

Supervised Learning

prediction

$$h: x \rightarrow y$$

Images cat

TEXT

Is hate speech?

House data

Price

Given Training Set

$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$$

$$x^{(i)} \in X, y^{(i)} \in Y$$

Do: find good $\star h: X \rightarrow y$ (hypothesis)

start with how accurate it is -

We care about new X values not in training set

if y is discrete \Rightarrow "classification"

continuous \Rightarrow "Regression"

First plot it, from the data set

Sales price	Lot in sqft
000 ~	000 ~

$\rightarrow h: \text{lot area} \rightarrow \text{price}$

How do we represent the h ?

$$h(x) = \theta_0 + \theta_1 x_1$$

size	bedroom	...	price
$x^{(1)}$	$x_1^{(1)} = 1$	$x_2^{(1)}$	400
2104	4		
$x^{(2)}$	2500	3	900

$$h(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d$$

$$= \sum_{y=0}^d \theta_y \cdot x_y \quad \underline{\text{NB}} \quad x_0 = 1$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

parameters

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \end{bmatrix}$$

features

$y^{(i)}$ is price

$(x^{(i)}, y^{(i)}) \leftarrow$ training graph

$$x = \begin{bmatrix} -x^{(1)}- \\ \vdots \\ -x^{(n)}- \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$$

matrix

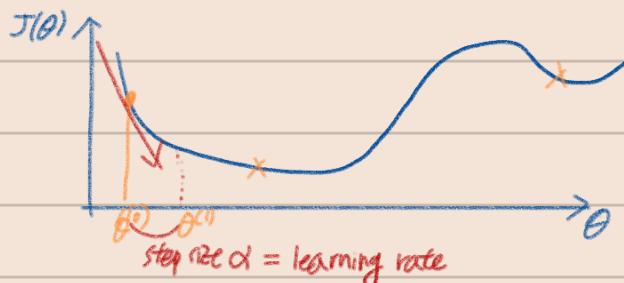
$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j$$

선형회귀 손실함수 Mean Squared Error

IDEA $J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 \Rightarrow \text{"optimization"}$

Minimize

Gradient Descent



선형회귀, 3차원 이상의 손실함수
convexfun + bounded
every local min = global min
→ 미지수 n개를 최적화 하기 어렵다
어려워서 최적화가 도달할 수 있다면 안정적

good clean model = bounded + convex

$$\theta^0 = 0$$

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \alpha \frac{\partial}{\partial \theta_j} J(\theta^{(t)}) \quad j = 0 \dots d$$

partial derivative w respect to theta

< too big: go to far
< too small: not enough progress

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \sum_{i=1}^n \frac{1}{2} \frac{\partial}{\partial \theta_j} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot \boxed{\frac{\partial}{\partial \theta_j} h_{\theta}(x^{(i)})} \end{aligned}$$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d$$

$$\frac{\partial}{\partial \theta_j} h_{\theta}(x) = x_j$$

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \alpha \cdot \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot \boxed{x_j^{(i)}}$$

$$\text{vector Equation} \Rightarrow \theta^{(t+1)} = \theta^{(t)} - \alpha \cdot \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) X^{(i)}$$

점화식.

수렴할 때까지 반복, 정해진 수만큼 진행

convex이므로 gradient는 항상 작아짐
but, for non-convex functions, we really don't know

(non-deterministic)

Batch vs. Stochastic mini batch

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \cdot \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

mini batch B (at random) $B \ll N$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha_B \sum_{i \in B}^n (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

come back when we look at different loss functions

Normal Equations ← 어떤 방식으로든 같은 결과를 얻는다: