

# LINEAR

## REGRESSION

→ In linear regression, we find a set of numbers called parameters from the training data and use it to predict values for testing data.

$$\rightarrow h_{\theta}(x) = \theta^T x = \sum_i^n \theta_i x_i ; \quad x_0 = 1$$

$m \rightarrow$  no. of training examples.

$n \rightarrow$  no. of features.

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

→ Cost function;

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

Has to be minimized.

→ So, which can be done by

Gradient-descent:

Start with some  $\theta$ ;

then learning rate

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

$$\Leftrightarrow \theta_j := \theta_j - \alpha [h_{\theta}(x) - y] x_j$$

For single training ex

→ For multiple training examples  $m$ ;

Repeat {  
     $\theta_j := \theta_j - \alpha \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] x_j^{(i)}$   
} for  $j = 0, 1, 2 \dots n$ .

→ Stochastic Gradient descent

Repeat {

For  $j=1$  to  $m$  {

$\theta_j := \theta_j - \alpha (h_{\theta}(x^{(j)}) - y^{(j)}) x_j^{(j)}$

}

→ Other Algorithms:

Normal Equation

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

→ Why least squares?

let us assume

$$y^{(i)} = \theta^T x^{(i)} + \underbrace{\epsilon^{(i)}}_{\text{Taken to be gaussian}}$$

$$\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$$

"from central limit theorem"

$$P(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

$\Rightarrow$  On writing  $\epsilon^{(i)}$  as  $y^{(i)} - \theta^T x^{(i)}$

$$P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$\text{i.e. } \boxed{y^{(i)} | x^{(i)}; \theta \sim \mathcal{N}(\theta^T x^{(i)}, \sigma^2)}$$

$$\mathcal{L}(\theta) = P(\bar{y} | x; \theta)$$

$$= \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$\ell(\theta) = \log \mathcal{L}(\theta)$$

$$= \sum_{i=1}^m \left[ \log\left(\frac{1}{\sqrt{2\pi}\sigma^2}\right) + \log\left(\exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)\right) \right]$$

$$= m \log \frac{1}{\sqrt{2\pi}\sigma^2} + \sum_{i=1}^m -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}$$

So,  $\mathcal{L}(\theta) \rightarrow \text{Maximised}$

$\Rightarrow \ell(\theta) \rightarrow \text{Maximised}$

$\Rightarrow \underbrace{\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2}_{\text{Cost function}} \rightarrow \text{Minimised}$