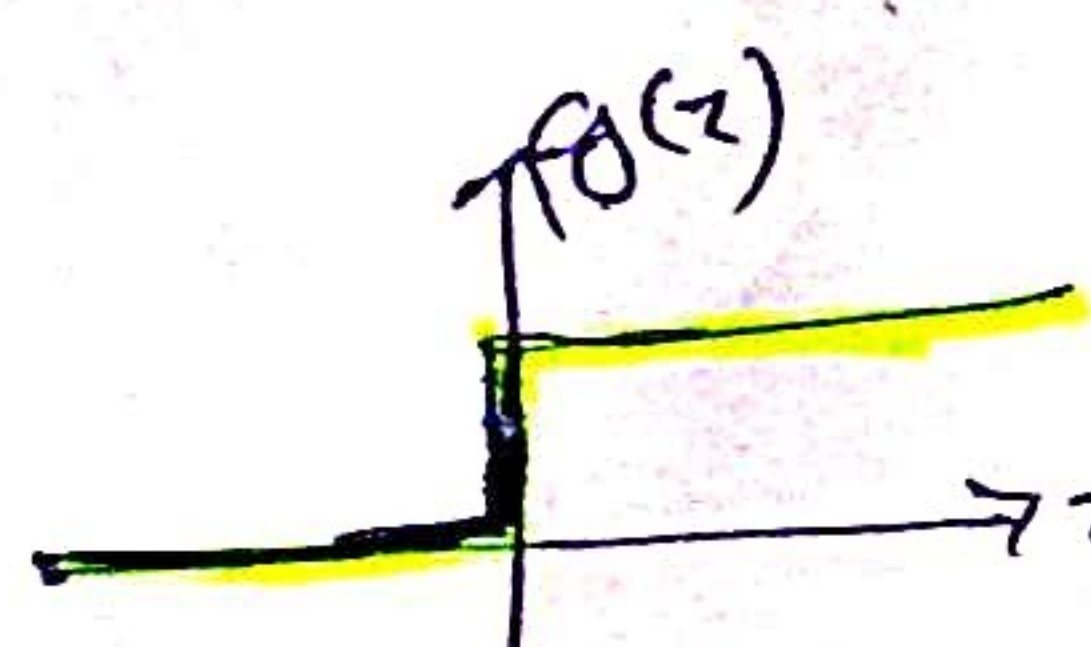


5. Kernelizing the Perceptron

⇒ Let there be a binary classification problem with $y \in \{0, 1\}$

⇒ Perceptron uses hypotheses of the form:

$$h_{\theta}(x) = g(\theta^T x)$$

 Where $g(z) = \text{Sign}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$

⇒ Here we consider Stochastic gradient descent-like implementation of the perceptron algorithm.

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

where $\theta^{(i)}$ is value of the parameters after the algorithm has seen i training examples.

$$\theta^0 = 0$$

Ⓐ Let k be a Mercer Kernel corresponding to some feature mapping ϕ .

(i) high-dimensional parameter vector can be represented as linear combination of input features $\phi(x)$

$$\theta^{(i)} = \sum_{j=1}^m \alpha_j^{(i)} \phi(x) \quad \left\{ \begin{array}{l} \text{This is a valid assumption} \\ \text{due to representation theorem} \end{array} \right\}$$

⇒ Here we only need to store $\alpha_1^{(i)}, \alpha_2^{(i)} \dots \alpha_m^{(i)}$ m parameter to represent probable infinite dimensional $\theta^{(i)}$.

(ii)

$$\begin{aligned}
 h_{\theta^{(i)}}(x^{(i+1)}) &= g\left(\theta^{(i)T} \phi(x^{(i+1)})\right) \\
 &= g\left(\left(\sum_{j=1}^m \alpha_j^{(i)} \phi(x^{(j)})\right)^T \phi(x^{(i+1)})\right) \\
 &= g\left(\sum_{j=1}^m \alpha_j^{(i)} \phi(x^{(j)})^T \phi(x^{(i+1)})\right)
 \end{aligned}$$

\Rightarrow Let \tilde{K} be the kernel matrix corresponding to kernel K and data set $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$

$\tilde{K} \in \mathbb{R}^{m \times m}$

$$h_{\theta^{(i)}}(x^{(i+1)}) = g\left(\tilde{K}(:, i+1)^T \alpha^{(i)}\right)$$

$\left\{ \begin{array}{l} \text{ } i+1\text{st Column of} \\ \text{Kernel matrix} \end{array} \right\}$

(iii)

$$\begin{aligned}
 \theta^{(i+1)} &:= \theta^{(i)} + \alpha (y^{(i+1)} - h_{\theta^{(i)}}(x^{(i+1)})) \alpha^{(i+1)} \\
 &\quad \downarrow \qquad \qquad \downarrow \\
 \sum_{j=1}^m \alpha_j^{(i+1)} \phi(x^{(j)}) &\quad \sum_{j=1}^m \alpha_j^{(i)} \phi(x^{(j)})
 \end{aligned}$$

$$\Rightarrow \text{Let } X = \begin{bmatrix} \phi(x^{(1)T}) \\ \phi(x^{(2)T}) \\ \vdots \\ \phi(x^{(m)T}) \end{bmatrix} \text{ be the design matrix.}$$

⇒ Multiplying both side of the equation with design matrix.

$$\sum_{j=1}^m \alpha_j^{(i+1)} \times \phi(x^{(j)}) = \sum_{j=1}^m \alpha_j^{(i)} \times \phi(x^{(j)})$$

$$+ \alpha \left(y^{(i+1)} - g(\tilde{K}(:, i+1)^T \alpha^{(i)}) \right) \times \phi(x^{(i+1)})$$

$$\tilde{K}(:, j)$$

$$\tilde{K} \alpha^{(i+1)}$$

$$\Rightarrow \tilde{K} \alpha^{(i+1)} = \tilde{K} \alpha^{(i)} + \alpha \left(y^{(i+1)} - g(\tilde{K}(:, i+1)^T \alpha^{(i)}) \right) \tilde{K}(:, i+1)$$

$$\alpha^{(i+1)} = \alpha^{(i)} + \alpha \left(y^{(i+1)} - g(\tilde{K}(:, i+1)^T \alpha^{(i)}) \right) \tilde{K}^{-1} \tilde{K}(:, i+1)$$