

10 Neural Networks + Deep Learning

$$Y = F(X) + \epsilon$$

$F(X)$ = multiple regression, tree,
bagged trees, random forest,
boosted trees ...

Neural networks specify $F(X)$ as a composition of functions.

10.1 Single layer neural networks

A feed-forward single layer neural network uses p variables in \underline{x} as input layer + feeds these into K hidden units (we pick K).

$$F(\underline{x}) = \beta_0 + \sum_{k=1}^K \beta_k g\left(\omega_{k0} + \sum_{j=1}^p \omega_{kj} x_j\right)$$

g is called an activation function to produce K activations:

$$A_k = g\left(w_{k0} + \sum_{j=1}^p w_{kj} x_j\right)$$

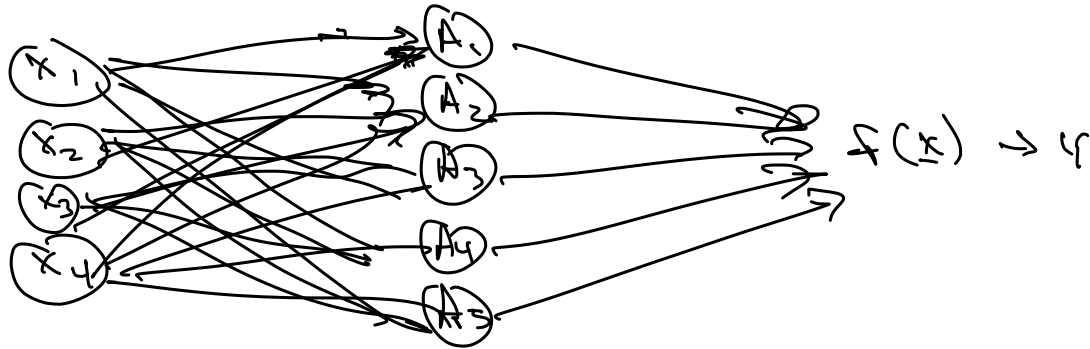
note these activations are a hidden layer in the regular regression

$$F(x) = \beta_0 + \sum_{k=1}^K p_k A_k$$

Input layer

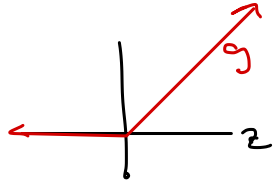
hidden layer

output layer



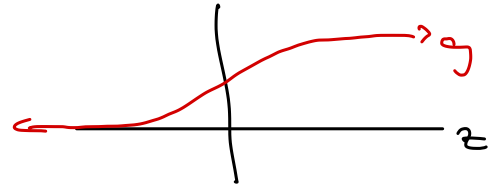
Activation function $g(\cdot)$ is usually (these days)
rectified linear unit (ReLU)

$$g(z) = (z)_+ = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}$$



Historically the sigmoid was used

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



Important for g to be nonlinear, otherwise $g(z) = z$

$$\begin{aligned} f(\underline{x}) &= \beta_0 + \sum_{k=1}^K \beta_k g\left(w_{k0} + \sum_{j=1}^p w_{kj} x_j\right) \\ &= \beta_0 + \sum_{k=1}^K \beta_k \left(w_{k0} + \sum_{j=1}^p w_{kj} x_j\right) \end{aligned}$$

$$\begin{aligned}
 &= \underbrace{\left(\beta_0 + \sum_{k=1}^K \beta_k \omega_{k0} \right)}_{\beta_0} + \sum_{j=1}^P \underbrace{\left(\sum_{k=1}^K \beta_k \omega_{kj} \right)}_{\beta_j} x_j \\
 &= \beta_0 + \sum_{j=1}^P \beta_j x_j
 \end{aligned}$$

= regular multiple regression.

Depending on g we can get interactive effects from \underline{x} :

$$p=2 \quad K=2 \quad g(z) = z^2$$

$$\beta_0 = 0 \quad \beta_1 = \frac{1}{4} \quad \beta_2 = -\frac{1}{4}$$

$$\omega_{10} = 0 \quad \omega_{11} = 1 \quad \omega_{12} = 1$$

$$\omega_{20} = 0 \quad \omega_{21} = 1 \quad \omega_{22} = -1$$

$$f(\underline{x}) = \beta_0 + \sum_{k=1}^2 \beta_k g\left(w_{k0} + \sum_{j=1}^2 w_{kj} x_j\right)$$

$$= \beta_0 + \beta_1 \left(w_{10} + w_{11} x_1 + w_{12} x_2\right)^2 \\ + \beta_2 \left(w_{20} + w_{21} x_1 + w_{22} x_2\right)^2$$

$$= \frac{1}{4} (x_1 + x_2)^2 - \frac{1}{4} (x_1 - x_2)^2$$

$$= x_1 x_2.$$

To estimate β_0, \dots, w_{kp} we could use

$$\sum_{i=1}^n (y_i - f(\underline{x}_i))^2$$

For classification problems, suppose $y \in \{1, \dots, M\}$,
model

$$P(Y=m|X) = \frac{e^{z_m}}{\sum_{l=1}^M e^{z_l}}$$

using the softmax function where

$$z_m = f_m(x) \quad m=1, \dots, M$$

↪
a neural network.

Estimation minimizes cross-entropy:

$$- \sum_{i=1}^n \sum_{m=1}^M y_{im} \log(f_m(x_i)).$$

10.2 Multilayer + Deep Networks

A multiple layer network gives rise to deep learning.

single layer

$$f(x) = \beta_0 + \sum_{k=1}^K \beta_k A_k$$

$$A_k = g\left(\omega_{k0} + \sum_{j=1}^p \omega_{kj} x_j\right)$$

two layer

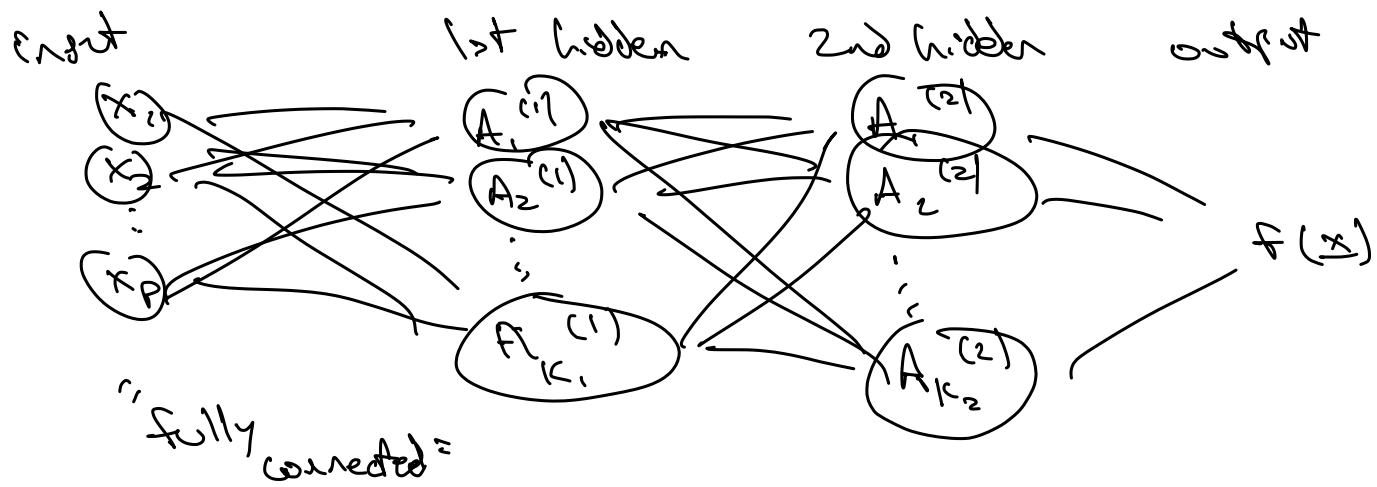
K_1 nodes in 1st layer, K_2 in 2nd layer

$$f(x) = \beta_0 + \sum_{k=1}^{K_2} \beta_k A_k^{(2)}$$

$$k=1, \dots, K_2 \rightarrow A_k^{(2)} = g\left(\omega_{k0}^{(2)} + \sum_{j=1}^{K_1} \omega_{kj}^{(2)} \underbrace{A_j^{(1)}}_{\text{not } x}\right)$$

$$\rightarrow A_k^{(1)} = g\left(\omega_{k0}^{(1)} + \sum_{j=1}^p \omega_{kj}^{(1)} x_j\right)$$

$k=1, \dots, K_1$



Deep network with L layers \rightarrow l th layer having K_l nodes

$$f(\underline{x}) = \beta_0 + \sum_{k=1}^{K_L} \beta_k A_k^{(L)}$$

$$A_k^{(L)} = g\left(w_{k0}^{(L)} + \sum_{j=1}^{K_{L-1}} w_{kj}^{(L)} A_j^{(L-1)}\right) \quad k=1, \dots, K_L$$

\vdots

$$A_k^{(2)} = g \left(\omega_{k0}^{(2)} + \sum_{j=1}^{K_1} \omega_{kj}^{(2)} A_j^{(1)} \right) \quad k=1, \dots, K_2$$

$$A_k^{(1)} = g \left(\omega_{k0}^{(1)} + \sum_{j=1}^P \omega_{kj}^{(1)} x_j \right) \quad k=1, \dots, K_1$$