**Q7** In conventional ~~ma~~ ML models involving linear regression, the ~~&~~ "hypothesis" ~~is~~ is a linear/affine function, ~~tha~~i.e of the form.

$$z = W^T X + b.$$

An activation layer introduces **non-linearity** in the model, e.g.

$$a = \sigma(z) = \sigma(W^T X + b)$$

$$\text{where } \sigma(x) = \frac{1}{1+e^{-x}}.$$

## Mathematical Requirements

**(1) Non-linearity**   $a(\lambda z_1 + \mu z_2) \neq \lambda a(z_1) + \mu a(z_2)$

**(2) Differentiability**

Although ReLU is not diff'ble, it is extensively used. ~~bc~~ At $0$, the issue is handled by employing **subgradient**.

**(3) Boundedness**   $\exists M > 0 \text{ s.t. } |f_i(x)| \leq M \quad \forall x \in \mathbb{R}.$
e.g $\sigma$, tanh etc.

However, there are activation fns that are unbdd. e.g ReLU.

Some of the

### Types

(1) **Sigmoid**, $\sigma(x) = \dfrac{1}{1+e^{-x}}$.

→ Used for $\underset{\wedge}{\text{binary}}$ classification problems.

→ $\sigma \in (0,1) \Rightarrow$ ~~can be returns~~

→ $\sigma(x)$ can be interpreted as probability.

(2) **ReLU**

$$\text{ReLU}(x) = \max(0, x) = \begin{cases} x, & x > 0 \\ 0, & x < 0 \end{cases}$$

→ Used to avoid vanishing gradients
→ Used in CNNs almost always $\left[\text{other alternatives} \rightarrow \begin{array}{l} \text{GELU} \\ \text{Leaky ReLU} \end{array}\right]$

(3) **Softmax**

$$\left[\text{softmax}(z)\right]_i = \frac{e^{z_i}}{\sum\limits_{j=1}^{} e^{z_j}} \qquad \hookrightarrow \text{no. of classes}$$

$$z = \begin{bmatrix} z_1 \\ z_2 \\ | \\ z_c \end{bmatrix}$$

→ Used for multi-class classification problems

(4) Leaky ReLU

$$f(x) = \max\{\alpha x, x\} \quad , \text{ where } \alpha > 0 \text{ is}$$
$$\text{fixed. } \in$$

→ Used in CNNs

→ Avoids dying ReLU issue, i.e. $f'(x) = \alpha \neq 0$

$$\left[ \text{as opposed to } \frac{d}{dx} ReLU(x) = 0 \, \forall \, x < 0 \right]$$