

**Supervised**  
 input  $\xrightarrow{\text{algorithm}}$  target  
 new input  $\nrightarrow$  ?  $\rightarrow$  if the new input does not meet some specific criteria, algorithm does not know what to do with the data  
 (Classification)  
 (Regression)

**unsupervised learning**  
 subset input  $\xrightarrow{\text{algorithm}}$  target  
 new input  $\xrightarrow{\text{algorithm}}$  target  
 Algorithm learns to make decisions after training (based on some criteria, it gets to target)  
 (E-M algorithm, PCA, NMF, etc.)  
 $\hookrightarrow$  latent variables algorithm

## Linear Regression

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$  (each  $x_i$  is a vector, see below to understand)  
 $y_1, \dots, y_n \in \mathbb{R}$  ( $y_1, \dots, y_n$  is a single vector)

$$y_i = \vec{x}_i^T \theta^* + \epsilon_i \quad \text{for } i = 1, \dots, n$$

is equivalent to:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix}_{n \times d} \begin{bmatrix} \theta_1^* \\ \theta_2^* \\ \vdots \\ \theta_d^* \end{bmatrix}_{d \times 1} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}_{n \times 1}$$

$\vec{x}_1^T \in \mathbb{R}^d$

this matrix can be also written as follows.

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} \xrightarrow{\text{transpose the original}} \begin{bmatrix} \vec{x}_1^T & \vec{x}_2^T & \dots & \vec{x}_n^T \end{bmatrix} \Leftrightarrow \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}_{n \times d}$$

$\vec{x}_1 \quad \vec{x}_2 \quad \dots \quad \vec{x}_n \in \mathbb{R}^d$

hence

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}_{n \times d} \begin{bmatrix} \theta_1^* \\ \theta_2^* \\ \vdots \\ \theta_d^* \end{bmatrix}_{d \times 1} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}_{n \times 1}$$

$= n \times 1$

✓ dimensions OK!

or equivalently:  $y_i = \vec{x}_i^T \theta^* + \epsilon_i \quad \text{for } i = 1, \dots, n$

# Ordinary Least Square

$$\hat{\theta}^{OLS} = \underset{\theta}{\operatorname{argmin}} \left( \frac{1}{2n} \|Y - \hat{Y}\|_2^2 \right)$$

sum of squared errors

(2)

but

this term here is actually optional, some literature includes, some not!

$$\|Y - \hat{Y}\|_2^2 = \|Y - X\hat{\theta}\|_2^2 = \sum_{i=1}^n (y_i - x_i^T \theta)^2$$

↔  
equivalent

recall from our Linear Algebra Lesson.  
the  $l_2$ -norm  $\|x\|_2 = \sqrt{\langle x, x \rangle}$   
Hence  $\|x\|_2^2 = \langle x, x \rangle$  (inner product)  
 $\|x\|_2^2 = \sum_{i=1}^n x_i^2$

Let's do both ways

① Using summation

$$\hat{\theta}^{OLS} = \underset{\theta}{\operatorname{argmin}} \frac{1}{2n} \left( \sum_{i=1}^n (y_i - x_i^T \theta)^2 \right) = \underset{\theta}{\operatorname{argmin}} \frac{1}{2n} \left( \sum_{i=1}^n (y_i^2 - 2y_i x_i^T \theta + (x_i^T \theta)^2) \right)$$

F.O.C.

First order condition

= derivative with respect to  $\theta$  and set to 0

$$\frac{\partial f(\theta)}{\partial \theta} = \frac{-2}{2n} \sum_{i=1}^n (y_i x_i^T)^T + \frac{2}{2n} \sum_{i=1}^n (x_i x_i^T \theta)$$

$$= -\frac{1}{n} \left( \sum_{i=1}^n x_i y_i - \sum_{i=1}^n (x_i x_i^T \theta) \right) \xrightarrow{\text{F.O.C.}} 0$$

$$\Leftrightarrow -\frac{1}{n} \sum x_i y_i = -\frac{1}{n} \sum x_i x_i^T \theta \Rightarrow \hat{\theta}^{OLS} = (\sum x_i x_i^T)^{-1} \sum x_i y_i$$

where:

$$\sum_{i=1}^n x_i x_i^T = \sum_{i=1}^n \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix} (x_{i1} \ x_{i2} \ \dots \ x_{id})_{1 \times d} = \begin{pmatrix} \sum x_{i1}^2 & \sum x_{i2} x_{i1} & \dots & \sum x_{id} x_{i1} \\ \vdots & \sum x_{i2}^2 & & \vdots \\ \sum x_{i1} x_{id} & \dots & \dots & \sum x_{id}^2 \end{pmatrix}_{d \times d}$$

must be invertible

And

$$\sum_{i=1}^n x_i y_i = \begin{pmatrix} \sum x_{i1} y_i \\ \sum x_{i2} y_i \\ \vdots \\ \sum x_{id} y_i \end{pmatrix}_{d \times 1}$$

using matrix notation.

$$Y \in \mathbb{R}^{n \times 1} \quad \theta \in \mathbb{R}^{d \times 1} \quad X \in \mathbb{R}^{n \times d} \quad (3)$$

$$\frac{1}{2n} \|Y - \hat{Y}\|_2^2 = \|Y - X\hat{\theta}\|_2^2 = (Y - X\hat{\theta})^T (Y - X\hat{\theta}) = Y^T Y - \underbrace{Y^T X \hat{\theta}}_{1 \times 1} - \underbrace{\hat{\theta}^T X^T Y}_{1 \times 1} + \underbrace{\hat{\theta}^T X^T X \hat{\theta}}_{n \times n}$$

$$f(\theta) = Y^T Y - 2 \theta^T X^T Y + \theta^T X^T X \theta$$

F.O.C.  $\frac{\partial f(\theta)}{\partial \theta} = \frac{1}{2n} (-2 X^T Y + 2 X^T X \theta) = 0$

$$\Leftrightarrow \frac{1}{n} X^T X \theta = \frac{1}{n} X^T Y \Rightarrow \hat{\theta}^{OLS} = (X^T X)^{-1} X^T Y$$

Equivalently

$$\hat{\theta}^{OLS} = \arg \min_{\theta} \underbrace{\frac{1}{2n} \|Y - X\theta\|_2^2}_{L(\theta)}$$

Gradient

$$\nabla L(\hat{\theta}) = \frac{1}{2n} (2)(-1) X^T (Y - X\theta)$$

$$\Leftrightarrow \frac{1}{n} X^T (X\theta - Y) = 0$$

$\Leftrightarrow$

$$X^T X \theta = X^T Y$$

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

Hessian

$$\nabla^2 L(\hat{\theta}) = \frac{\partial^2 L(\theta)}{\partial \theta^2} = \frac{1}{n} X^T X > 0$$

positive definite

(strictly convex)

global minimum

Likelihood function (since  $\varepsilon \sim N(0, \sigma^2 I_n) \Rightarrow \hat{\theta} \sim MN(\bar{\theta}, \sigma^2 (X^T X)^{-1})$   $\bar{\theta}$  is a vector

$$L(\theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2} \frac{(Y - X\theta)^T (Y - X\theta)}{\|Y - X\theta\|_2^2}}$$

take log:  $\log L(\theta) = + \log 1 - \frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|Y - X\theta\|_2^2$

1<sup>st</sup> derivative  $\frac{\partial \log L(\theta)}{\partial \theta} \Rightarrow$  for some  $\sigma^2$  constant, with  $\varepsilon \sim N(0, \sigma^2 I_n)$   
OLS provides same solution as MLE.