

# ① Training Data Set

$$D = \{(x^{(i)}, y^{(i)}), i = 1, 2, \dots, m\}$$

$x$  = Number of features

The cost function for stochastic gradient descent.

$$J_{\text{train}}(\vec{\theta}) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\vec{\theta}, x^{(i)}, y^{(i)})$$

where the cost  $(\vec{\theta}, (x^{(i)}, y^{(i)}))$  is the cost associated with  $i$ -th training example.

$$\begin{aligned} \text{cost}(\vec{\theta}, (x^{(i)}, y^{(i)})) &= \frac{1}{2} (h_{\vec{\theta}}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2} [\vec{\theta}^T x^{(i)} - y^{(i)}]^2 \\ &= \frac{1}{2} [\theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)}]^2 \end{aligned}$$

The update rule for the stochastic gradient descent,

$$i = 1 \text{ to } m$$

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} \text{cost}(\vec{\theta}, (x^{(i)}, y^{(i)}))$$

which turns out

$$\frac{\partial}{\partial \theta_j} \text{cost}(\vec{\theta}, (x^{(i)}, y^{(i)}))$$

$$= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2} [\theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_j x_j^{(i)} + \theta_n x_n^{(i)} - y^{(i)}]^2$$

$$= \frac{1}{2} \times 2 [\theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)}] x_j^{(i)}$$

1 continued

Stochastic Gradient Descent update rule,

for  $i = 1$  to  $m$

$$\theta_j = \theta_j - \alpha (\theta_0 + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)} - y \times x_j^{(i)})$$

for  $i = 1$  to  $m$

$$\theta_j = \theta_j - \alpha (\vec{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_j^{(i)}$$

②

first, we need to be clear about the probabilistic assumptions. You must, when using the results of the analysis, be aware that the results may not be relevant in cases where the assumptions do not hold.

Next, we must use the maximum likelihood estimation formula which is defined as  $x_1, x_2, \dots, x_n$  will be observations, of  $n$  independent and identically distributed random variables, drawn from a probability Distribution  $f_{\theta}(x)$ , where  $f_{\theta}(x)$  is known to be from a distribution  $f_{\theta}$  that depend on some parameters  $\theta$ .

Let's assume that the target variables and the inputs are related using this equation:  $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$ , where  $\epsilon^{(i)}$  is an error term that captures unmodelled effects, or random noise. We can then further assume that  $\epsilon^{(i)}$  are distributed IID according to Gaussian distributions with mean zero and some variance  $\sigma^2$ . We then write this assumption as  $\epsilon^{(i)} \sim N(0, \sigma^2)$

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

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

$P(y^{(i)} | x^{(i)}; \theta)$  indicates that  $y^{(i)}$  given  $x^{(i)}$  is distribution and parameterized by  $\theta$ . You may also write the distribution of  $y^{(i)}$  as  $y^{(i)}; \theta \sim N(\theta^T x^{(i)}, \sigma^2)$ .

Given  $X$  and  $\theta$ , the probability of the data given by  $P(\vec{y} | X; \theta)$ . This quantity is typically viewed a function of  $\vec{y}$  for a fixed value of  $\theta$ . When you wish to view this as a function of  $\theta$ , we will call it a likelihood function.



2 continue

$$L(\theta) = L(\theta; X, \hat{Y}) = P(\hat{Y} | X; \theta)$$

The independent assumption on the  $x^{(i)}$  can be written as

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

The principle of maximum likelihood says that we should choose  $\theta$  to make the data as high probability as possible. Instead of maximizing  $L(\theta)$ , we can also maximize any strictly increasing function of  $L(\theta)$ . In particular, the derivations will be simpler if we maximize the log likelihood  $l(\theta)$

$$\begin{aligned} l(\theta) &= \log L(\theta) = \log \left( \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \right) \\ &= m \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{\sigma^2} \times \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \end{aligned}$$

maximizing  $l(\theta)$  gives the same answer as minimizing  $\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$

Maximizing Likelihood Estimation isn't necessarily optimal, due to other estimation algorithms that achieve better results, it does have its own attractive properties. The most important being consistency. A sequence of Maximizing Likelihood Estimation on an increasing number of operations will converge to the true value of the parameters.

Using the above two definitions we can say  $\hat{\theta}$  is a reasonable choice when using  $\theta$ .

③

$$\log h_0(x^i) = \log \frac{1}{1 + e^{-ox^i}} = -\log(1 + e^{-ox^i})$$

$$\begin{aligned} \log(1 - h_0(x^i)) &= \log\left(1 - \frac{1}{1 + e^{-ox^i}}\right) \\ &= \log(e^{-ox^i}) - \log(1 + e^{-ox^i}) \\ &= -ox^i - \log(1 + e^{-ox^i}) \end{aligned}$$

The original cost function is in the form of

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^m y^i \log(h_0(x^i)) + (1 - y^i) \log(1 - h_0(x^i))$$

We then plug in the two simplified expressions

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^m [-y^i (\log(1 + e^{-ox^i})) + (1 - y^i) (-ox^i - \log(1 + e^{-ox^i}))]$$

$$= -\frac{1}{n} \sum_{i=1}^m [y^i ox^i - ox^i - \log(1 + e^{-ox^i})]$$

$$= -\frac{1}{n} \sum_{i=1}^m [y^i ox^i - \log(1 + e^{-ox^i})]$$

The second equality

$$= -ox^i - \log(1 + e^{-ox^i})$$

$$= -[\log e^{-ox^i} + \log(1 + e^{-ox^i})]$$

$$= -\log(1 + e^{-ox^i})$$

And to compute the partial derivatives

$$\frac{\partial}{\partial \theta_j} y^i ox^i = y^i x_j^i$$

$$\frac{\partial}{\partial \theta_j} \log(1 + e^{-ox^i}) = \frac{x_j^i e^{-ox^i}}{1 + e^{-ox^i}} = x_j^i h_0(x^i)$$