

Weighted least Square, Logistics Regression Newton's method

⇒ ~~Linear Regression~~, you

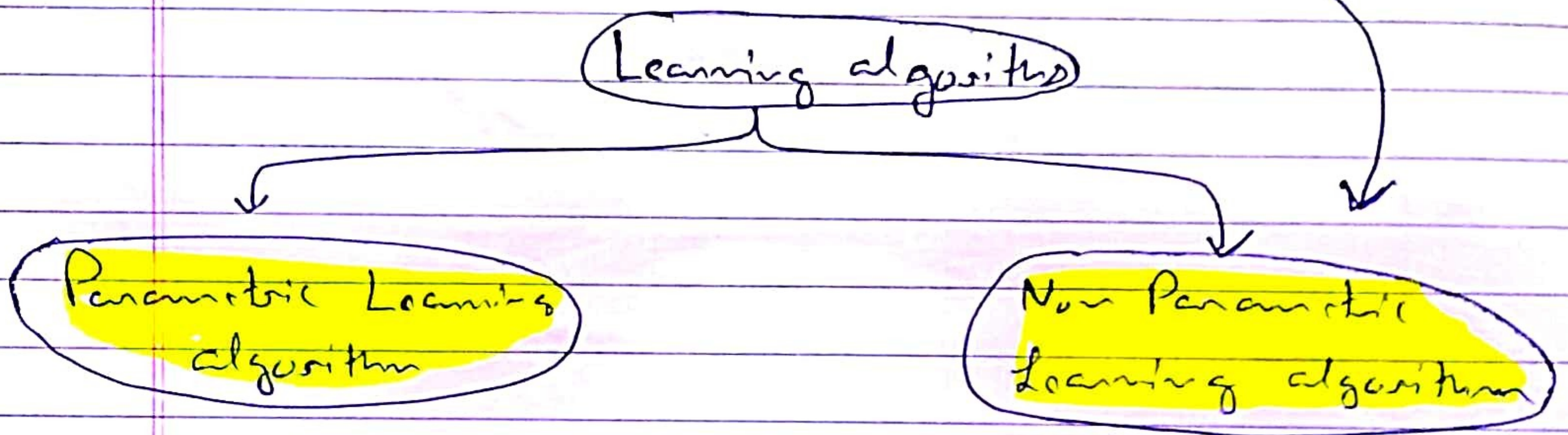
⇒ Using Linear Regression you can also train a non linear model by using a simple trick.

$$Y = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 \quad \left\{ \text{Non linear model in } x_1 \right\}$$

Let new feature $x_2 = x_1^2$

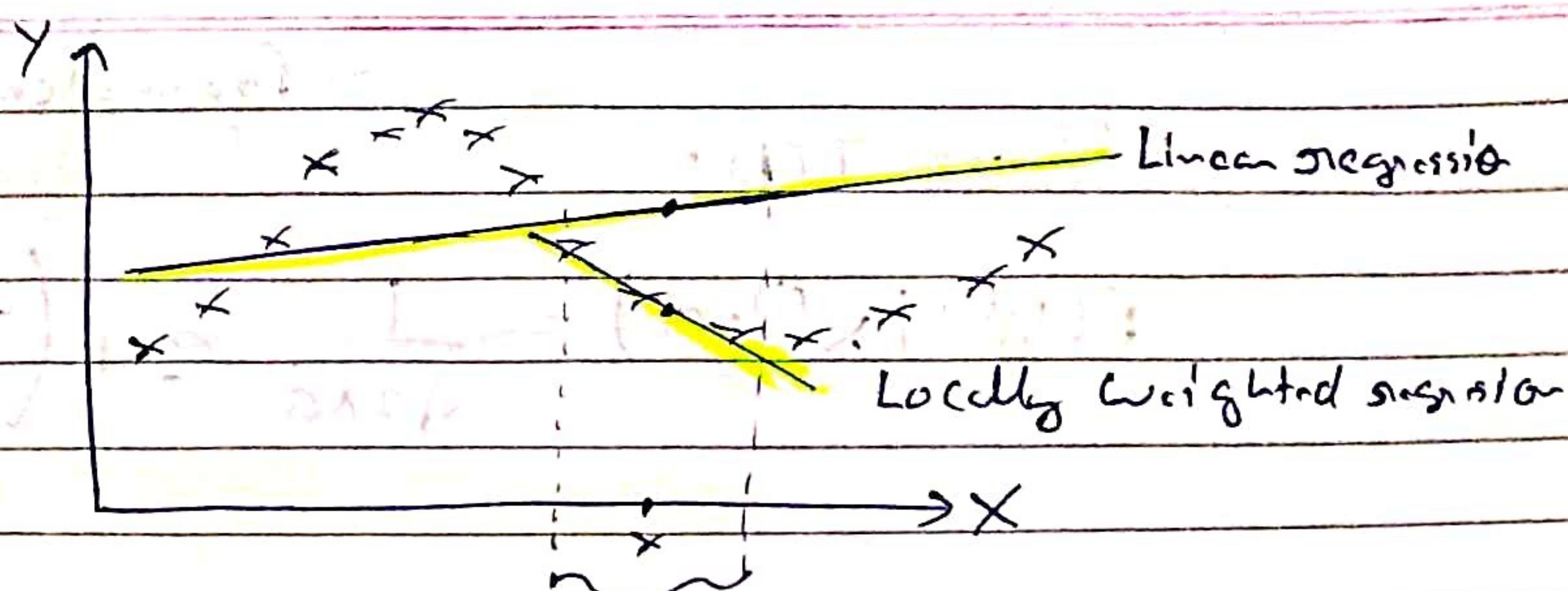
$$\text{So } Y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \quad \left\{ \begin{array}{l} \text{Now the model} \\ \text{is linear in two} \\ \text{features } x_1 \text{ \& } x_2 \end{array} \right\}$$

* Locally weighted regression



$\left\{ \begin{array}{l} \text{Fit a fixed set of parameters} \\ (\theta_i) \text{ to data} \end{array} \right\}$

$\left\{ \begin{array}{l} \text{Amount of data (Parameters)} \\ \text{you need to keep grows} \\ \text{with size of data} \end{array} \right\}$



Fit Θ to minimize

$$J = \sum_{i=1}^m \omega^{(i)} (y^{(i)} - \Theta^T x^{(i)})^2$$

→ weight function

$$\omega^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right) \quad \left\{ \begin{array}{l} \text{Common choice} \\ \text{of weighting} \\ \text{function} \end{array} \right\}$$

⇒ Generally used when you have a relatively low dimensional dataset and we have a lot of data.

★ Probabilistic interpretation

$$y^{(i)} = \Theta^T x^{(i)} + \epsilon^{(i)}$$

→ {Unmodeled effect
random noise etc}

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

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

$\epsilon^{(i)}$ are IID. → Parameterized by θ

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

$$y^{(i)} | x^{(i)}; \theta \sim N(\theta^T x^{(i)}, \sigma^2)$$

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

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

⇒ With MAP estimation:

~~$$\theta = \arg \max_{\theta} \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$~~

$$\theta = \min_{\theta} \underbrace{\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2}_{J(\theta)}$$

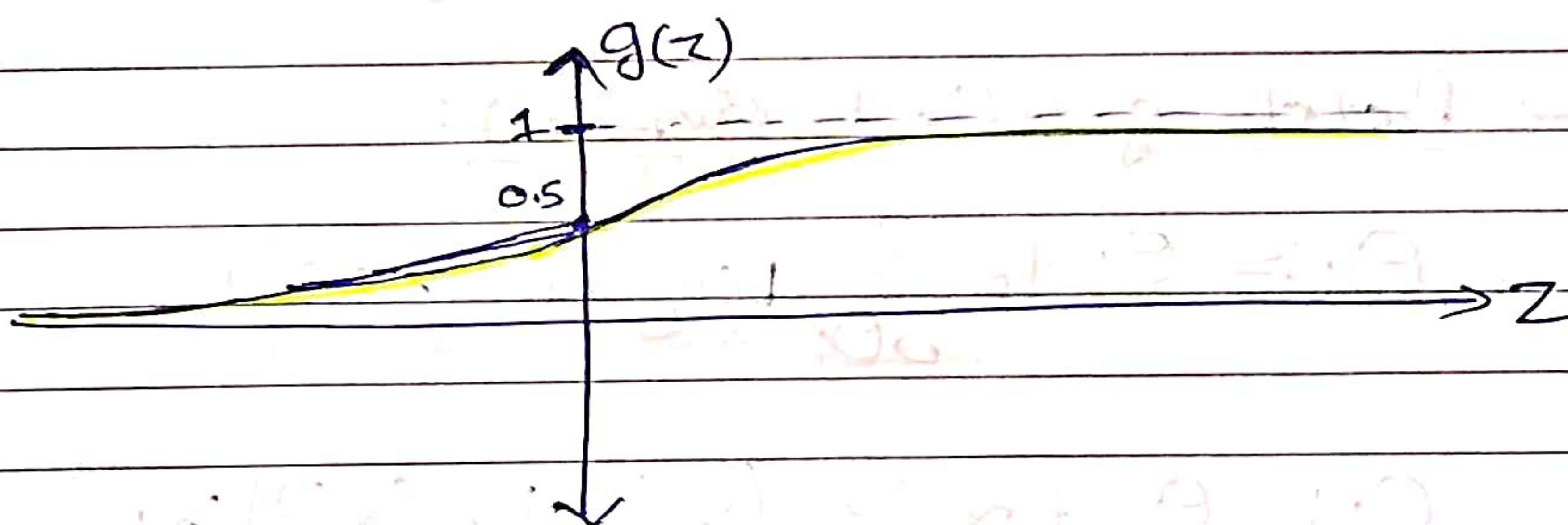
★ Classification

① Binary Classification $\{y \in \{0, 1\}\}$

Logistic Regression

$$g(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid or Logistic function



Let
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$
 Hypothesis function

$$P(Y=1 | X; \theta) = h_{\theta}(x)$$

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

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

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

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

$$\log P(Y|X;\theta) = \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)}))$$

⇒ Choose θ to maximize $\log P(Y|X;\theta)$.

⇒ Batch gradient ^{ascent} ~~descent~~:

$$\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} \log P(Y|X;\theta)$$

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

⇒ The above function has only one optimum, no local optimum.

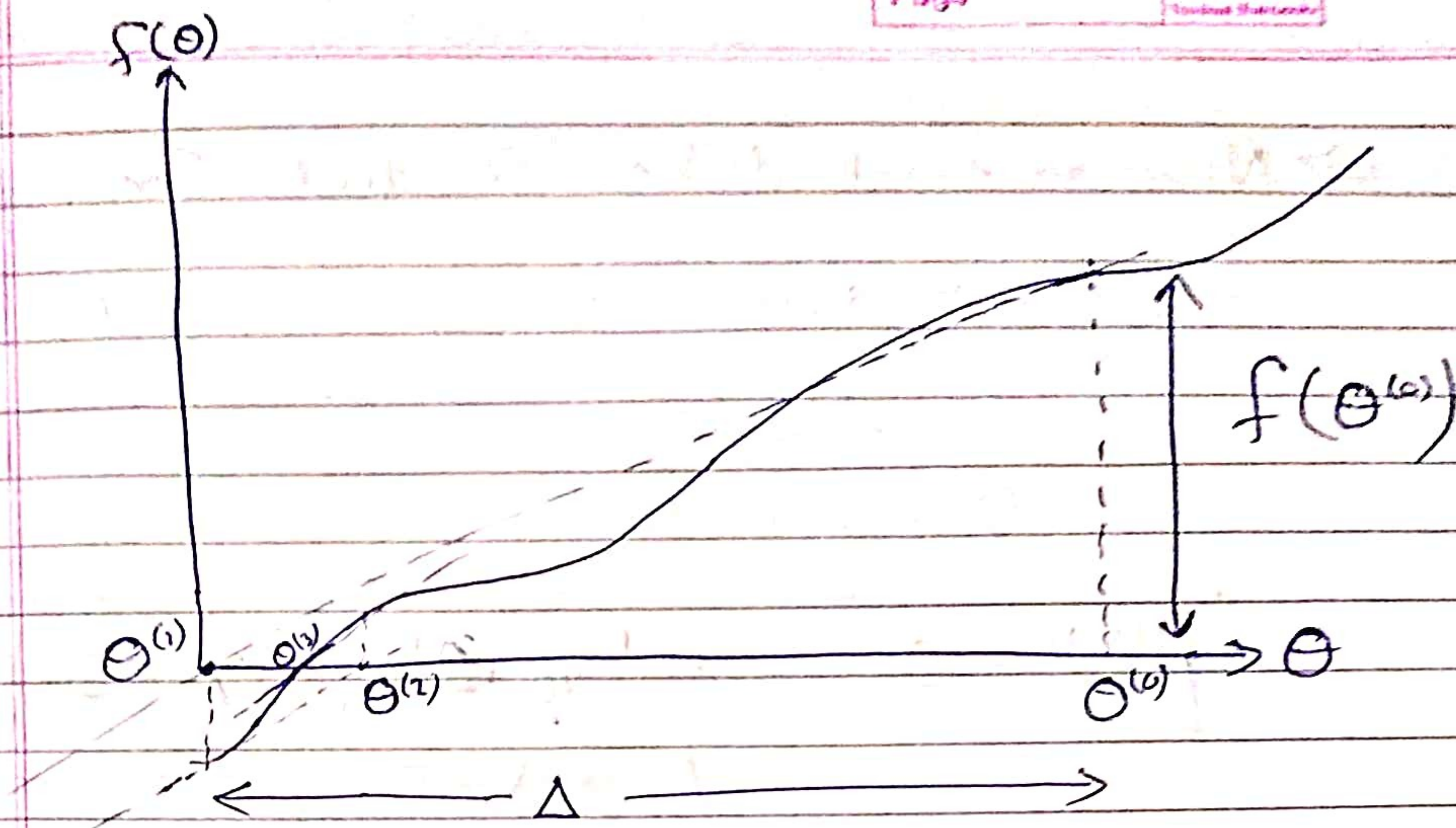
* Newton's Method

"Given f , find θ , st $f(\theta) = 0$ "

Approach: (1) Start with a guess $\theta^{(0)}$;

$$(2) \theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

(3) Iterate



$$\theta^{(1)} = \theta^{(0)} - \Delta$$

$$f'(\theta^{(0)}) = \frac{f(\theta^{(0)})}{\Delta}$$

$$\theta^{(1)} = \theta^{(0)} - \frac{f(\theta^{(0)})}{f'(\theta^{(0)})}$$

Am general,

$$\theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

$$\Rightarrow \text{Let } f(\theta) = \log(P(y|x; \theta))$$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\frac{\partial}{\partial \theta} \log(P(y|x; \theta))}{\frac{\partial^2}{\partial \theta^2} \log(P(y|x; \theta))}$$

⇒ Newton method has "Quadratic Convergence".

0.01 → 0.0001 → 0.00000001

Error

⇒ When θ is a vector:

$$\theta^{(t+1)} = \theta^{(t)} + H^{-1} \nabla_{\theta} \log(P(y|x; \theta))$$

\nwarrow Hessian matrix \nearrow Gradient

$$H_{ij} = \frac{\partial^2 f}{\partial \theta_i \partial \theta_j}$$

————— X ————— X —————