

## LINEAR REGRESSION

GIVEN A PAIR OF RV  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$ , in linear regression analysis we are interested to relate  $Y$  with the observations  $X$  using a REGRESSION FUNCTION

$$m: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$m(x) = E(Y|X=x)$$

Why do we choose this form? Given a generic function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , in general

$$\begin{aligned} E[(f(x) - m(x))(m(x) - Y)] &\stackrel{LIE}{=} E[E[(f(x) - m(x))(m(x) - Y)|X]] \\ \stackrel{LIE \text{ (LOW VARIANCE EXPECTATIONS)}}{=} E[f(x) - m(x)] \cdot E[m(x) - Y] &= E[f(x) - m(x)] \cdot E[m(x) - E(Y|X)] \\ E[X] = E[E[XY]] &= E[(f(x) - m(x)) \cdot (m(x) - E(Y|X))] \\ = \sum_{i=1}^n E[X|A_i] P(A_i) &= E[(f(x) - m(x)) \cdot (m(x) - E(Y|X))] \\ &= 0 \end{aligned}$$

$$\begin{aligned} E[(f(x) - Y)^2] &= E[(f(x) - m(x) + m(x) - Y)^2] = \\ &E[(f(x) - m(x))^2 + 2(f(x) - m(x))(m(x) - Y) + (m(x) - Y)^2] \\ &= E[(f(x) - m(x))^2] + 2E[(f(x) - m(x))(m(x) - Y)] + E[(m(x) - Y)^2] \\ &= E[(f(x) - m(x))^2] + E[(m(x) - Y)^2] \end{aligned}$$

This quantity is minimized when  $f(x) = m(x)$  because the quantity  $E[(m(x) - Y)^2]$  is unavoidable, is always present

So  $m(x)$  satisfies the following  $E[(m(x) - Y)^2] = \min_{f: \mathbb{R}^d \rightarrow \mathbb{R}} E[(f(x) - Y)^2]$  and we can say that the OPTIMAL APPROXIMATION (wrt the MEAN SQUARE) of  $Y$  by a function  $f$  is the regression function  $m(x)$

$X, Y$  are RVs so they are UNKNOWN  $\xrightarrow{\text{BUT}}$  we can estimate  $m(x)$  using data

## LINEAR REGRESSION MODEL

$D = \{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_m, y_m)\}$   $(\bar{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$  and they are i.i.d.

$$y_i = m(x_i) + \varepsilon_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2} + \dots + \theta_d x_{id} + \varepsilon_i$$

$i = 1, \dots, m$

!  $m$  is the dimension of the data set  
 $d$  is the dimension of the vector  $\rightarrow$  i.e.  $d$  coordinates

$$\bar{x}_i = (x_{i1}, \dots, x_{id})^d$$

$\theta_0$  is the intercept,  $\epsilon_i$  is the noise such that  $\mathbb{E}[\epsilon_i] = 0$   
 $\text{Var}[\epsilon_i] = \sigma^2$

$$\bar{x}_i \rightarrow \bar{x}_i = \begin{pmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{id} \end{pmatrix} \in \mathbb{R}^{d+1} \quad \bar{\theta} = \begin{pmatrix} \theta_0 \\ \vdots \\ \theta_d \end{pmatrix} \in \mathbb{R}^{d+1}$$

$$y_i = \begin{pmatrix} 1, x_{i1}, x_{i2}, \dots, x_{id} \end{pmatrix} \begin{pmatrix} \theta_0 \\ \vdots \\ \theta_d \end{pmatrix} + \epsilon_i = \bar{x}_i^T \cdot \bar{\theta} + \epsilon_i \quad (i=1, \dots, m)$$

Since the previous is valid for each point of our dataset. We can do

$$\bar{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m, \quad \bar{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_m \end{pmatrix} \in \mathbb{R}^m$$

$$X \equiv \begin{pmatrix} - \bar{x}_1^T - \\ \vdots \\ - \bar{x}_m^T - \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & \dots & x_{md} \end{pmatrix}$$

$[m \times (d+1)]$

$$\begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & \dots & x_{md} \end{pmatrix} \begin{pmatrix} \theta_0 \\ \vdots \\ \theta_d \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_m \end{pmatrix} = \bar{y} = X\bar{\theta} + \bar{\epsilon}$$

$$y_i = \sum_{j=0}^d X_{ij} \theta_j + \epsilon_i$$

Residuals

LINEAR REGRESSION MODEL

→ LEAST SQUARES ⇒ METHOD FOR ESTIMATING THE REGRESSION MODEL

💡 MINIMIZING THE RESIDUALS SUM OF SQUARES (RSS)

$$RSS(\bar{\theta}) = \sum_{i=1}^m \varepsilon_i^2 = \sum_{i=1}^m \left( y_i - \sum_{j=1}^d X_{ij} \theta_j \right)^2$$

In general,  $\bar{u}^T u = \sum_{i=1}^k u_i^2 = \|u\|_2^2 = \langle u, u \rangle$

$$\begin{aligned} RSS(\bar{\theta}) &= (\bar{y} - X\bar{\theta})^T (\bar{y} - X\bar{\theta}) = (\bar{y}^T - \bar{\theta}^T X^T) (\bar{y} - X\bar{\theta}) \\ &= \bar{y}^T \bar{y} - \bar{y}^T X \bar{\theta} - \bar{\theta}^T X^T \bar{y} + \bar{\theta}^T X^T X \bar{\theta} \quad (\bar{y}^T X \bar{\theta})^T \\ &= \bar{y}^T \bar{y} - 2 \bar{\theta}^T X^T \bar{y} + \bar{\theta}^T X^T X \bar{\theta} \end{aligned}$$

m=1

d=2

$$\frac{1}{2} \begin{bmatrix} 1 & 1 & 2 \\ -1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} = 8$$

$$\begin{bmatrix} 3 & 3 & 4 \\ 2 & 2 & 2 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} = 20$$

6 + 14 = 20

$$\begin{bmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = 8$$

$$\begin{bmatrix} 0 & 6 \\ 0 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = 10$$

$$\min_{\bar{\theta}} RSS(\bar{\theta}) \Rightarrow \left. \frac{\partial RSS(\bar{\theta})}{\partial \bar{\theta}} \right|_{\bar{\theta}=\hat{\theta}} = 0 - 2 X^T \bar{y} + 2 X^T X \hat{\theta} = 0$$

$\hat{\theta} = (\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_d)^T$  is the estimator that minimize RSS

$$\hat{\beta} = \arg \min_{\beta} RSS(\beta) = (X^T X)^{-1} X^T \bar{y}$$

EXAMPLE

m=3, d=1

{(1, 3.1), (3, 10.5), (4, 17.6)}

$$\bar{x} = \begin{pmatrix} 1 \\ 3 \\ 4 \end{pmatrix}$$

$$\bar{y} = \begin{pmatrix} 3.1 \\ 10.5 \\ 17.6 \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 4 \end{pmatrix}$$

.....

SINGULAR MATRIX  $\rightarrow$  NOT INVERTIBLE BECAUSE OF THE DETERMINANT  
determinant =  $\phi$

$\rightarrow$  if  $X^T X$  is not invertible?

### REGULARIZATION

- RIDGE  $\rightarrow$  it introduces a penalty term depending on the square ( $L^2$ -norm) of the vector of the regression coefficients

$$\hat{\theta} = \arg \min_{\theta} \left[ \sum_{i=1}^m \left( y_i - \sum_{j=1}^d x_{ij} \theta_j \right)^2 + \lambda \sum_{j=1}^d \theta_j^2 \right]$$

$$RSS = (\bar{y} - X\bar{\theta})^T (\bar{y} - X\bar{\theta}) + \lambda \bar{\theta}^T \bar{\theta}$$

$$\bar{\theta}^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T \bar{y}$$

- LASSO

$\rightarrow$  PENALIZATION TERM IS  $\lambda \sum_{j=1}^d |\bar{\theta}_j|$

This regularization has not CLOSED FORM  $L^1$ -norm  
we need something to compute it