}\right)^2}$$

This expression can be further simplified since:

$$\begin{aligned} & \arg \max_{\alpha_1, \dots, \alpha_k} \prod_n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y - (\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k)}{\sigma}\right)^2} \\ &= \arg \max_{\alpha_1, \dots, \alpha_k} \prod_n e^{-\frac{1}{2} \left(\frac{y - (\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k)}{\sigma}\right)^2} \\ &= \arg \max_{\alpha_1, \dots, \alpha_k} e^{-\frac{1}{2} \sum_n \left(\frac{y - (\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k)}{\sigma}\right)^2} \end{aligned}$$

(since  $e^a e^b = e^{a+b}$ )

$$\begin{aligned}
 &= \underset{\theta_1, \dots, \theta_n}{\operatorname{argmax}} \ell \sim n \\
 &= \underset{\theta_1, \dots, \theta_n}{\operatorname{argmin}} \sum_n \left( \frac{y - (\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n)}{\sigma} \right)^2 \\
 &= \underset{\theta_1, \dots, \theta_n}{\operatorname{argmin}} \sum_n (y - (\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n))^2
 \end{aligned}$$

( $e^u e^v = e^{u+v}$ )  
 (since  $e^{-x}$   
 is decreasing)

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

This method often inspires the loss functions  
 for various neural networks and other models