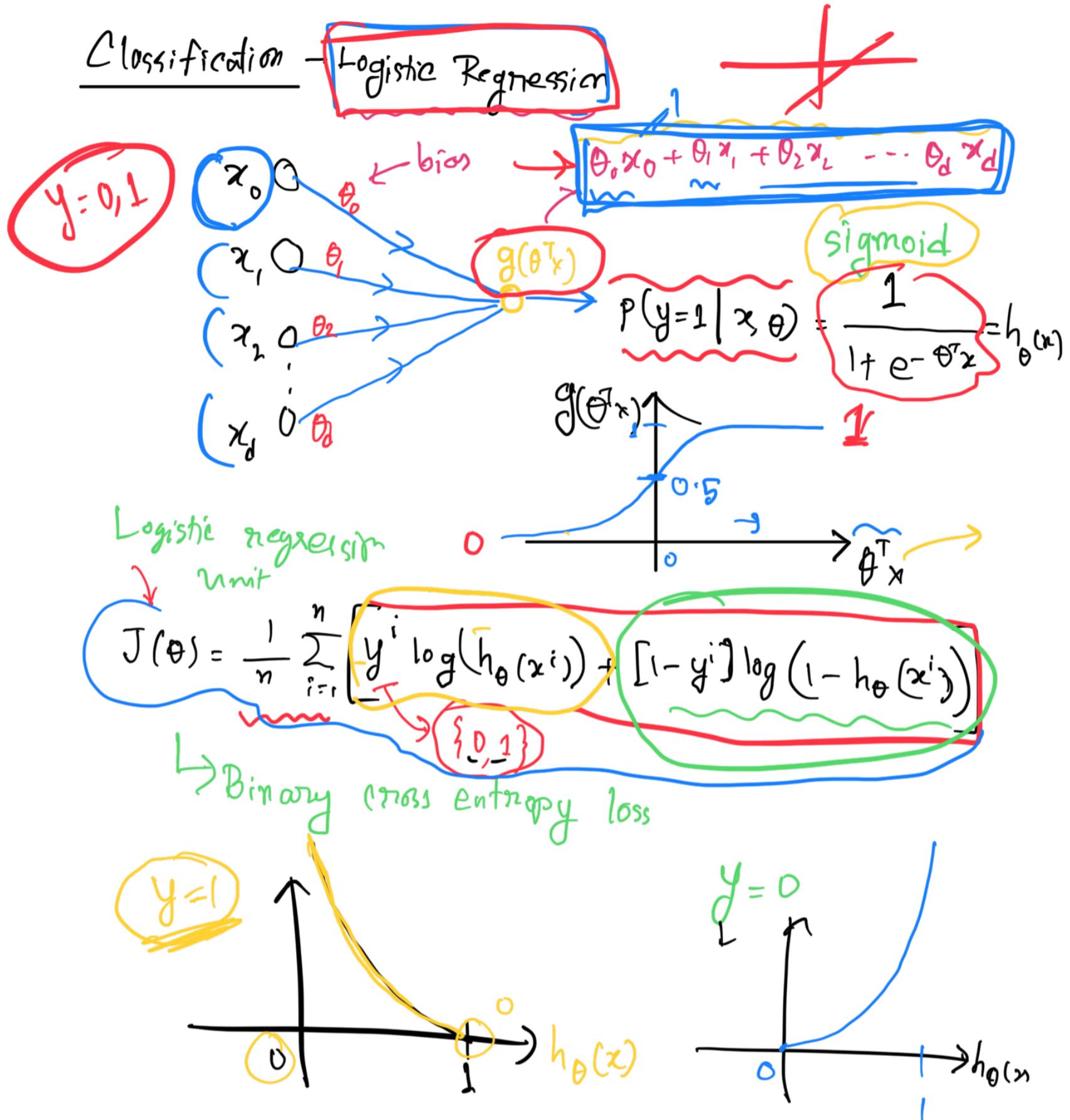
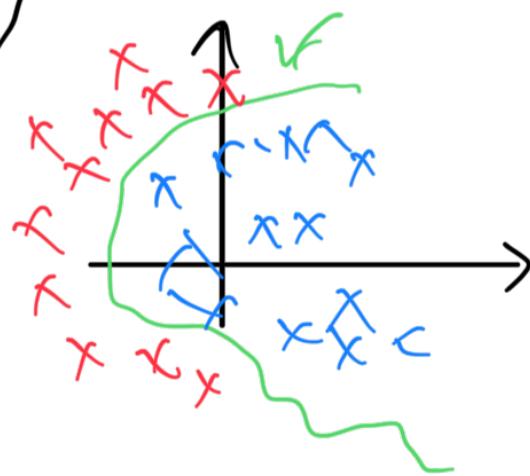
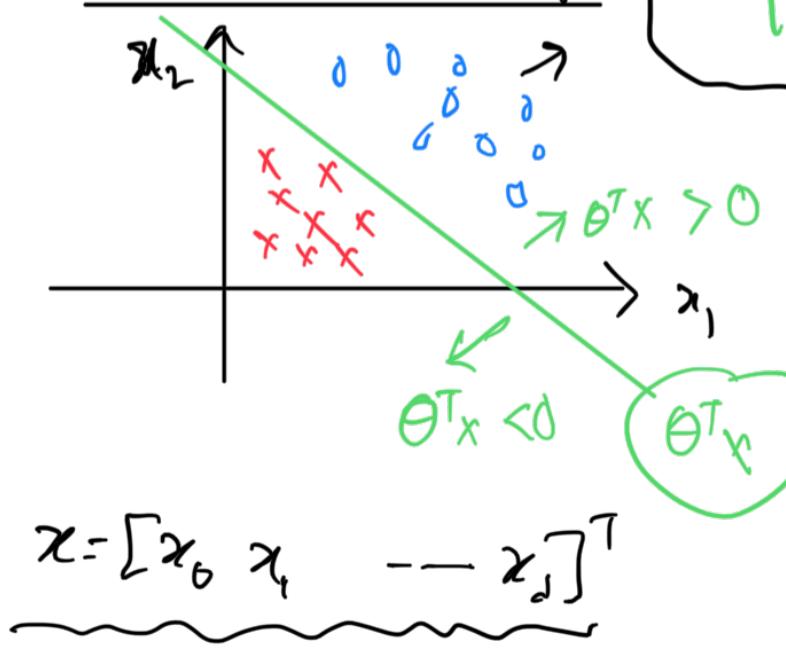


Classification - Logistic Regression

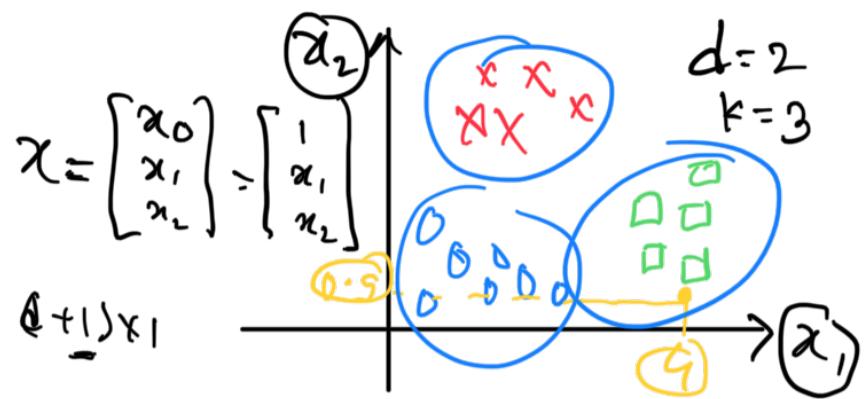


(Decision Boundary)



$$x = [x_0, x_1, \dots, x_d, x_0^2, x_1^2, \dots, x_d^2, x_0 x_1, x_0 x_2, \dots]^T$$

Multiclass classification



Email: Friends/fun/spam



Alg 0-1: One-vs-All

idea: Train a logistic regression classifier for each

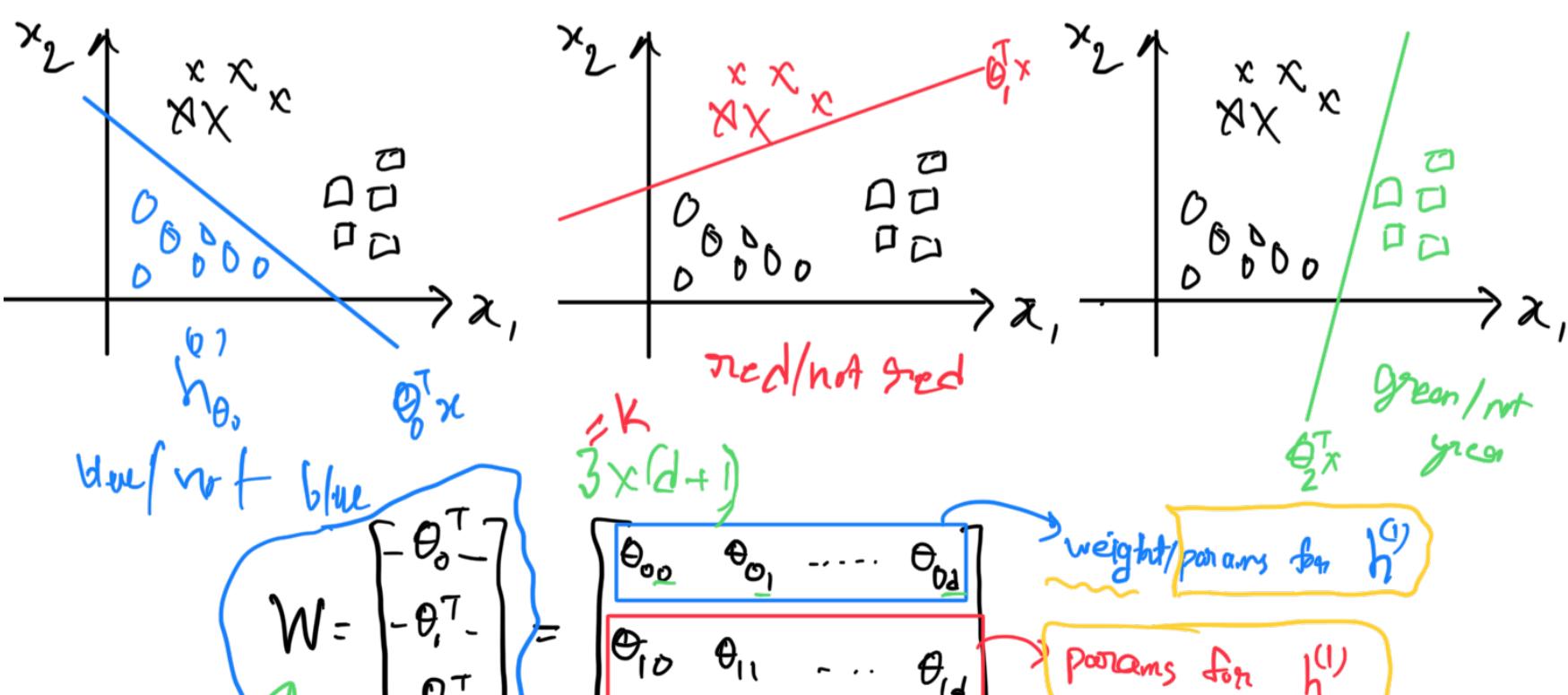
class i to predict the probability of $y=i$

$$y_{\text{pred}} = \arg \max_i h^{(i)}(x)$$

$$(1) = h_{\theta_0}^{(0)}(x) = \frac{1}{1 + e^{-\theta_0^T x}} = P(y=0 | x, \theta_0), \quad \theta_0 = [\theta_{00}, \theta_{01}, \theta_{02}, \dots, \theta_{0d}]^T$$

$$(2) = h_{\theta_1}^{(1)}(x) = \frac{1}{1 + e^{-\theta_1^T x}} = P(y=1 | x, \theta_1), \quad \theta_1 = [\theta_{10}, \theta_{11}, \theta_{12}, \dots, \theta_{1d}]^T$$

$$(3) = h_{\theta_2}^{(2)}(x) = \frac{1}{1 + e^{-\theta_2^T x}} = P(y=2 | x, \theta_2), \quad \theta_2 = [\theta_{20}, \theta_{21}, \theta_{22}, \dots, \theta_{2d}]^T$$



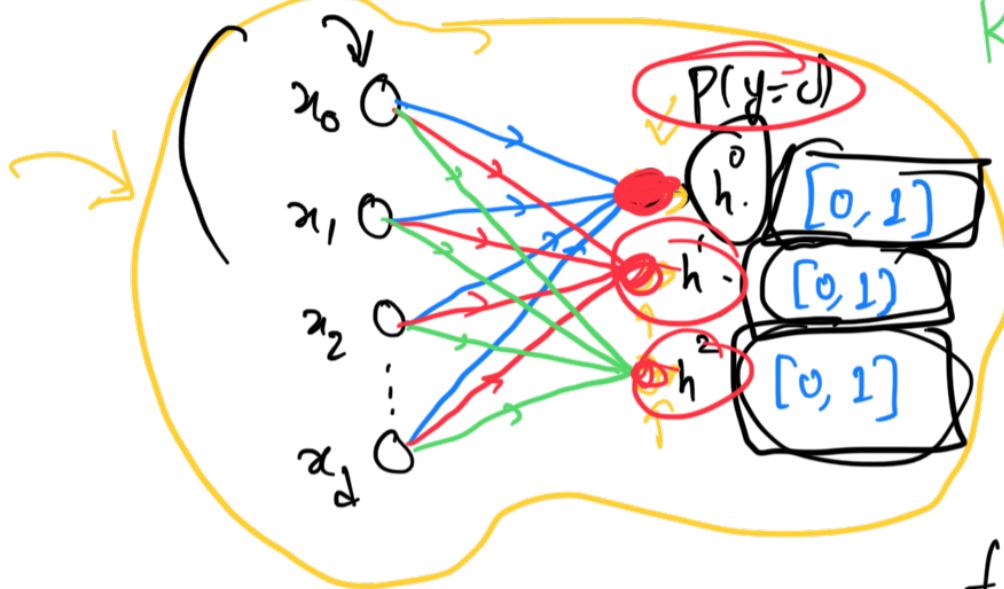
Weight

$$\begin{bmatrix} h_0 \\ h_1 \\ h_2 \end{bmatrix} = h = f(Wx) \quad K \times 1$$

of class $K \times (d+1)$

$$= f\left(\begin{bmatrix} \theta_0^T x \\ \theta_1^T x \\ \vdots \\ \theta_K^T x \end{bmatrix}\right) = \begin{bmatrix} f(\theta_0^T x) \\ f(\theta_1^T x) \\ \vdots \\ f(\theta_K^T x) \end{bmatrix} = \begin{bmatrix} \frac{1}{1+e^{-\theta_0^T x}} \\ \frac{1}{1+e^{-\theta_1^T x}} \\ \vdots \\ \frac{1}{1+e^{-\theta_K^T x}} \end{bmatrix}$$

K: number of class.



$$h = f(Wx) = \begin{bmatrix} h_0 \\ h_1 \\ h_2 \end{bmatrix} \quad \begin{array}{l} \text{Pr} = h_0 \\ \text{Pr} = h_1 \\ \text{Pr} = h_2 \end{array}$$

$$K \times (d+1)$$

$$f = \begin{bmatrix} \dots \end{bmatrix} \rightarrow \text{Prob. dist. of } y.$$

Numerical example

$$\theta_0 = [1, -1, -1]^T, \quad \theta_1 = [-1, -0.5, 1]^T, \quad \theta_2 = [-4, 2, -1]^T$$

$$\rightarrow x = [1, 4, 0.5]^T$$

$$W = \begin{bmatrix} 1 & -1 & -1 \\ -1 & -0.5 & 1 \\ -4 & 2 & -1 \end{bmatrix}$$

$$\Rightarrow Wx = \begin{bmatrix} 1 \times 1 + (-1) \times 4 + (-1) \times 0.5 \\ (-1) \times 1 + (-0.5) \times 4 + 1 \times 0.5 \\ (-4) \times 1 + 2 \times 4 + (-1) \times 0.5 \end{bmatrix} = \begin{bmatrix} -3.5 \\ -2.5 \\ 3.5 \end{bmatrix}$$

↓ weight

$$f(Wx) = \begin{bmatrix} \frac{1}{1+e^{(-3.5)}} \\ \frac{1}{1+e^{(-2.5)}} \\ \frac{1}{1+e^{(3.5)}} \end{bmatrix}$$

Activation function

$$= \begin{bmatrix} h_0 \\ h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} 0.03 \\ 0.08 \\ 0.92 \end{bmatrix}$$

$$0.92 + 0.08 \rightarrow 0.01$$

$$\therefore y_{\text{pred}} = 2 (\square)$$

$$\approx 1.08$$

$$\sum h = 2$$

Learning

$$J^{(0)}(\theta_0)$$

$$J^{(1)}(\theta_1)$$

$$J^{(2)}(\theta_2)$$

Binary Cross Entropy Loss

Training/Optimizing - separately [Gr. D. / Adam / RMSprop]

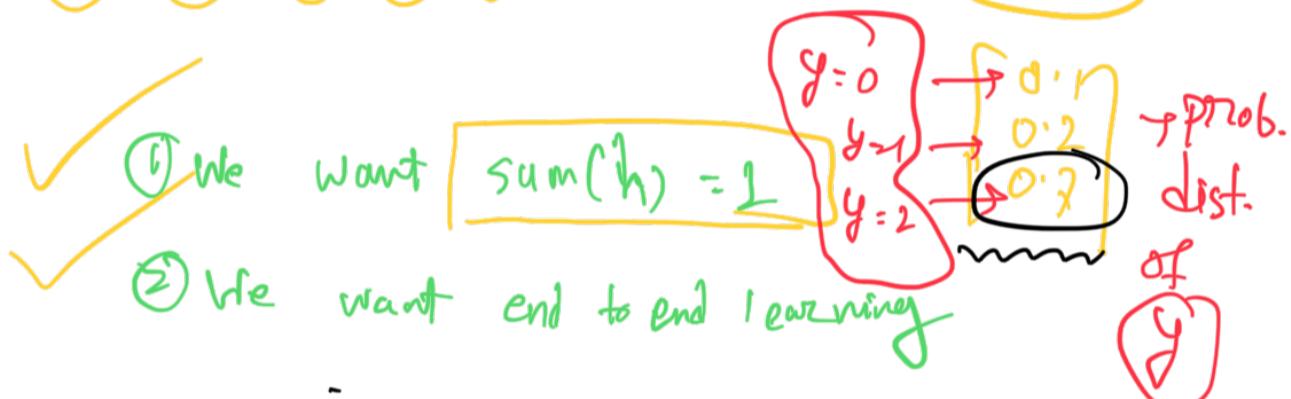
Prob:

① Have to train separately

ob $y = L_{\text{ImageNet}} \rightarrow K = 1000$

② No interpretation for h [$\sum(h) > 1$]

Sol^m: Softmax Classifier - Extension of L.R.



⇒ modifications

→ Activation/non-l.

① f_{non} ($f = \text{sigmoid}$ to $f = \text{softmax}$) $\sum h = 1$

② Change y to one-hot vector

$$\begin{aligned} \textcircled{1} \Rightarrow & \quad \textcircled{2} \quad \textcircled{3} \quad f_{\text{non}} \quad J(\theta) = \text{BCE} \quad \text{to CE} \quad [\text{cross entropy}] \\ (K \times (d+1)) \Rightarrow & W_x = \left[z_0, z_1, \dots, z_{K-1} \right]^T \rightarrow f(W_x) = \left[\frac{1}{1+e^{-z_0}}, \frac{1}{1+e^{-z_1}}, \dots \right]^T \end{aligned}$$

$$f(W_x) = \left[f_0, f_1, \dots, f_K \right]^T \quad K \times 1$$

$K \times 1$

$$f_1 = P(y=1)$$

$$f_2 = P(y=2)$$

$$\sum f_i = 1$$

$$f_i = \frac{e^{z_i}}{\sum_{j=0}^{K-1} e^{z_j}}$$

$$f_2 = \frac{e^{z_2}}{e^{z_0} + e^{z_1} + e^{z_2}}$$

$$\textcircled{0} \quad \sum_{i=0}^{K-1} f_i = 1,$$

$$\textcircled{11} \quad 0 \leq f_i \leq 1$$

$$P(y=i) = f_i = \frac{e^{z_i}}{\sum_{j=0}^{K-1} e^{z_j}}$$

$f \rightarrow$ probability distribution of y

$$(2) \Rightarrow y = [0, 1, 2] \rightarrow y = \left\{ \begin{matrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, & 1 \\ \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, & 2 \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} & 0 \end{matrix} \right\}$$

Avg. over all training sample

* One-hot encoding

$$(3) J(W) = - \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{K-1} y^{(i)}_k \log f_k$$

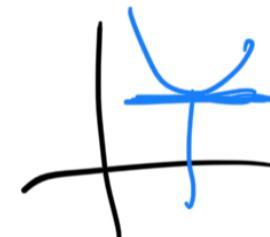
Cross Entropy Loss

$W \rightarrow$ weight matrix
 \Rightarrow best

Optimizer \rightarrow same as before (Gradient Descent)

$$W^{(0)} = \text{random}$$

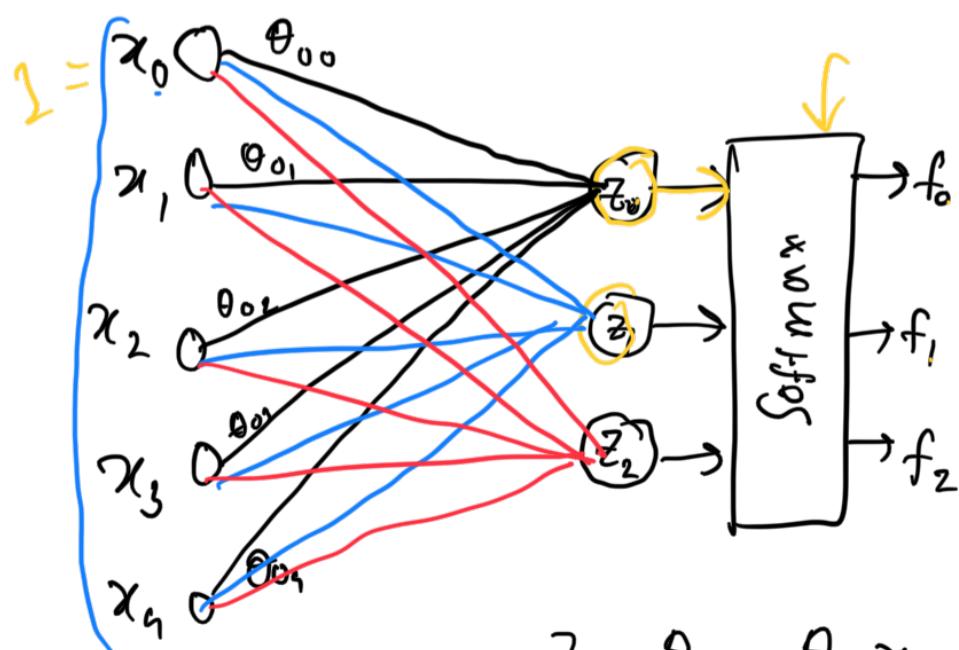
while not converged:



$$\tilde{W}^{(k+1)} = \tilde{W}^{(k)} - \nabla_{\tilde{W}} \tilde{J}$$

if $\nabla_{\tilde{W}} \tilde{J} \approx 0$:

converged = True.



$$w_x z_i = []$$

$$f(z) =$$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \end{bmatrix}$$

$$f = \text{Softmax}(w_x)$$

$$z_0 = \theta_{00} + \theta_{01}x_1 + \dots + \theta_{04}x_4$$

$$\text{for } \theta_0 \rightarrow \theta_0 = [0, 0.01, -0.05, 0.1, 0.05]^T$$

$$\text{for } \theta_1 \rightarrow \theta_1 = [0.2, 0.2, 0.2, 0.05, 0.16]^T$$

$$\text{for } \theta_2 \rightarrow \theta_2 = [0.3, 0.0, -0.45, -0.2, 0.33]^T$$

$$x^{(1)} = [1, -15, 22, -44, 56]$$

One-hot

$$y^{(1)} = 2 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$W = \begin{bmatrix} 0 & 0.01 & -0.05 & 0.1 & 0.05 \\ 0.2 & 0.2 & 0.2 & 0.05 & 0.16 \\ -0.3 & 0.0 & -0.45 & -0.2 & 0.33 \end{bmatrix}$$

$$z = \begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix} = Wx = \begin{bmatrix} -2.58 \\ 0.86 \\ 0.28 \end{bmatrix} \rightarrow z_0, z_1, z_2$$

$$f(Wx) = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \end{bmatrix}$$

$$f_0 = \frac{e^{z_0}}{e^{z_0} + e^{z_1} + e^{z_2}} = \frac{e^{-2.58}}{e^{-2.58} + e^{0.86} + e^{0.28}}$$

$$= \frac{0.058}{0.058 + 2.36 + 1.32} = 0.016$$

$$= \begin{bmatrix} 0.016 \\ 0.631 \\ 0.353 \end{bmatrix}$$

$$0.016 + 0.631 + 0.353 = 1$$

$$y_{\text{prob}} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = 1$$

$$f_1 = \frac{e^{z_1}}{e^{z_0} + e^{z_1} + e^{z_2}} = \frac{2.36}{0.058 + 2.36 + 1.32}$$

$$f_2 = \frac{e^{z_2}}{e^{z_0} + e^{z_1} + e^{z_2}} = \frac{1.32}{0.058 + 2.36 + 1.32} = 0.353$$

$$J_i = - \sum_{k=0}^{K-1} y_k \log f_k = -$$

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \rightarrow y_0, y_1, y_2$$

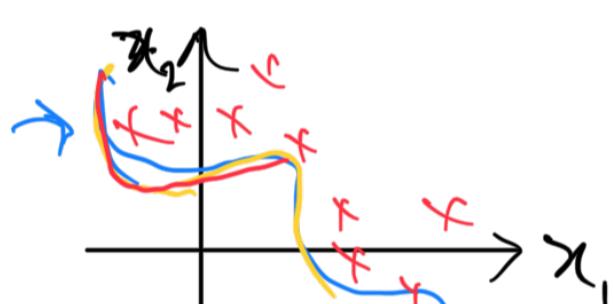
$$= - [y_0 \times \log f_0 + y_1 \times \log f_1 + y_2 \times \log f_2]$$

$$= - [0 \times \log 0.016 + 0 \times \log (0.631) + 1 \times \log (0.353)]$$

$$1.34, 2.06, 0.08 \Rightarrow J(w) = \frac{1}{4} [1.04 + 1.34 + 2.06 + 0.08]$$

Neural Networks (NN)

Why?



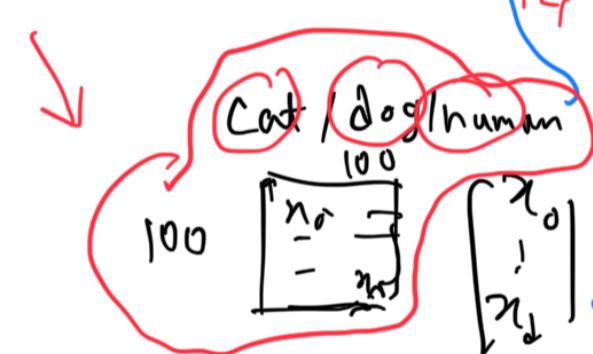
$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$\sin(x_1 + x_2)$

$$\begin{aligned} &x_0, x_1, x_2 \\ &x_0^2, x_1^2, x_2^2 \\ &x_0 x_1, x_0 x_2, x_1 x_2 \\ &\sin(x_1 + x_2) \end{aligned}$$

$$x_0^3, x_0^2 x_1, \dots$$

~~$x_0^3, x_0^2 x_1, \dots$~~



$$\begin{aligned} &x_0^2, x_0 x_1, x_1 x_2, \dots \\ &(\frac{n^2}{2}) \rightarrow 50 \times 10^6 !! \end{aligned}$$

→ NN are universal function approximators

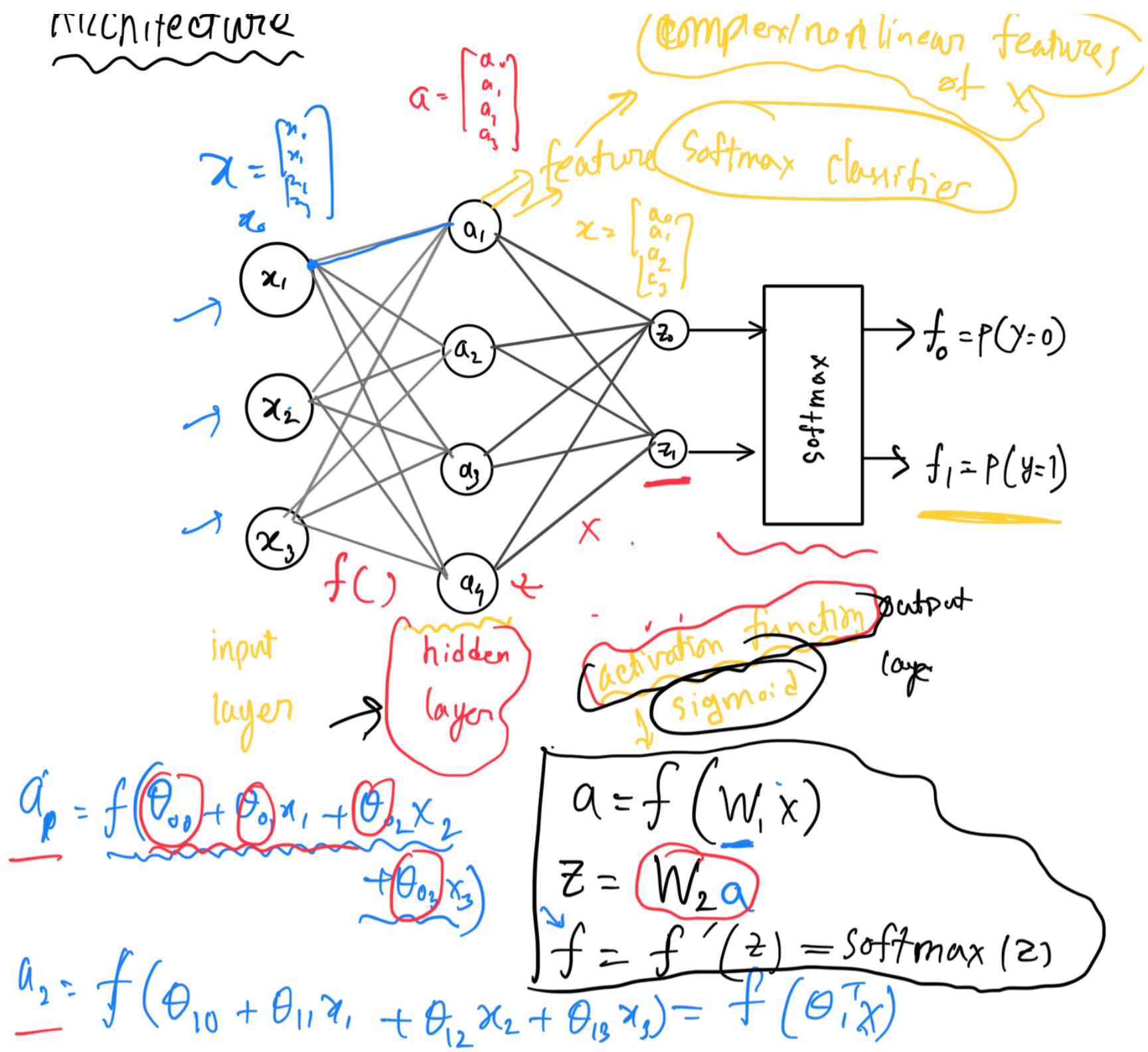
→ instead of increasing length, we add depth

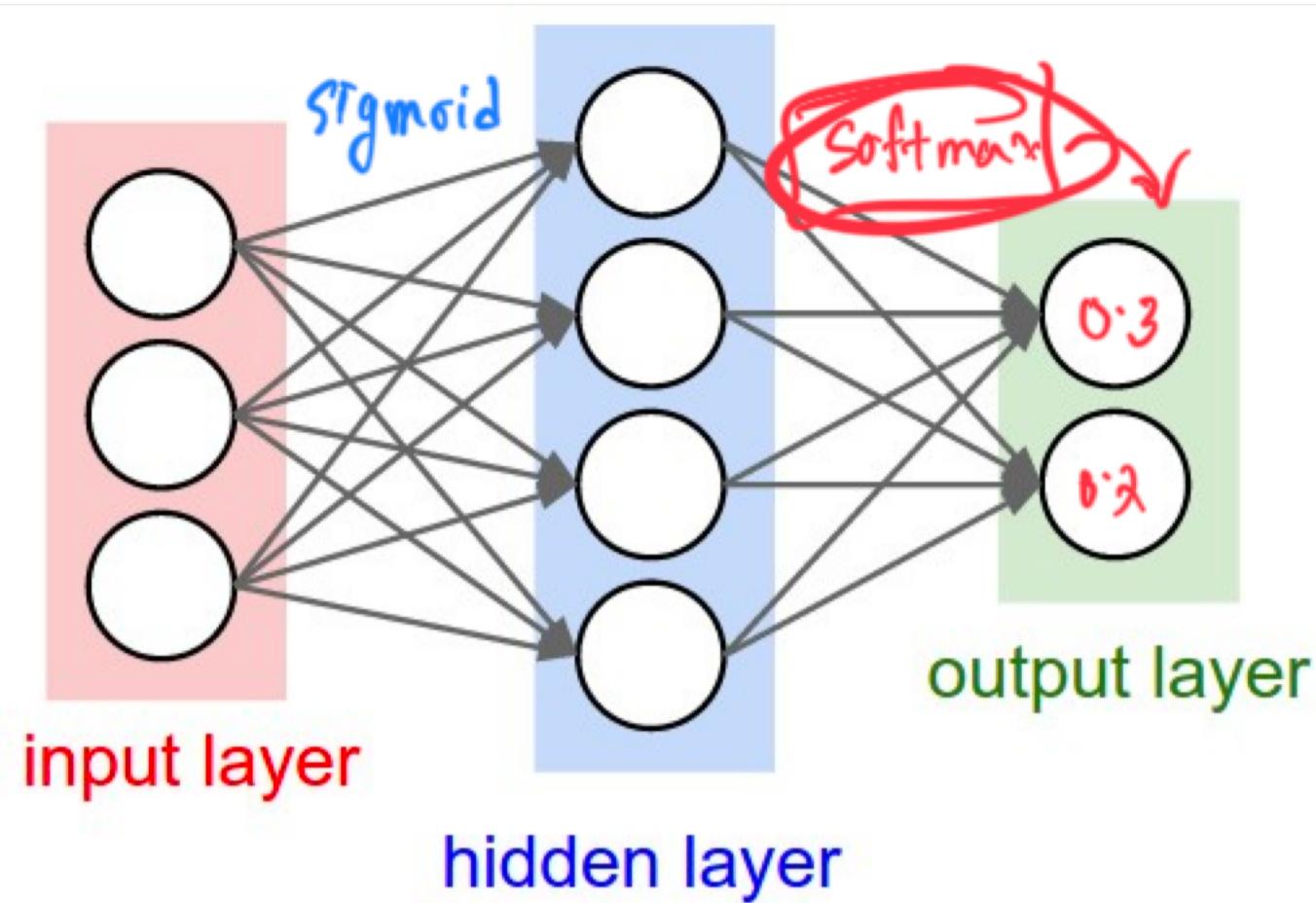
→ Deep neural networks can learn ~ more complex

functions, more acc.

→ Learn features themselves!

Architecture





More architectures

a

$$\rightarrow a = g(W_1 x)$$

$$g(z) = \text{non-linearity} \\ = \text{activation} \\ = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} F &= f(W_2 a) \\ &= \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_k \end{bmatrix} \rightarrow \end{aligned}$$

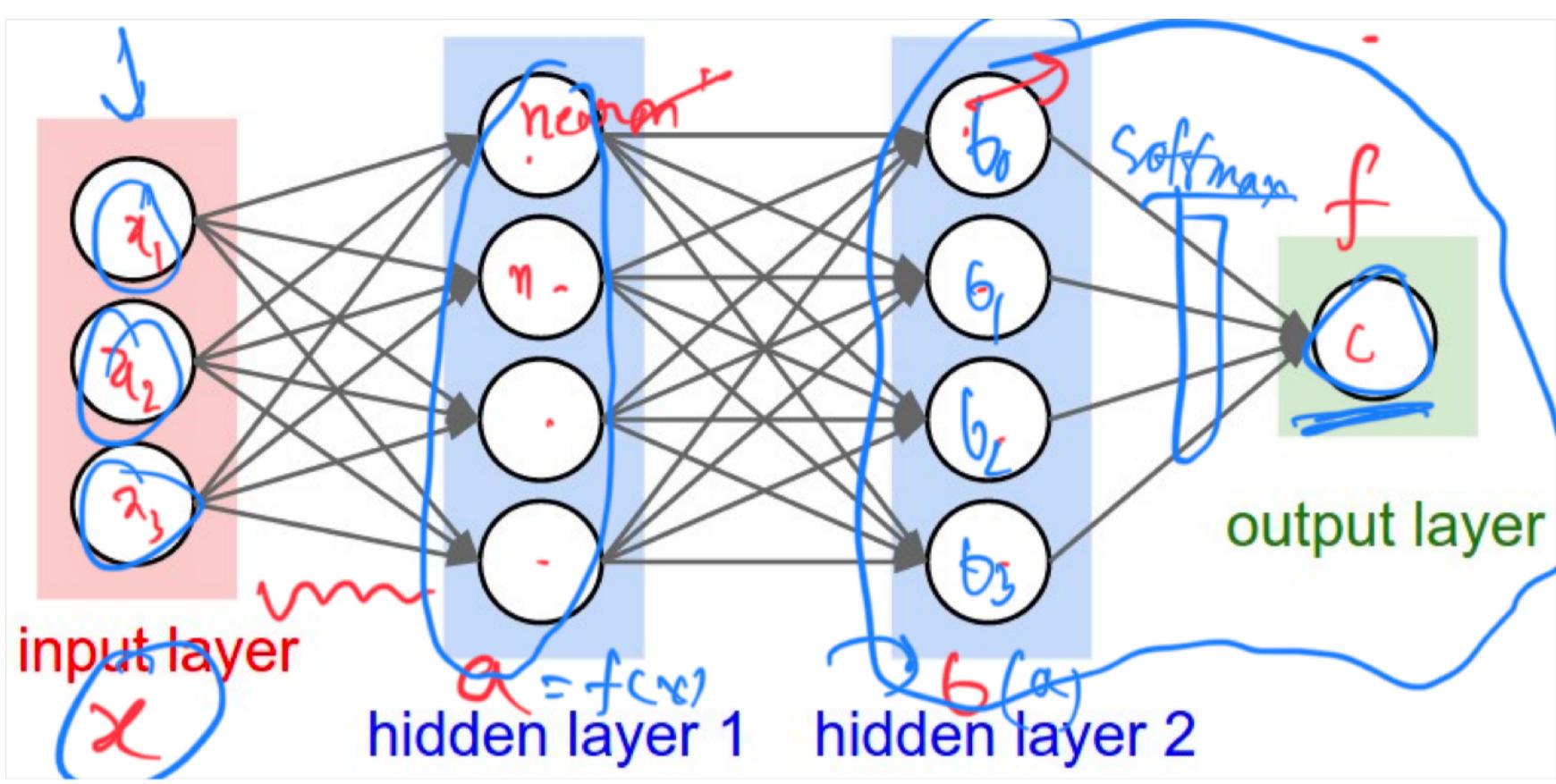
$$f = \text{Softmax}$$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \left[- \sum_{k=0}^{k_1} y_k \log f_k \right]$$

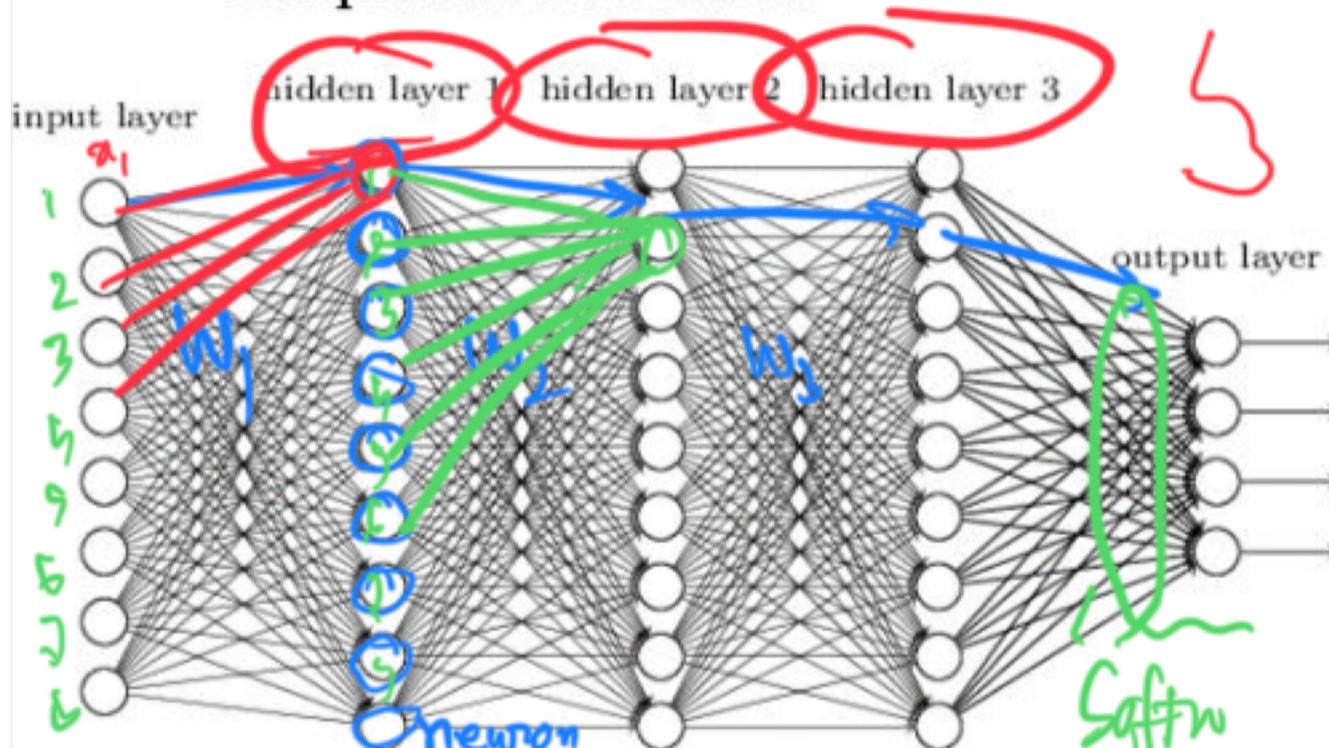


$$\begin{cases} W_1^{k+1} = W_1^k - \nabla_{W_1} J \\ W_2^{k+1} = W_2^k - \nabla_{W_2} J \end{cases}$$

Chain rule



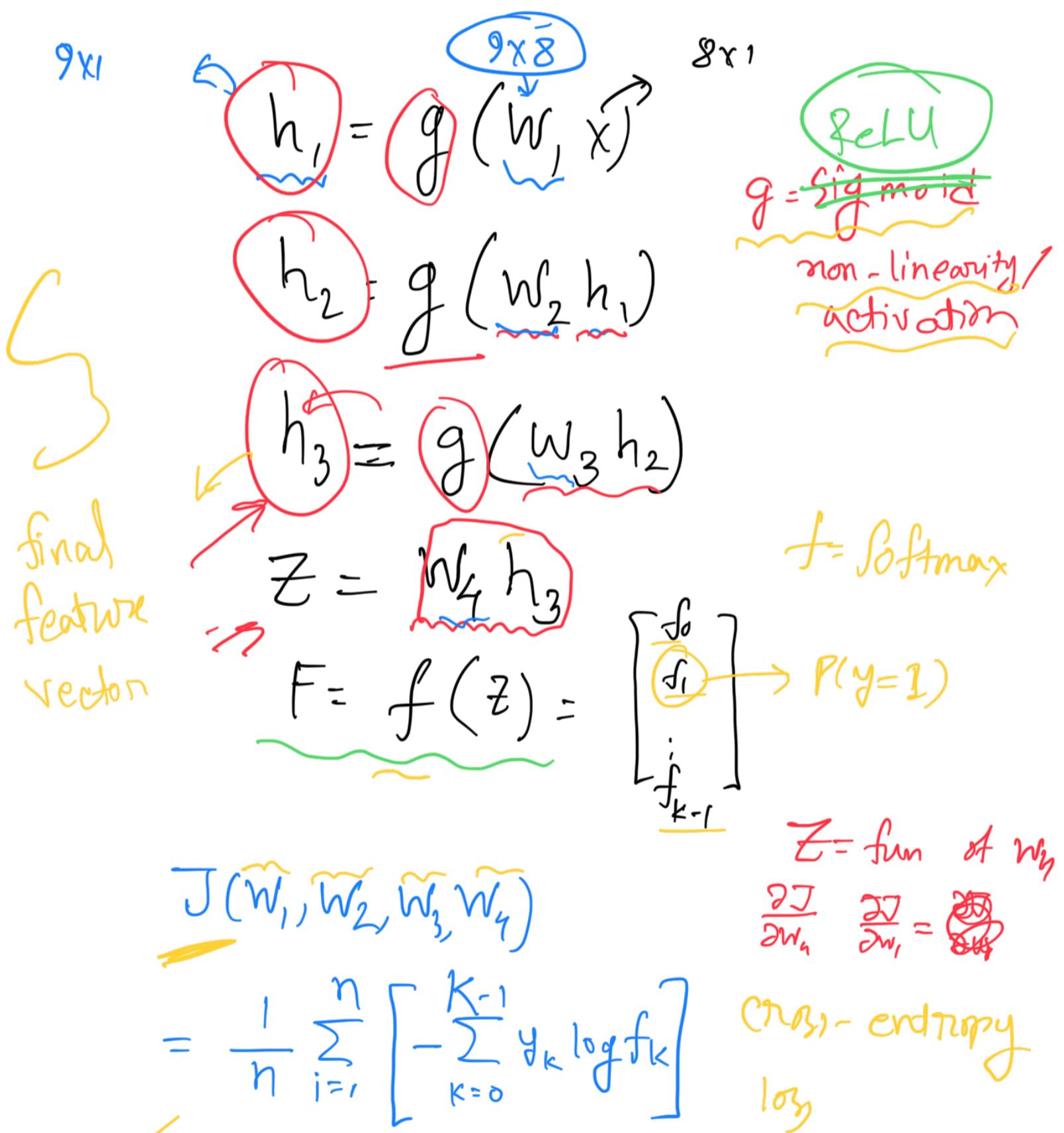
Deep neural network



M
N
M
M

Training neural networks

- artificial neural networks
- feed forward n
- fully connected n
- Dense n



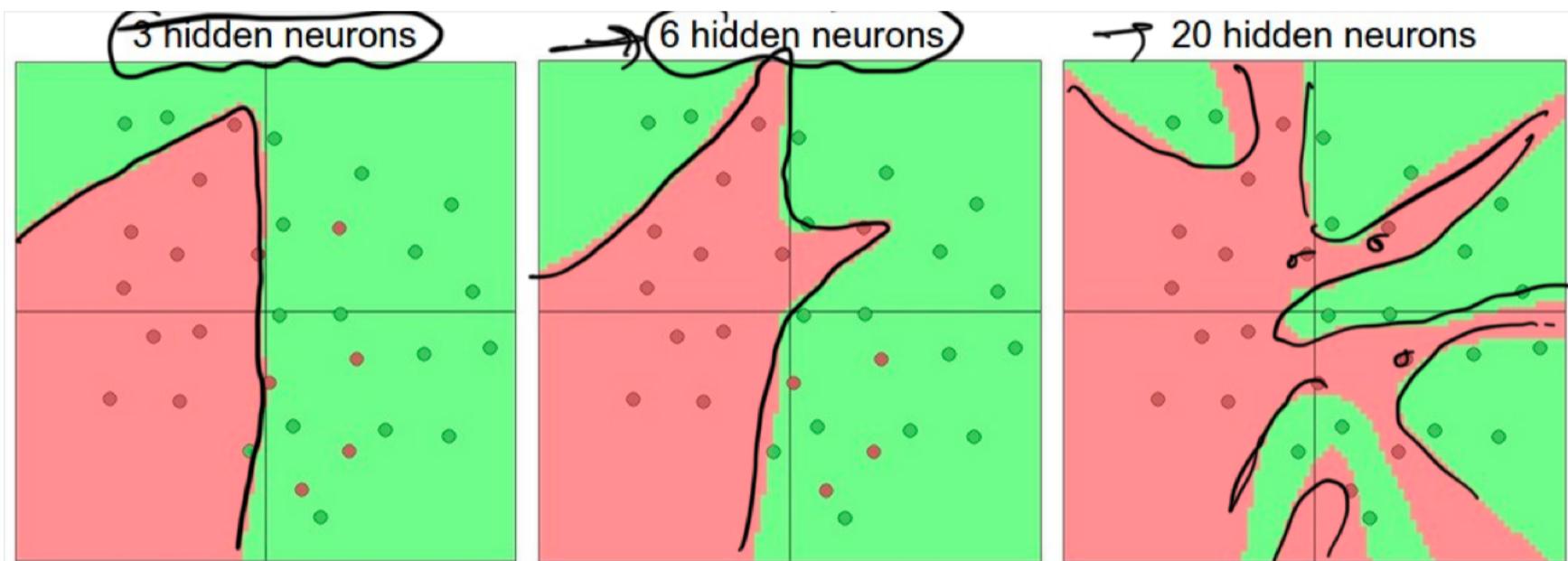
Gradient Descent

$$w_i^{(k+1)} = w_i^{(k)} - \nabla_{w_i} J$$

$$w_2 = w_2 - \nabla_{w_2} J$$

Practical consideration

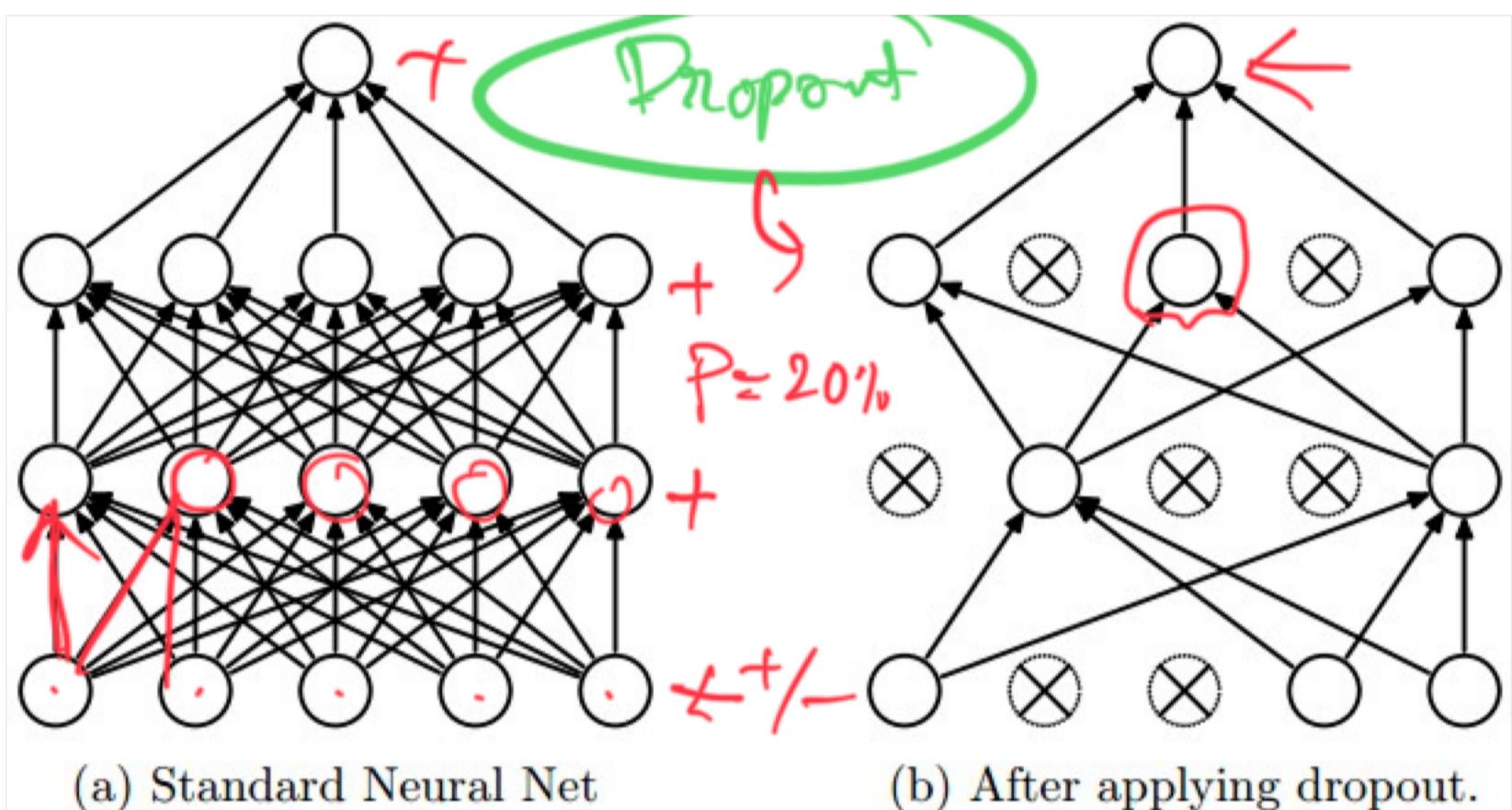
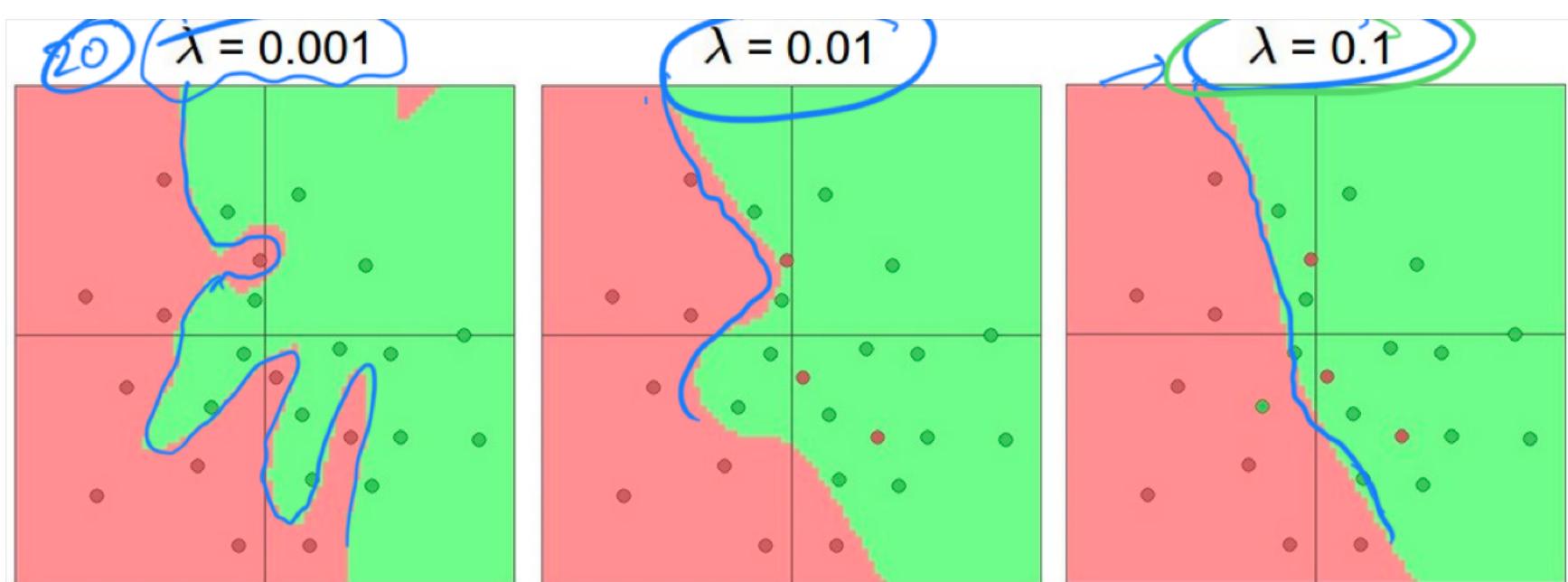
① Number of neurons



validation set! must

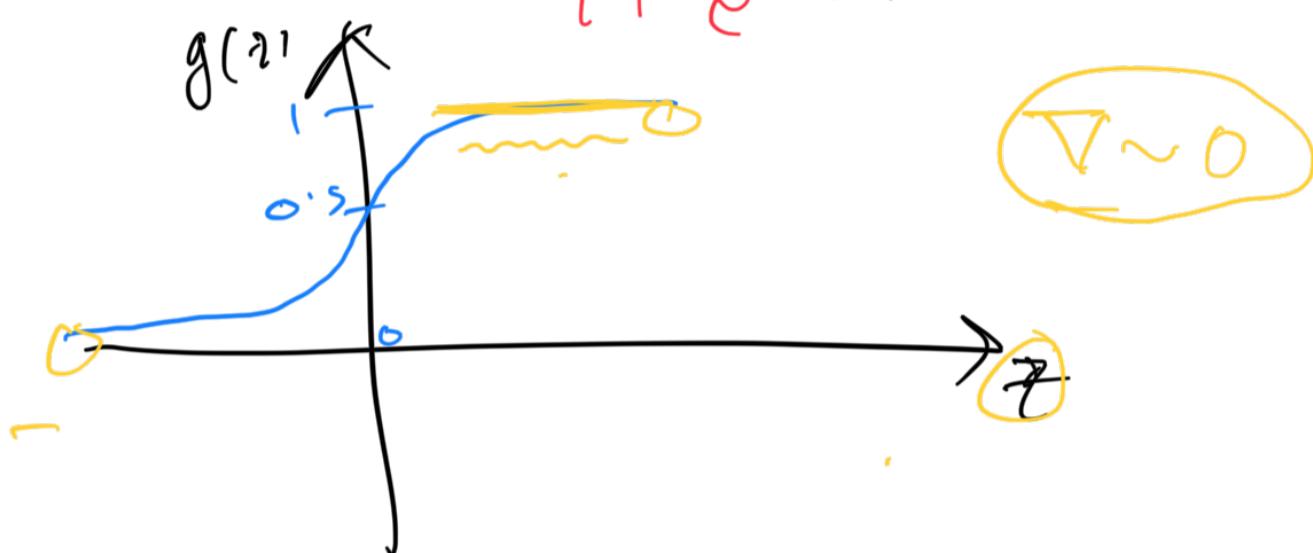
② Overfitting \rightarrow Regularization

$$J(\theta) + \alpha \sum |\theta_j|^2$$

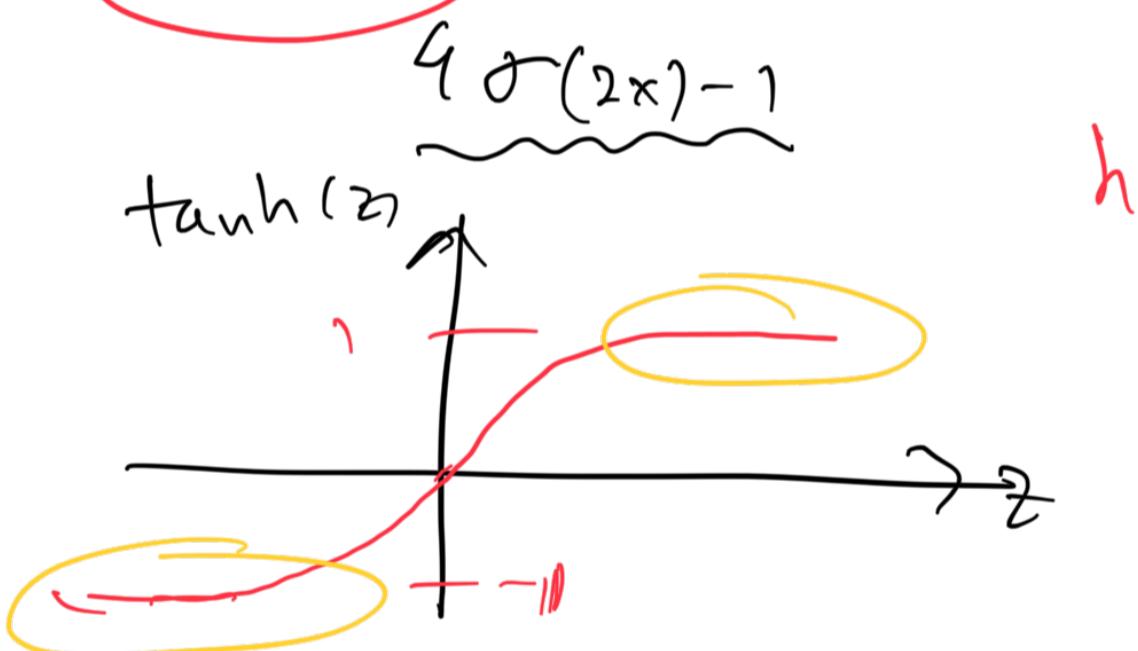


Activation function

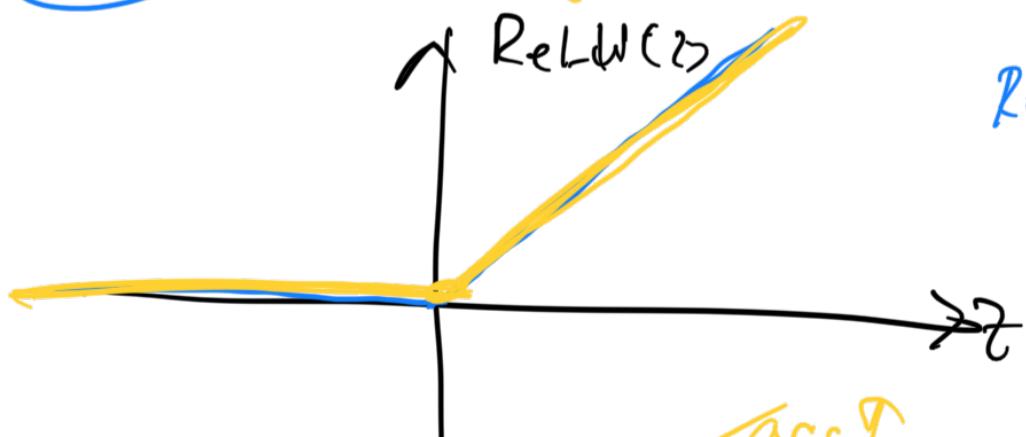
$$g(z) = \frac{1}{1 + e^{-\theta^T x}}$$



tanh

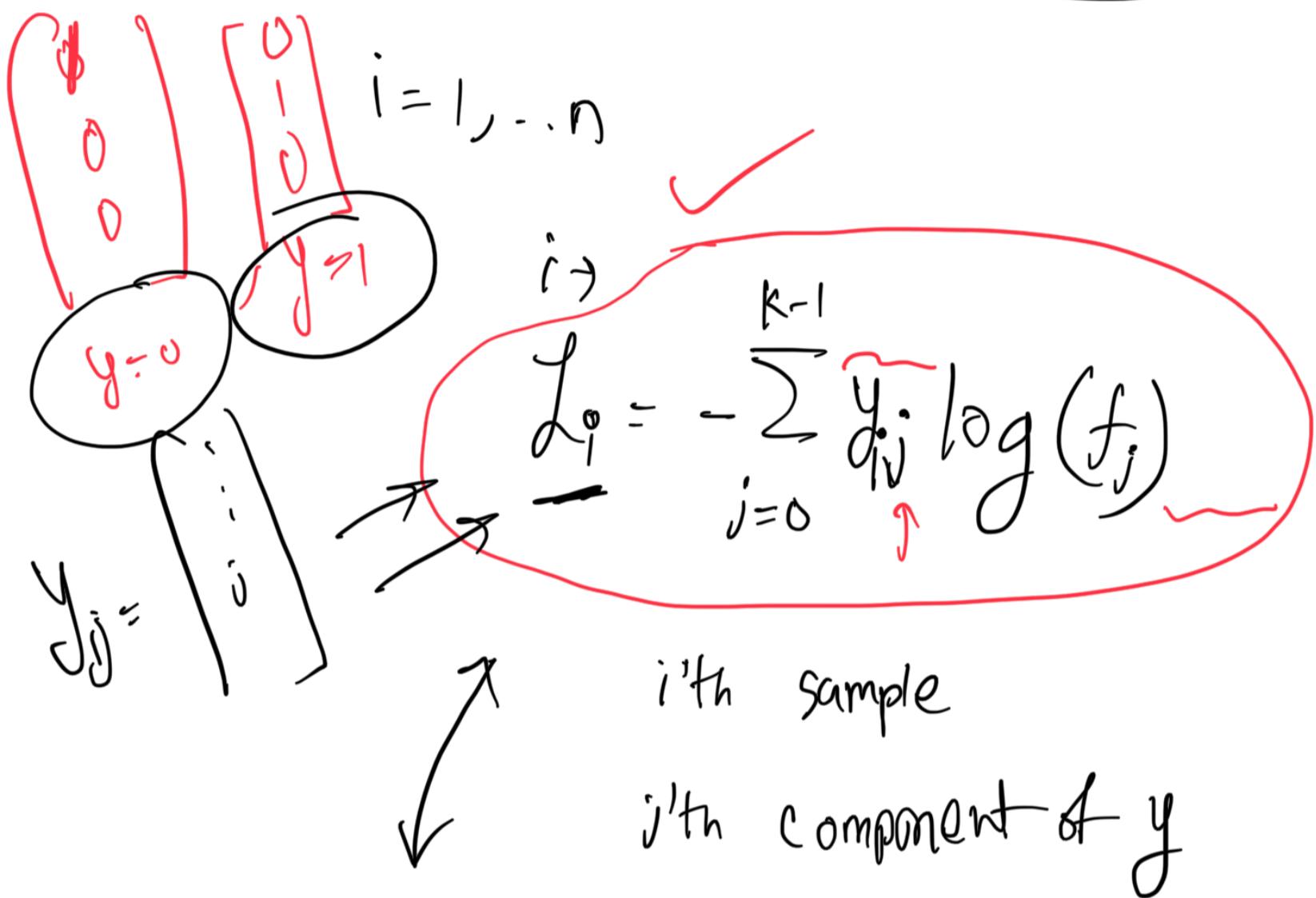
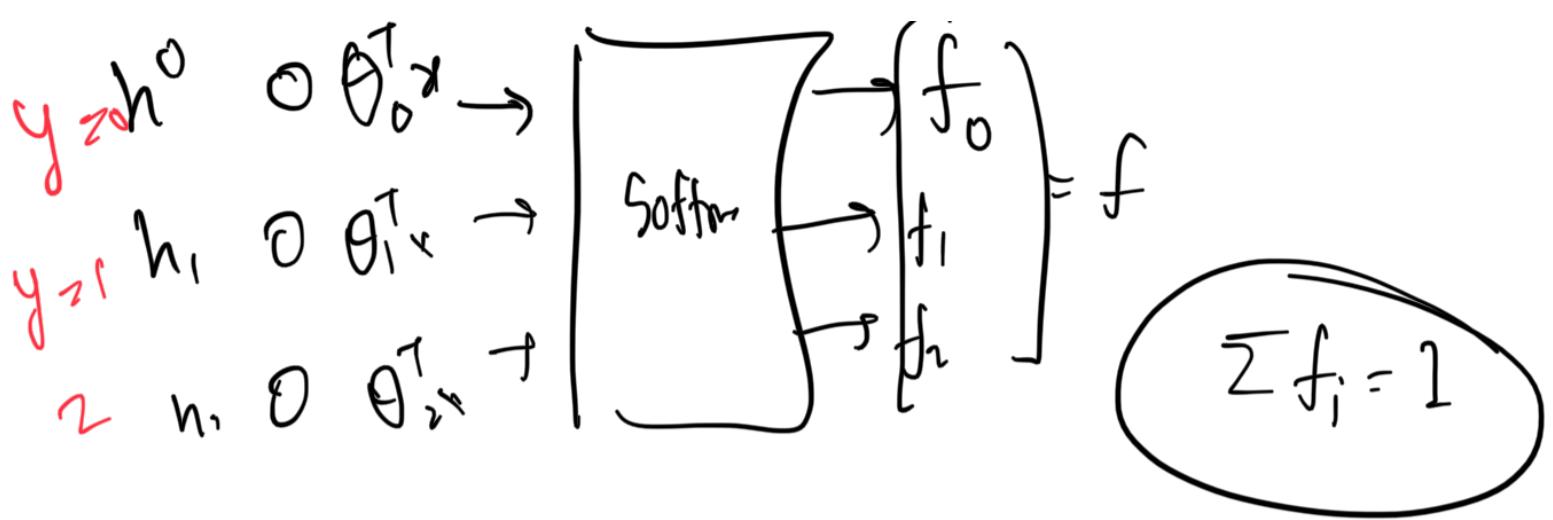


ReLU / Leaky ReLU



ReLU(z)
Leaky ReLU(z) if $z \geq 0$

acc ↑
speed ↑



$L_i = -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$ $y_i = \arg \max y$
 or equivalently $L_i = -f_{y_i} + \log \sum_j e^{f_j}$