

***CS434a/541a: Pattern Recognition***  
***Prof. Olga Veksler***

**Lecture 13**

# *Today*

---

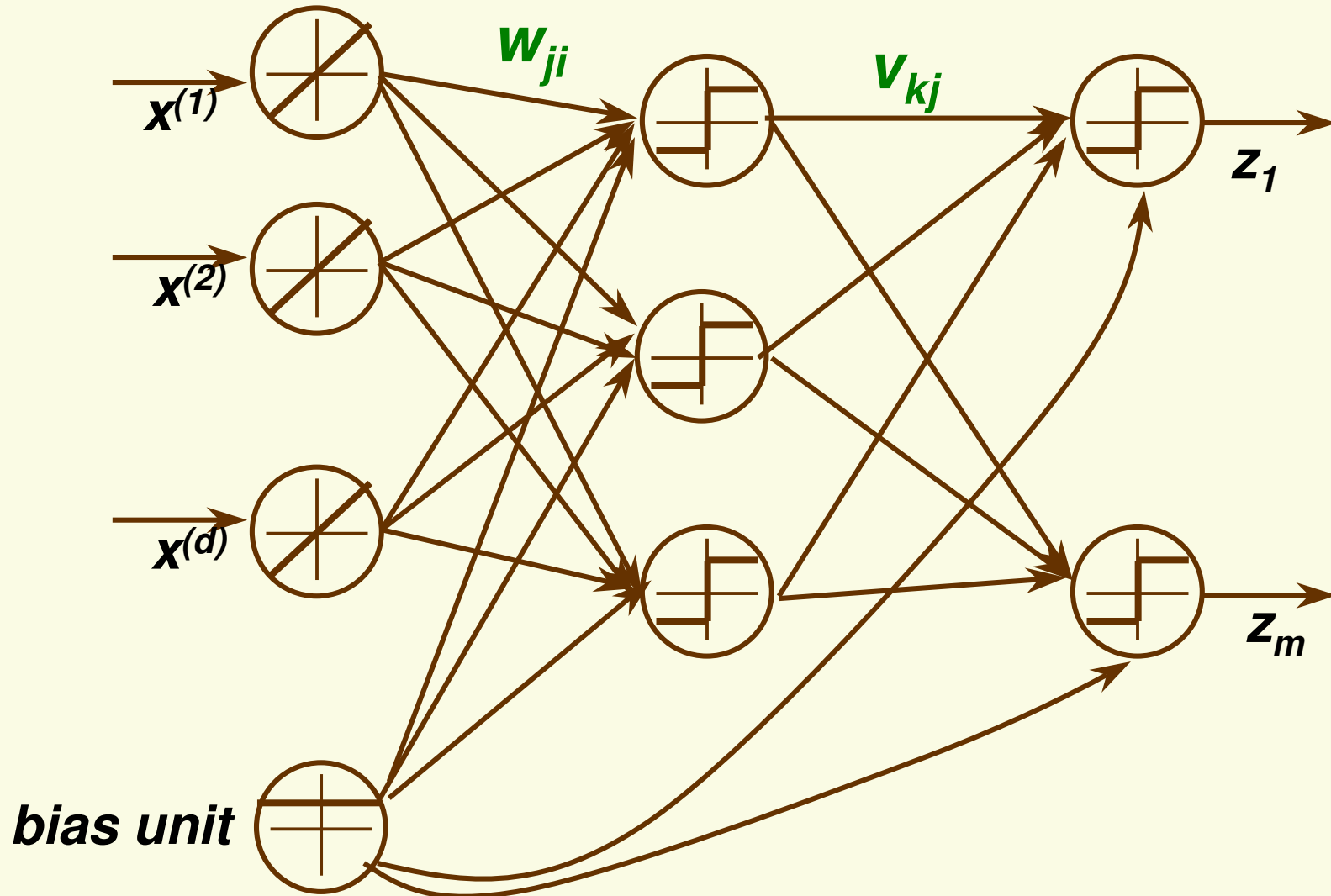
- Continue Multilayer Neural Networks (MNN)
  - Review MNN structure
  - Backpropagation
  - Training Protocols

# MNN: Feed Forward Operation

**input layer:**  
 $d$  features

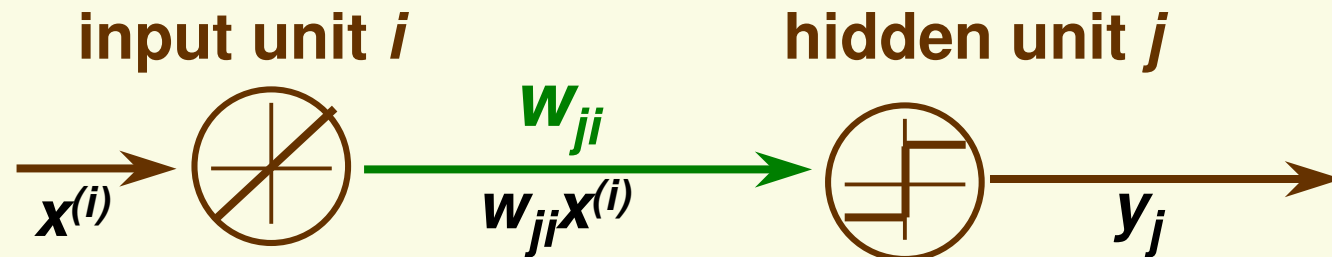
**hidden layer:**

**output layer:**  
 $m$  outputs, one for each class

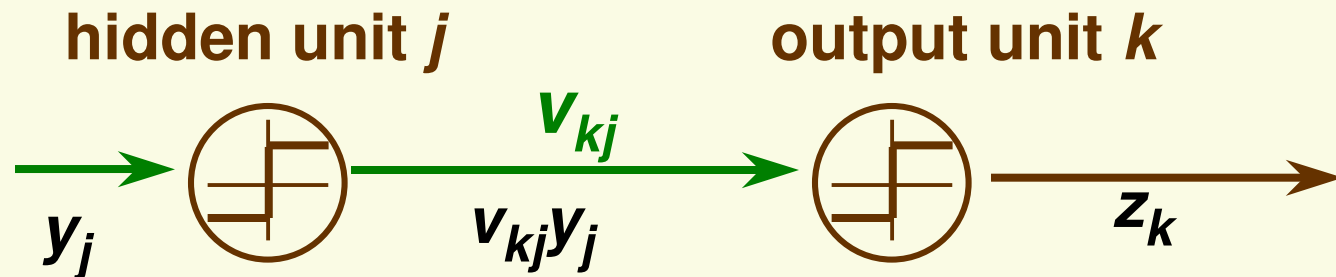


# MNN: Notation for Weights

- Use  $w_{ji}$  to denote the weight between input unit  $i$  and hidden unit  $j$



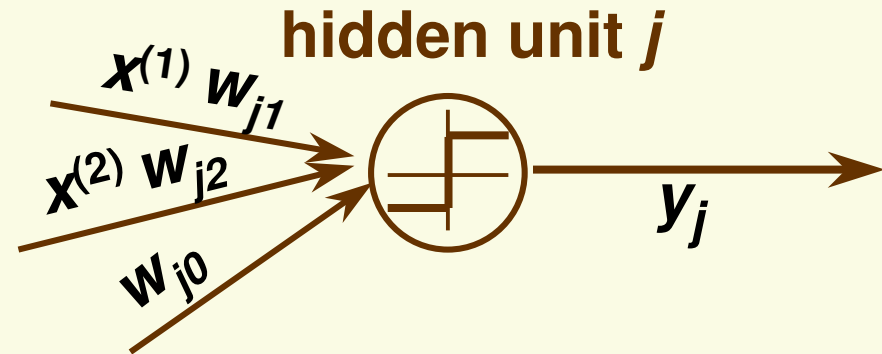
- Use  $v_{kj}$  to denote the weight between hidden unit  $j$  and output unit  $k$



# MNN: Notation for Activation

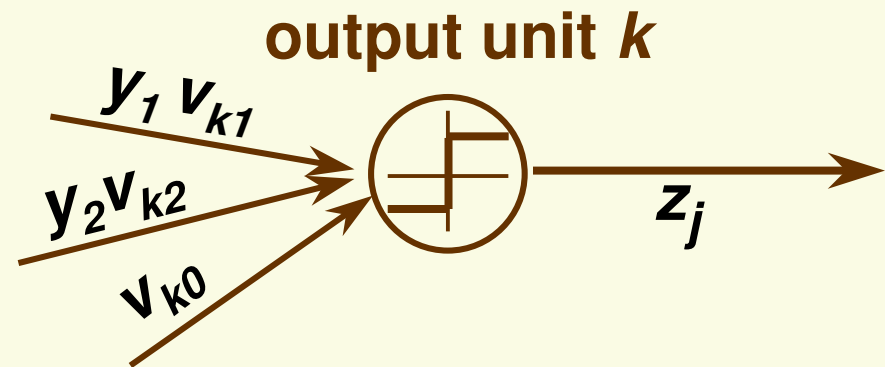
- Use  $net_j$  to denote the activation and hidden unit  $j$

$$net_j = \sum_{i=1}^d x^{(i)} w_{ji} + w_{j0}$$



- Use  $net_k^*$  to denote the activation at output unit  $k$

$$net_k^* = \sum_{j=1}^{N_H} y_j v_{kj} + v_{k0}$$



# Discriminant Function

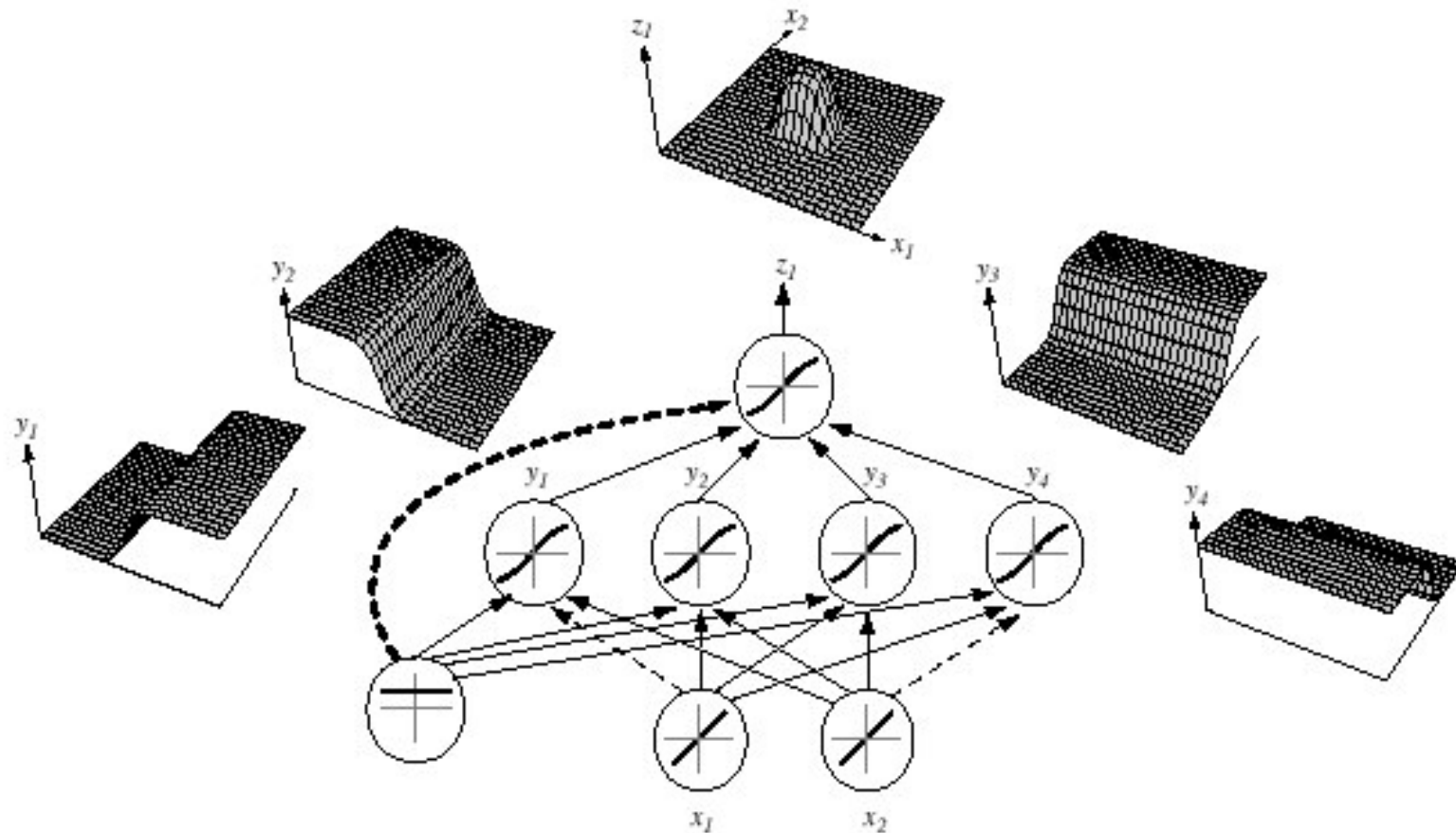
- Discriminant function for class  $k$  (the output of the  $k$ th output unit)

$$\begin{aligned} g_k(\mathbf{x}) = \mathbf{z}_k = & \\ = & \underbrace{f\left(\sum_{j=1}^{N_H} \mathbf{v}_{kj} f\left(\sum_{i=1}^d \mathbf{w}_{ji} \mathbf{x}^{(i)} + \mathbf{w}_{j0}\right) + \mathbf{v}_{k0}\right)}_{\text{activation at } k\text{th output unit}} \end{aligned}$$

activation at  $j$ th hidden unit

- Rich expressive power: every **continuous** discriminant function can be implemented with enough hidden units, 1 hidden layer, and proper nonlinear activation functions

# Expressive Power



**FIGURE 6.2.** A 2-4-1 network (with bias) along with the response functions at different units; each hidden output unit has sigmoidal activation function  $f(\cdot)$ . In the case shown, the hidden unit outputs are paired in opposition thereby producing a “bump” at the output unit. Given a sufficiently large number of hidden units, any continuous function from input to output can be approximated arbitrarily well by such a network. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

# MNN Activation function

- Must be nonlinear for expressive power larger than that of perceptron
  - If use linear activation function at hidden layer, can only deal with linearly separable classes
  - Suppose at hidden unit  $j$ ,  $h(\mathbf{u}) = \mathbf{a}_j \mathbf{u}$

$$\begin{aligned}
 g_k(\mathbf{x}) &= f \left( \sum_{j=1}^{N_H} \mathbf{v}_{kj} h \left( \sum_{i=1}^d \mathbf{w}_{ji} \mathbf{x}^{(i)} + \mathbf{w}_{j0} \right) + \mathbf{v}_{k0} \right) \\
 &= f \left( \sum_{j=1}^{N_H} \mathbf{v}_{kj} \mathbf{a}_j \left( \sum_{i=1}^d \mathbf{w}_{ji} \mathbf{x}^{(i)} + \mathbf{w}_{j0} \right) + \mathbf{v}_{k0} \right) \\
 &= f \left( \sum_{i=1}^d \sum_{j=1}^{N_H} (\mathbf{v}_{kj} \mathbf{a}_j \mathbf{w}_{ji} \mathbf{x}^{(i)} + \mathbf{w}_{j0}) + \mathbf{v}_{k0} \right) \\
 &= f \left( \sum_{i=1}^d \mathbf{x}^{(i)} \sum_{j=1}^{N_H} \mathbf{v}_{kj} \mathbf{a}_j \mathbf{w}_{ji} + \left( \sum_{j=1}^{N_H} \mathbf{w}_{j0} + \mathbf{v}_{k0} \right) \right)
 \end{aligned}$$

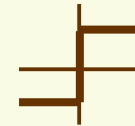


# ***MNN Activation function***

---

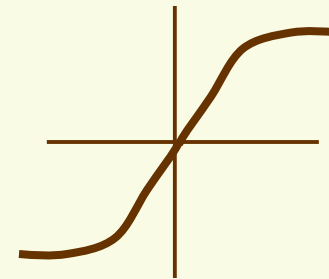
- In previous example, used discontinuous activation function

$$f(\mathit{net}_k) = \begin{cases} 1 & \text{if } \mathit{net}_k \geq 0 \\ -1 & \text{if } \mathit{net}_k < 0 \end{cases}$$



- We will use gradient descent for learning, so we need to use continuous activation function

***sigmoid*** function



- From now on, assume  $f$  is a differentiable function

# ***MNN: Modes of Operation***

---

- Network have two modes of operation:
  - ***Feedforward***

The feedforward operations consists of presenting a pattern to the input units and passing (or feeding) the signals through the network in order to get outputs units (no cycles!)
  - ***Learning***

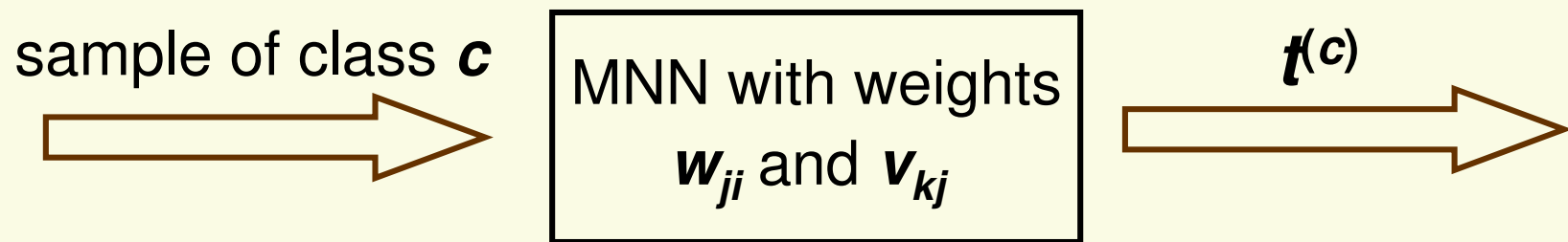
The supervised learning consists of presenting an input pattern and modifying the network parameters (weights) to reduce distances between the computed output and the desired output

# MNN: Class Representation

- Training samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$  each of class  $1, \dots, m$
- Let network output  $\mathbf{z}$  represent class  $\mathbf{c}$  as **target**  $\mathbf{t}^{(c)}$

$$\mathbf{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_c \\ \vdots \\ z_m \end{bmatrix} = \mathbf{t}^{(c)} = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow \text{c-th row}$$

## Our Ultimate Goal For FeedForward Operation

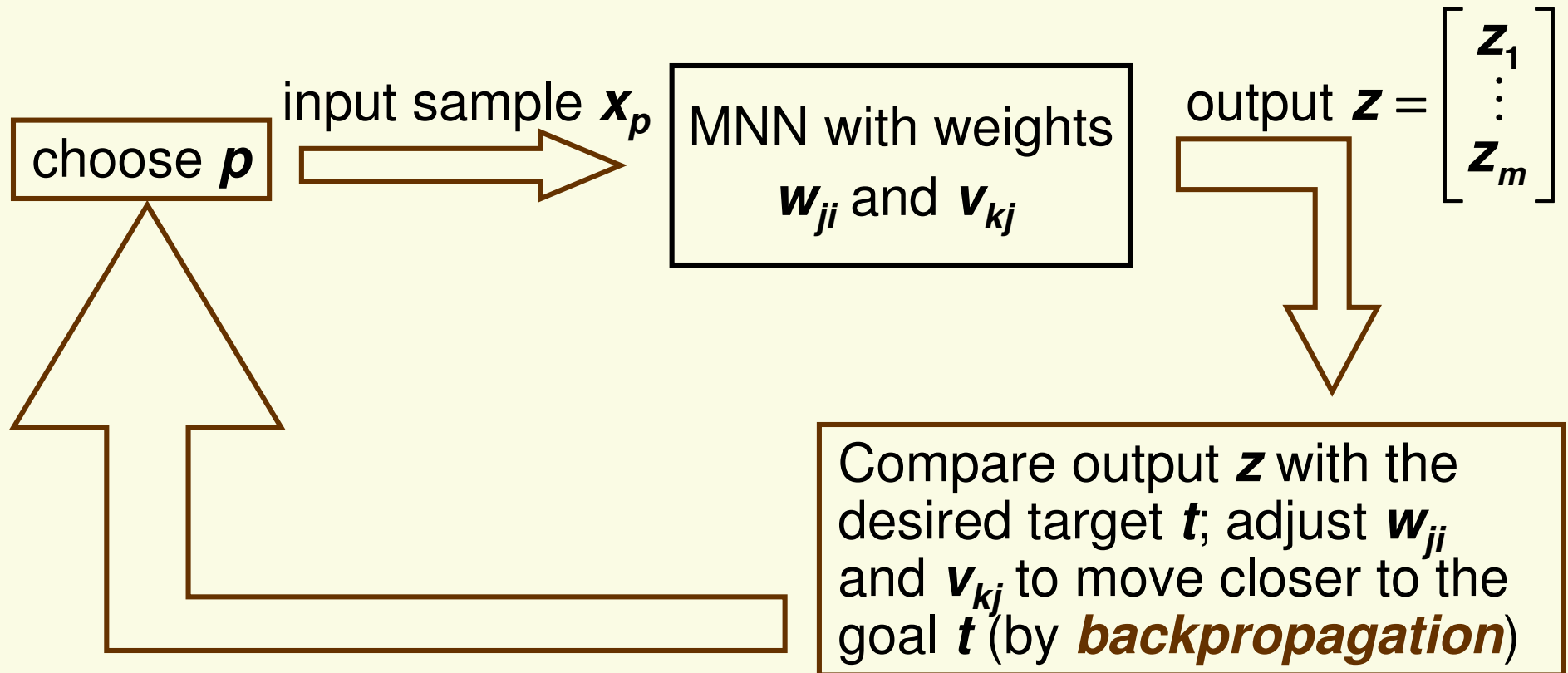


## MNN training to achieve the Ultimate Goal

Modify (learn) MNN parameters  $\mathbf{w}_{ji}$  and  $\mathbf{v}_{kj}$  so that for each **training** sample of class  $\mathbf{c}$  MNN output  $\mathbf{z} = \mathbf{t}^{(c)}$

# Network Training (learning)

1. Initialize weights  $\mathbf{w}_{ji}$  and  $\mathbf{v}_{kj}$  randomly
2. Iterate until a stopping criterion is reached



# BackPropagation

- Learn  $\mathbf{w}_{ji}$  and  $\mathbf{v}_{kj}$  by minimizing the training error
- What is the training error?
- Suppose the output of MNN for sample  $\mathbf{x}$  is  $\mathbf{z}$  and the target (desired output for  $\mathbf{x}$ ) is  $\mathbf{t}$
- Error on one sample:  $\mathcal{J}(\mathbf{w}, \mathbf{v}) = \frac{1}{2} \sum_{c=1}^m (t_c - z_c)^2$
- Training error:  $\mathcal{J}(\mathbf{w}, \mathbf{v}) = \frac{1}{2} \sum_{i=1}^n \sum_{c=1}^m (t_c^{(i)} - z_c^{(i)})^2$

- Use gradient descent:

$\mathbf{v}^{(0)}, \mathbf{w}^{(0)} = \text{random}$

*repeat until convergence:*

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} \mathcal{J}(\mathbf{w}^{(t)})$$

$$\mathbf{v}^{(t+1)} = \mathbf{v}^{(t)} - \eta \nabla_{\mathbf{v}} \mathcal{J}(\mathbf{v}^{(t)})$$

# BackPropagation

- For simplicity, first take training error for one sample  $\mathbf{x}_i$

$$J(\mathbf{w}, \mathbf{v}) = \frac{1}{2} \sum_{c=1}^m (t_c - \mathbf{z}_c)^2$$

*function of w,v*

*fixed constant*

$$\mathbf{z}_k = f \left( \sum_{j=1}^{N_H} \mathbf{v}_{kj} f \left( \sum_{i=1}^d \mathbf{w}_{ji} \mathbf{x}^{(i)} + \mathbf{w}_{j0} \right) + \mathbf{v}_{k0} \right)$$

- Need to compute
  - partial derivative w.r.t. hidden-to-output weights  $\frac{\partial J}{\partial \mathbf{v}_{kj}}$
  - partial derivative w.r.t. input-to-hidden weights  $\frac{\partial J}{\partial \mathbf{w}_{ji}}$

# BackPropagation: Layered Model

activation at  
hidden unit  $j$

$$net_j = \sum_{i=1}^d x^{(i)} w_{ji} + w_{j0}$$



output at  
hidden unit  $j$

$$y_j = f(net_j)$$



activation at  
output unit  $k$

$$net_k^* = \sum_{j=1}^{N_H} y_j v_{kj} + v_{k0}$$



activation at  
output unit  $k$

$$z_k = f(net_k^*)$$



objective function

$$J(\mathbf{w}, \mathbf{v}) = \frac{1}{2} \sum_{c=1}^m (t_c - z_c)^2$$

chain rule  $\frac{\partial J}{\partial v_{kj}}$

chain rule  $\frac{\partial J}{\partial w_{ji}}$

# BackPropagation

$$net_k = \sum_{j=1}^{N_H} y_j v_{kj} + v_{k0} \Rightarrow z_k = f(net_k^*) \Rightarrow J(w, v) = \frac{1}{2} \sum_{c=1}^m (t_c - z_c)^2$$

- First compute hidden-to-output derivatives  $\frac{\partial J}{\partial v_{kj}}$

$$\frac{\partial J}{\partial v_{kj}} = \frac{1}{2} \sum_{c=1}^m \frac{\partial}{\partial v_{kj}} (t_c - z_c)^2 = \sum_{c=1}^m (t_c - z_c) \frac{\partial}{\partial v_{kj}} (t_c - z_c)$$

$$= (t_k - z_k) \frac{\partial}{\partial v_{kj}} (t_k - z_k) = -(t_k - z_k) \frac{\partial}{\partial v_{kj}} (z_k)$$

$$= -(t_k - z_k) \frac{\partial z_k}{\partial net_k^*} \frac{\partial net_k^*}{\partial v_{kj}}$$

$$= \begin{cases} -(t_k - z_k) f'(net_k^*) y_j & \text{if } j \neq 0 \\ -(t_k - z_k) f'(net_k^*) & \text{if } j = 0 \end{cases}$$



# BackPropagation

---

Gradient Descent *Single Sample* Update Rule for hidden-to-output weights  $v_{kj}$

$$j > 0: \mathbf{v}_{kj}^{(t+1)} = \mathbf{v}_{kj}^{(t)} + \eta(t_k - z_k) f'(net_k^*) y_j$$

$$j = 0 \text{ (bias weight): } \mathbf{v}_{k0}^{(t+1)} = \mathbf{v}_{k0}^{(t)} + \eta(t_k - z_k) f'(net_k^*)$$

# BackPropagation

- Now compute input-to-hidden  $\frac{\partial J}{\partial w_{ji}}$

$$\begin{aligned}
 \frac{\partial J}{\partial w_{ji}} &= \sum_{k=1}^m (t_k - z_k) \frac{\partial}{\partial w_{ji}} (t_k - z_k) \\
 &= - \sum_{k=1}^m (t_k - z_k) \frac{\partial z_k}{\partial w_{ji}} = - \sum_{k=1}^m (t_k - z_k) \frac{\partial z_k}{\partial net_k^*} \frac{\partial net_k^*}{\partial w_{ji}} \\
 &= - \sum_{k=1}^m (t_k - z_k) f'(net_k^*) \frac{\partial net_k^*}{\partial y_j} \frac{\partial y_j}{\partial w_{ji}} \\
 &= - \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj} \frac{\partial y_j}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} \\
 &= - \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj} \frac{\partial y_j}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} \\
 &= \begin{cases} - \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj} f'(net_j) x^{(i)} & \text{if } i \neq 0 \\ - \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj} f'(net_j) & \text{if } i = 0 \end{cases}
 \end{aligned}$$

$$\begin{aligned}
 &\Downarrow \\
 net_h &= \sum_{h=1}^d x^{(i)} w_{hi} + w_{h0} \\
 &\Downarrow \\
 y_j &= f(net_j) \\
 &\Downarrow \\
 net_k^* &= \sum_{s=1}^{N_H} y_s v_{ks} + v_{k0} \\
 &\Downarrow \\
 z_k &= f(net_k^*) \\
 &\Downarrow \\
 J(w, v) &= \frac{1}{2} \sum_{c=1}^m (t_c - z_c)^2
 \end{aligned}$$

# BackPropagation

---

$$\frac{\partial J}{\partial w_{ji}} = \begin{cases} -f'(net_j) x^{(i)} \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj} & \text{if } i \neq 0 \\ -f'(net_j) \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj} & \text{if } i = 0 \end{cases}$$

Gradient Descent **Single Sample** Update Rule for input-to-hidden weights  $w_{ji}$

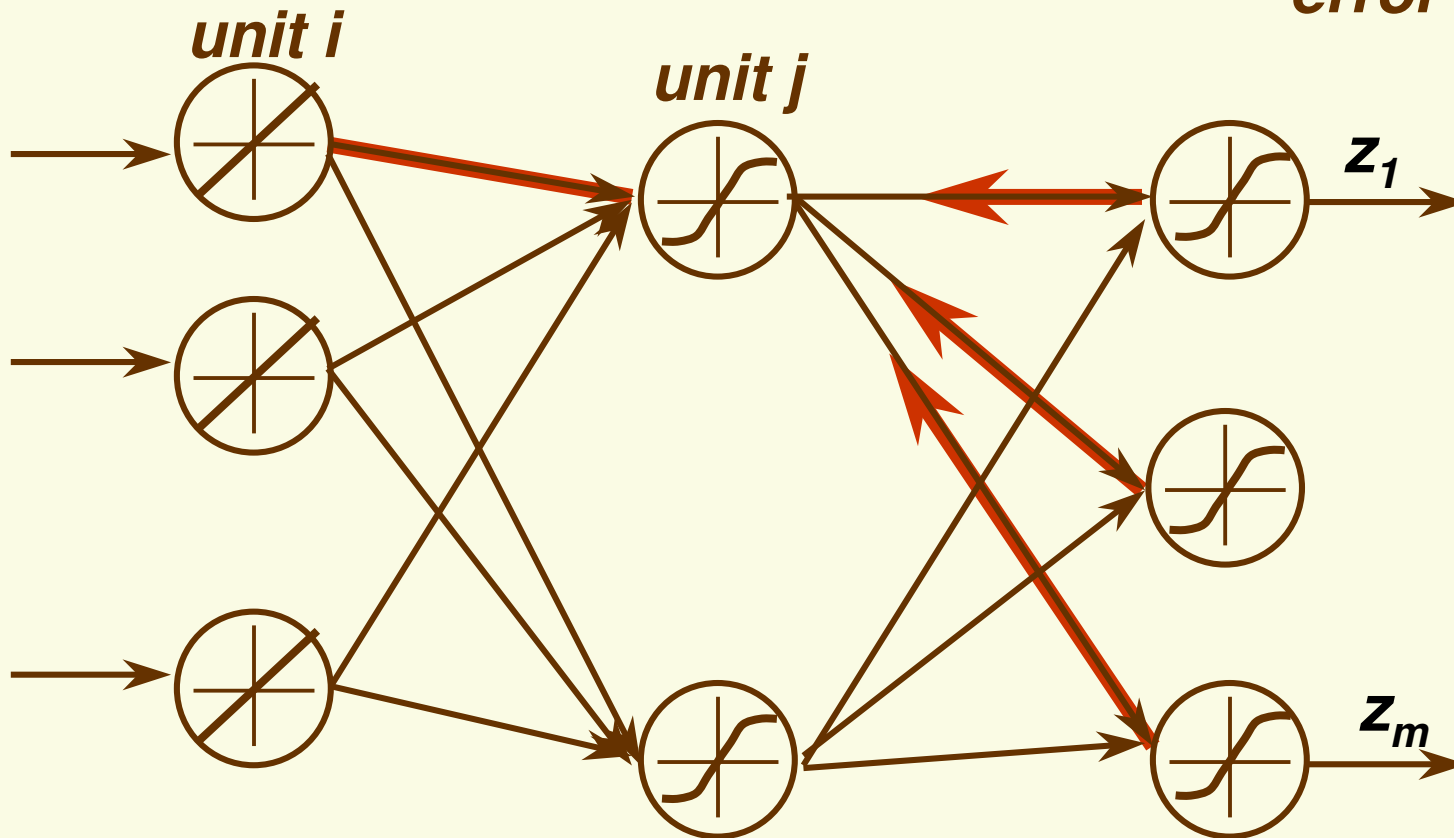
$$i > 0: w_{ji}^{(t+1)} = w_{ji}^{(t)} + \eta f'(net_j) x^{(i)} \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj}$$

$$i = 0 \text{ (bias weight): } w_{j0}^{(t+1)} = w_{j0}^{(t)} + \eta f'(net_j) \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj}$$

# BackPropagation of Errors

$$\frac{\partial J}{\partial w_{ji}} = -f'(net_j) x^{(i)} \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj}$$

$$\frac{\partial J}{\partial v_{kj}} = -\underbrace{(t_k - z_k)}_{\text{error}} f'(net_k^*) y_j$$



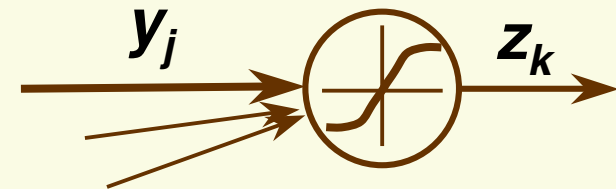
- Name “backpropagation” because during training, errors propagated back from output to hidden layer

# BackPropagation

- Consider update rule for hidden-to-output weights:

$$\mathbf{v}_{kj}^{(t+1)} = \mathbf{v}_{kj}^{(t)} + \eta(\mathbf{t}_k - \mathbf{z}_k) f'(\mathbf{net}_k^*) \mathbf{y}_j$$

- Suppose  $\mathbf{t}_k - \mathbf{z}_k > 0$
- Then output of the  $k$ th hidden unit is too small:  $\mathbf{t}_k > \mathbf{z}_k$
- Typically activation function  $f$  is s.t.  $f' > 0$
- Thus  $(\mathbf{t}_k - \mathbf{z}_k) f'(\mathbf{net}_k^*) > 0$
- There are 2 cases:



1.  $\mathbf{y}_j > 0$ , then to increase  $\mathbf{z}_k$ , should increase weight  $\mathbf{v}_{kj}$  which is exactly what we do since  $\eta(\mathbf{t}_k - \mathbf{z}_k) f'(\mathbf{net}_k^*) \mathbf{y}_j > 0$
2.  $\mathbf{y}_j < 0$ , then to increase  $\mathbf{z}_k$ , should decrease weight  $\mathbf{v}_{kj}$  which is exactly what we do since  $\eta(\mathbf{t}_k - \mathbf{z}_k) f'(\mathbf{net}_k^*) \mathbf{y}_j < 0$

# BackPropagation

---

- The case  $t_k - z_k < 0$  is analogous
- Similarly, can show that input-to-hidden weights make sense
- Important: weights should be initialized to random **nonzero** numbers

$$\frac{\partial J}{\partial w_{ji}} = -f'(net_j) x^{(i)} \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj}$$

- if  $v_{kj} = 0$ , input-to-hidden weights  $w_{ji}$  never updated

# *Training Protocols*

---

- How to present samples in training set and update the weights?
- Three major training protocols:
  1. Stochastic
    - Patterns are chosen randomly from the training set, and network weights are updated after every sample presentation
  2. Batch
    - weights are update based on all samples; iterate weight update
  3. Online
    - each sample is presented only once, weight update after each sample presentation

# Stochastic Back Propagation

## 1. Initialize

- number of hidden layers  $n_H$
- weights  $\mathbf{w}$ ,  $\mathbf{v}$
- convergence criterion  $\theta$  and learning rate  $\eta$
- time  $t = 0$

## 2. do

$\mathbf{x} \leftarrow$  randomly chosen training pattern

for all  $0 \leq i \leq d$ ,  $0 \leq j \leq n_H$ ,  $0 \leq k \leq m$

$$\mathbf{w}_{ji} = \mathbf{w}_{ji} + \eta f'(\mathbf{net}_j) \mathbf{x}^{(i)} \sum_{k=1}^m (t_k - z_k) f'(\mathbf{net}_k^*) \mathbf{v}_{kj}$$

$$\mathbf{w}_{j0} = \mathbf{w}_{j0} + \eta f'(\mathbf{net}_j) \sum_{k=1}^m (t_k - z_k) f'(\mathbf{net}_k^*) \mathbf{v}_{kj}$$

$$\mathbf{v}_{kj} = \mathbf{v}_{kj} + \eta (t_k - z_k) f'(\mathbf{net}_k^*) \mathbf{y}_j$$

$$\mathbf{v}_{k0} = \mathbf{v}_{k0} + \eta (t_k - z_k) f'(\mathbf{net}_k^*)$$

$$t = t + 1$$

until  $\|\mathbf{J}\| < \theta$

## 3. return $\mathbf{v}$ , $\mathbf{w}$



# Batch Back Propagation

- This is the **true** gradient descent, (unlike stochastic propagation)
- For simplicity, derived backpropagation for a single sample objective function:

$$J(\mathbf{w}, \mathbf{v}) = \frac{1}{2} \sum_{c=1}^m (t_c - z_c)^2$$

- The full objective function:

$$J(\mathbf{w}, \mathbf{v}) = \frac{1}{2} \sum_{i=1}^n \sum_{c=1}^m (t_c^{(i)} - z_c^{(i)})^2$$

- Derivative of full objective function is just a sum of derivatives for each sample:

$$\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}, \mathbf{v}) = \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \mathbf{w}} \left( \sum_{c=1}^m (t_c^{(i)} - z_c^{(i)})^2 \right)$$

*already derived this*

# Batch Back Propagation

---

- For example,

$$\frac{\partial \mathbf{J}}{\partial \mathbf{w}_{ji}} = \sum_{p=1}^n -f'(\mathit{net}_j) \mathbf{x}_p^{(i)} \sum_{k=1}^m (t_k - \mathbf{z}_k) f'(\mathit{net}_k^*) \mathbf{v}_{kj}$$

# Batch Back Propagation

1. Initialize  $n_H$ ,  $w$ ,  $v$ ,  $\theta$ ,  $\eta$ ,  $t = 0$

2. do

$$\Delta v_{kj} = \Delta v_{k0} = \Delta w_{ji} = \Delta w_{j0} = 0$$

one epoch

for all  $1 \leq p \leq n$

for all  $0 \leq i \leq d$ ,  $0 \leq j \leq n_H$ ,  $0 \leq k \leq m$

$$\Delta v_{kj} = \Delta v_{kj} + \eta(t_k - z_k) f'(net_k^*) y_j$$

$$\Delta v_{k0} = \Delta v_{k0} + \eta(t_k - z_k) f'(net_k^*)$$

$$\Delta w_{ji} = \Delta w_{ji} + \eta f'(net_j) x_p^{(i)} \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj}$$

$$\Delta w_{j0} = \Delta w_{j0} + \eta f'(net_j) \sum_{k=1}^m (t_k - z_k) f'(net_k^*) v_{kj}$$

$$v_{kj} = v_{kj} + \Delta v_{kj}; v_{k0} = v_{k0} + \Delta v_{k0}; w_{ji} = w_{ji} + \Delta w_{ji}; w_{j0} = w_{j0} + \Delta w_{j0}$$

$$t = t + 1$$

until  $\|J\| < \theta$

3. return  $v$ ,  $w$

# *Training Protocols*

---

## 1. Batch

- True gradient descent

## 2. Stochastic

- Faster than batch method
- Usually the recommended way

## 3. Online

- Used when number of samples is so large it does not fit in the memory
- Dependent on the order of sample presentation
- Should be avoided when possible