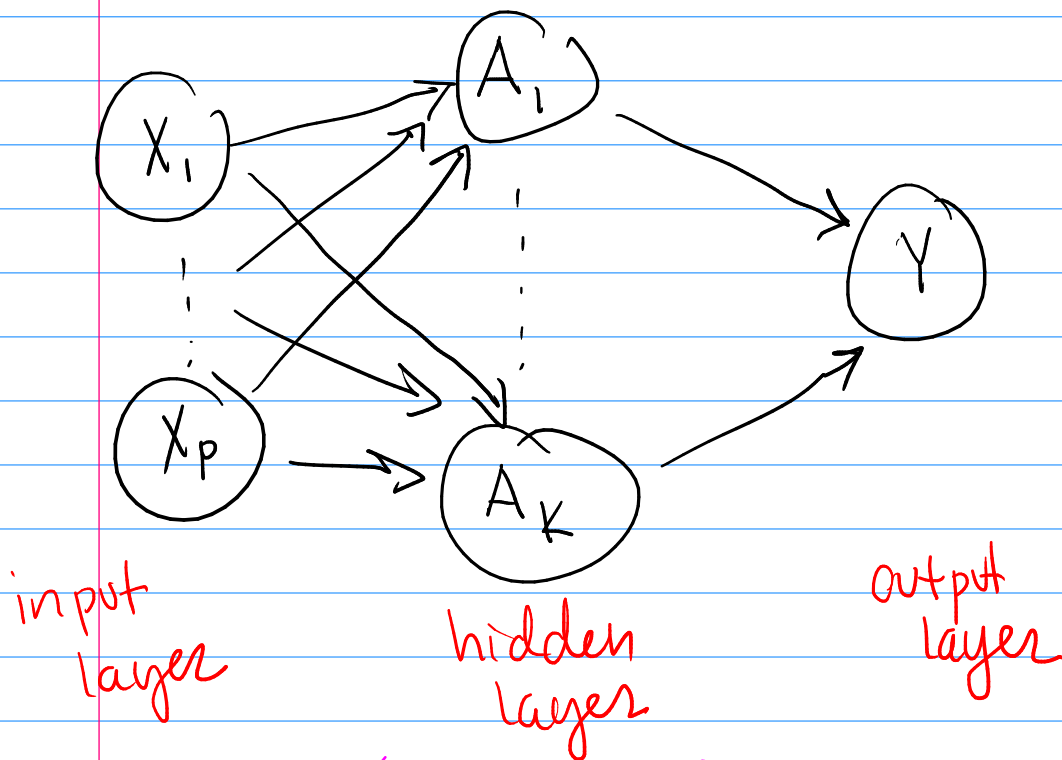# Neural Networks

Basically: build a prediction function $\hat{f}$
as a composition of functions

$$\hat{f} = \hat{f}_K(\hat{f}_{K-1}(\hat{f}_{K-2}(\hat{f}_{K-3}(\cdots\hat{f}_1(x)))) \cdots)$$

## Single hidden layer (feed-forward network)



input
layer

hidden
layer

output
layer

① $A_k = \sigma(W_k^T X + b_k)$ ⟵ $b_k \in \mathbb{R}$

$(X_1, \cdots, X_p)$

$W_k \in \mathbb{R}^p$

"activation function", non-linear

$\textcircled{2}$ $Y = u^T A + c$

$c \in \mathbb{R}$

$u \in \mathbb{R}^k$
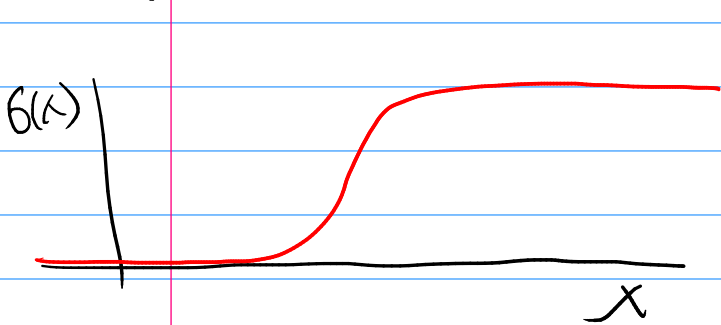
The activation function $\sigma$ can be of many different forms,

popular choices:

$-\ \sigma(x) = \dfrac{1}{1 + e^{-x}}$ (sigmoid)

$-\ \sigma(x) = \max(0, x)$ (ReLU)

sigmoid

relu

$\sigma(x)$

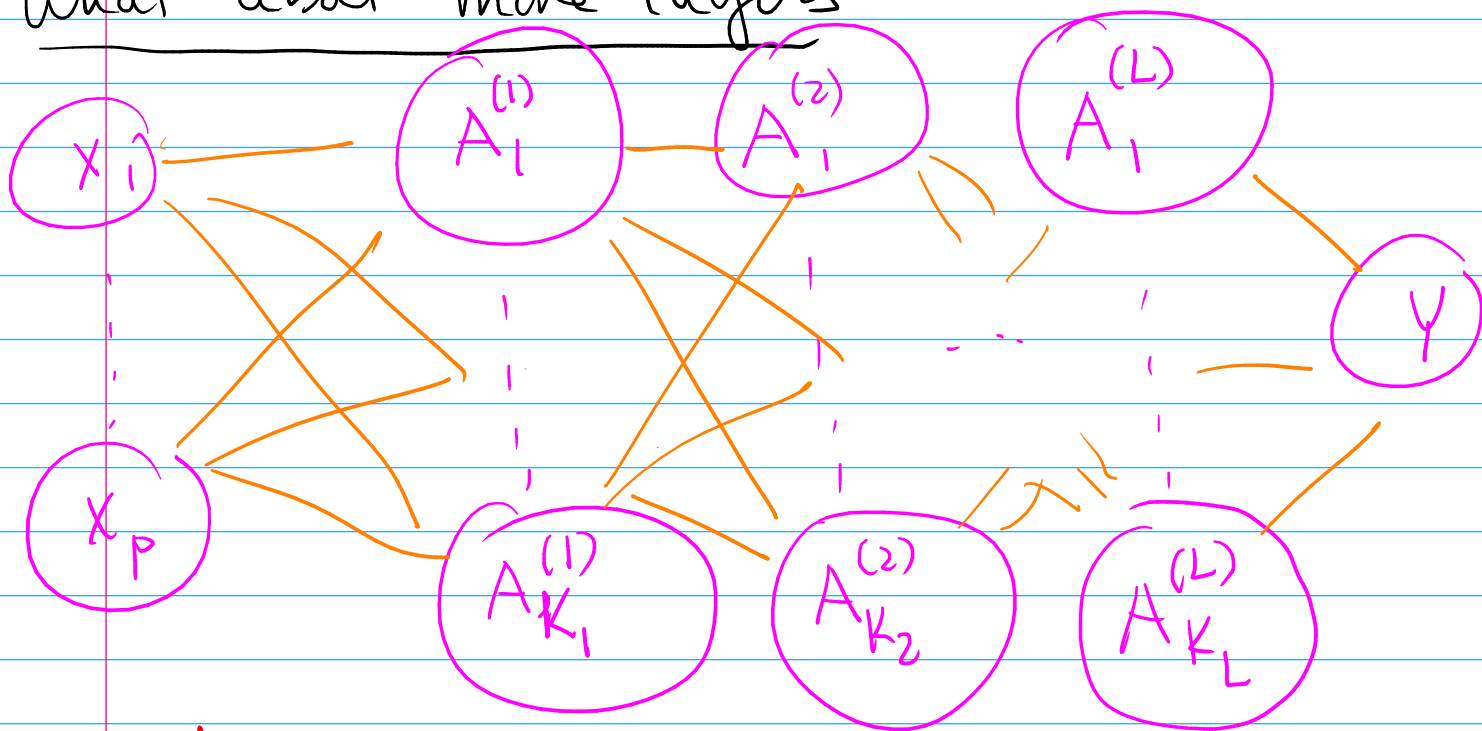$\sigma(x)$

$x$

$0$

applying element-wise

$\hat{f}(x) = u^T \sigma(W^T x + b) + c$

mtx of weights

# What abat more layers



**First hidden**

$$A_k^{(1)} = \sigma\left(W_k^{(1)T} X + b_k^{(1)}\right)$$

**or**

$$A^{(1)} = \sigma_l\left(W^{(1)T} X + b^{(1)}\right)$$

$$\left(A_1^{(1)}, \ldots, A_{K_1}^{(1)}\right)$$

$$mtx \quad W^{(1)} = \begin{bmatrix} - W_1^{(1)} - \\ - W_{K_1}^{(1)} - \end{bmatrix}$$

$$b = \begin{bmatrix} b_1^{(1)} \\ \vdots \\ b_{K_1}^{(1)} \end{bmatrix}$$

**Subsequent hidden layers**

$$A^{(\ell)} = \sigma_\ell\left(W^{(\ell)T} A^{(\ell-1)} + b^{(\ell)}\right)$$

**Output layer**

$$Y = u^T A^{(L)} + c$$

$$\hat{f}(x) = \ldots \ldots \sigma_3\left(W^{(3)T} \sigma_2\left(W^{(2)T} \sigma_1\left(\boxed{W^{(1)T} x + b^{(1)}}\right) + b^{(2)}\right) + b^{(3)}\right)$$

To fit, I need to specify all my params:

$$b^{(1)}, \ldots, b^{(L)}, W^{(1)}, \ldots, W^{(L)}$$

## For classification!

All the same, except last layer typically

Y has M values (for each of M classes)



← the prob. of class m

$$Y = \text{SoftMax}\left(U^T A^{(k)} + C\right)$$ ← btwn $[0,1]$ sum to 1

$$\text{Softmax}(z)_i = e^{z_i} / \sum_i e^{z_i}$$

How to learn params?

Let $\Theta$ = vec. of all params

then we "learn" $\Theta$ as $\hat{\Theta}$ where

$$\hat{\Theta} = \arg\min_{\Theta} \sum_n L(y_n, f_\Theta(x_n))$$

$f_\Theta = NN$

Problem: - super-high-dim'le opt. problem
- likely to over-fit

Soln:

① learn "slowly" via gradient descent

② regularize

Gradient descent:

⓪ initial guess $\Theta^{(0)}$

① For $t = 1, \dots,$

$$\Theta^{(t)} = \Theta^{(t-1)} - \alpha \nabla_\Theta L \big|_{\Theta^{(t-1)}}$$

② Stop when $L$ isn't decreasing w/ further steps

Calculate $\nabla_\theta L$ via back-propagation

<span style="color:red">(applying Chain rule from calc I)</span>

$$\frac{\partial}{\partial x} f(g(x)) = f'(g(x)) g'(x)$$

GD can be slow!

Speed up! Stochastic Grad. Descent.

$$\nabla_\theta \sum_n L(y_n, f_\theta(x_n)) = \sum_n \nabla_\theta L(y_n, f_\theta(x_n))$$

<span style="color:magenta">don't sum over all training data – use a subset.</span>

Also regularize!

(1) penalize loss e.g. $L + \|\theta\|_2^2$

$$L + \|\theta\|_1$$

(2) early stopping – stop SGD when val. perf. stop getting better

(3) drop out: randomly set some weights to zero during training