


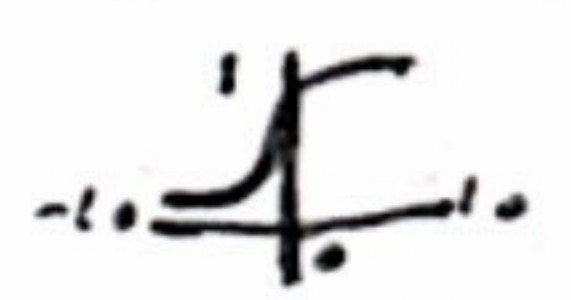
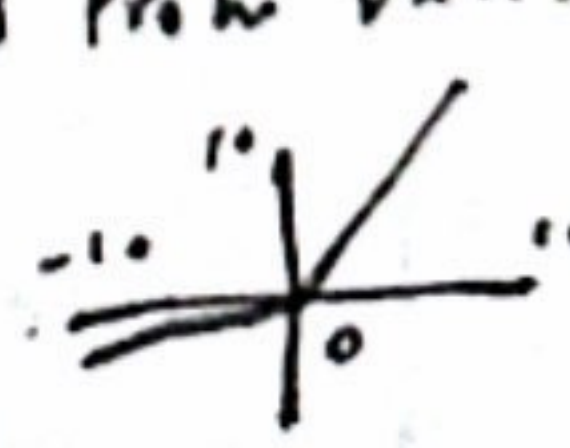
Activation Functions (and gradient vanishing)

- as seen in the NN notes, an activation function in a NN determines whether a neuron should be activated (ie pass info forward) it decides how much signal a neuron passes forward & has without it a NN would become an giant linear equation and become useless. Activation function introduce non-linearity to the network allowing it to learn.

Ex: after a layer computes $z = wx + b$ the activation function transforms it into something non-linear like $\text{Relu}(z)$ so the next layer can learn these patterns (in short if a result is over some threshold the neuron fires and passes information in form of a signal it on/off and how much is it on by, to the next layer)

Numeric Ex: input: $x = [0.5, 0.3]$ weights: $w = [0.4, 0.6]$ bias: $b = 0.1$
Weighted sum of Neuron: $z = wx + b = 0.5 \times 0.4 + 0.3 \times 0.6 + 0.1 = 0.32$
apply activation func of choice $\text{ReLU}(0.32) = \max(0, 0.32) = 0.32$

Types of activation functions (common ones)

- **ReLU**: $\max(0, x)$ Simple fast, and works well in practice, but can kill neurons when inputs go negative for too long.

- **Sigmoid**: squashes values to $[0, 1]$ good for probabilities, bad for deep layers as it makes gradients vanish meaning \Rightarrow During backprop (method used in NN to learn from mistakes by calculating how much each weight contributed to the error) then adjusting weights in direction that reduces error more from output layer to input. The gradient you need to update early layers shrinks so much it basically becomes 0. When that happens layers stop learning they freeze while later layers do the training. Because some functions squash values into tiny ranges where the func's derivatives are close to zero or multiply enough near zero derivatives across many layers and signal dies.

- **Tanh**: like sigmoid but centered at 0; $\sigma(0,0)$, still suffers from vanishing grad.
- **Leaky ReLU**: $\max(0.5x, x)$ fixes dead neurons by allowing a small slope for negatives (a band-aid for ReLU's dying neuron prob).

- **Gelu**: used in transformers, smooth, non linear, better gradient flow than ReLU works better since it reflects the idea that small neg values should not be treated as strictly dead.
