

Lesson 4: Feedforward Neural Network (FFNN)

Definition: 1) A FFNN consists of layers of computational units (i.e., neurons), usually interconnected in a feedforward way.

b) A FFNN is an artificial Neural Network (ANN) wherein connections between the nodes (neurons) do not form a cycle.

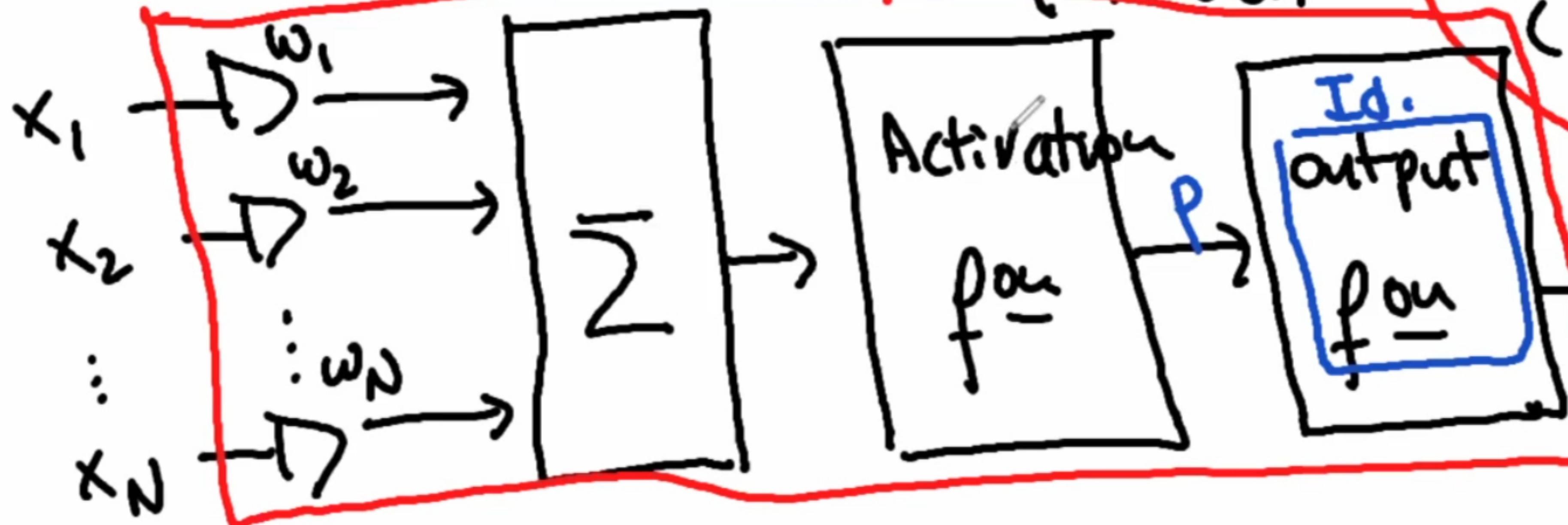
Context:

Biological NN (BNN)

Artificial NN (ANN)

Artificial Neural Network

Def: Artificial Neuron (with a model)



(i) FFNN

(Feedforward NN)

(ii) FBNN / RNN

(Feedback NN / Recurrent NN)

with Memory

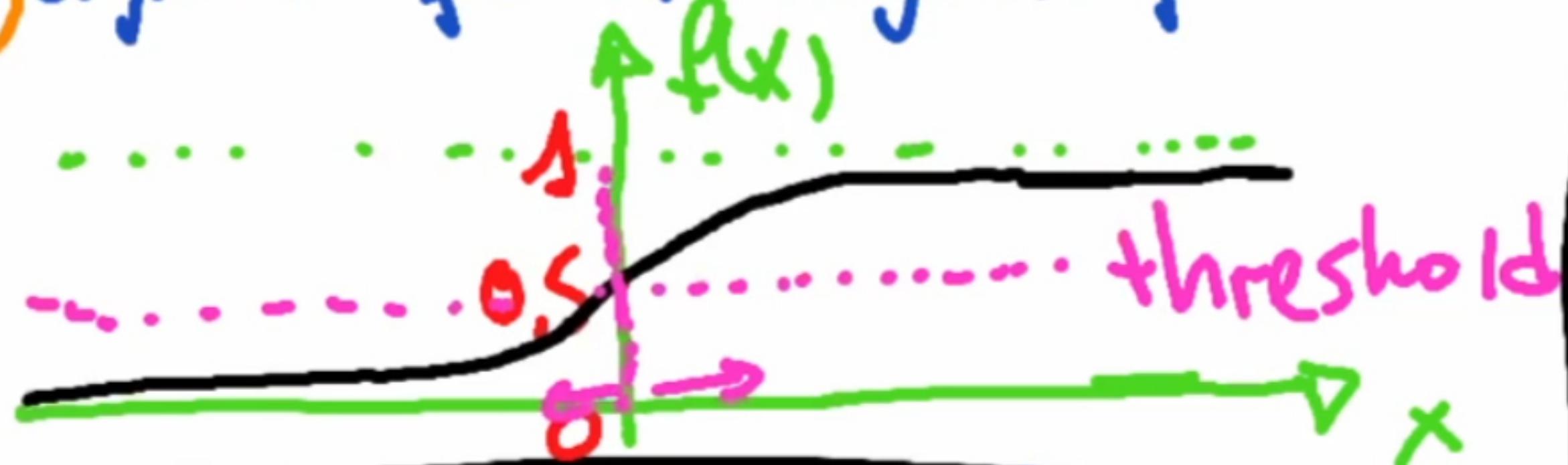
In this course,
output f_{out} is the
identity f_{id}

Types of activation functions (Basic)

For this course!

Graph

1) Sigmoid function or logistic func



2) Hyperbolic tangent func



3) Rectified Linear Unit (ReLU)



Equation

$$f(x) = \frac{1}{1 + e^{-x}}$$

$$f'(x) = f(x)(1 - f(x))$$

Derivative

$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f'(x) = (-f(x))^2$$

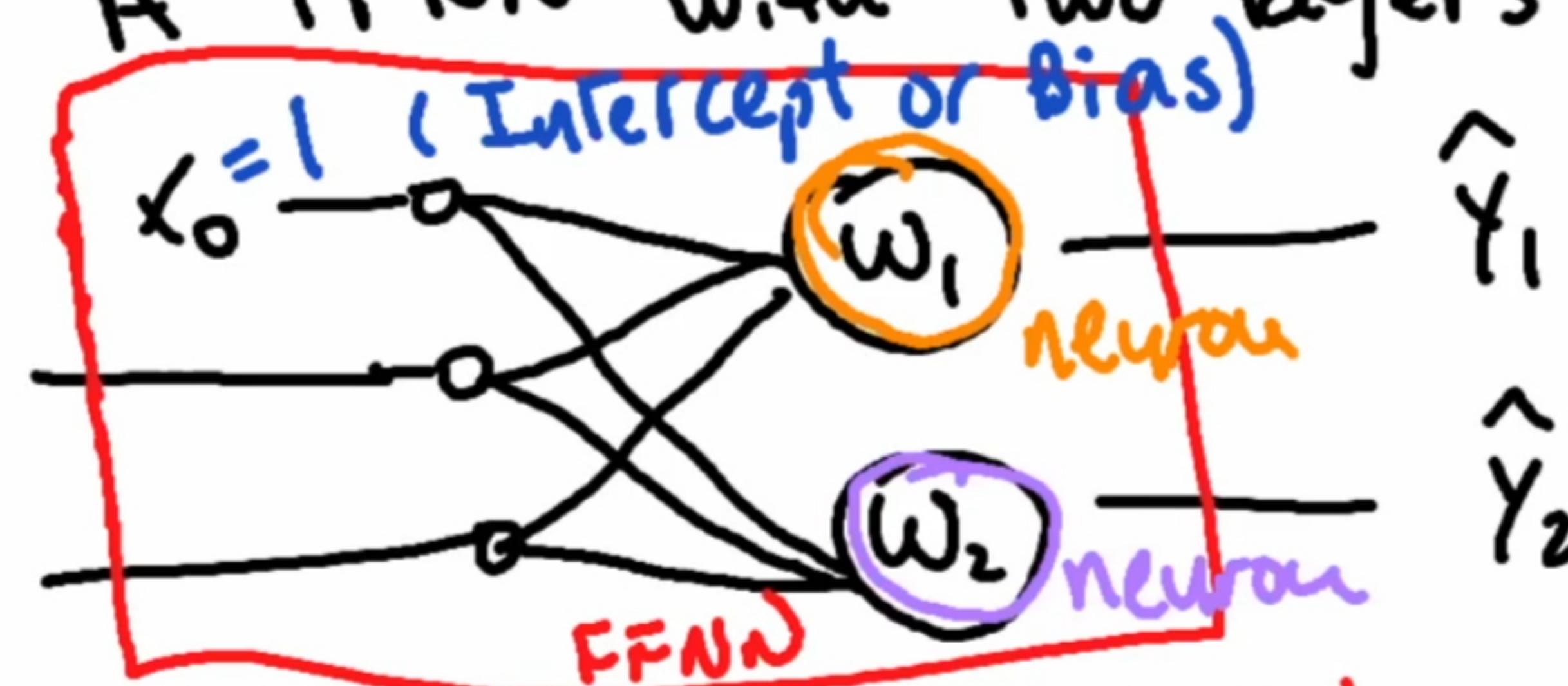
$$f(x) = \begin{cases} 0, & \text{if } x \leq 0 \\ x, & \text{if } x > 0 \end{cases} = \max\{0, x\} = x \cdot \mathbb{I}\{x > 0\}$$

$$f'(x) = \begin{cases} 0, & x \leq 0 \\ 1, & x > 0 \end{cases}$$

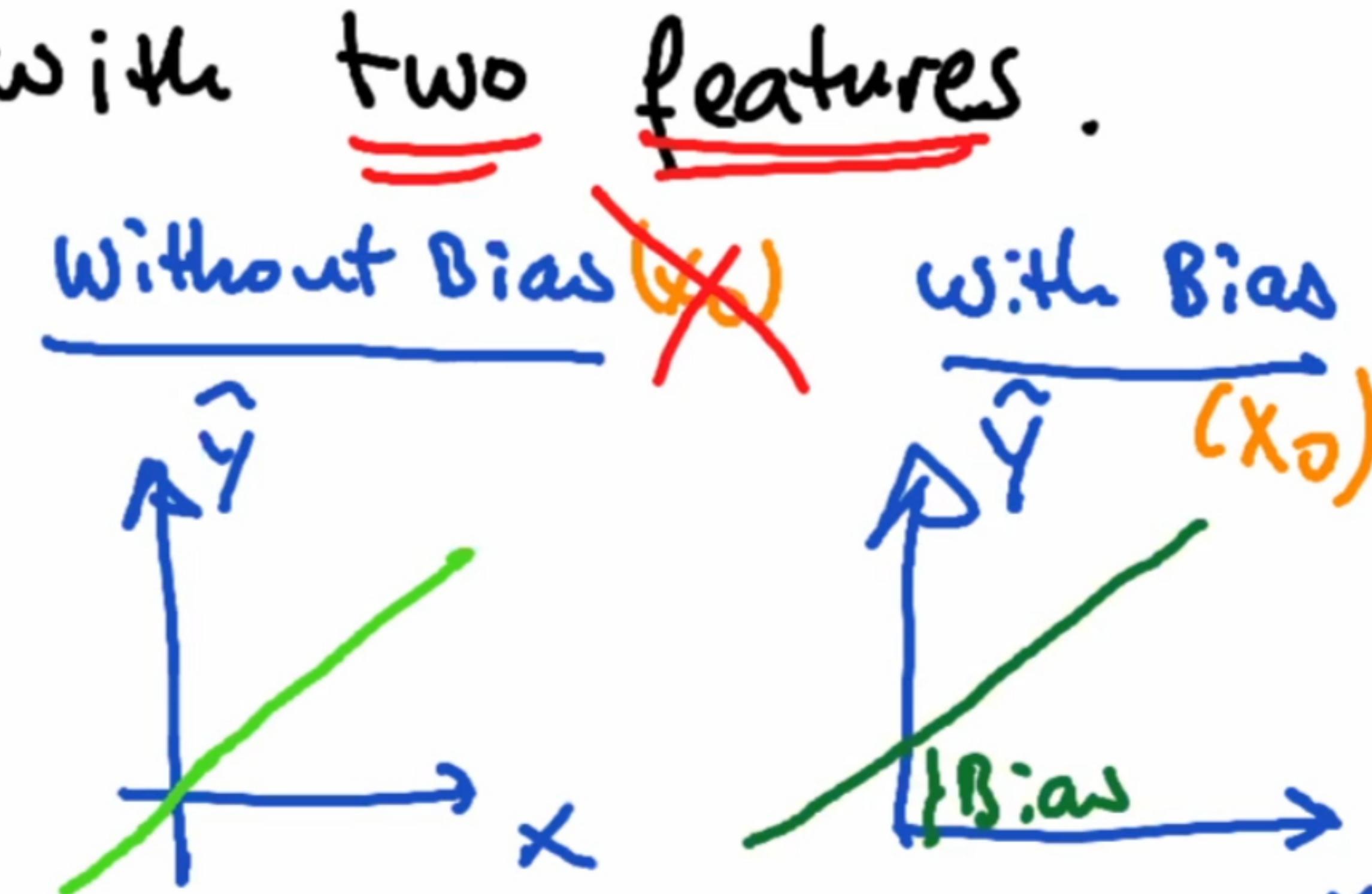
Indicator func: $\mathbb{I}\{x > 0\} = \begin{cases} 0, & x \leq 0 \\ 1, & x > 0 \end{cases}$

Example: A FFNN with two layers with two features.

Features



Input layer Output layer



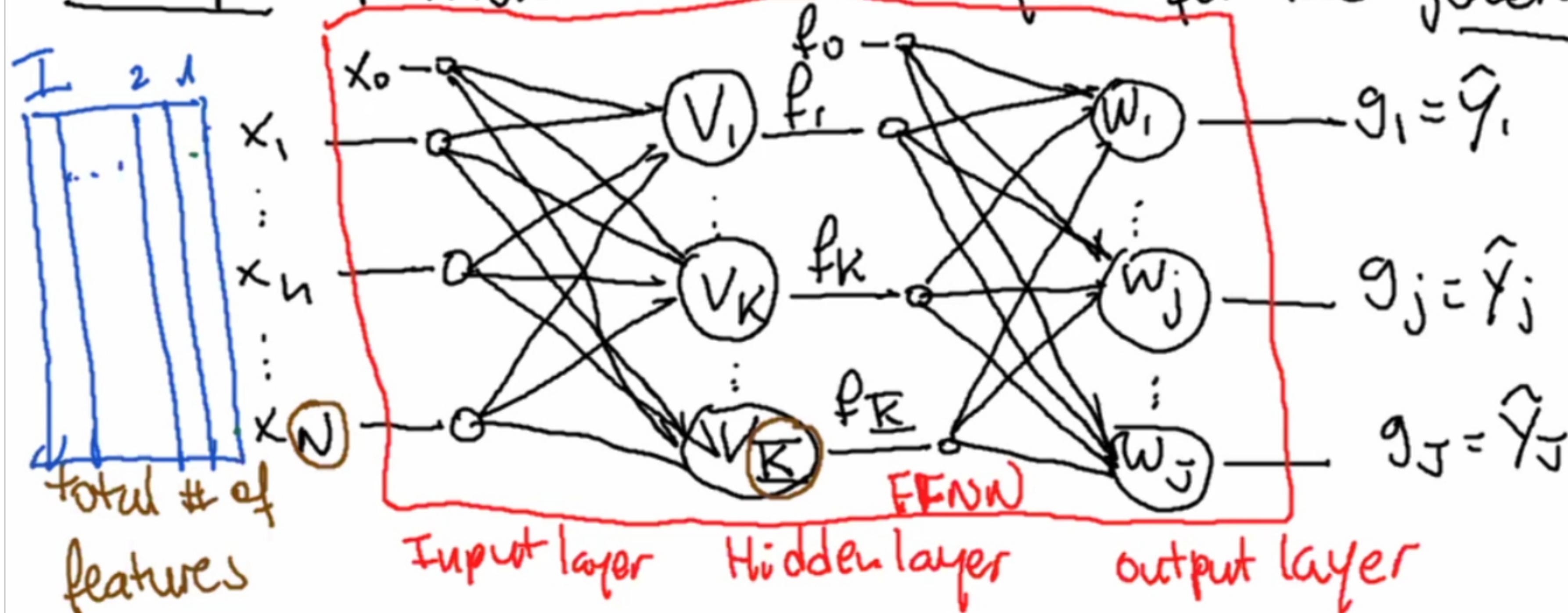
$$W = \begin{bmatrix} 1 & 1 \\ w_1 & w_2 \\ 1 & 1 \end{bmatrix}$$

total # of features
 $N+1$
 3×2

$$\hat{y}_1 = \begin{bmatrix} w_{01} & w_{02} \\ w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$$

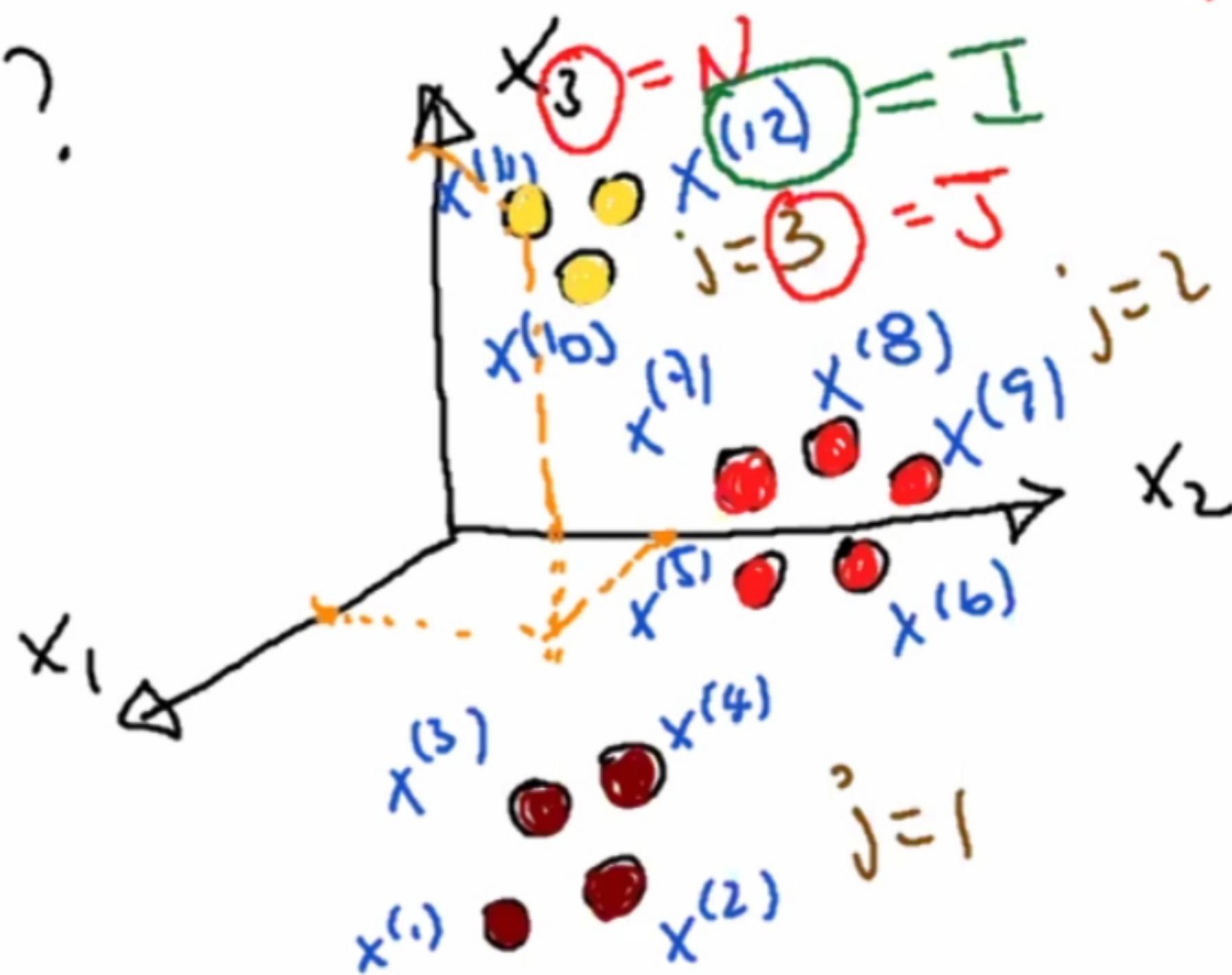
output-variable dimension

Example: A FFNN with three layers for the general case:



- ✓ $N \triangleq$ total # of features
- ✓ $K \triangleq$ total # of hidden neurons
- ✓ $J \triangleq$ output-variable dimension
- ✓ $I \triangleq$ total # of training examples
or By definition

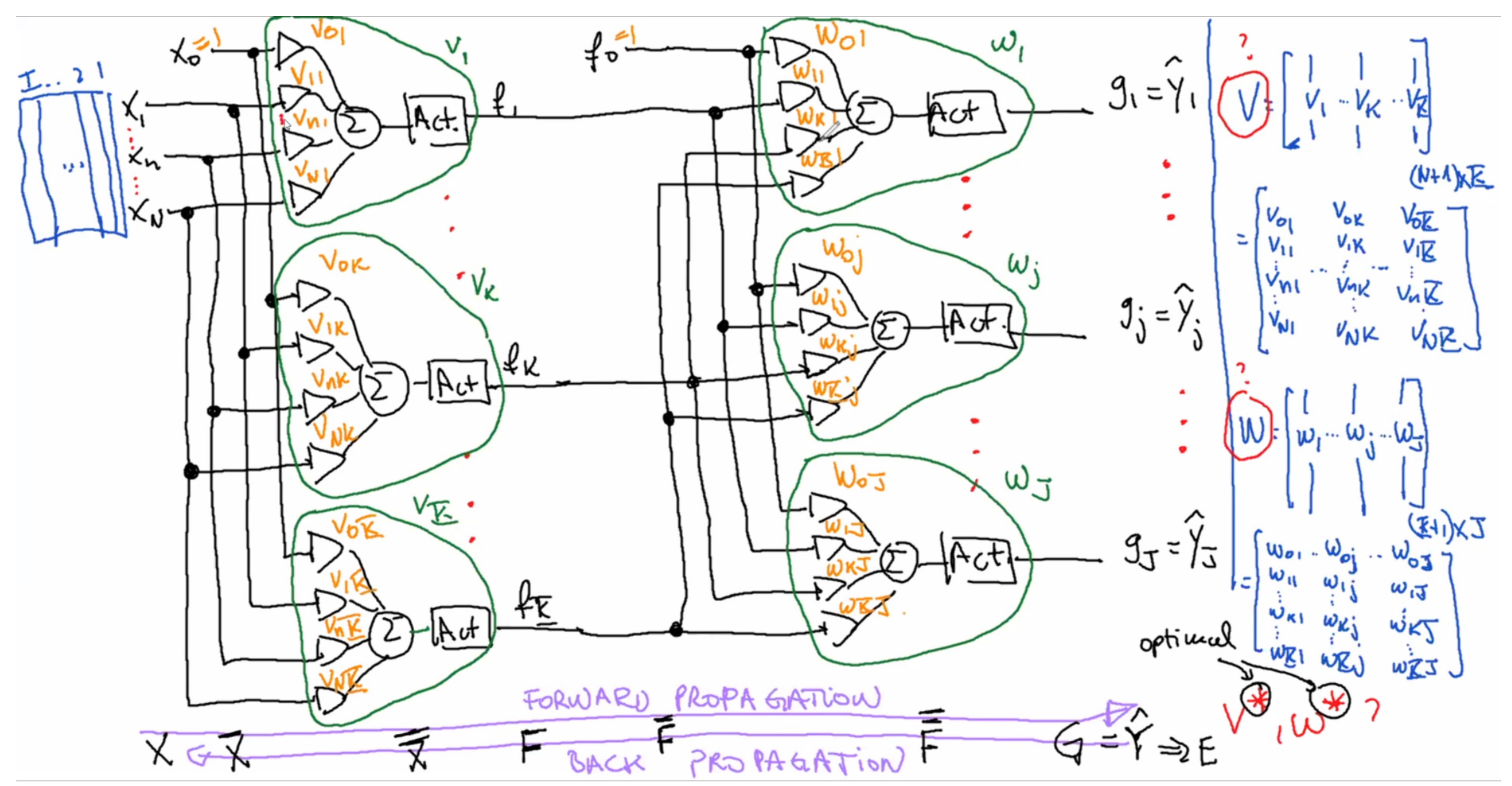
$I ?$



$$x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \end{bmatrix}$$

$$i \in \{1, 2, \dots, I\}$$

$$I = 12$$



Goal of FFNN:

Define $\theta = (V, W)$. Then,

$$\theta^* = (V^*, W^*) = \underset{\theta}{\operatorname{argmin}} E(x, y; \theta) ?$$

Example:

$$Z = \begin{bmatrix} \text{index} & \text{value} \\ 1 & 3 \\ 2 & 2 \\ 3 & 1 \end{bmatrix} \quad \min Z = 1$$

$$\operatorname{argmin} Z = 3$$

Def: Error $E = \frac{1}{2} \sum_{i=1}^I (\|\hat{y}^{(i)} - y^{(i)}\|_2^2) \triangleq \text{sum of squared errors (SSE)}$

Example: 2-norm $\hat{y}^{(i)} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, y^{(i)} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \|\hat{y}^{(i)} - y^{(i)}\|_2 = \sqrt{(1-3)^2 + (2-4)^2}$

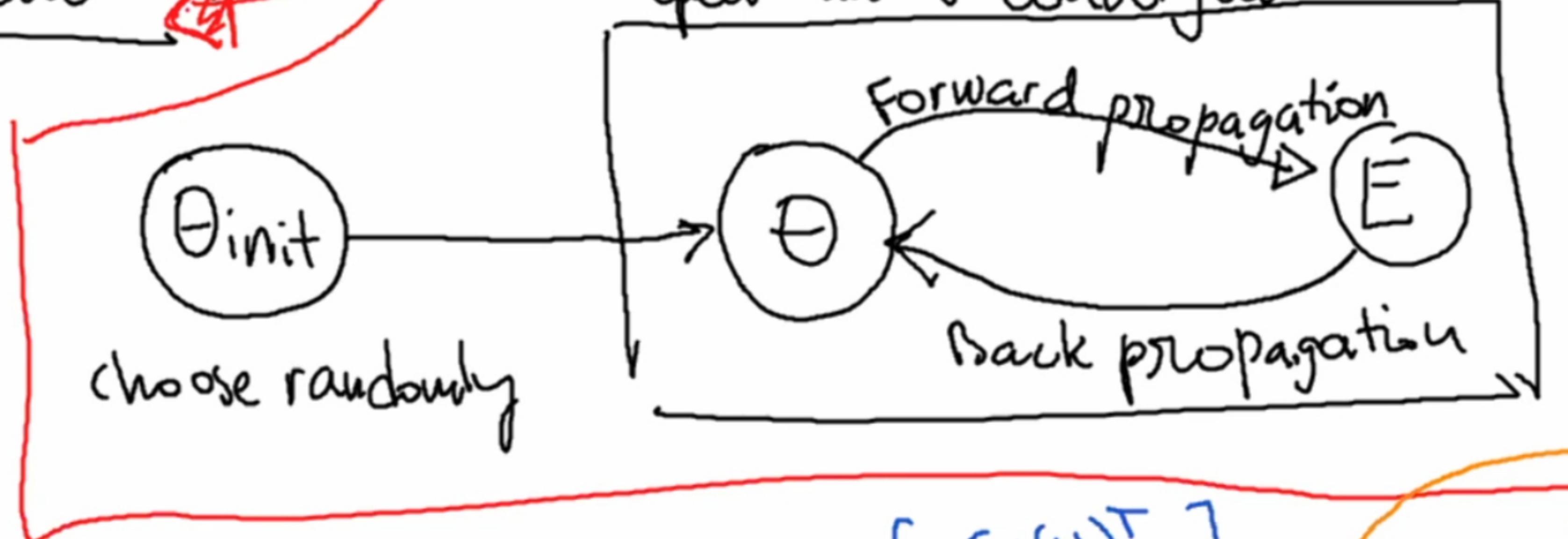
$$= \frac{1}{2} \sum_{i=1}^I \sum_{j=1}^J (\hat{y}_j^{(i)} - y_j^{(i)})^2$$

$\hat{y} \triangleq \text{Estimated output (FFNN)}$

$y \triangleq \text{Actual output}$

~~MXE~~

Approach



$\theta^* \Rightarrow \hat{Y}$

elementwise inversion

Ex. $A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}; A^{-1} = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ -\frac{1}{4} & \frac{1}{5} & \frac{1}{6} \end{bmatrix}$

$X^{(i)} = \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_N^{(i)} \end{bmatrix}; (X^{(i)})^T = \begin{bmatrix} x_1^{(i)} & \dots & x_N^{(i)} \end{bmatrix}$

Forward propagation:

$\bar{X} = \begin{bmatrix} X \\ \vdots \\ 1 \end{bmatrix}; \bar{X} = \begin{bmatrix} (\bar{x}^{(1)})^T \\ \vdots \\ (\bar{x}^{(N+1)})^T \end{bmatrix}; \bar{X} = \bar{X}_{I \times K} \cdot V_{(N+1) \times B} = \begin{bmatrix} (\bar{x}^{(1)})^T v \\ \vdots \\ (\bar{x}^{(B+1)})^T v \end{bmatrix}; F = \begin{bmatrix} 1 & F \\ \vdots & \vdots \\ 1 & F \end{bmatrix} = \begin{bmatrix} (\bar{f}^{(1)})^T \\ \vdots \\ (\bar{f}^{(B+1)})^T \end{bmatrix}; \bar{F} = \bar{F}_{I \times J} \cdot W_{(B+1) \times J} = \begin{bmatrix} (\bar{f}^{(1)})^T w \\ \vdots \\ (\bar{f}^{(B+1)})^T w \end{bmatrix}; G = \left(1 + \exp(-\bar{F}) \right)^{-1} = \hat{Y}$

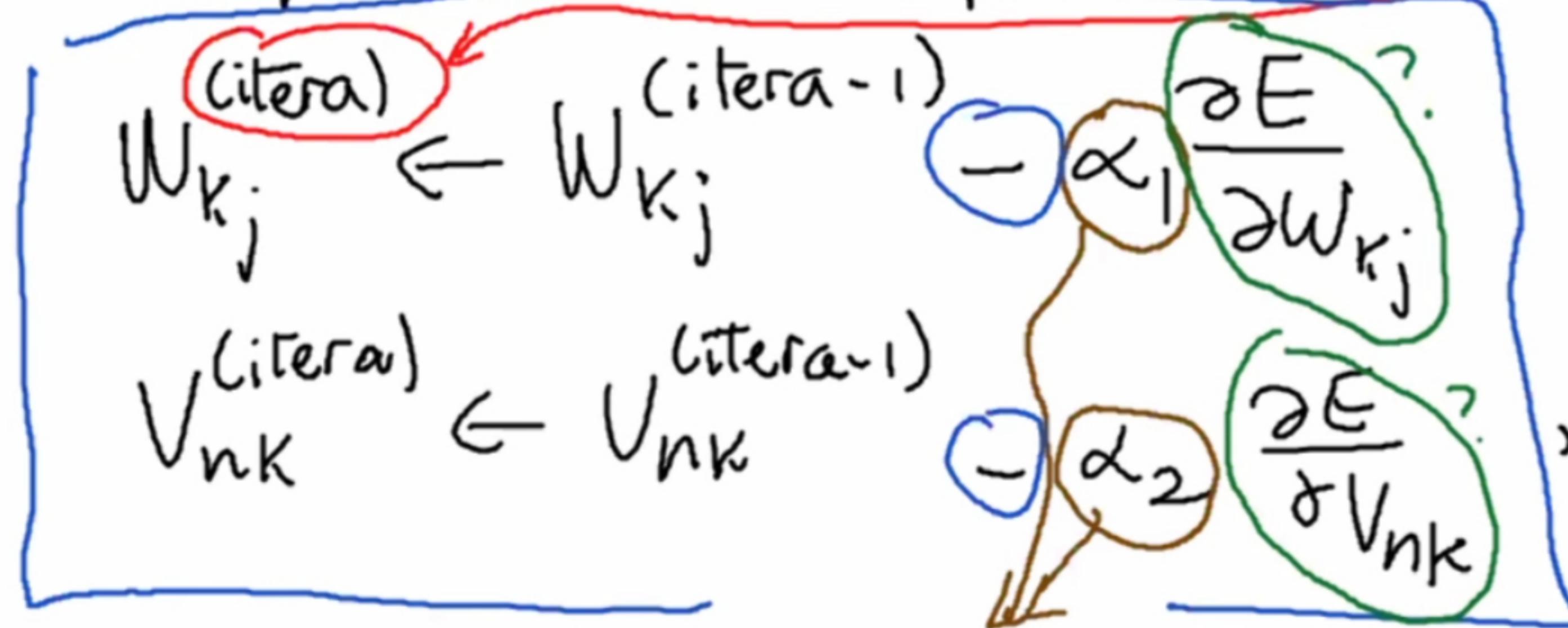
$E = \frac{1}{2} \sum_{i=1}^I \sum_{j=1}^J (\hat{y}_j^{(i)} - y_j^{(i)})^2$

sigmoid $\{ \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}_{I \times K} \cdot -1 \}$

$F = (1 + \exp(-\bar{F}))^{-1}$

Back propagation

Goal: Update the model parameters



Batch Gradient Descent learning rate
"Group" (BGD)
"0,01" or "0,001"

$i_{\text{itera}} = \{1, 2, \dots, I_{\text{itera-max}}\}$

$k = \{0, 1, \dots, K\}, j = \{1, 2, \dots, J\}$

~~$n = \{0, 1, \dots, N\}, k = \{1, 2, \dots, K\}$~~

$$\begin{aligned} \frac{\partial E}{\partial w_{kj}} &= \frac{\partial}{\partial w_{kj}} \frac{1}{2} \sum_{i=1}^I \sum_{j'=1}^J (\hat{y}_{j'}^{(i)} - y_{j'}^{(i)})^2 = \frac{\partial}{\partial w_{kj}} \frac{1}{2} \sum_{i=1}^I \sum_{j'=1}^J (g_{j'}^{(i)} - y_{j'}^{(i)})^2 \\ &= \sum_{i=1}^I \frac{1}{2} \frac{\partial}{\partial w_{kj}} \sum_{j'=1}^J (g_{j'}^{(i)} - y_{j'}^{(i)})^2 = \sum_{i=1}^I \cancel{\frac{1}{2}} \cancel{x} (g_j^{(i)} - y_j^{(i)}) \frac{\partial g_j^{(i)}}{\partial w_{kj}} = \sum_{i=1}^I (g_j^{(i)} - y_j^{(i)}) g_j^{(i)} (1 - g_j^{(i)}) \bar{f}_k^{(i)} \end{aligned}$$

$$\frac{\partial}{\partial w_{kj}} [(g_1^{(i)} - y_1^{(i)})^2 + \dots + (g_j^{(i)} - y_j^{(i)})^2 + \dots + (g_J^{(i)} - y_J^{(i)})^2]$$

$$\frac{\partial g_j^{(i)}}{\partial w_{kj}} = \frac{\partial}{\partial w_{kj}} \left[\frac{e^{-\bar{f}_k^{(i)} w_j}}{1 + e^{-\bar{f}_k^{(i)} w_j}} \right] = -e^{-\bar{f}_k^{(i)} w_j} \frac{\frac{\partial}{\partial w_{kj}} (-\bar{f}_k^{(i)} w_j)}{(1 + e^{-\bar{f}_k^{(i)} w_j})^2}$$

$$(\gamma')' = \frac{-\gamma'}{\gamma'^2}$$

$$= \frac{1}{1 + e^{-(\bar{f}_k^{(i)})^\top w_j}} \left(1 - \frac{1}{1 + e^{-(\bar{f}_k^{(i)})^\top w_j}} \right)$$

$$\bar{f}_k^{(i)} = g_j^{(i)} (1 - g_j^{(i)}) \bar{f}_k^{(i)}$$

$$\frac{\partial E}{\partial v_{nk}} = \frac{\partial}{\partial v_{nk}} \left[\frac{1}{2} \sum_{i=1}^I \sum_{j=1}^J (g_j^{(i)} - y_j^{(i)})^2 \right] = \sum_{i=1}^I \sum_{j=1}^J \frac{1}{2} \frac{\partial}{\partial v_{nk}} (g_j^{(i)} - y_j^{(i)})^2 = \sum_{i=1}^I \sum_{j=1}^J (g_j^{(i)} - y_j^{(i)}) \frac{\partial g_j^{(i)}}{\partial v_{nk}}$$

$$\frac{\partial g_j^{(i)}}{\partial v_{nk}} = \frac{\partial g_j^{(i)}}{\partial \bar{f}_k^{(i)}} \frac{\partial \bar{f}_k^{(i)}}{\partial v_{nk}}$$

⊗ ⊗

$$\begin{aligned} \textcircled{R} \quad \frac{\partial g_j^{(i)}}{\partial \bar{f}_k^{(i)}} &= \frac{\partial}{\partial \bar{f}_k^{(i)}} \left[\left(1 + e^{-\bar{f}_k^{(i) T} w_j} \right)^{-1} \right] \\ &= -e^{-\bar{f}_k^{(i) T} w_j} (-w_{kj}) \\ &= -\frac{e^{-\bar{f}_k^{(i) T} w_j} (-w_{kj})}{(1 + e^{-\bar{f}_k^{(i) T} w_j})^2} \\ &= \frac{1}{1 + e^{-\bar{f}_k^{(i) T} w_j}} \left(1 - \frac{1}{1 + e^{-\bar{f}_k^{(i) T} w_j}} \right) w_{kj} \\ &= g_j^{(i)} (1 - g_j^{(i)}) w_{kj} = f_k^{(i)} \left(1 - f_k^{(i)} \right) \bar{x}_n^{(i)} \end{aligned}$$

Pseudo-code

$$\alpha_1 = \alpha_2 = 0.001$$

Initialize V, W randomly

for $\text{itera} = \{1, 2, \dots, \text{Itera_max}\}$:

FWP : $X, \bar{X}, \tilde{X}, F, \bar{F}, \tilde{F}, G, E$

BP :

$$w_{kj}^{(\text{itera})} \leftarrow w_{kj}^{(\text{itera})} - \alpha_1 \frac{\partial E}{\partial w_{kj}}$$

$$v_{nk}^{(\text{itera})} \leftarrow v_{nk}^{(\text{itera})} - \alpha_2 \frac{\partial E}{\partial v_{nk}}$$

$$k = \{0, 1, \dots, K\}, j = \{1, 2, \dots, J\}$$

$$n = \{0, 1, \dots, N\}, k = \{1, 2, \dots, K\}$$

$$\Rightarrow V^*, W^*$$