

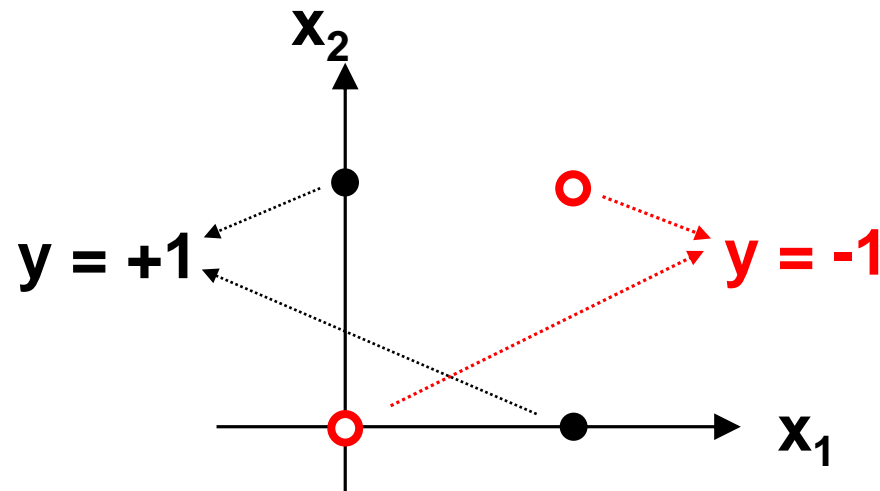
Artificial Neural Network for Classification (2)

Hanwool Jeong

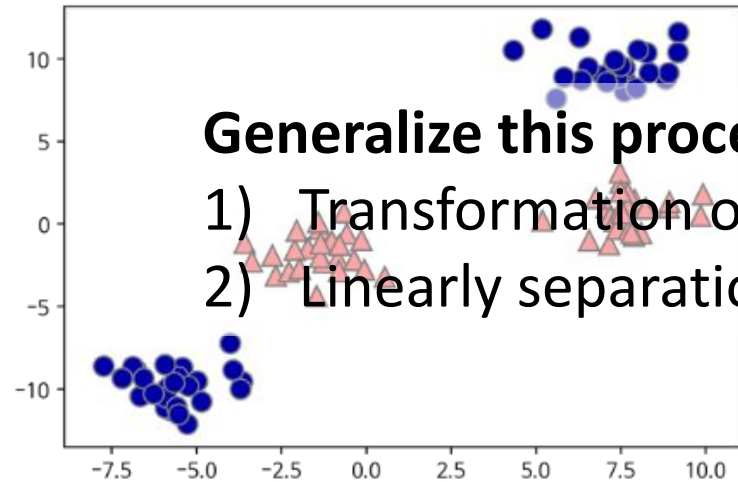
hwjeong@kw.ac.kr

Can Perceptron Solve XOR Problem?

- Linearly non-separable dataset cannot be classified by the perceptron [Minsky, 1969].
- However, in 1974, Werbos's dissertation says linearly non-separable dataset can be classified by the perceptron through multi layers.
➔ That is multi layer perceptron can solve XOR problem.

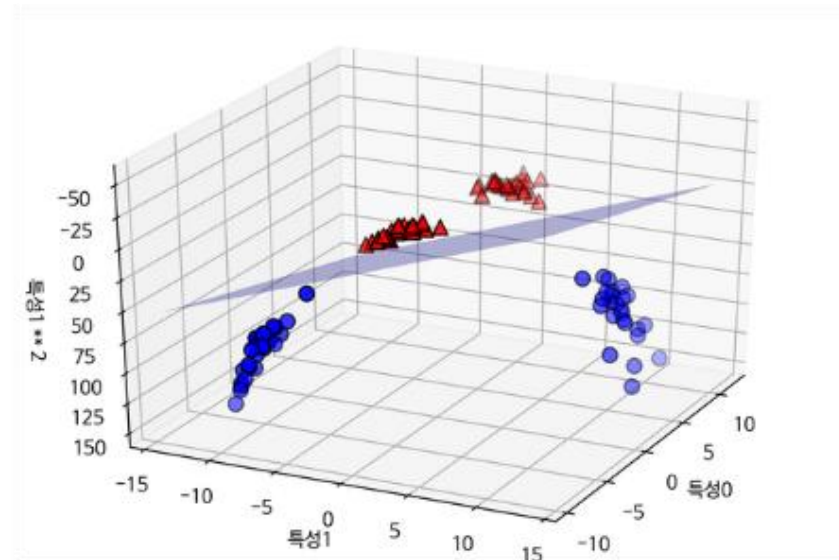
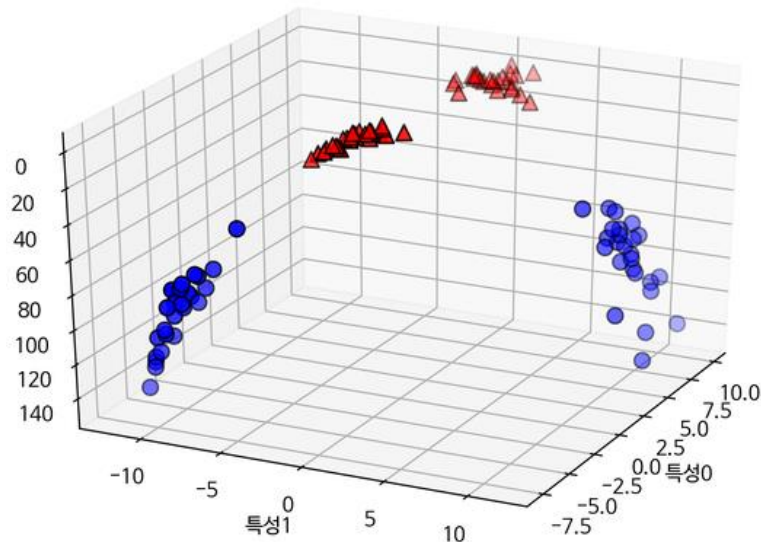
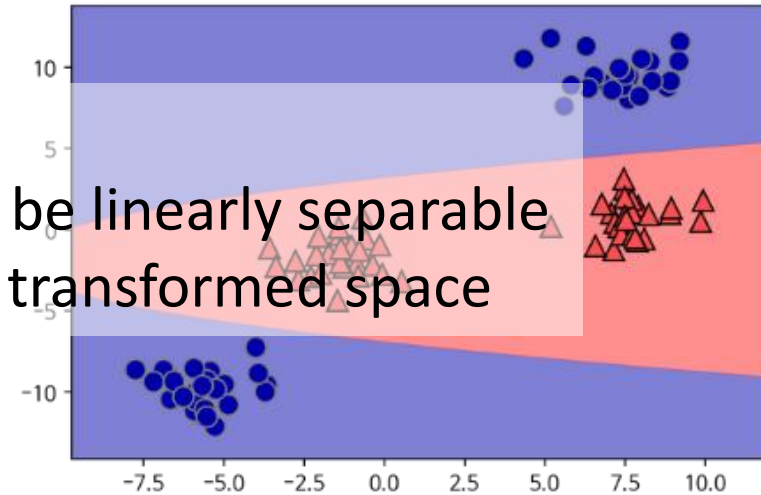


Revisit: How Do We Handle Linearly inseparable Data in SVM?



Generalize this procedure!

- 1) Transformation of data to be linearly separable
- 2) Linearly separation in the transformed space



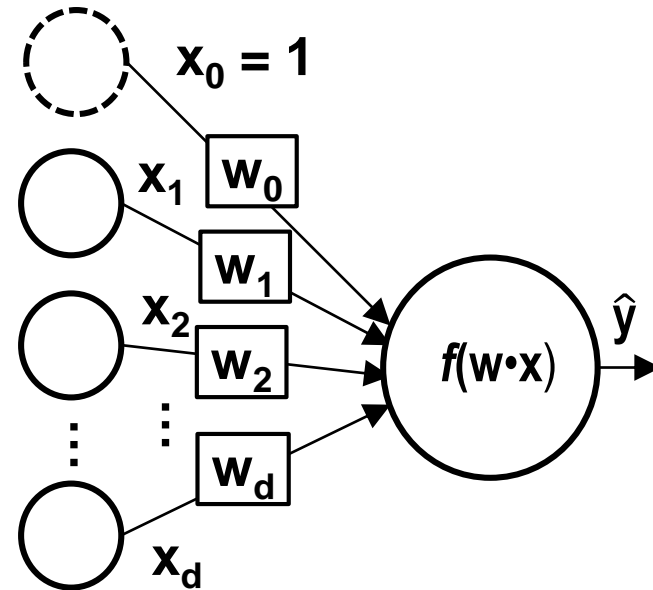
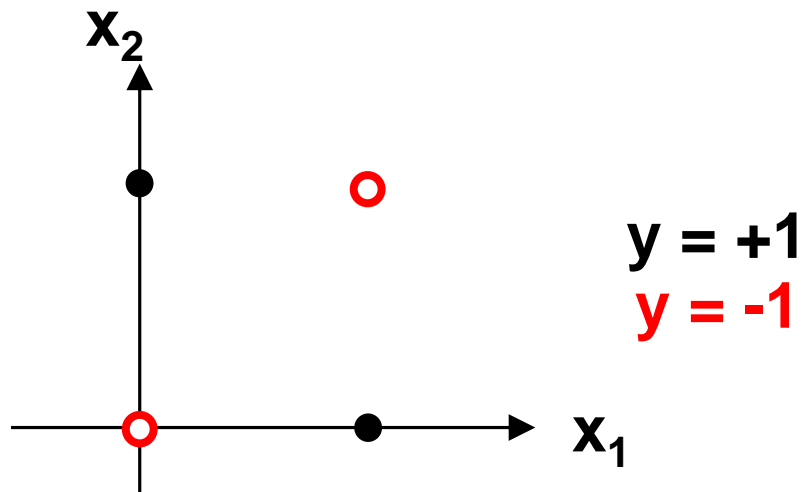
Ideation!

Can data be **transformed** using linear operation (=inner product) followed by non-linear sign function?

→ Let's think over how perceptron classifies data

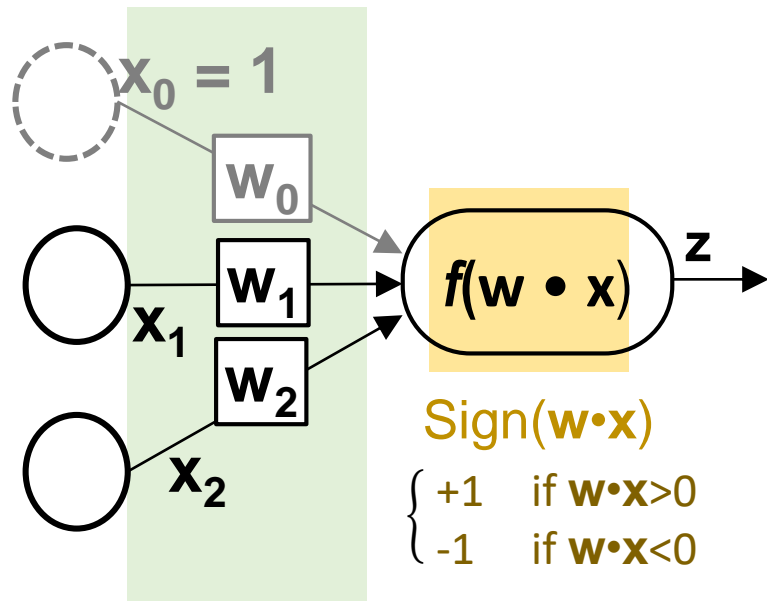
Generalize this procedure!

- 1) Transformation of data to be linearly separable
- 2) Linearly separation in the transformed space



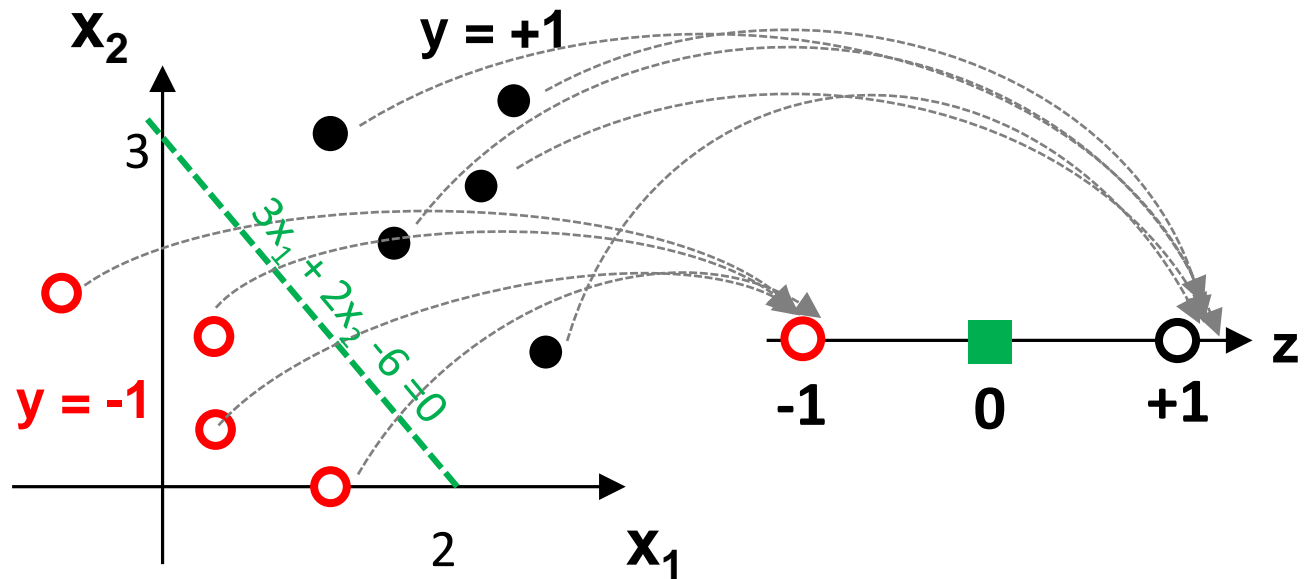
You Should Feel What Perceptron Can Do

- It is nothing but **data transformation!**



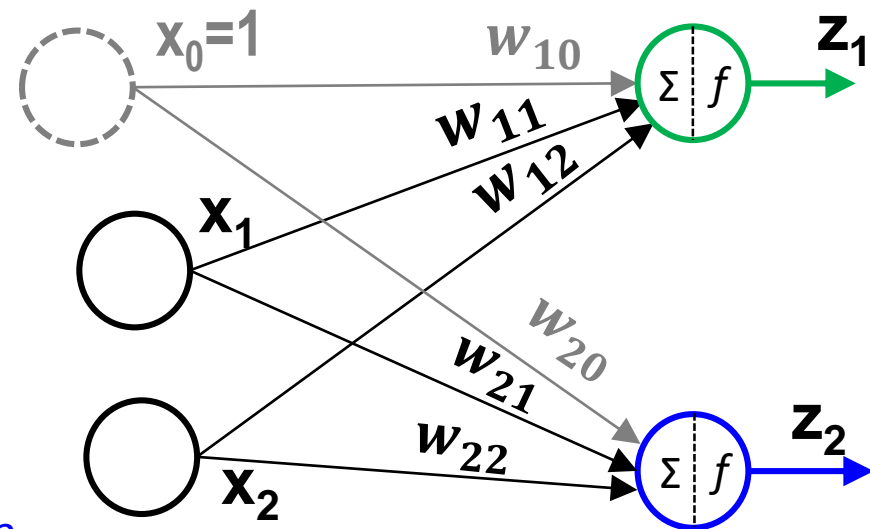
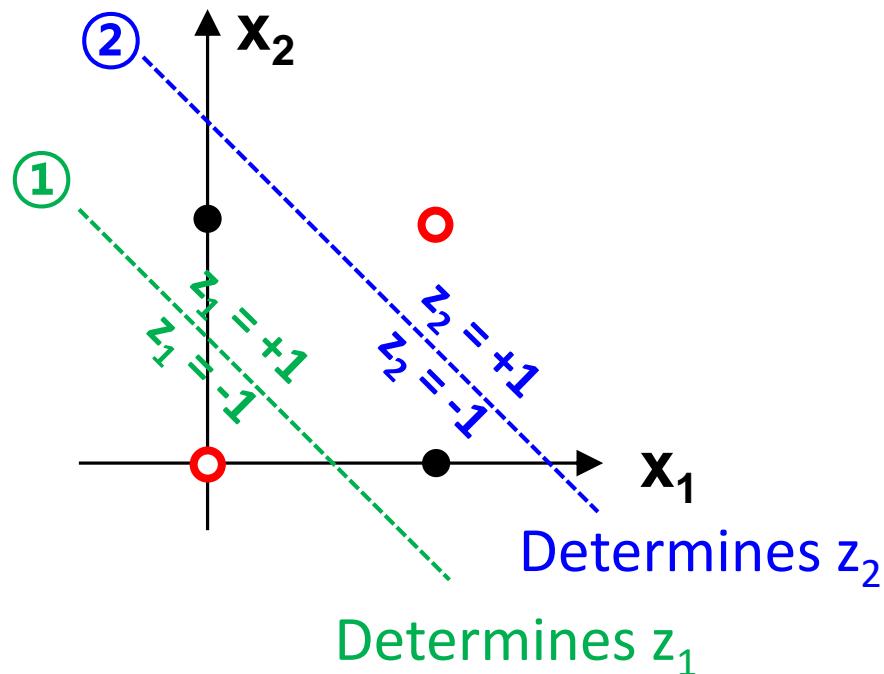
$$[w_0, w_1, w_2] = [-6, 3, 2]$$

$$w \cdot x = 3x_1 + 2x_2 - 6$$

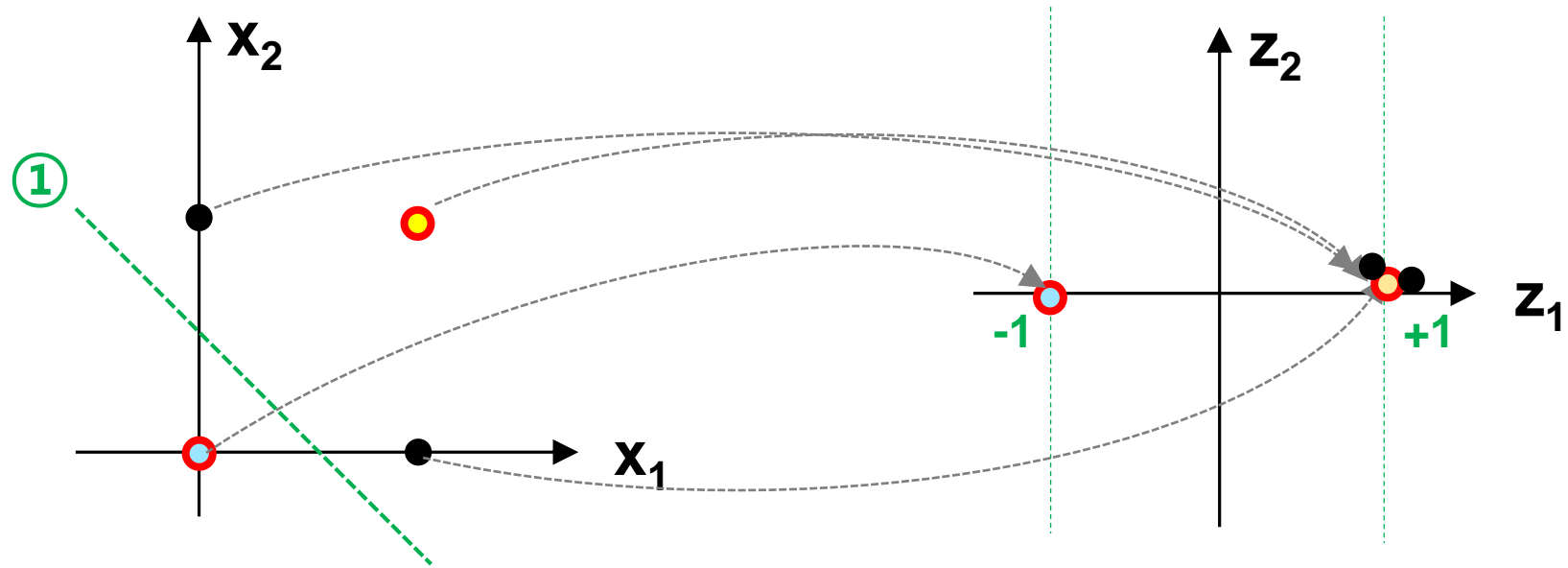


Transformation of Dataset

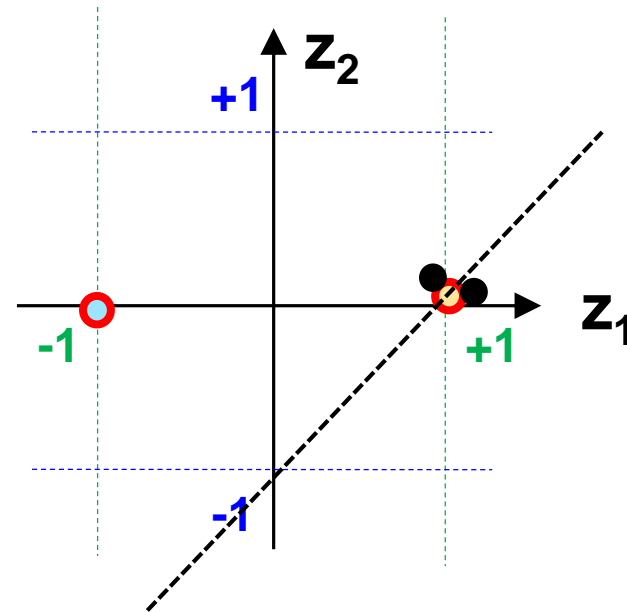
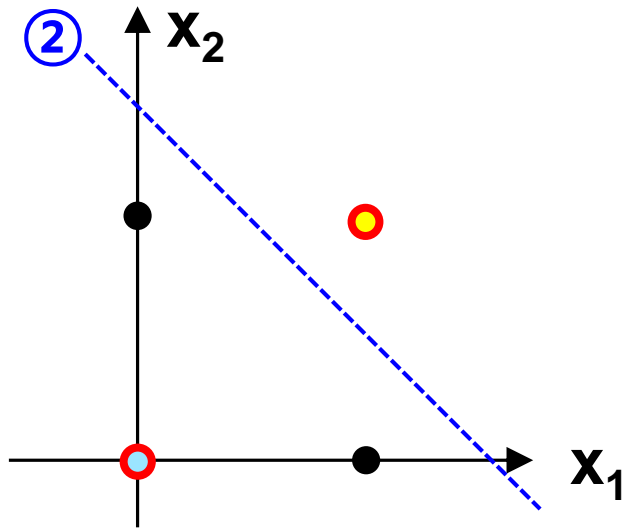
- Key idea is to transfer the data into new space.
- x_1 - x_2 space \rightarrow z_1 - z_2 space
- Linearly non-separable \rightarrow Linearly separable?



Step1 : Determines z_1

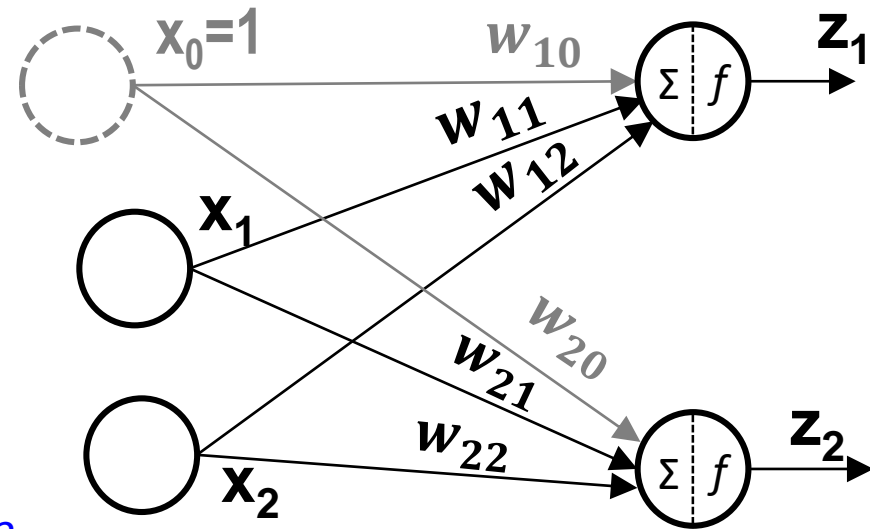
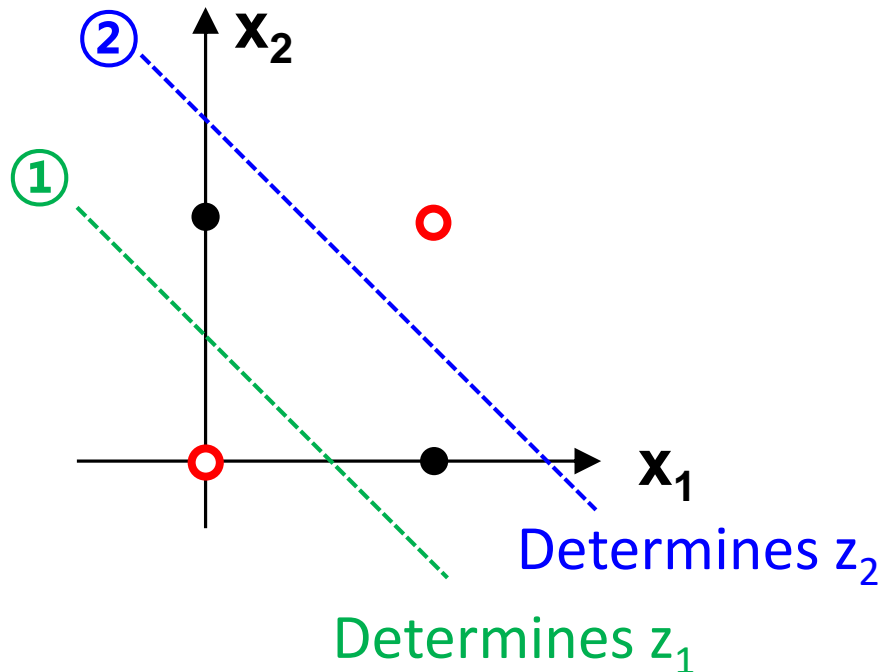


Step2 : Determines z_2



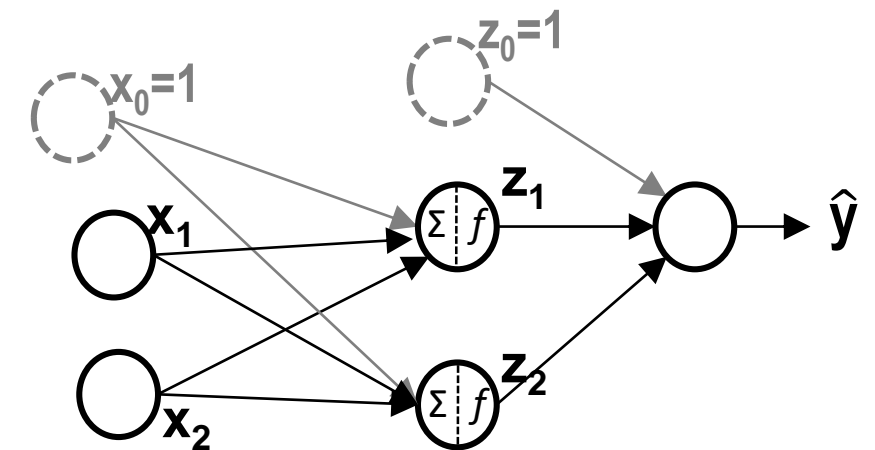
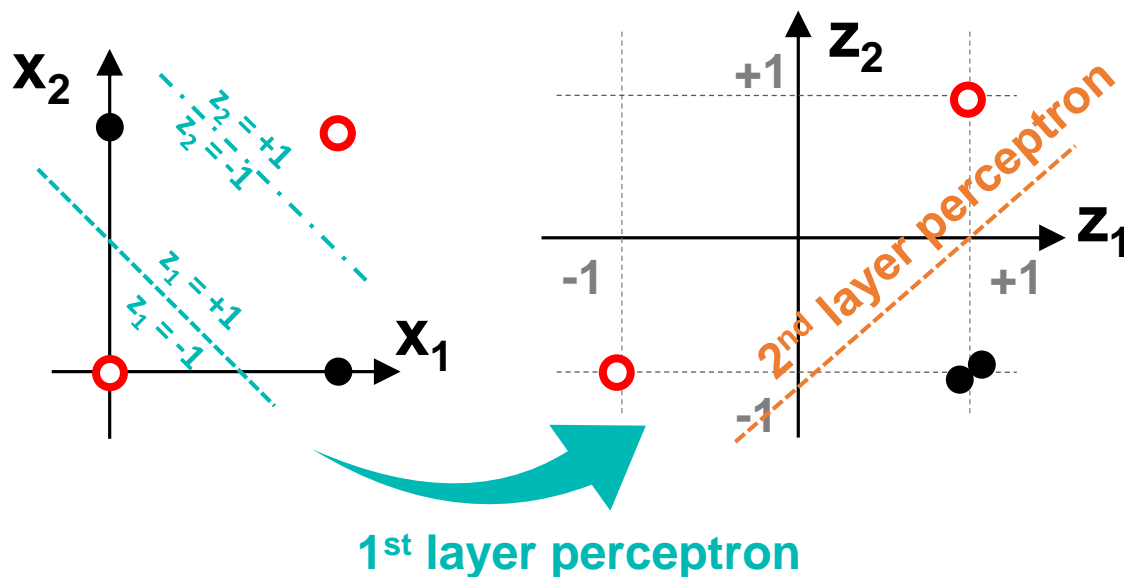
Can you Determine W Matrix?

- Saying $W = \begin{bmatrix} W_{10} & W_{11} & W_{12} \\ W_{20} & W_{21} & W_{22} \end{bmatrix}$,
Bias



Multi Layer Perceptron (MLP)

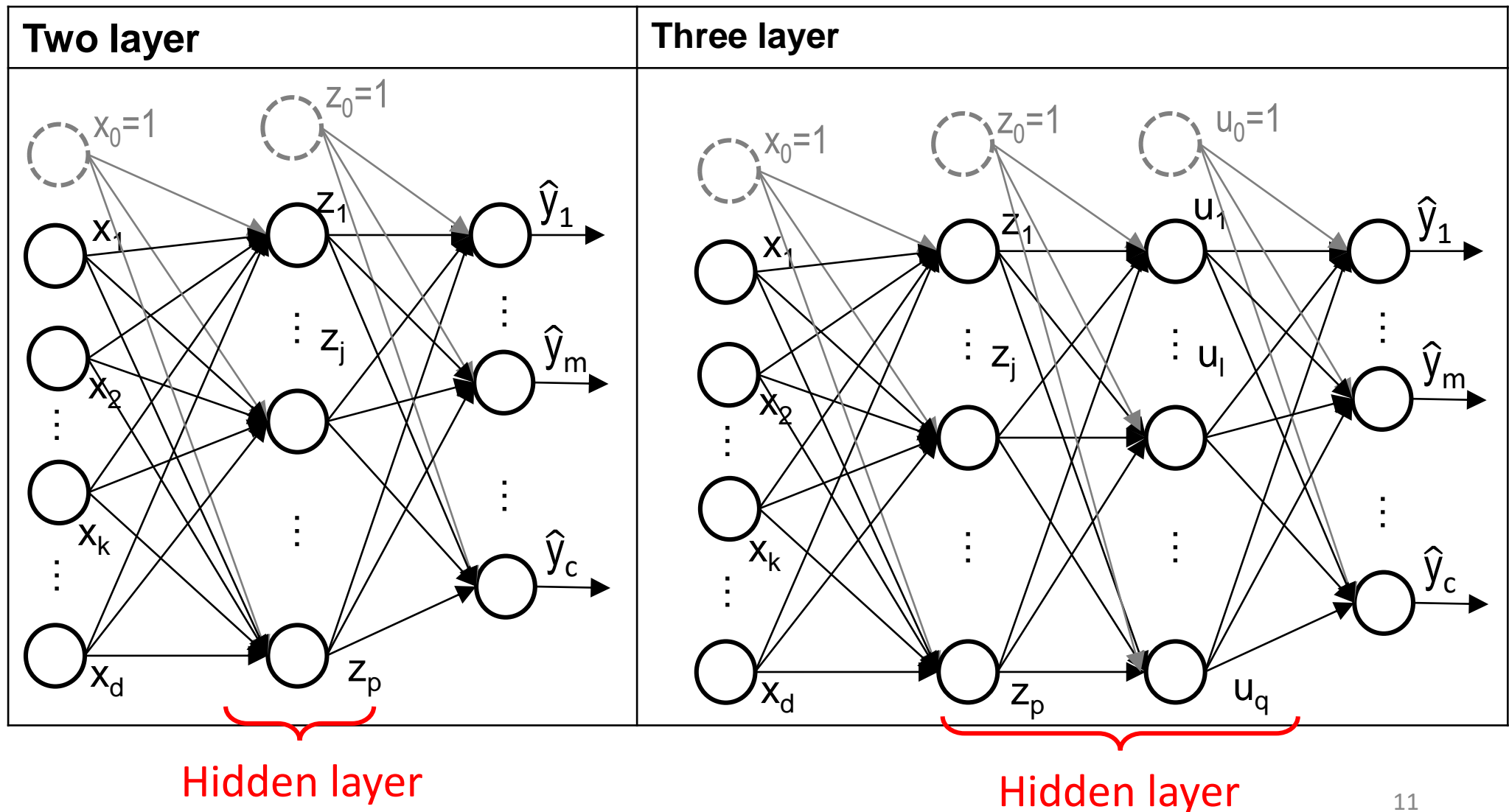
- After became linearly separable in z_1 - z_2 plane, we can do the same what we did previously, forming **MLP**



→ Successfully classifying XOR data

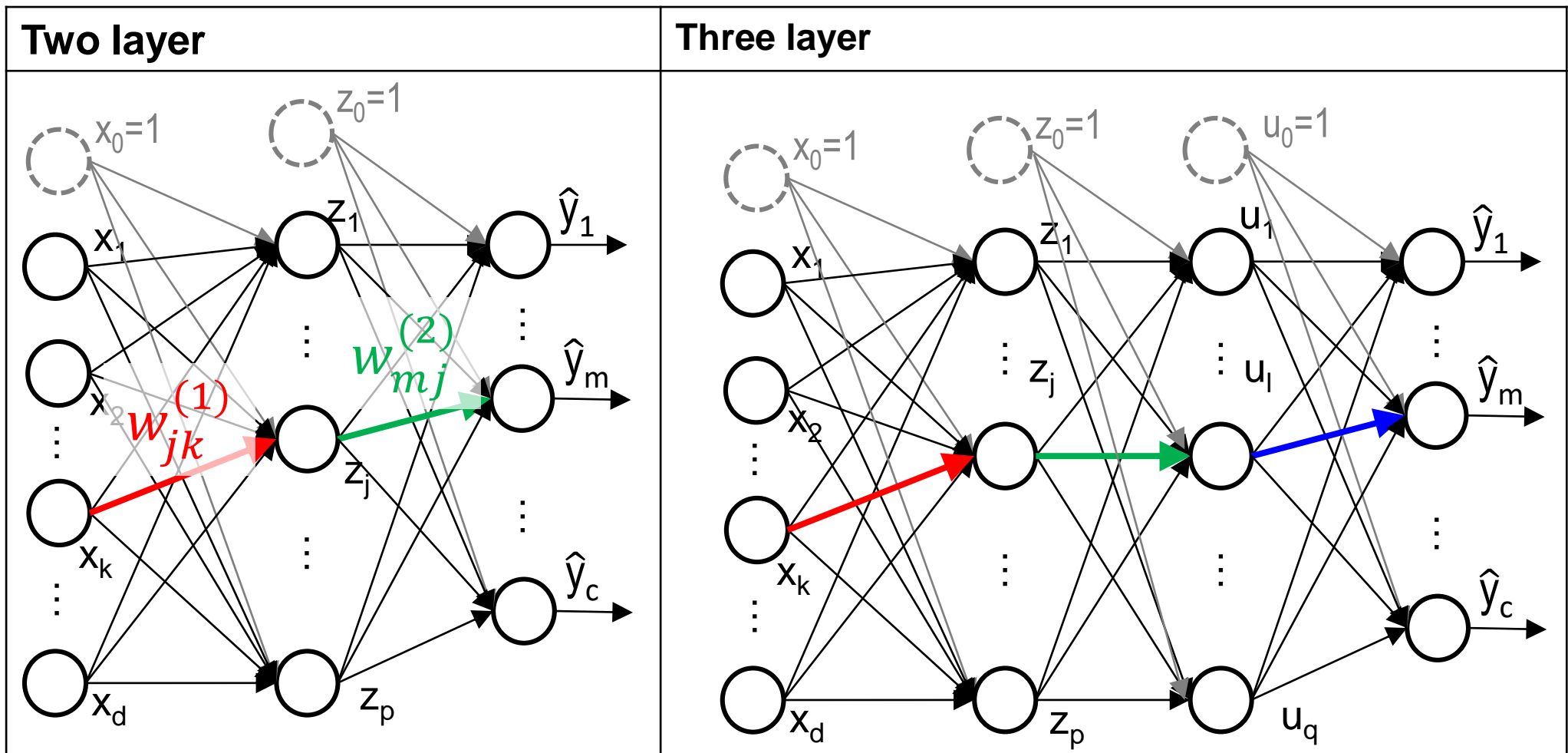
General Structure of MLP

- Please note how the final output configuration looks like



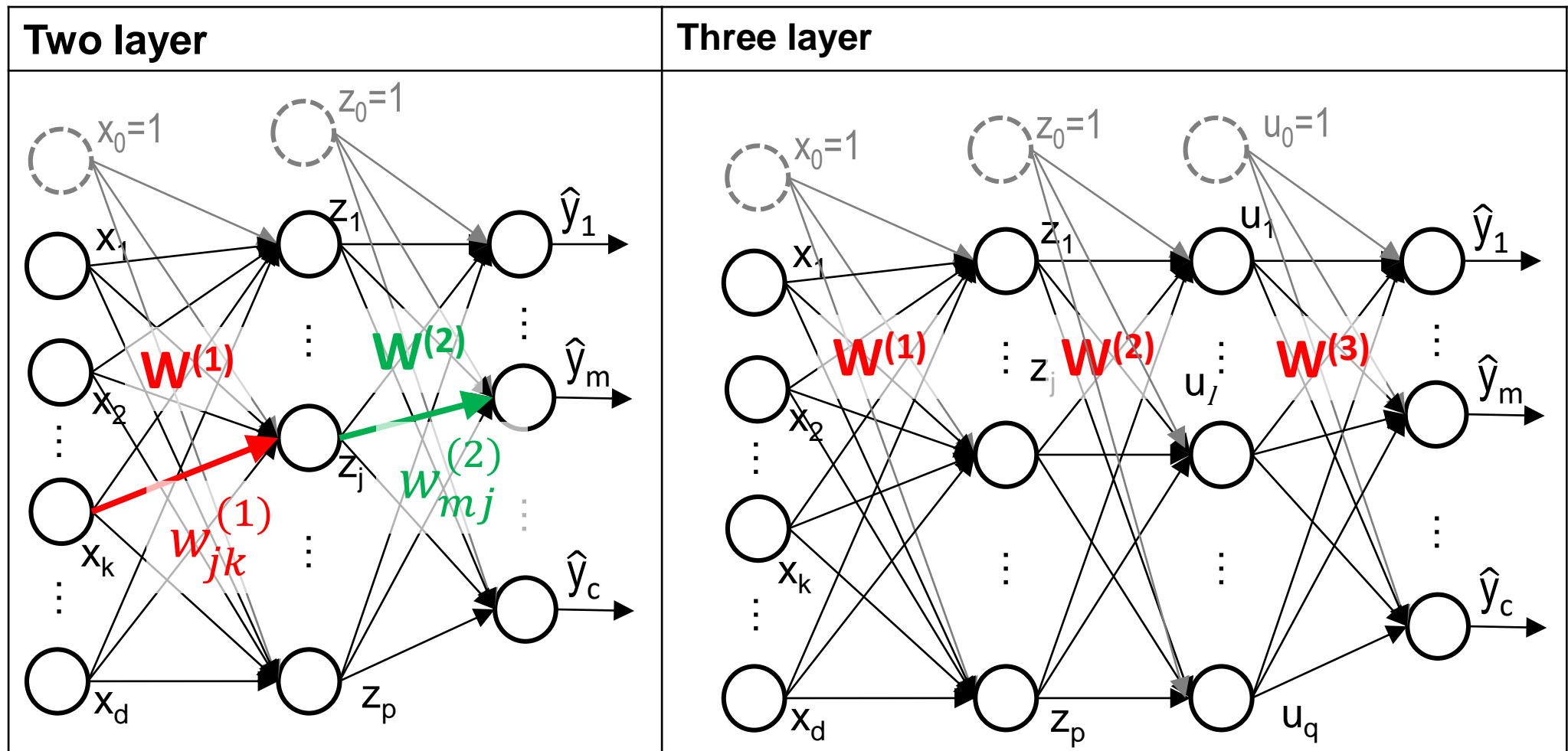
Weight Component Representation

- Please note how the final output configuration looks like

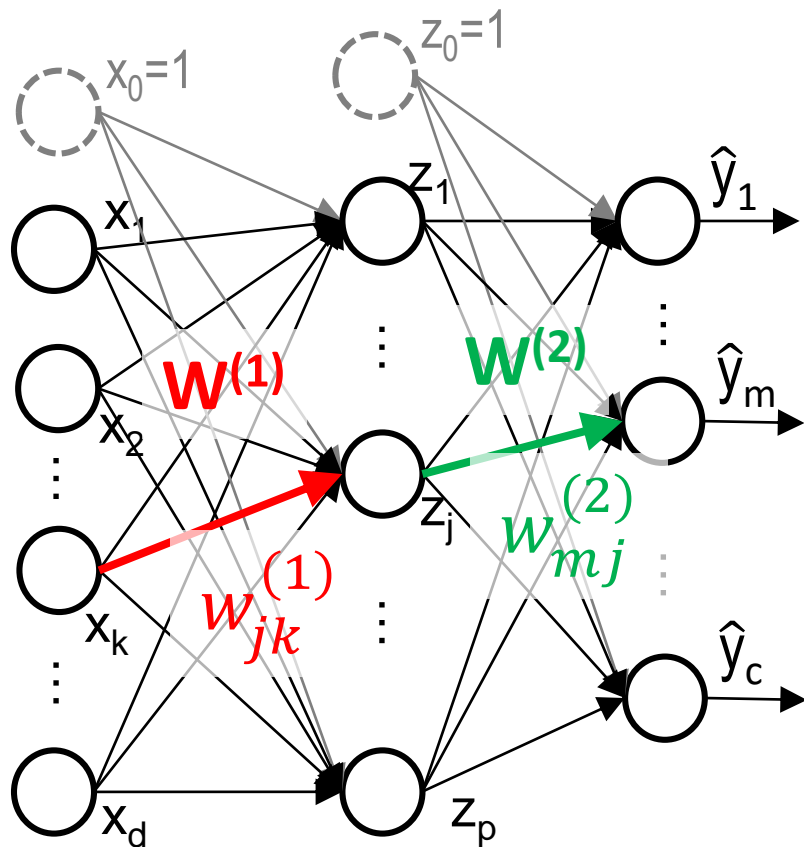


Weight Matrix Representation

- Please note how the final output configuration looks like



Form of W Matrix



[

$$\begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_k \\ \dots \\ x_d \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ \dots \\ z_j \\ \dots \\ z_p \end{bmatrix}$$

From now on, you just simply think, “Oh, There is W matrix which can transform \vec{x} into \vec{z} .”

Operation Representation in MLP

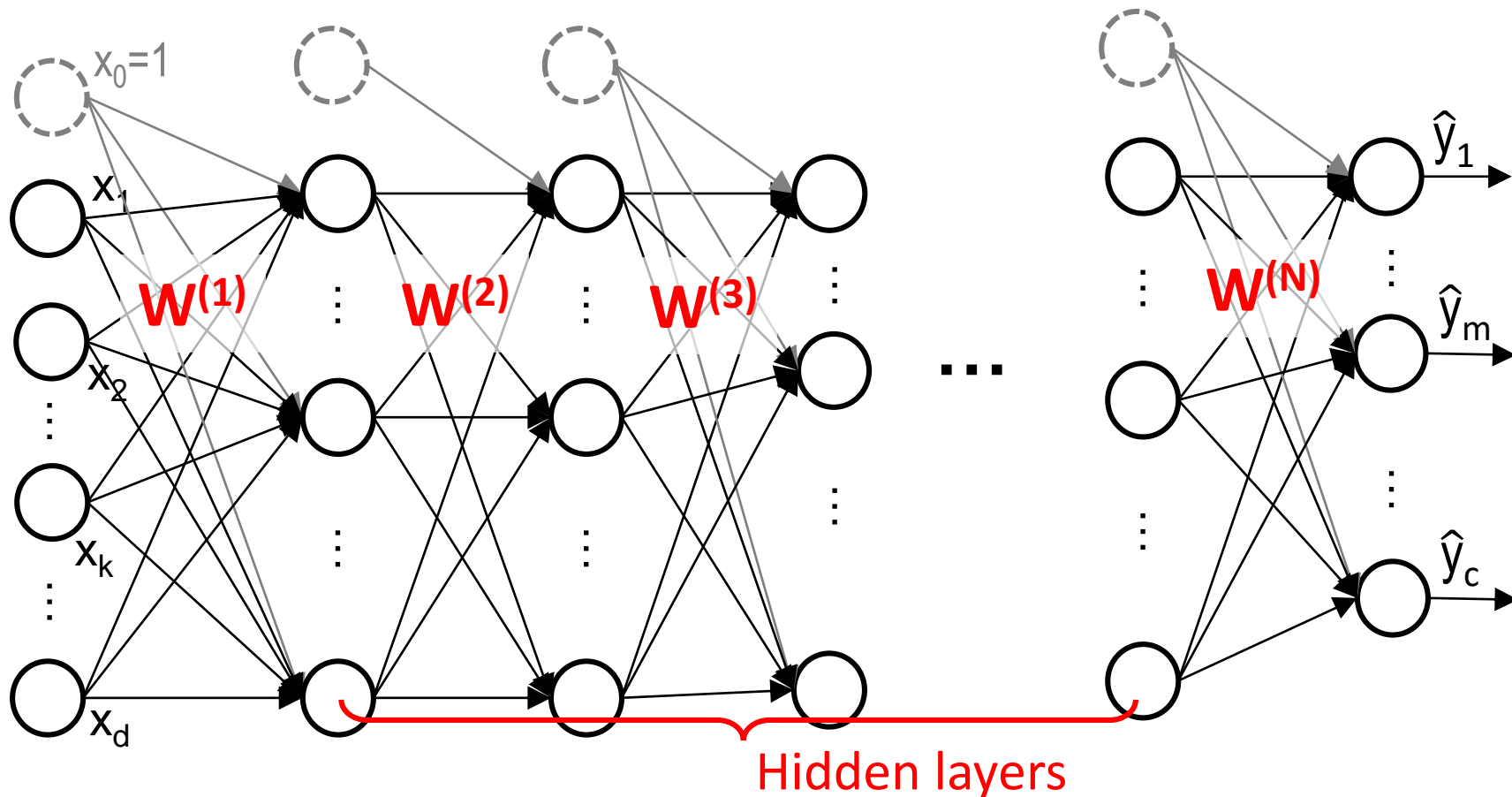
- For example, $W^{(1)}$ looks like

$$W^{(1)} = \begin{bmatrix} w_{10}^{(1)} & w_{11}^{(1)} & \dots & w_{1d}^{(1)} \\ w_{20}^{(1)} & w_{21}^{(1)} & \dots & w_{2d}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{p0}^{(1)} & w_{p1}^{(1)} & \dots & w_{pd}^{(1)} \end{bmatrix}$$

- For two layer perceptron with activation function of f
1st layer : $z_j = f(w_{j0}x_0 + w_{j1}x_1 + \dots + w_{jk}x_k \dots + w_{jd}x_d) = f(\mathbf{w}_j \bullet \mathbf{x})$
That is, $\mathbf{z} = f(\mathbf{W}^{(1)}\mathbf{x})$
Output (=2nd) layer : $\hat{\mathbf{y}} = f(W^{(2)}\mathbf{z}) = f(W^{(2)}f(\mathbf{W}^{(1)}\mathbf{x}))$
- How about for three layer perceptron?

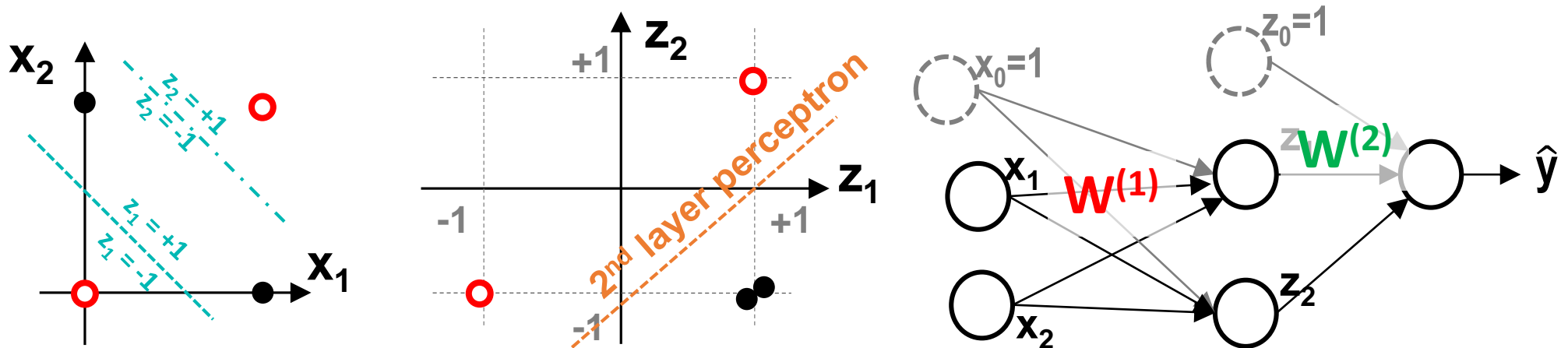
Deep MLP

- More than two layers (3, 4, ...)
 - Not necessarily transform data being linearly separable then classify.
- ➔ **Arbitrary operations are performed in hidden layers.**



Checkpoints

- ✓ Applying perceptron XOR data through MLP
- ✓ Matrix representation for MLP: $\hat{y} = f(W^{(2)}f(W^{(1)}\mathbf{x}))$
- ✓ Deep MLP
- ✓ Coming up next: how can we determine W ? = training!



General Procedure for Training

- 1) Defining cost/loss function $J(\theta)$
- 2) Then minimizing cost function (MLE vs. MAP)
 - MLE cost function is derived based on $p(\mathbf{D}/\theta)$.
 - Some transform is applied for mathematical convenience.
 - For example, $NLL = -\log p(\mathbf{D}/\theta)$ or $(\frac{1}{2})NLL^2, \dots$
 - Gaussian distribution is used for PDF, minimizing the square of Euclidean distance, $\|\hat{\mathbf{y}} - \mathbf{y}\|^2$, becomes the cost function.
 - $\hat{\mathbf{y}}$ =output predicted value by model
 - \mathbf{y} =output given by the training set (answer)
 - After the cost function is given, differentiation or gradient descent can be used.

$$\theta_{\text{next}} = \theta_{\text{present}} - \eta \nabla J(\theta)$$

Cost Function Development in Training for MLP

- Cost function? For each sample, we can consider the following as cost function.

$$\|\hat{\mathbf{y}}_i - \mathbf{y}_i\|^2$$

- How about for all samples? We have the design matrix \mathbf{X} .
- Saying $\hat{\mathbf{Y}} = F(\mathbf{X})$ where F is overall MLP classifier for given \mathbf{Y} , then the cost becomes

$$\|\hat{\mathbf{Y}} - \mathbf{Y}\|^2 = \|F(\mathbf{X}) - \mathbf{Y}\|^2$$

- The training is to find \mathbf{W}^* that minimizes the above. That is,

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \|F(\mathbf{X}) - \mathbf{Y}\|^2$$

- It should be noted that $F(\mathbf{X}) = F(\mathbf{X}|\mathbf{W})$

Applying GD Method

- Defining a cost function $J(\mathbf{W})$ as follows for mathematical convenience.

$$J(\mathbf{W}) = \frac{1}{2} \|F(\mathbf{X}) - \mathbf{Y}\|^2 = \frac{1}{2} \|\hat{\mathbf{Y}} - \mathbf{Y}\|^2$$

- For MLP, the gradient descent to find \mathbf{W}^* is

$$\mathbf{W}^{(1)} = \mathbf{W}^{(1)} - \eta \frac{\partial J}{\partial \mathbf{W}^{(1)}}$$

$$\mathbf{W}^{(2)} = \mathbf{W}^{(2)} - \eta \frac{\partial J}{\partial \mathbf{W}^{(2)}}$$

$$\mathbf{W}^{(3)} = \mathbf{W}^{(3)} - \eta \frac{\partial J}{\partial \mathbf{W}^{(3)}}$$

....

- We should find the expression for the above.

Revisit the Location of $W^{(l)}$

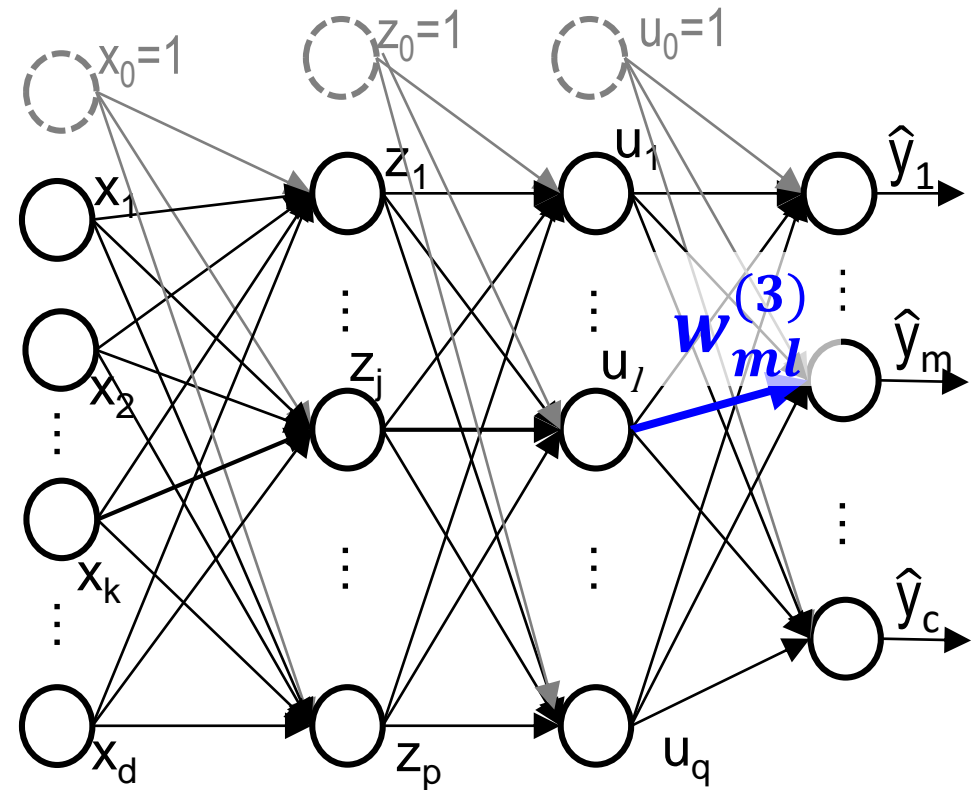
- Which one is easier to find out?

$$W^{(1)} = W^{(1)} - \eta \frac{\partial J}{\partial W^{(1)}}$$

$$W^{(2)} = W^{(2)} - \eta \frac{\partial J}{\partial W^{(2)}}$$

$$W^{(3)} = W^{(3)} - \eta \frac{\partial J}{\partial W^{(3)}}$$

$$J(\mathbf{W}) = \frac{1}{2} \|F(\mathbf{X}) - \mathbf{Y}\|^2 = \frac{1}{2} \|\hat{\mathbf{Y}} - \mathbf{Y}\|^2$$



Starting From a Component of $W^{(3)}$

$$w_{ml}^{(3)} = w_{ml}^{(3)} - \eta \frac{\partial J}{\partial w_{ml}^{(3)}}$$

- To focus on $\frac{\partial J}{\partial w_{ml}^{(3)}}$,

$$\frac{\partial J}{\partial w_{ml}^{(3)}} = \frac{1}{2} \frac{\partial (\|\hat{\mathbf{Y}} - \mathbf{Y}\|^2)}{\partial w_{ml}^{(3)}} = \frac{1}{2} \frac{\partial (\hat{y}_m - y_m)^2}{\partial w_{ml}^{(3)}}$$

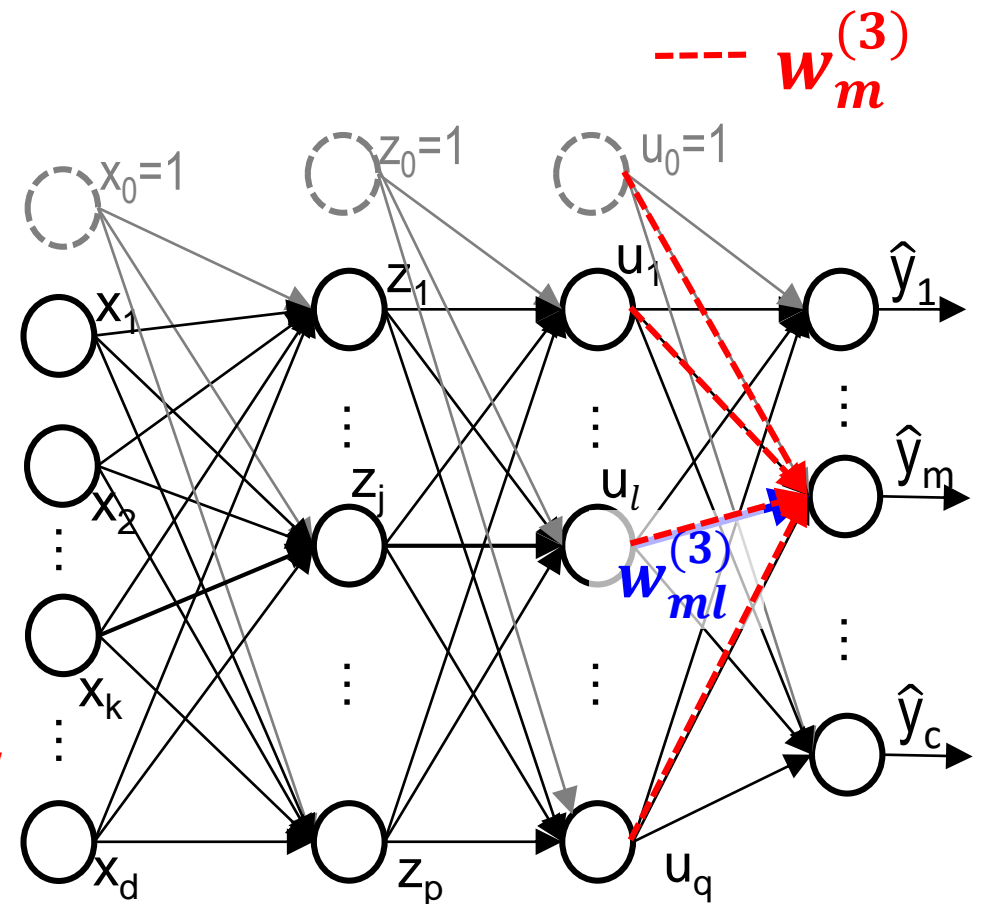
$$= (\hat{y}_m - y_m) \frac{\partial \hat{y}_m}{\partial w_{ml}^{(3)}}$$

$w_{ml}^{(3)}$ vs. \hat{y}_m ?
It is in $f(\mathbf{w}_m^{(3)} \cdot \mathbf{u})$

$$= (\hat{y}_m - y_m) \frac{\partial f(\mathbf{w}_m^{(3)} \cdot \mathbf{u})}{\partial w_{ml}^{(3)}}$$

We don't know f yet

$$= (\hat{y}_m - y_m) f'(\mathbf{w}_m^{(3)} \cdot \mathbf{u}) u_l$$



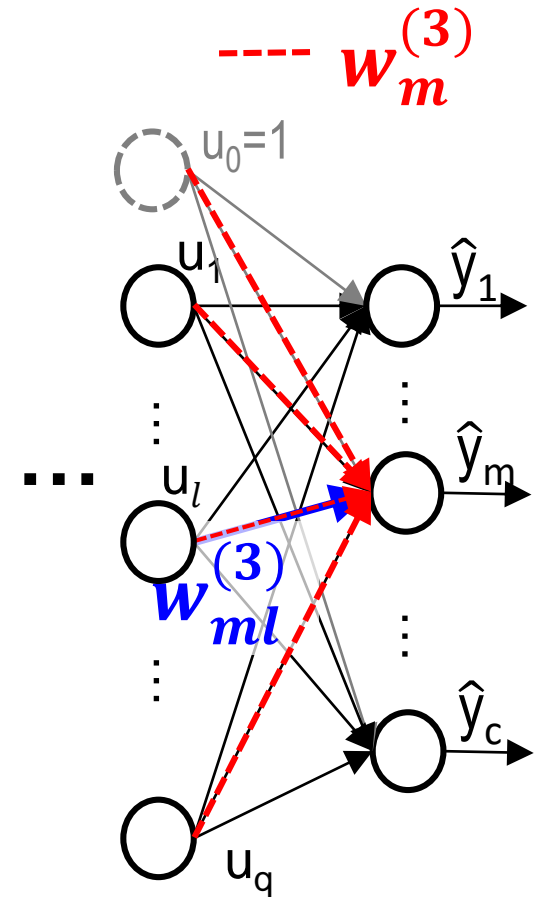
GD for $\mathbf{W}^{(3)}$

- For each m , δ_m can be defined, thus

$$\frac{\partial \mathbf{J}}{\partial w_{ml}^{(3)}} = (\hat{y}_m - y_m) f'(\mathbf{w}_m^{(3)} \cdot \mathbf{u}) u_l = \delta_m u_l$$

- Defining $\boldsymbol{\delta} = [\delta_1, \dots, \delta_c]^\top$ and $\mathbf{u} = [u_1, \dots, u_q]^\top$

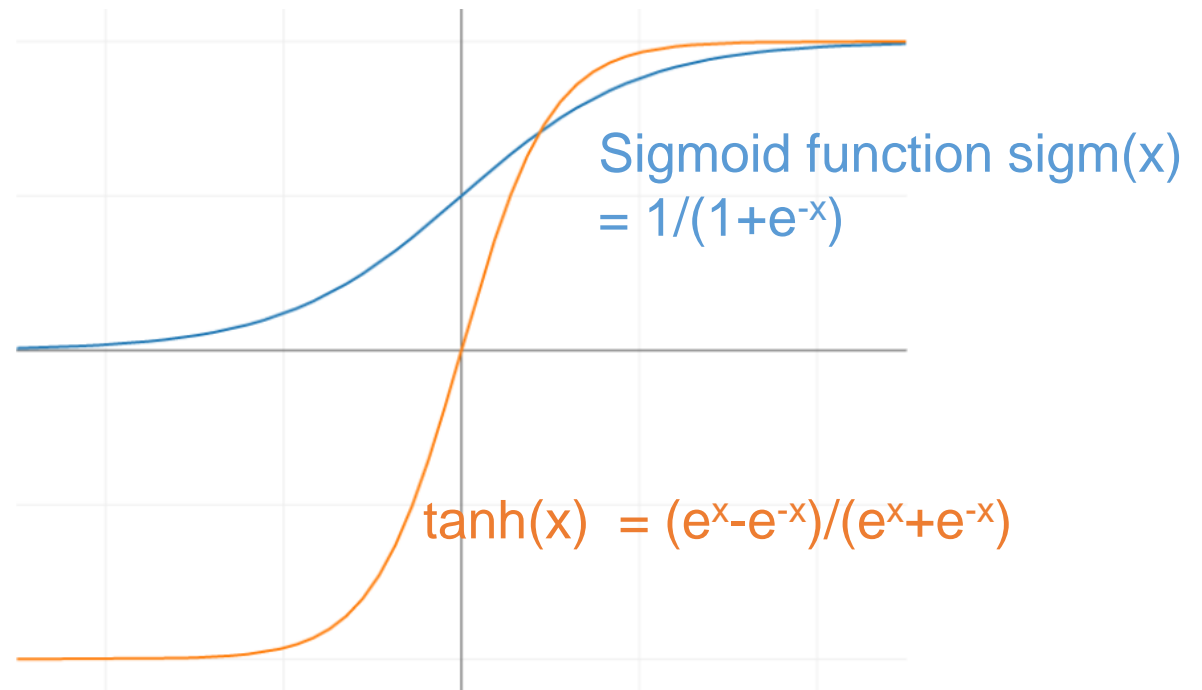
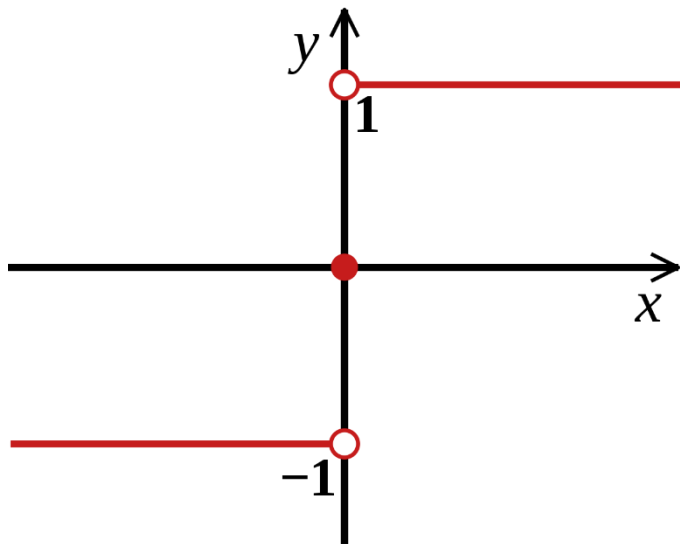
$$\mathbf{W}^{(3)} = \mathbf{W}^{(3)} - \eta \frac{\partial \mathbf{J}}{\partial \mathbf{W}^{(3)}} = \mathbf{W}^{(3)} - \eta \boldsymbol{\delta} \mathbf{u}^\top$$



We Need $f'(\bullet)$; Choice of Activation Function

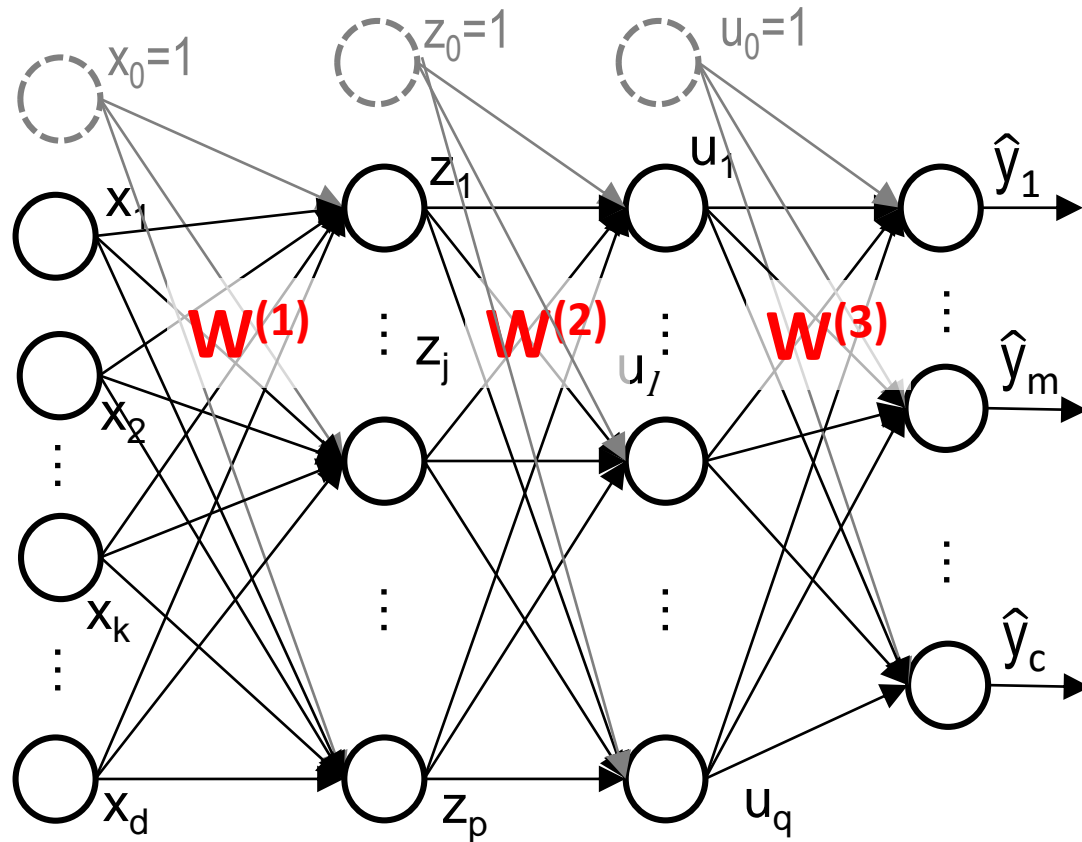
- We need differentiable activation function such as sigmoid or tanh, instead of sign function.

$f(x)$ as sign function or $\text{sgn}(x)$



What We Did is

$$\mathbf{W}^{(3)} = \mathbf{W}^{(3)} - \eta \frac{\partial J}{\partial \mathbf{W}^{(3)}} = \mathbf{W}^{(3)} - \eta \delta \mathbf{u}^T$$



For prediction,
 $\hat{\mathbf{y}} = f(\mathbf{W}^{(3)} f(\mathbf{W}^{(2)} f(\mathbf{W}^{(1)} \mathbf{x})))$

How About $W^{(2)}$?

It may affect all outputs

$$\frac{\partial J}{\partial w_{lj}^{(2)}} = \frac{1}{2} \frac{\partial (\|\hat{\mathbf{Y}} - \mathbf{Y}\|^2)}{\partial w_{lj}^{(2)}} = \frac{1}{2} \frac{\partial \sum_{m=1}^c (\hat{y}_m - y_m)^2}{\partial w_{lj}^{(2)}}$$

$$= \sum_{m=1}^c (\hat{y}_m - y_m) \frac{\partial \hat{y}_m}{\partial w_{lj}^{(2)}} \quad \mathbf{w}_{lj}^{(2)} \text{ vs. } \hat{y}_m?$$

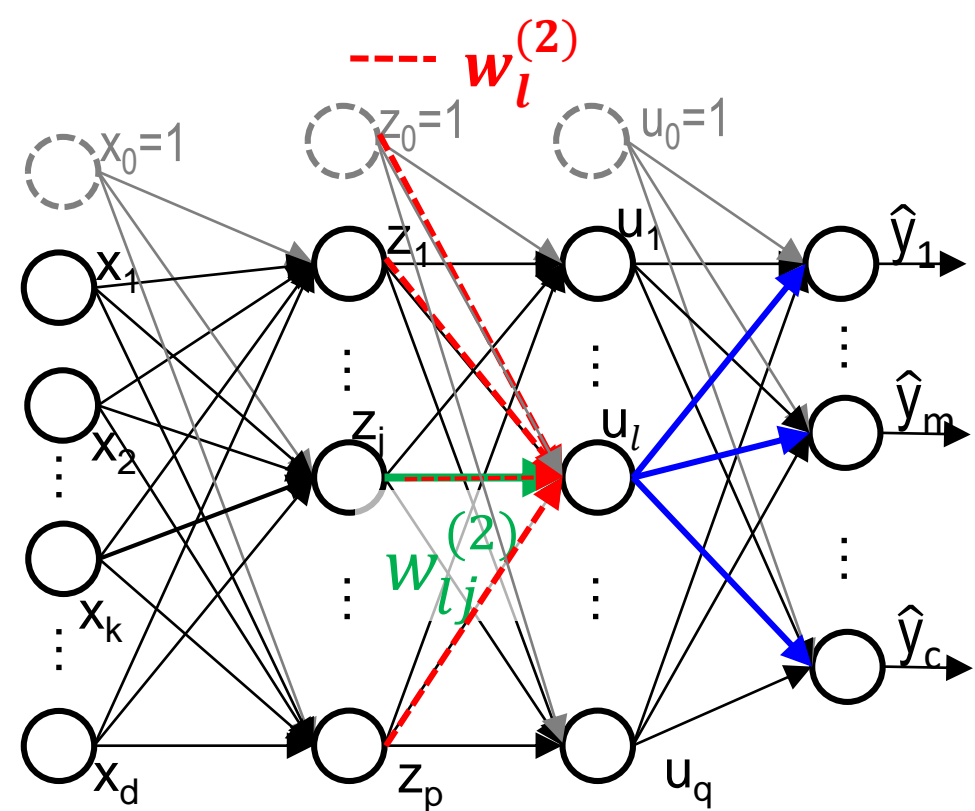
$$= \sum_{m=1}^c (\hat{y}_m - y_m) \frac{\partial f(\mathbf{w}_m^{(3)} \cdot \mathbf{u})}{\partial w_{lj}^{(2)}} \quad \text{It is in } f(\mathbf{w}_m^{(3)} \cdot \mathbf{u})$$

$$= \sum_{m=1}^c (\hat{y}_m - y_m) f'(\mathbf{w}_m^{(3)} \cdot \mathbf{u}) \frac{\partial (\mathbf{w}_m^{(3)} \cdot \mathbf{u})}{\partial w_{lj}^{(2)}}$$

$$= \sum_{m=1}^c \left\{ (\hat{y}_m - y_m) f'(\mathbf{w}_m^{(3)} \cdot \mathbf{u}) w_{ml}^{(3)} \right\} \frac{\partial u_l}{\partial w_{lj}^{(2)}}$$

$$= \frac{\partial f(\mathbf{w}_l^{(2)} \cdot \mathbf{z})}{\partial w_{lj}^{(2)}} \sum_{m=1}^c \left\{ (\hat{y}_m - y_m) f'(\mathbf{w}_m^{(3)} \cdot \mathbf{u}) w_{ml}^{(3)} \right\}$$

$$= z_j f'(\mathbf{w}_l^{(2)} \cdot \mathbf{z}) \sum_{m=1}^c \left\{ (\hat{y}_m - y_m) f'(\mathbf{w}_m^{(3)} \cdot \mathbf{u}) w_{ml}^{(3)} \right\}$$



$w_{lj}^{(2)}$ vs. $w_m^{(3)} \cdot u$?

$w_{lj}^{(2)} \rightarrow f(\mathbf{w}_l^{(2)} \cdot \mathbf{z}) = u_l \rightarrow w_m^{(3)} \cdot \mathbf{u}$

GD for $\mathbf{W}^{(2)}$

$$\frac{\partial J}{\partial w_{lj}^{(2)}} = z_j f'(\mathbf{w}_l^{(2)} \cdot \mathbf{z}) \sum_{m=1}^c \left\{ (\hat{y}_m - y_m) f'(\mathbf{w}_m^{(3)} \cdot \mathbf{u}) w_{ml}^{(3)} \right\}$$

- For each l , γ_l can be defined as

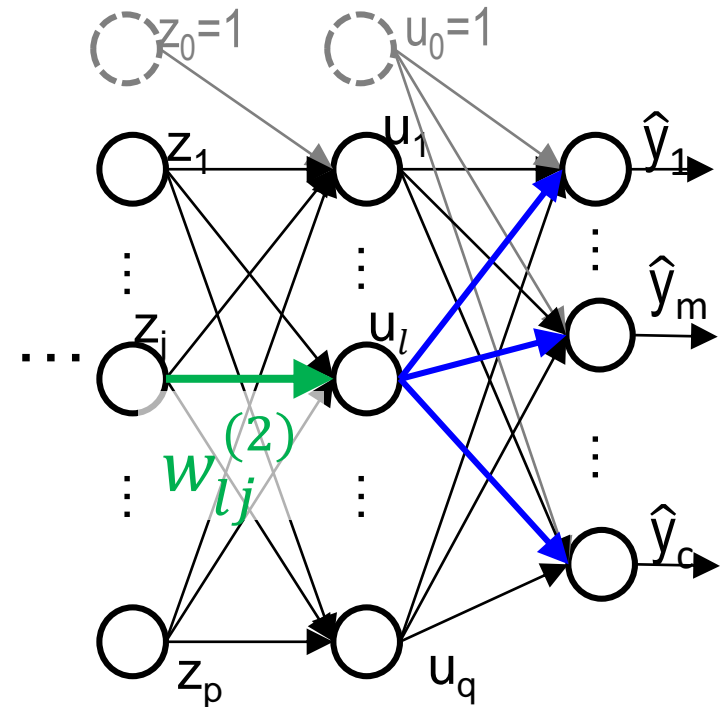
$$\begin{aligned} \gamma_l &= f'(\mathbf{w}_l^{(2)} \cdot \mathbf{z}) \sum_{m=1}^c \left\{ (\hat{y}_m - y_m) f'(\mathbf{w}_m^{(3)} \cdot \mathbf{u}) w_{ml}^{(3)} \right\} \\ &= f'(\mathbf{w}_l^{(2)} \cdot \mathbf{z}) \sum_{m=1}^c \left\{ \delta_m w_{ml}^{(3)} \right\} \end{aligned}$$

- Then,

$$\frac{\partial J}{\partial w_{lj}^{(2)}} = \gamma_l z_j$$

- Defining $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_q]^T$ & $\mathbf{z} = [z_1, \dots, z_p]^T$

$$\mathbf{W}^{(2)} = \mathbf{W}^{(2)} - \eta \frac{\partial J}{\partial \mathbf{W}^{(2)}} = \mathbf{W}^{(2)} - \eta \boldsymbol{\gamma} \mathbf{z}^T$$



You should get a feel for the terms

We Can Extend it to $\mathbf{W}^{(1)}$

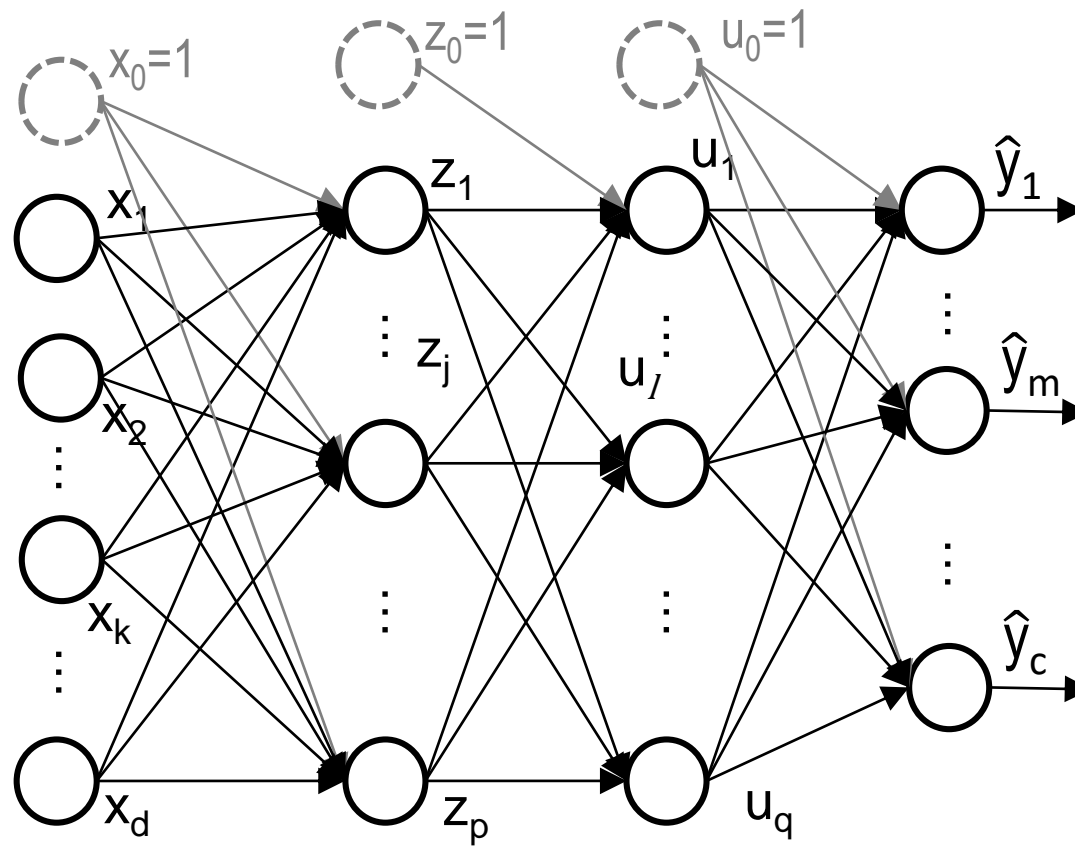
- Calculating β for each j ,

$$\beta_j = f'(\mathbf{w}_j^{(1)} \cdot \mathbf{x}) \sum_{l=1}^q \{y_l w_{lj}^{(2)}\}$$

- Then, with a vector definition $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_p)$

$$\mathbf{W}^{(1)} = \mathbf{W}^{(1)} - \eta \frac{\partial \mathbf{J}}{\partial \mathbf{W}^{(1)}} = \mathbf{W}^{(1)} - \eta \boldsymbol{\beta} \mathbf{x}^T$$

Direction?



Training Algorithm (Mini-batch version)

Initialize $\mathbf{W}^{(1)}$, $\mathbf{W}^{(2)}$, $\mathbf{W}^{(3)}$ properly

Until (No change):

Extracting n' samples from \mathbf{X} to form \mathbf{X}'

$$\Delta\mathbf{W}^{(1)} = 0, \Delta\mathbf{W}^{(2)} = 0, \Delta\mathbf{W}^{(3)} = 0$$

for each comp of \mathbf{X}' :

$$\mathbf{z} = f(\mathbf{W}^{(1)}\mathbf{x})$$

$$\mathbf{u} = f(\mathbf{W}^{(2)}\mathbf{z})$$

$$\mathbf{y}_h = f(\mathbf{W}^{(3)}\mathbf{u})$$

$$\boldsymbol{\delta} = (\mathbf{y}_h - \mathbf{y}) \times f'(\mathbf{W}^{(3)}\mathbf{u}) \quad \#c \times 1$$

$$\boldsymbol{\gamma} = \boldsymbol{\delta} \bullet \widetilde{\mathbf{W}}^{(3)} \times f'(\mathbf{W}^{(2)}\mathbf{z}) \quad \#q \times 1$$

$$\boldsymbol{\beta} = \boldsymbol{\gamma} \bullet \widetilde{\mathbf{W}}^{(2)} \times f'(\mathbf{W}^{(1)}\mathbf{z}) \quad \#p \times 1$$

$$\mathbf{W}^{(3)} = \mathbf{W}^{(3)} - (1/n')\eta\boldsymbol{\delta}\mathbf{u}^T$$

$$\mathbf{W}^{(2)} = \mathbf{W}^{(2)} - (1/n')\eta\boldsymbol{\gamma}\mathbf{z}^T$$

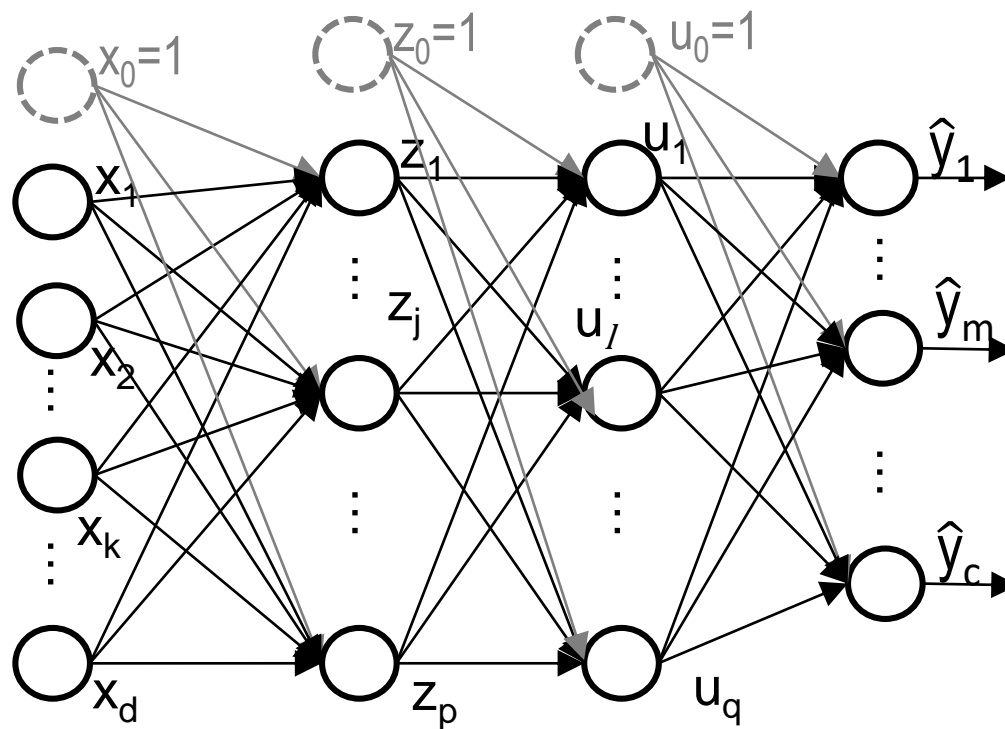
$$\mathbf{W}^{(1)} = \mathbf{W}^{(1)} - (1/n')\eta\boldsymbol{\beta}\mathbf{x}^T$$

You can see the \mathbf{W} is
determined in backward

→ Called error back(ward) propagation

Prediction

- At the output stage, the softmax function is utilized for deriving probability (soft decision)
- For hard decision, $\mathbf{m}^* = \underset{m}{\operatorname{argmax}} \hat{y}_m$



Softmax function

$$\sigma(\hat{y}_m) = \frac{\exp(\hat{y}_m)}{\sum_{k=1}^c \exp(\hat{y}_k)}$$

[2, 0.4, 0.2]

→ [0.73, 0.15, 0.12]

→ Called error forward propagation

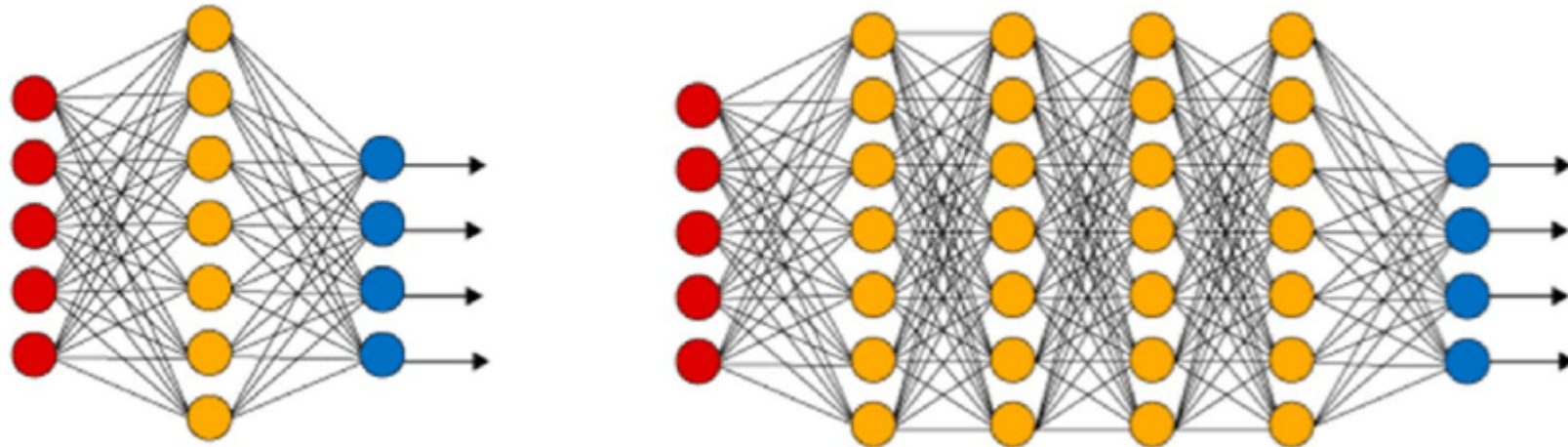
Limitation of Deep Layer

- Large computation
- Vanishing/Exploding gradient
- Overfitting risks

Enable of Deep Layer: Deep Learning

- Large computation
 - ➔ Development of high performance hardware
 - ➔ Use of more efficient algorithm such as CNN
- Vanishing/Exploding gradient
 - ➔ Use of improved activation function such as RELU or its variants
 - ➔ Use of cross-entropy or LLE for cost function
- Overfitting risks
 - ➔ Development of various regulation techniques

Deep Learning from Deep MLP (Deep MLP)



Package for Deep Learning in Python; Tensorflow

- Using tensorflow developed by google
 - Compatible with sklearn (tf.learn)
 - Various high-level API such as Keras or Pretty Tensor
 - Visualization is easy with TensorBoard
- <https://www.tensorflow.org/>
- Install tensorflow package by typing in Anaconda prompt:
conda install tensorflow

Checkpoints

- ✓ Backward propagation for training MLP
- ✓ Deep learning and its limitation
- ✓ Coming up next: Tensorflow for python coding MLP