

Supervised learning setup

Notation. M : # samples

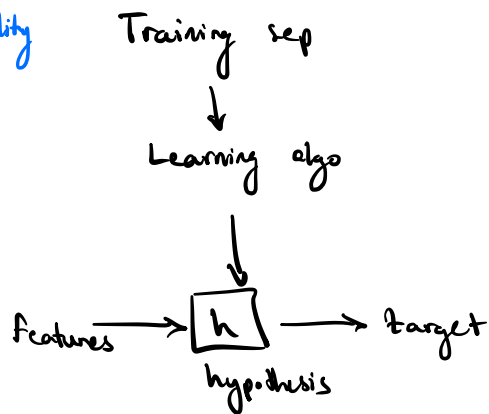
X : features (input)

Y : target (output)

(X, Y) : training data

$(x^{(i)}, y^{(i)})$: i th training example

Functionality



Formalization Initially (linear regression): $h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots +$

Define $x_0 = 1$:

$$h(x) = \sum_{j=0}^n \theta_j x_j \quad (n = \text{\#features})$$

$$\text{Define } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$\Rightarrow h(x) = \theta^T x$$

Cost function Consider $h_\theta(x) - y$
↑ prediction
↑ observed true value

Traditional learning problem:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\theta_f = \underset{\theta}{\operatorname{argmin}} J(\theta)$$

Note h maps x to y
 J maps θ to \mathbb{R}

Gradient Descent



Objective: find θ_j 's that minimize $J(\theta)$ over (x, y)

- Idea:
1. choose random starting location ($\theta = \vec{0}$)
 2. find direction of steepest descent
 3. take step of predetermined size α in that direction (change θ to reduce $J(\theta)$)

$$\theta_j := \theta_i - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad j=0, \dots, n$$

↑
learning rate

$$\boxed{\theta_j := \theta - \alpha \nabla J(\theta)}$$

4. if not reached local min, goto 2
5. profit!

Further derivation

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Simplified: $J(\theta) = \frac{1}{2} (h_{\theta}(x) - y)^2$

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

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\Rightarrow \theta_j = \theta_j - \alpha (h_{\theta}(x) - y) x_j$$

Final definition

Repeatedly!

$$\boxed{\theta_j = \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad j=1 \dots n}$$

Batch Gradient Descent - GD wrt entire training data

Stochastic Gradient Descent - much faster

Repeat: $\{$

For $i=1 \dots m \{$

$$\theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad j=1 \dots n$$

$\}$
 $\}$

Minibatch Gradient Descent - look at 64 or 128 examples,
do SGD step per minibatch

Explicitly solving linear regression

$$X = \begin{bmatrix} -x^{(1)\top} \\ -x^{(2)\top} \\ \vdots \\ -x^{(n)\top} \end{bmatrix} \in \mathbb{R}^{n \times n+1} \quad (\text{design matrix})$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

Normal equation

θ^* = optimal value for θ

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

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