# Regression — Closed Form Solution

$$\theta^* = (X^T X)^{-1} X^T y$$

$$X = \begin{bmatrix} & \vdots & \\ - & x^i & - \\ & \vdots & \end{bmatrix}_{m \times (n+1)}$$

$\theta^T x$

$m$ - no. of ex

$n$ = no. of feature

$x_j^i$  $i^{th}$ ex  $j^{th}$ feature

$$\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}_{(n+1 \times 1)}$$

$$y = \theta_1 x + \theta_0$$

## Proof

$$loss = \frac{1}{2} \sum_{i=1}^{m} (y_{pred} - y)^2$$

$$X = \begin{bmatrix} x^1 & 1 \\ \vdots & 1 \\ - x^i & 1 \\ x^m & 1 \end{bmatrix}_{m \times (n+1)} \qquad y = \begin{bmatrix} y^1 \\ y^2 \\ y^m \end{bmatrix}_{m \times 1}$$

$$y_{pred} = h_\theta(x) = \sum_{i=0}^{n} \theta_i x_i^j$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}_{\left(\frac{n}{+} \atop 1\right) \times 1}$$

we introduce one more column of 1's in $x$ as

$$h_\theta(x) = \theta_0 \boxed{x_0} + \theta_1 x_1 + \theta_2 x_2 + \dots \theta_i x_i$$

$$x_0 = 1$$

$$= \theta^T x = \sum_{i=0}^{m} \theta_i x_i$$

## ~~Proof~~

$$\underbrace{loss = \frac{1}{2} (X\theta - y)^2}_{\text{matrix notation}} = \frac{1}{2} \sum_{i=1}^{m} (y_{pred}^i - y^i)^2$$

$X_0^1$ = 1ˢᵗ example 0ᵗʰ feature

$$\begin{bmatrix} X_0^1 & X_1^1 & X_2^1 \\ X_0^2 & X_1^2 & X_2^2 \\ X_0^3 & X_1^3 & X_2^3 \\ \vdots \\ X_0^n & X_1^m & X_2^m \end{bmatrix} = X \qquad \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} = \theta$$

$\underset{1}{\underbrace{\phantom{xx}}}$

$$X\theta = \begin{bmatrix} \theta_0 X_0^1 + \theta_1 X_1^1 + \theta_2 X_2^1 \\ \theta_0 X_0^2 + \theta_1 X_1^2 + \theta_2 X_2^2 \\ \vdots \end{bmatrix}$$

$$= \begin{bmatrix} \theta_0 + \theta_1 X_1^1 + \theta_2 X_2^1 \\ \theta_0 * \theta_1 X_1^2 + \theta_2 X_2^2 \\ \vdots \end{bmatrix}$$

$$X\theta = \begin{bmatrix} y_{pred}^1 \\ y_{pred}^2 \\ \vdots \\ y_{pred}^m \end{bmatrix}$$

$$(\underbrace{X\theta}_{} - \underbrace{y}_{})^2 = loss$$

$y_{pred}$     $y_{actual}$

$$\left( \begin{bmatrix} \\ \\ \\ \end{bmatrix}_{M \times 1} - \begin{bmatrix} \\ \\ \\ \end{bmatrix}_{M \times 1} \right)^2 = \text{scalar}$$

square of a vector is scalar

$$J(\theta) = (X\theta - y)^2 \qquad z^2 = z^T z$$
$$= (X\theta - y)^T (X\theta - y) \quad \text{loss in matrix}$$
$$\qquad \qquad \qquad \qquad \qquad \text{notation}$$
$$: (\theta^T X^T - y^T)(X\theta - y)$$
$$\qquad \qquad \qquad \qquad \qquad (AB)^T = B^T A^T$$
$$J(\theta) = \theta^T X^T X \theta - \theta^T X^T y - y^T X \theta + y^T y$$

minimize $J(\theta)$

$$\nabla_\theta J(\theta) = 0$$

computing first derivative of $J(\theta)$

$$\nabla_\theta J(\theta) = \nabla_\theta (\theta^T X^T X \theta \overset{③rule}{} - \theta^T X^T y - y^T X \theta + y^T y)^0$$

$$= 2 X^T X \theta - \nabla_\theta (\theta^T X^T y + \theta^T X^T y)$$

$$= 2 X^T X \theta - \nabla_\theta (2\theta^T X^T y)^{②rule}$$

$$= 2 X^T X \theta - 2 X^T y = 0 \text{ for minima}$$

$$= X^T X \theta = X^T y$$

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

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

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

Rules of matrix calculus

① $\nabla_\theta a^T \theta = a$ ② $\nabla_\theta \theta^T a = a$

③ $\nabla_\theta \theta^T A \theta = 2 A \theta$

## Gradient Descent
→ iterative approach
→ figure out learning rate
→ slow if n is small
→ useful for big dataset
→ minibatch is useful for very large dataset

## Closed form solution
→ small dataset
→ get directly solution
→ expensive if dataset is large