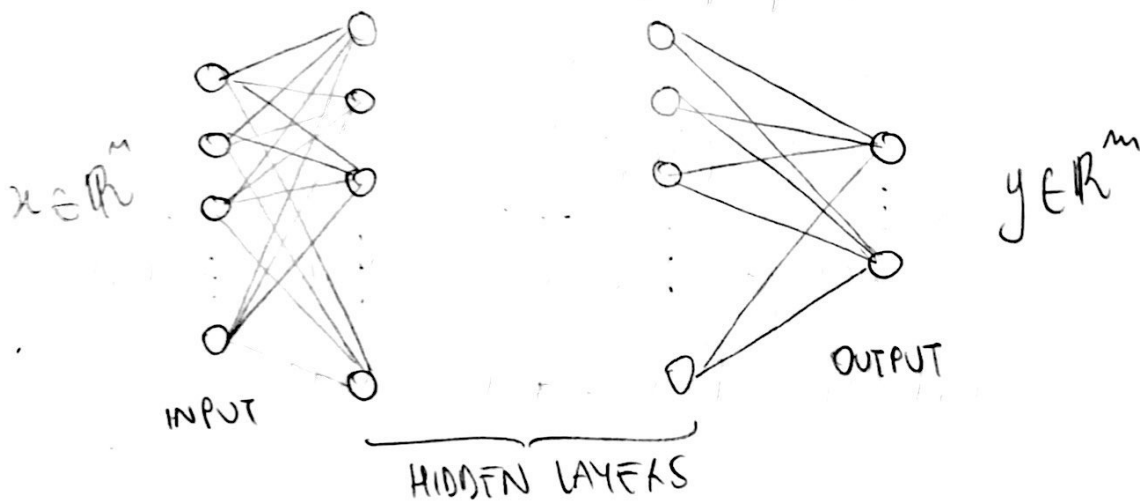# ARTIFICIAL NEURAL NETWORKS

HUBEL + WIESEL $\rightarrow$ NOBEL PRIZE.

1980 - NEOCOGNITRON : FIRST MODEL TO SIMULATE THE HIERARCHICAL STRUCTURE OF THE CAT VISUAL CORTEX

2012 - ImageNet : REVOLUTIONIZE WHAT WE THINK ABOUT OF NEURAL NETWORKS. $\rightarrow$ HUGE DATASET

THEY INVOLVE OPTIMIZATION AT VERY LARGE SCALE.



$x \in \mathbb{R}^n$    INPUT    HIDDEN LAYERS    OUTPUT    $y \in \mathbb{R}^m$

$$y = f(x, \beta)$$

COMPUTER VISION $\rightarrow$ INPUT IS AN IMAGE

WE CAN SEE THE MAP BETWEEN INPUTS AND FIRST LAYER AS $x^{(1)} = A_1 x$ AND SIMILARLY FOR ALL THE NEXT LAYERS.

$$\left.\begin{array}{l} x^{(1)} = A_1 x \\ x^{(2)} = A_2 x^{(1)} \\ \quad \vdots \\ y = A_m x^{(m-1)} \end{array}\right\} \quad \text{NESTED STRUCTURE}$$

OR MORE COMPACTLY :

$$y = A_m \ldots A_2 A_1 x$$

31

IF ALL LINEAR:

$$\bar{A} = A_m \dots A_2 A_1$$

$$y = \bar{A} x$$

THE BEST WE CAN DO IS JUST A LINEAR MAPPING. WE NEED TO INTRODUCE NONLINEARITY:

$$x^{(1)} = f_1 (A_1, x)$$
$$x^{(2)} = f_2 (A_2, x^{(1)})$$
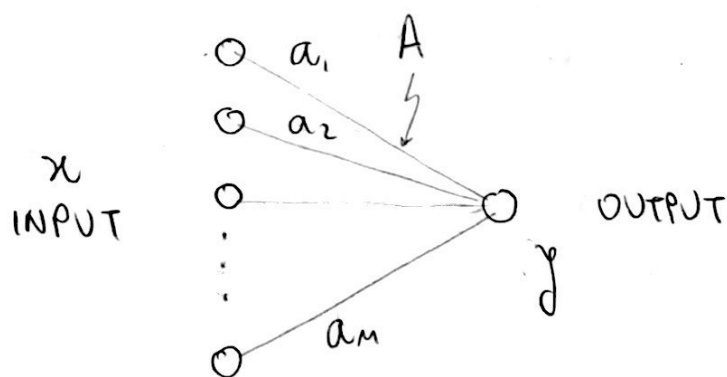
$$\vdots$$

$$y = f_m (A_m, x^{(m-1)})$$

THESE OPERATIONS WILL NOT COLLAPSE TO A SINGLE MATRIX:

$$y = f_m (A_m, \dots f_2 (A_2, f_1 (A_1, x)))$$

LET'S CONSIDER AN EXAMPLE

$x$ - dog/cat picture $64 \times 64$ images

$y$ - $\{dog, cat\} = \{+1, -1\}$



$$x = \begin{bmatrix} x_1 & x_2 & \dots & x_m \end{bmatrix} \quad y = \begin{bmatrix} +1, -1, \dots, +1 \end{bmatrix}$$

32

$$Y = AX$$

$$[+1 \ldots, \quad ,-1] = [a_1 \, a_2 \ldots a_m] \begin{bmatrix} | & & | \\ x_1 & \ldots & x_m \\ | & & | \end{bmatrix}$$

THE EASIEST WAY TO COMPUTE THE WAY IS WITH PSEUDOINVERSE.

THE SYSTEM IS OVERDETERMINED

IN THIS EXAMPLE THE ACTIVATION FUNCTION WAS THE IDENTITY:

$$F(x) = x$$

NONLINEAR ACTIVATION FUNCTIONS ARE MORE POWERFUL:

CONTINUOUS VERSION OF

- $F(x) = \begin{cases} 0 & x \leq 0 \\ 1 & x > 0 \end{cases}$  BINARY STEP

- $F(x) = \dfrac{1}{1 + e^{-x}}$  LOGISTIC

- $F(x) = \tanh(x)$  HYPERBOLIC TANGENT

- $F(x) = \begin{cases} 0 & x \leq 0 \\ x & x > 0 \end{cases}$  RELU

RELU HELP TO PROPAGATE GRADIENT WHEN $x \to \mp \infty$ SINCE LOGISTIC AND TANH TENDS TO OF FLAT IN THESE REGIONS.

AN ANN IS A NONLINEAR MAPPING FROM INPUTS TO OUTPUTS WITH A BUNCH OF PARAMETERS AND A COMPOSITIONAL STRUCTURE.

▷ GRADIENT DESCENT & BACKPROP.

LET'S START CONSIDERING THE SIMPLEST NETWORK POSSIBLE

$$x \; \bigcirc \xrightarrow{\quad a \quad} \bigcirc_{z} \xrightarrow{\quad b \quad} \bigcirc \; y$$

33

$$x \to f(x,a) \to z \to g(z,b) \to y$$

BACKPROPAGATION IS ITERATIVELY APPLY CHAINRULE WHICH WORKS SINCE THE NET IS COMPOSITIONAL:

$$y = g(z,b) = g(f(x,a);b) \leftarrow \text{COMPOSITION}$$

WE HAVE TRAINING DATA X THAT MAPS INTO OUTPUT DATA. OUR JOB IS TO FIGURE OUT $a$ AND $b$. THERE IS NO UPPERBOUND IN THE ERROR SO WE CAN OBTAIN A MINIMUM WITH THE DERIVATIVE.

$$E = \frac{1}{2}(y - y_0)^2$$

$$\underset{\text{MODEL}}{\nearrow} \qquad \underset{\text{TRUTH}}{\nwarrow}$$

$$\frac{\partial E}{\partial a} = 0 \to (y - y_0)\frac{dy}{dz}\frac{dz}{da} = 0$$

$$\frac{\partial E}{\partial b} = 0 \to (y - y_0)\frac{dy}{db} = 0$$

NOW WE UPDATE THE NET' PARAMETERS WITH GRADIENT DESCENT:

$$a_{k+1} = a_k + \delta \left.\frac{\partial E}{\partial a}\right|_k$$

$$b_{k+1} = b_k + \delta \left.\frac{\partial E}{\partial b}\right|_k$$

$\delta$ IS CALLED THE LEARNING RATE.

LET'S CONSIDER THE LINEAR CASE

$$z = ax$$
$$y = bz$$

$$\frac{\partial E}{\partial a} = -(y_0 - y)bx$$

$$\frac{\partial E}{\partial b} = -(y_0 - y)z = -(y_0 - y)ax$$

34

$$y = f(x, \beta) = f(x, A_1, A_2, \ldots, A_m)$$

$$\arg\min_{A_j} E(A_1, A_2, \ldots, A_m) = \arg\min_{A_j} \sum_{k=1}^{m} \underbrace{\left( f(x_k, \beta) - y_k \right)^2}_{E_k}$$

LARGE SCALE OPTIMIZATION

$$\underline{x}_{j+1} = \underline{x}_j - \delta \nabla f(\bar{x}_j)$$

THE BEST WAY TO UPDATE WEIGHTS IS WITH STOCHASTIC GRADIENT DESCENT. INSTEAD OF USING ALL THE SAMPLES IN A TIME WE USE JUST A RANDOM BATCH