

Page \_\_\_\_\_

$i^{\text{th}}$  training example (instance)

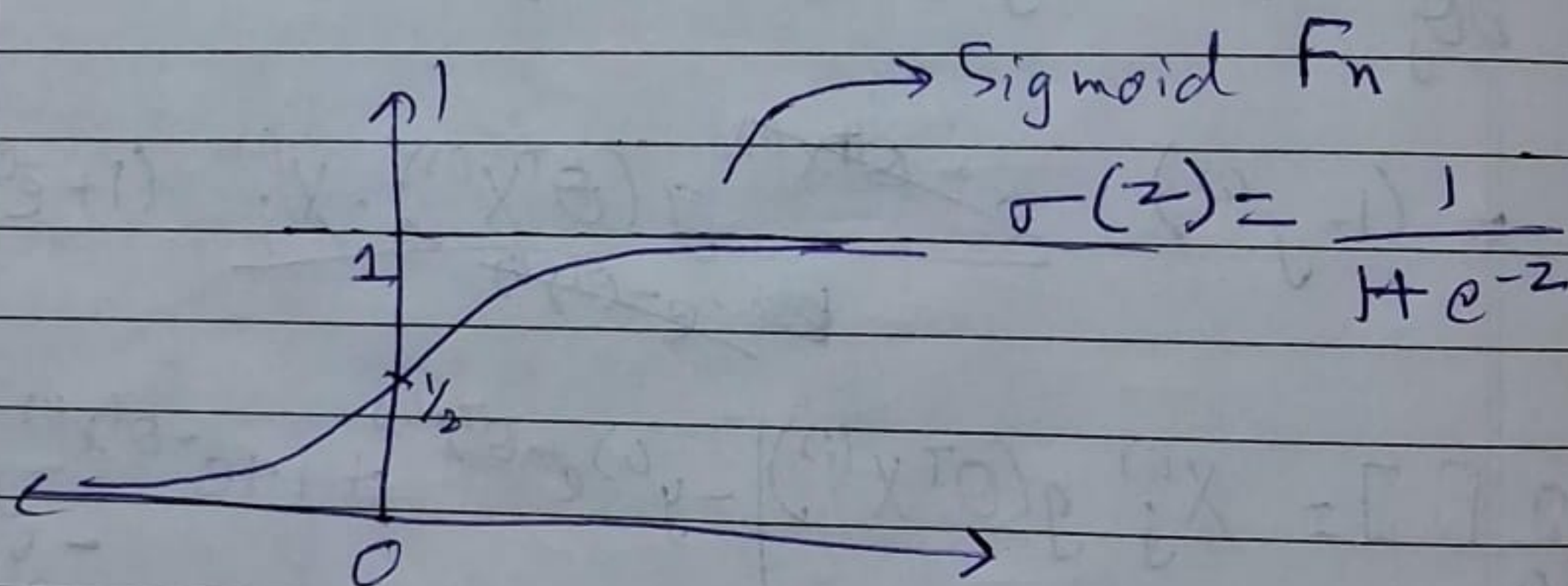
$$X^{(i)} = \begin{bmatrix} 20 \\ 45 \\ \vdots \end{bmatrix}$$

$X_j^{(i)}$

$m$  instances having  $n$  features

$$\Theta^T X^{(i)} = z^{(i)} \rightarrow \text{"Logits"}$$

$(n+1)$  parameters       $(n+1)$  features



$$\underbrace{h_{\Theta}^{(i)}(X^{(i)})}_{\text{estimated Probability}} = \sigma(z^{(i)})$$



Assume

$$P(y=1|x; \theta) = h_{\theta}(x).$$

$$P(y=0|x; \theta) = 1 - h_{\theta}(x).$$

$$\because y \in \{0, 1\}$$

$$P(y|x; \theta) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{1-y}.$$

(Likelihood)

$$L(\theta) = P(\vec{y} | X; \theta)$$

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

$$= \prod_{i=1}^m \left[ h_{\theta}(x^{(i)})^{y^{(i)}} \cdot (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}} \right].$$



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

$$= \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))$$

Goal: maximize  $l(\theta)$

Define

$$J(\theta) = -\frac{1}{m} l(\theta) \rightarrow \text{Negative Avg log-Likelihood}$$

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

$\rightarrow$  Cross-Entropy Loss

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

$$\frac{\partial J}{\partial l} = -\frac{1}{m}$$

$\rightarrow$  because  $h_{\theta}(x^{(i)}) = \sigma(z^{(i)})$

$$\frac{\partial l}{\partial \Theta_j} = \frac{\partial l}{\partial h} \cdot \frac{\partial h}{\partial z} \cdot \frac{\partial z}{\partial \Theta_j}$$

$$= \left[ \frac{y}{h} - \frac{(1-y)}{1-h} \right] \left[ \sigma'(z) (1-\sigma(z)) \right] X_j$$

$$= \left[ \frac{y}{h} - \frac{(1-y)}{(1-h)} \right] [h(1-h)] X_j^{(i)}$$

$$= \frac{(y - yh - h + yh)}{h(1-h)} \times h(1-h) \times X_j^{(i)} = (y-h) \cdot X_j^{(i)}$$



$$\Rightarrow \frac{\partial J}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

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

~~$$\frac{1}{m} X^T (y - h_{\theta}(x))$$~~

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta)$$

Learning  
Rate