

Supervised Learning

$x \rightarrow y$
(input) (output)

Learn Hypothesis: $h(x) \approx y$

Terminology

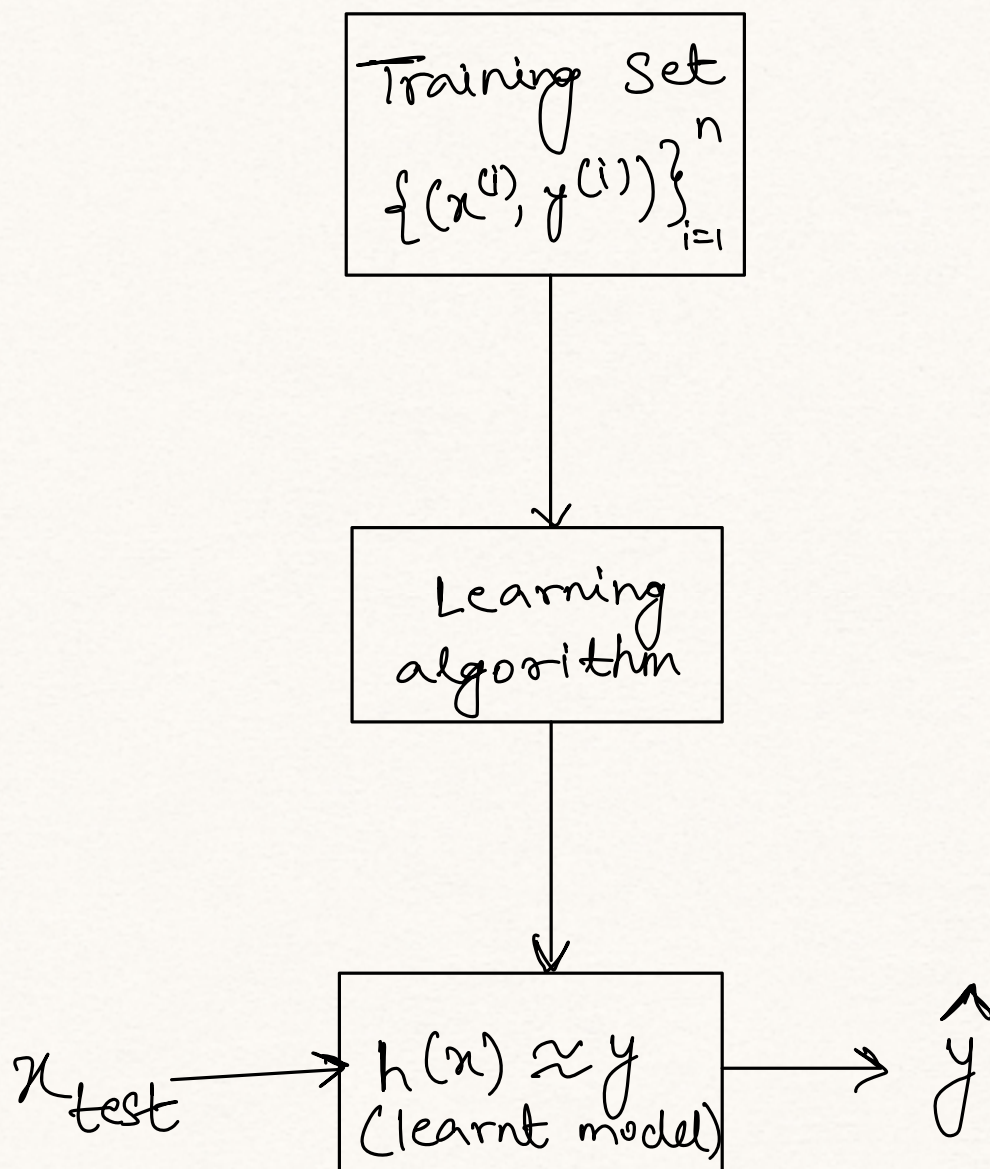
n : number of (x, y) examples in training set.

d : $x \in \mathbb{R}^d \rightarrow$ dimension of input

$x^{(i)}$: i th example input

$y^{(i)}$: i th example output / label / Ground truth

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



Linear Regression

$x \in \mathbb{R}^d$, $y \in \mathbb{R}$, n such examples.

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

$$\theta \in \mathbb{R}^{d+1}$$

$$h_{\theta}(x) = \left(\sum_{i=1}^d \theta_i x_i \right) + \theta_0$$

$\boxed{x_0 = 1}$ — intercept term

$$h_{\theta}(x) = \sum_{i=0}^d \theta_i x_i = \theta^T x$$

Cost Function / Loss function

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

Squared cost error funcⁿ.

$$\hat{\theta} = \arg \min_{\theta} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Gradient Descent

$\theta^{(0)} := \text{Initialization}$

$$\theta_j^{(1)} = \theta_j^{(0)} - \alpha \frac{\partial}{\partial \theta_j} J(\theta^{(0)})$$

α : learning rate

Repeat till convergence:

$$\theta^{(1)} = \theta^{(0)} - \alpha \nabla_{\theta} J(\theta^{(0)}) \quad \text{In vector form}$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} J(\theta^{(t)})$$

- To check convergence we pre-define an $\epsilon < 10^{-5}$ & check either $\|\nabla_{\theta} (J(\theta^{(t)}))\|$ or $\|\theta^{(t)} - \theta^{(t-1)}\|$ or $\|J(\theta^{(t)}) - J(\theta^{(t-1)})\|$ becomes $< \epsilon$.

Gradient descent on Linear Regression

$\theta^{(0)} = \text{initialisation}$

Repeat until convergence:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} J(\theta^{(t)})$$

$$\begin{aligned}
&= \theta^{(t)} - \alpha \nabla_{\theta} \left[\frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right] \\
&= \theta^{(t)} - \alpha \nabla_{\theta} \left[\frac{1}{2} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2 \right] \\
&= \theta^{(t)} - \alpha \left[\frac{1}{2} \sum_{i=1}^n 2 (\theta^T x^{(i)} - y^{(i)}) x^{(i)} \right] \\
&= \theta^{(t)} - \alpha \left[\sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)}) x^{(i)} \right]
\end{aligned}$$

\uparrow vector of dim. $d+1$. \uparrow scalar \uparrow vector of dim. $d+1$

Stochastic Gradient Descent (SGD)

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \cdot \nabla_{\theta} \tilde{J}(\theta)$$

$$\tilde{J}(\theta) = \frac{1}{2} (\theta^T x^{(k)} - y^{(k)})^2$$

k : uniformly at random sampled from training set.

(Random updates)

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

Design matrix

$$X = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(i)} \\ \vdots \\ x^{(n)} \end{bmatrix} \quad \begin{matrix} \xleftarrow{d+1} \\ \xrightarrow{n} \end{matrix}$$

$$\vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

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

$$\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

$$\begin{matrix} X\theta - y \\ \mathbb{R}^{n \times d+1} \cdot \mathbb{R}^{d+1} - \mathbb{R}^n \\ \mathbb{R}^n - \mathbb{R}^n \\ \mathbb{R}^n \end{matrix}$$

$$= \begin{bmatrix} x^{(1)T} \theta - y^{(1)} \\ \vdots \\ x^{(i)T} \theta - y^{(i)} \\ \vdots \\ x^{(n)T} \theta - y^{(n)} \end{bmatrix}$$

$$J(\theta) = \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

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

$$\nabla_{\theta} \frac{1}{2} (x\theta - y)^T (x\theta - y)$$

$$\nabla_{\theta} \frac{1}{2} [(x\theta)^T (x\theta) - (x\theta)^T y - y^T (x\theta) + y^T y]$$

$$\nabla_{\theta} \frac{1}{2} [\theta^T (x^T x) \theta - 2 \theta^T (x^T y) + y^T y]$$

$$= \frac{1}{2} [2(x^T x) \theta - 2x^T y] = 0$$

$$\Rightarrow (x^T x) \theta = x^T y \rightarrow \text{Normal Eq'n}$$

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

Given $x^T x$ is invertible.