

# Neural Networks

## Part 3 – Training with Backpropagation

CSE 4309 – Machine Learning  
Vassilis Athitsos  
Computer Science and Engineering Department  
University of Texas at Arlington

# Review: A Multiclass Example

- Suppose we have this training set:
  - $\mathbf{x}_1 = (0.5, 2.4, 8.3, 1.2, 4.5)^T$ ,  $q_1 = \text{dog}$
  - $\mathbf{x}_2 = (3.4, 0.6, 4.4, 6.2, 1.0)^T$ ,  $q_2 = \text{dog}$
  - $\mathbf{x}_3 = (4.7, 1.9, 6.7, 1.2, 3.9)^T$ ,  $q_3 = \text{cat}$
  - $\mathbf{x}_4 = (2.6, 1.3, 9.4, 0.7, 5.1)^T$ ,  $q_4 = \text{fox}$
  - $\mathbf{x}_5 = (8.5, 4.6, 3.6, 2.0, 6.2)^T$ ,  $q_5 = \text{cat}$
  - $\mathbf{x}_6 = (5.2, 8.1, 7.3, 4.2, 1.6)^T$ ,  $q_6 = \text{fox}$
- In this training set:
  - We have three classes.
  - Each training input is a five-dimensional vector.

# Review: Generating One-Hot Vectors

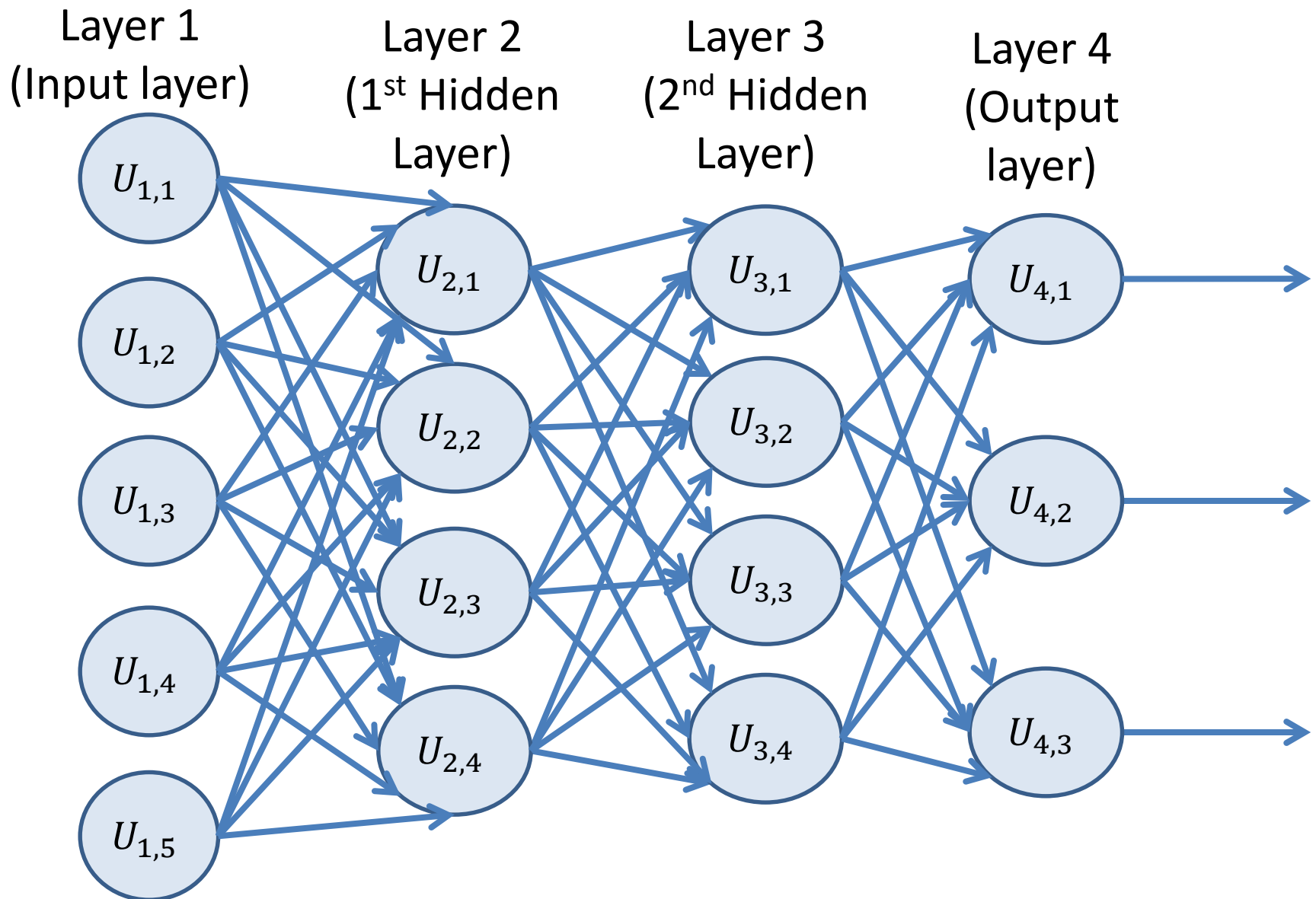
- Before we train a neural network, we must convert class labels to one-hot-vectors.
- Step 1: convert class labels  $q_n$  to new class labels  $s_n$ , which are integers between 1 and  $K$ .
  - $K$  is the number of classes,  $K = 3$  in our example.
- Step 2: convert labels  $s_n$  to  $K$ -dimensional one-hot vectors  $\mathbf{t}_n$ .
- Result: training set with new class labels  $s_n$  and one-hot vectors  $\mathbf{t}_n$ :

– $\mathbf{x}_1 = (0.5, 2.4, 8.3, 1.2, 4.5)^T$ ,	$s_1 = 1$	$\mathbf{t}_1 = (1, 0, 0)^T$
– $\mathbf{x}_2 = (3.4, 0.6, 4.4, 6.2, 1.0)^T$ ,	$s_2 = 1$	$\mathbf{t}_2 = (1, 0, 0)^T$
– $\mathbf{x}_3 = (4.7, 1.9, 6.7, 1.2, 3.9)^T$ ,	$s_3 = 2$	$\mathbf{t}_3 = (0, 1, 0)^T$
– $\mathbf{x}_4 = (2.6, 1.3, 9.4, 0.7, 5.1)^T$ ,	$s_4 = 3$	$\mathbf{t}_4 = (0, 0, 1)^T$
– $\mathbf{x}_5 = (8.5, 4.6, 3.6, 2.0, 6.2)^T$ ,	$s_5 = 2$	$\mathbf{t}_5 = (0, 1, 0)^T$
– $\mathbf{x}_6 = (5.2, 8.1, 7.3, 4.2, 1.6)^T$ ,	$s_6 = 3$	$\mathbf{t}_6 = (0, 0, 1)^T$

# Review: Multiclass Neural Networks

- For perceptrons, we saw that we can perform multiclass (i.e., for more than two classes) classification by training one perceptron for each class.
- For neural networks, we will train a SINGLE neural network, with MULTIPLE output units.
  - The number of output units will be equal to the number of classes.

# Review: A Network for Our Example



# Neural Network Notation

- $L$  is the total number of layers in the neural network.
- $D$  is the number of dimensions of the input.
- $K$  is the number of classes we want to recognize.
- Each unit, is denoted as  $U_{l,i}$ , where :
  - $1 \leq l \leq L$ , and  $l$  is the layer index
  - $1 \leq i$ , and  $i$  is the index of the unit within layer  $l$ .
  - Layer 1 is the input layer. Units  $U_{1,1}, \dots, U_{1,D}$  are the **input units**.
  - Layer  $L$  is the output layer. Units  $U_{L,1}, \dots, U_{L,K}$  are the **output units**.
- We denote by  $w_{l,i,j}$  the weight of the edge connecting the output of unit  $U_{l-1,j}$  to an input of unit  $U_{l,i}$ .
  - $j$  is the index of the unit in layer  $l - 1$ .
  - $i$  is the index of the unit in layer  $l$ .
- We denote by  $b_{l,i}$  the bias weight of  $U_{l,i}$ .

# Neural Network Notation

- We denote by  $J_l$  the number of units in layer  $l$ .
  - For the input layer,  $J_1 = D$ .
  - For the output layer,  $J_L = K$  (the number of classes).
  - For each hidden layer  $l$ ,  $J_l$  is a hyperparameter.
- We denote by  $a_{l,i}$  the weighted sum that is calculated at  $U_{l,i}$ .

$$a_{l,i} = b_{l,i} + \sum_{j=1}^{J_{l-1}} (w_{l,i,j} * z_{l-1,i})$$

- Note that this weighted sum is NOT applicable when  $l = 1$ , it only starts getting calculated for  $l \geq 2$ . So,  $a_{1,i}$  is not defined.
- We denote by  $z_{l,i}$  the output of unit  $U_{l,i}$ .
  - If  $l = 1$ , then  $z_{1,i} = x_i$ .
  - If  $l \geq 2$ , then  $z_{l,i} = \sigma(a_{l,i}) = \frac{1}{1+e^{-a_{l,i}}}$

# Neural Network Notation

- We denote by  $J_l$  the number of units in layer  $l$ .
  - For the input layer,  $J_1 = D$ .
  - For the output layer,  $J_L = K$  (the number of classes).
  - For each hidden layer  $l$ ,  $J_l$  is a hyperparameter.
- We denote by  $a_{l,i}$  the weighted sum that is calculated at  $U_{l,i}$ .

$$a_{l,i} = b_{l,i} + \sum_{j=1}^{J_{l-1}} (w_{l,i,j} * z_{l-1,i})$$

- Note that this weighted sum is NOT applicable when  $l = 1$ , it only starts getting calculated for  $l \geq 2$ . So,  $a_{1,i}$  is not defined.
- We denote by  $z_{l,i}$  the output of unit  $U_{l,i}$ .
  - If  $l = 1$ , then  $z_{1,i} = x_i$ .
  - If  $l \geq 2$ , then  $z_{l,i} = \sigma(a_{l,i}) = \frac{1}{1+e^{-a_{l,i}}}$

In these slides, we assume that we are using the sigmoid as activation function.



# Squared Error for Neural Networks

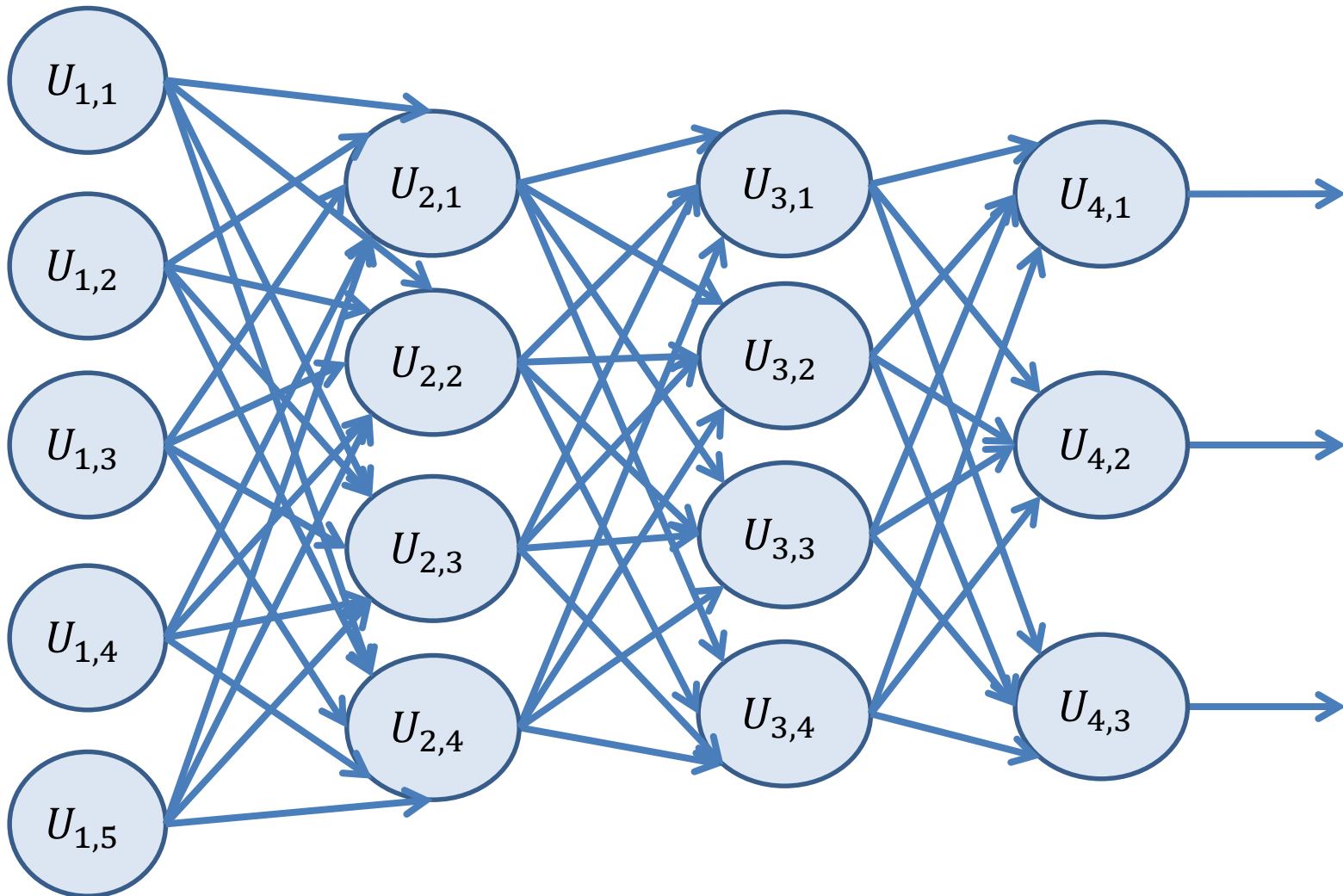
- We denote by  $\mathbf{b}$  the vector of all bias weights  $b_{l,i}$ .
- We denote by  $\mathbf{w}$  the vector of all weights  $w_{l,i,j}$ .
- We denote by  $E_n(\mathbf{b}, \mathbf{w})$  the contribution that training input  $\mathbf{x}_n$  makes to the overall error, given  $\mathbf{b}, \mathbf{w}$ .

$$E_n(\mathbf{b}, \mathbf{w}) = \frac{1}{2} \sum_{c=1}^K \left\{ (t_{n,c} - z_{L,c})^2 \right\}$$

- Remember:
  - Target output  $\mathbf{t}_n$  is a one-hot vector.
  - We denote by  $t_{n,c}$  the c-th dimension of target output  $\mathbf{t}_n$ .

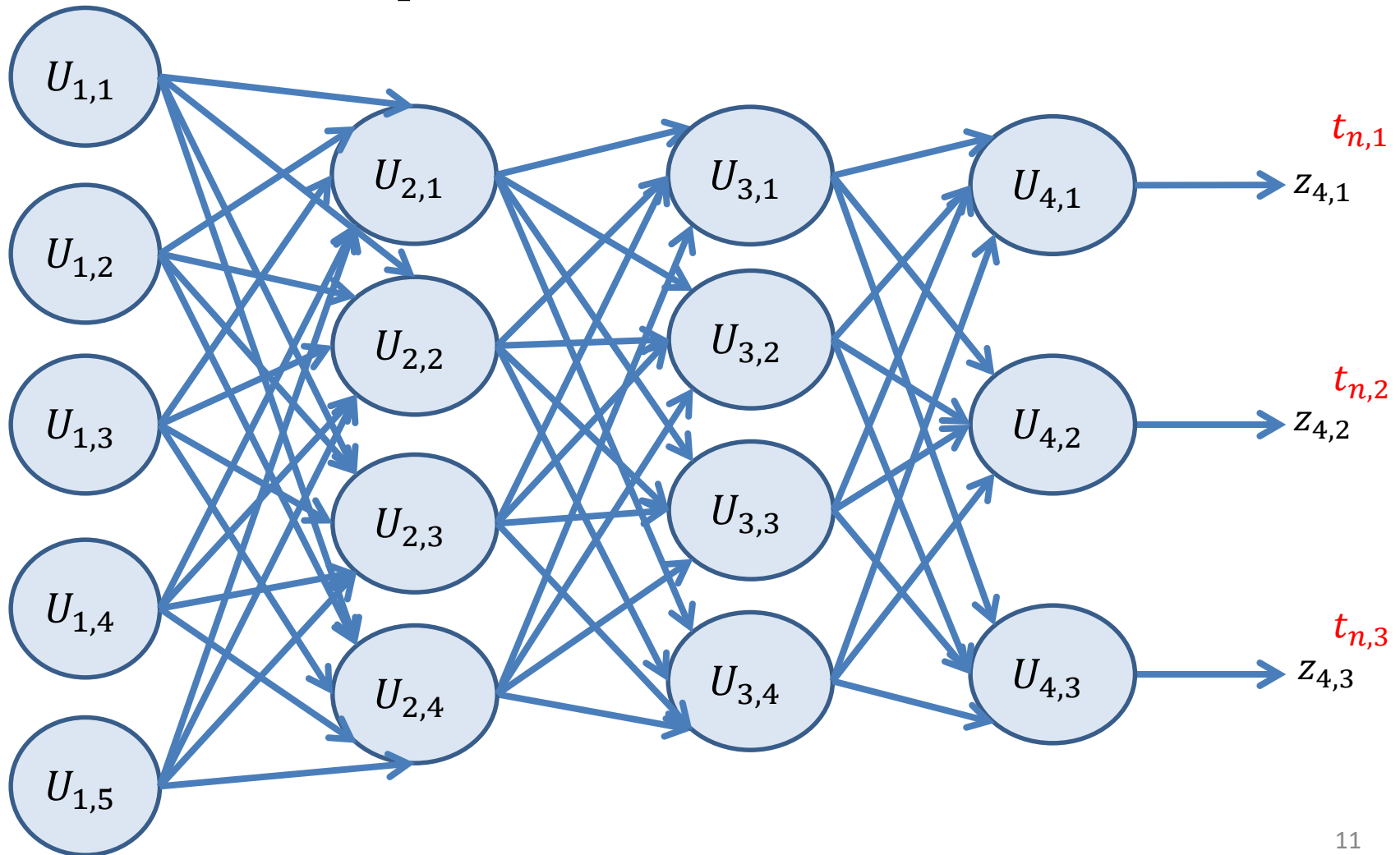
In our three-class example:

- $E_n = \frac{1}{2} \sum_{c=1}^K \left\{ (t_{n,c} - z_{L,c})^2 \right\}$
- What does each  $z_{L,c}$  correspond to?



In our three-class example:

- $E_n = \frac{1}{2} \sum_{c=1}^K \left\{ (t_{n,c} - z_{L,c})^2 \right\}$
- $E_n = \frac{1}{2} \left[ (t_{n,1} - z_{4,1})^2 + (t_{n,2} - z_{4,2})^2 + (t_{n,3} - z_{4,3})^2 \right]$



# Squared Error for Neural Networks

- We denote by  $E(\mathbf{b}, \mathbf{w})$  the overall error over all training examples for the network specified by  $\mathbf{b}, \mathbf{w}$ .

$$E(\mathbf{b}, \mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{b}, \mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \sum_{c=1}^K \left\{ (t_{n,c} - z_{L,c})^2 \right\}$$

- This is now a double summation.
  - We sum over all training examples  $\mathbf{x}_n$ .
  - For each  $\mathbf{x}_n$ , we sum over all perceptrons in the output layer.

# Training Neural Networks

- To train a neural network, we use gradient descent.
  - We follow the same approach of sequential learning that we followed for training single perceptrons.
- Given a training example  $\mathbf{x}_n$  and target output  $\mathbf{t}_n$ :
  - Compute the training error  $E_n(\mathbf{b}, \mathbf{w})$ .
  - Compute the gradients  $\frac{\partial E_n}{\partial \mathbf{b}}$  and  $\frac{\partial E_n}{\partial \mathbf{w}}$ .
  - Based on the gradients, we can update all weights  $\mathbf{b}$  and  $\mathbf{w}$ .
- The process of computing the gradient and updating neural network weights is called **backpropagation**.
- We will see the solution when we use the sigmoidal function as activation function  $h$ .

# Computing the Gradient

- Overall, we want to compute  $\frac{\partial E_n}{\partial \mathbf{b}}$  and  $\frac{\partial E_n}{\partial \mathbf{w}}$ .
- This is the same as computing:
  - For each bias weight  $b_{l,i}$ , the partial derivative  $\frac{\partial E_n}{\partial b_{l,i}}$
  - For each  $w_{l,i,j}$ , the partial derivative  $\frac{\partial E_n}{\partial w_{l,i,j}}$ .
- To compute  $\frac{\partial E_n}{\partial b_{l,i}}$  and  $\frac{\partial E_n}{\partial w_{l,i,j}}$ , we will use this strategy:
  - Decompose  $E_n$  into a composition of simpler functions.
  - Compute the derivative of each of those simpler functions.
  - Apply the chain rule to obtain  $\frac{\partial E_n}{\partial b_{l,i}}$  and  $\frac{\partial E_n}{\partial w_{l,i,j}}$ .

# Decomposing the Error Function

- Let  $U_{l,i}$  be a perceptron in the neural network.
- Define  $a_{l,i}(\mathbf{x}_n, \mathbf{b}, \mathbf{w})$  to be the weighted sum of the inputs of  $U_{l,i}$ , given input  $\mathbf{x}_n$  and given the current values of  $\mathbf{b}$  and  $\mathbf{w}$ .

$$a_{l,i}(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) = b_{l,i} + \sum_{j=1}^{J_{l-1}} \left( w_{l,i,j} * z_{l-1,j}(\mathbf{x}_n, \mathbf{w}) \right)$$

- Remember that  $z_{l,i}$  is the output of unit  $U_{l,i}$ , and it is obtained by applying the sigmoid function on top of  $a_{l,i}(\mathbf{x}_n, \mathbf{b}, \mathbf{w})$ .

$$z_{l,i} = \sigma \left( a_{l,i}(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) \right) = \frac{1}{1 + e^{-a_{l,i}(\mathbf{x}_n, \mathbf{b}, \mathbf{w})}}$$

# Decomposing the Error Function

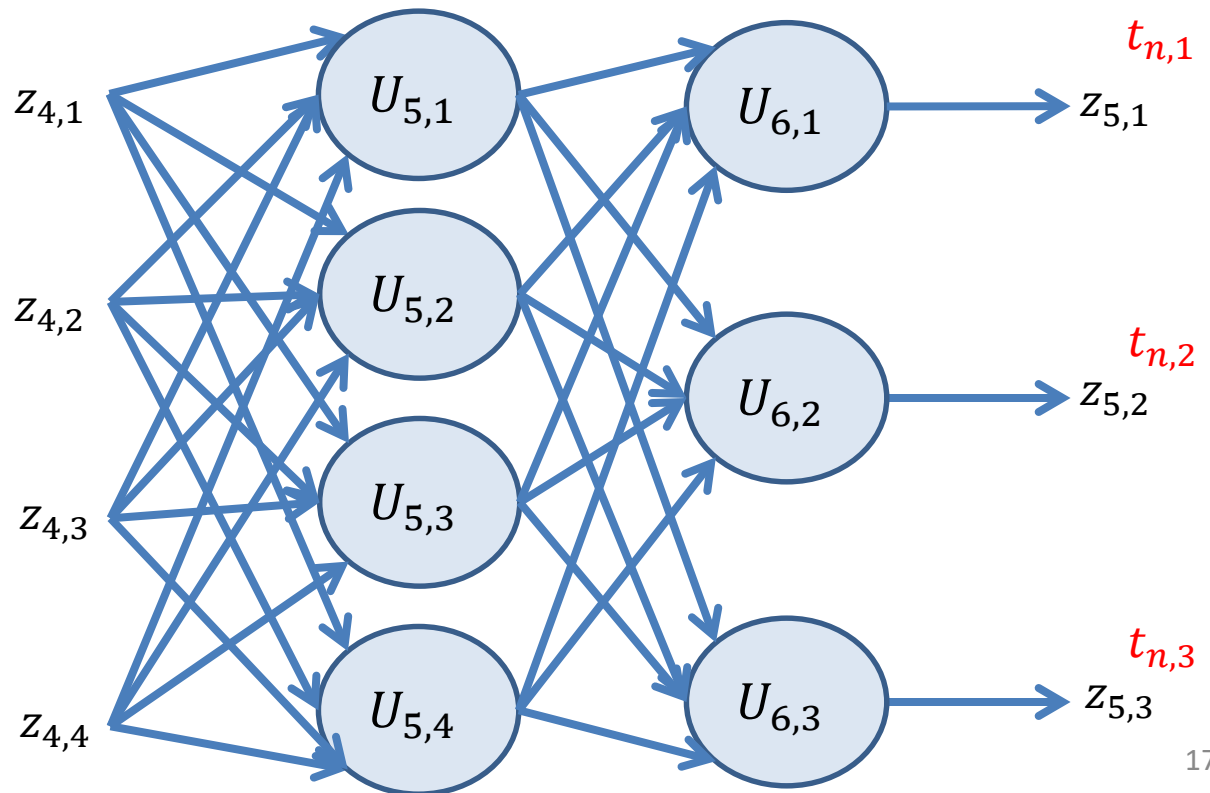
- Define  $\mathbf{z}_l$  to be a vector containing the outputs of all units at layer  $l$ .
  - Using our notation,  $\mathbf{z}_l = (z_{l,1}, z_{l,2}, \dots, z_{l,J_l})^T$ , where  $J_l$  is the number of units at layer  $l$ .
- Define function  $E_{n,l}(\mathbf{z}_l, \mathbf{b}, \mathbf{w})$  to be the error of the network given outputs  $\mathbf{z}_l$ .
- Intuition for  $E_{n,l}(\mathbf{z}_l, \mathbf{b}, \mathbf{w})$ :
  - Suppose that you know  $\mathbf{z}_l$ , and the weights for all layers after layer  $l$ .
  - Then, you can still compute the output of the network, and the error  $E_n(\mathbf{b}, \mathbf{w})$ .



# Visualizing Function $E_{n,l}$

- Suppose we know the target output  $\mathbf{t}_n$ , and all weights  $\mathbf{b}$  and  $\mathbf{w}$ .
- If we know the output  $\mathbf{z}_l$  of layer  $l$ :
  - Can we compute the output of the network?
  - Can we compute the error  $E_n$ ?

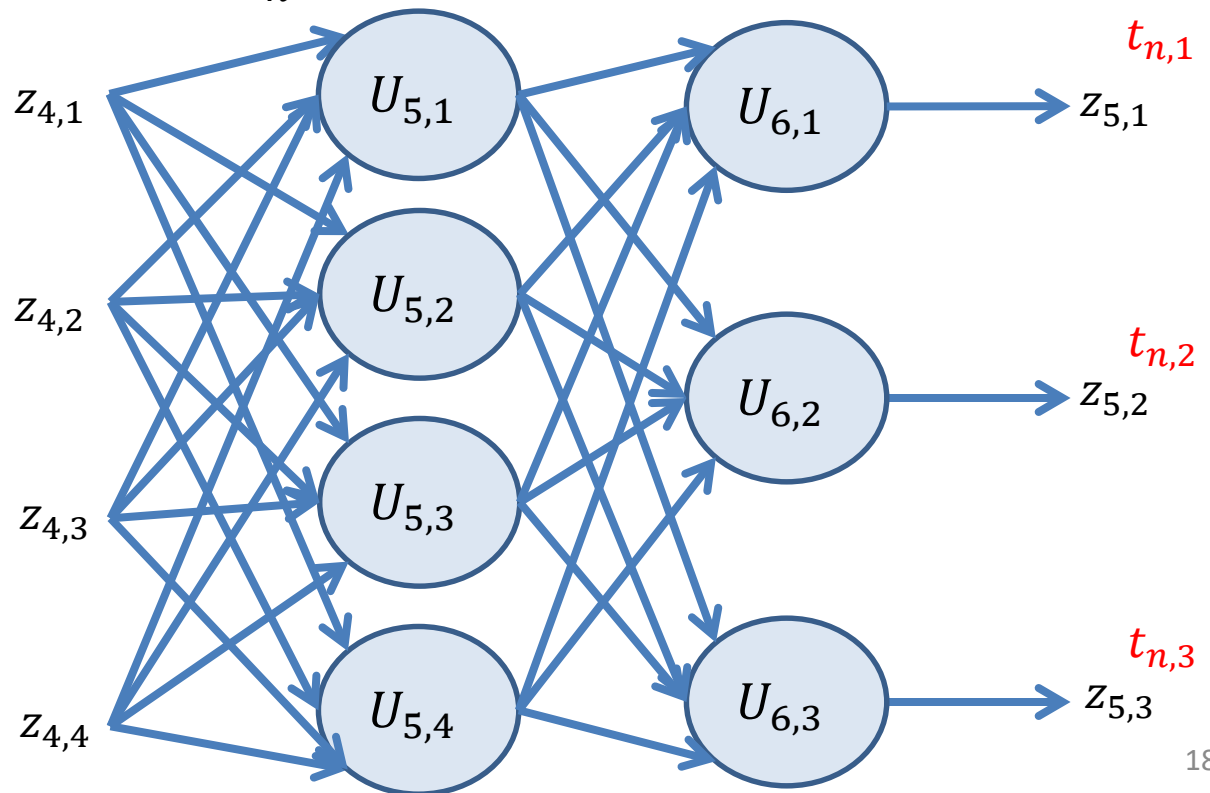
Layers 1 to 4,  
not shown.



# Visualizing Function $E_{n,l}$

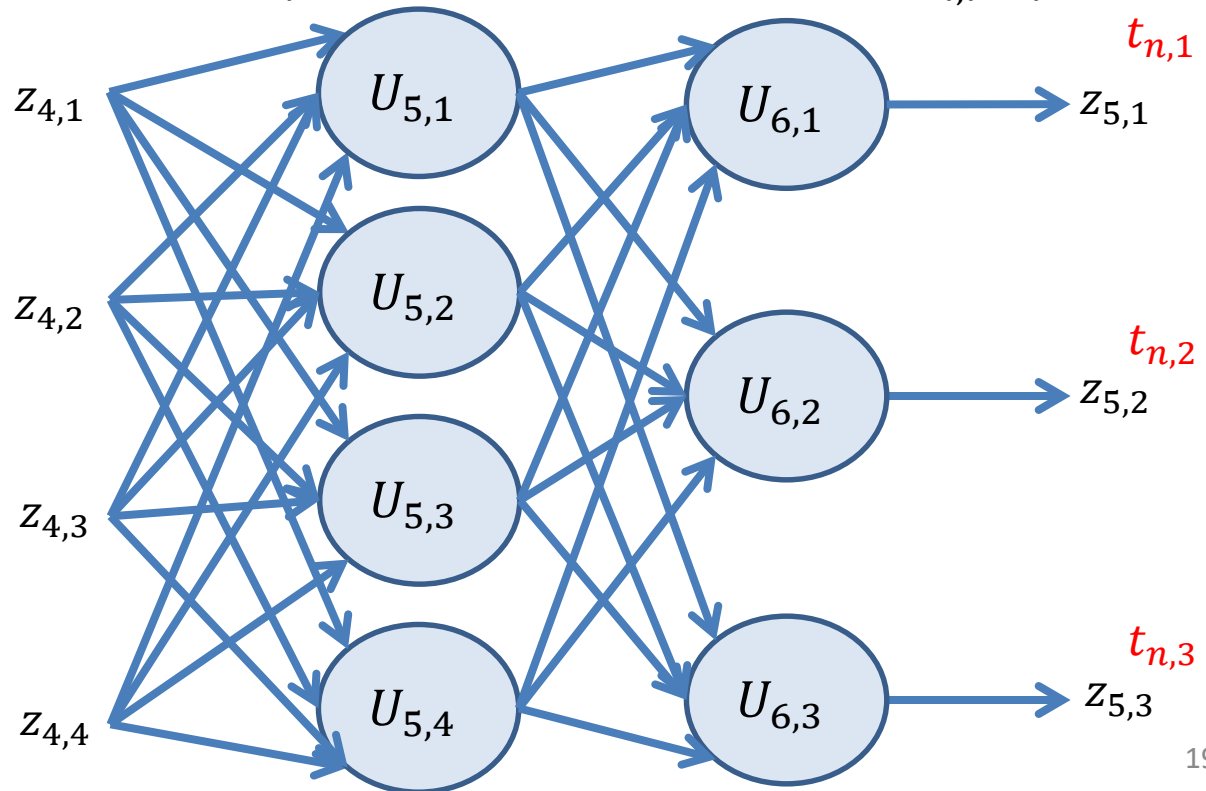
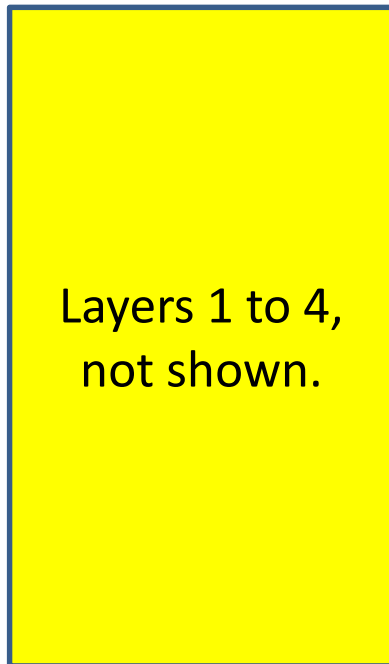
- In this example, the network has six layers.
  - We have no idea what happens in layers 1 to 4.
  - However, we are given the output  $\mathbf{z}_4$  of layer 4.
  - Can we compute the output of the network?
  - Can we compute the error  $E_n$ ?

Layers 1 to 4,  
not shown.



# Visualizing Function $E_{n,l}$

- In this example, given the output  $\mathbf{z}_4$  of layer 4, if we know all weights  $\mathbf{b}$  and  $\mathbf{w}$ , we can compute the final output and the error.
  - Given  $\mathbf{z}_4$ , we can compute the output  $\mathbf{z}_5$  of layer 5.
  - Given  $\mathbf{z}_5$ , we can compute the output  $\mathbf{z}_6$  of layer 6, (the output layer).
  - Given  $\mathbf{z}_6$  and target output  $\mathbf{t}_n$ , we can compute the error  $E_{n,l}(\mathbf{z}_l, \mathbf{b}, \mathbf{w})$ .



# Decomposing the Error Function

- We have three auxiliary functions:
  - $a_{l,i}(\mathbf{x}_n, \mathbf{b}, \mathbf{w})$
  - $\sigma(\alpha)$
  - $E_{n,l}(\mathbf{z}_l, \mathbf{b}, \mathbf{w})$
- Then,  $E_n$  is a composition of functions  $E_{n,l}$ ,  $\sigma$ ,  $a_{l,i}$ .

$$E_n(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) = E_{n,l}(\mathbf{z}_l, \mathbf{b}, \mathbf{w})$$

$$= E_{n,l} \left( \underbrace{\left( \sigma \left( a_{l,1}(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) \right), \dots, \sigma \left( a_{l,J_l}(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) \right) \right)}_{\mathbf{z}_l}, \mathbf{b}, \mathbf{w} \right)$$

# Computing the Gradient of $E_n$

$$\begin{aligned} E_n(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) &= E_{n,l} \left( (z_{l,1}, z_{l,2}, \dots, z_{l,J_l}), \mathbf{b}, \mathbf{w} \right) \\ &= E_{n,l} \left( \left( \sigma \left( a_{l,1}(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) \right), \dots, \sigma \left( a_{l,J_l}(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) \right) \right), \mathbf{b}, \mathbf{w} \right) \end{aligned}$$

- So,  $E_n$  is a composition of function  $E_{n,l}$ , function  $\sigma$ , and functions  $a_{l,i}$ .
- This allows us to compute  $\frac{\partial E_n}{\partial b_{l,i}}$  and  $\frac{\partial E_n}{\partial w_{l,i,j}}$  by applying the chain rule.

# Computing the Gradient of $E_n$

$$\begin{aligned} E_n(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) &= E_{n,l} \left( (z_{l,1}, z_{l,2}, \dots, z_{l,J_l}), \mathbf{b}, \mathbf{w} \right) \\ &= E_{n,l} \left( \left( \sigma \left( a_{l,1}(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) \right), \dots, \sigma \left( a_{l,J_l}(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) \right) \right), \mathbf{b}, \mathbf{w} \right) \end{aligned}$$

- Applying the chain rule:

$$\frac{\partial E_n}{\partial b_{l,i}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial b_{l,i}}$$

$$\frac{\partial E_n}{\partial w_{l,i,j}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial w_{l,i,j}}$$

# Computing $\frac{\partial a_{l,i}}{\partial b_{l,i}}$ and $\frac{\partial a_{l,i}}{\partial w_{l,i,j}}$

$$\frac{\partial E_n}{\partial b_{l,i}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial b_{l,i}}$$

$$\frac{\partial a_{l,i}}{\partial b_{l,i}} = \frac{\partial \left( b_{l,i} + \sum_{k=1}^{J_{l-1}} (w_{l,i,k} * z_{l-1,k}) \right)}{\partial b_{l,i}}$$

=???

This is actually very simple. It is of this form:

$$\frac{\partial (x + \text{stuff that is independent of } x)}{\partial x} = ???$$

# Computing $\frac{\partial a_{l,i}}{\partial b_{l,i}}$ and $\frac{\partial a_{l,i}}{\partial w_{l,i,j}}$

$$\frac{\partial E_n}{\partial b_{l,i}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial b_{l,i}}$$

$$\frac{\partial a_{l,i}}{\partial b_{l,i}} = \frac{\partial \left( b_{l,i} + \sum_{k=1}^{J_{l-1}} (w_{l,i,k} * z_{l-1,k}) \right)}{\partial b_{l,i}}$$

$$= 1$$

This is actually very simple. It is of this form:

$$\frac{\partial (x + \text{stuff that is independent of } x)}{\partial x} = 1$$



Computing  $\frac{\partial a_{l,i}}{\partial b_{l,i}}$  and  $\frac{\partial a_{l,i}}{\partial w_{l,i,j}}$

$$\frac{\partial E_n}{\partial w_{l,i,j}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial w_{l,i,j}}$$

$$\frac{\partial a_{l,i}}{\partial w_{l,i,j}} = \frac{\partial (b_{l,i} + \sum_{k=1}^{J_{l-1}} (w_{l,i,k} * z_{l-1,k}))}{\partial w_{l,i,j}}$$

= ???

How does  $w_{l,i,j}$  influence  $\sum_{k=1}^{J_{l-1}} (w_{l,i,k} * z_{l-1,k})$ ?

# Computing $\frac{\partial a_{l,i}}{\partial b_{l,i}}$ and $\frac{\partial a_{l,i}}{\partial w_{l,i,j}}$

$$\frac{\partial E_n}{\partial w_{l,i,j}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial w_{l,i,j}}$$

$$\frac{\partial a_{l,i}}{\partial w_{l,i,j}} = \frac{\partial (b_{l,i} + \sum_{k=1}^{J_{l-1}} (w_{l,i,k} * z_{l-1,k}))}{\partial w_{l,i,j}}$$

= ???

In  $\sum_{k=1}^{J_{l-1}} (w_{l,i,k} * z_{l-1,k})$ , at some point  $k = j$ . Then,  $w_{l,i,k}$  is multiplied by  $z_{l-1,j}$ . Based on that,  
 $\frac{\partial a_{l,i}}{\partial w_{l,i,j}} = ???$

Computing  $\frac{\partial a_{l,i}}{\partial b_{l,i}}$  and  $\frac{\partial a_{l,i}}{\partial w_{l,i,j}}$

$$\frac{\partial E_n}{\partial w_{l,i,j}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial w_{l,i,j}}$$

$$\frac{\partial a_{l,i}}{\partial w_{l,i,j}} = \frac{\partial (b_{l,i} + \sum_{k=1}^{J_{l-1}} (w_{l,i,k} * z_{l-1,k}))}{\partial w_{l,i,j}}$$

$$= z_{l-1,j}$$

# Computing $\frac{\partial z_{l,i}}{\partial a_{l,i}}$

$$\frac{\partial E_n}{\partial w_{l,i,j}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial w_{l,i,j}}$$

$$\frac{\partial z_{l,i}}{\partial a_{l,i}} = \frac{\partial (\sigma(a_{l,i}))}{\partial a_{l,i}} = \sigma(a_{l,i}) * (1 - \sigma(a_{l,i})) = z_{l,i} * (1 - z_{l,i})$$

- We just use the known formula for the derivative of  $\sigma$ .
  - One of the reasons we like using the sigmoidal function for activation is that its derivative has such a simple form.

# Computing $\frac{\partial E_{n,l}}{\partial z_{l,i}}$ , Case 1:

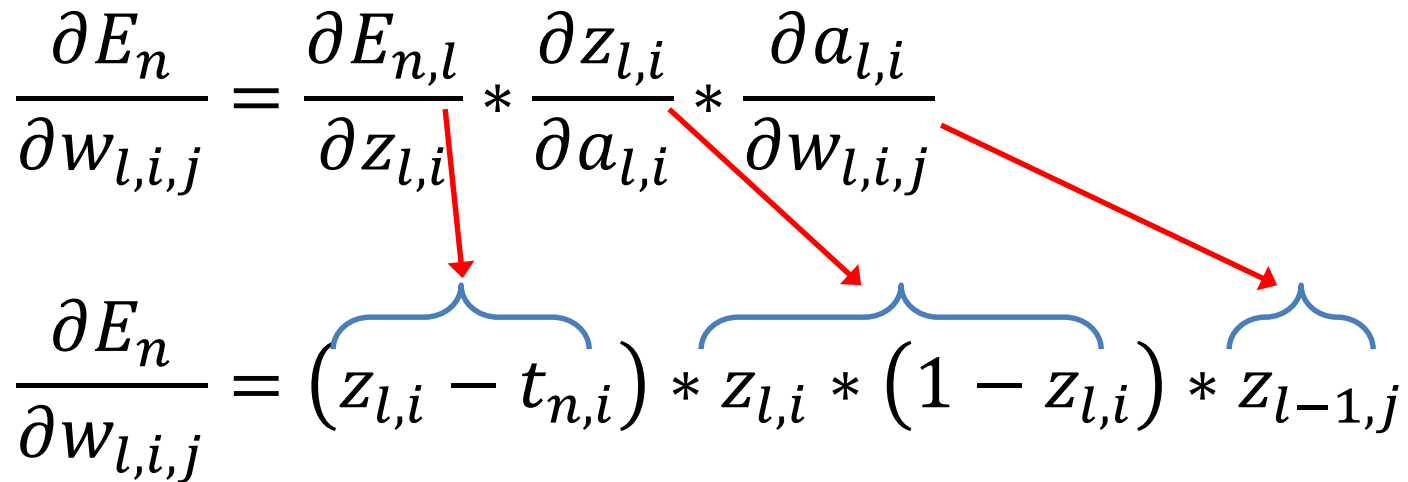
## If Unit $U_{l,i}$ Is an Output Unit

- If  $U_{l,i}$  is an output unit, then  $z_{l,i}$  is an output of the entire network.
- $z_{l,i}$  contributes to the error the term  $\frac{1}{2}(t_{n,i} - z_{l,i})^2$ .
- Therefore:

$$\frac{\partial E_{n,l}}{\partial z_{l,i}} = \frac{\partial \frac{1}{2}(t_{n,i} - z_{l,i})^2}{\partial z_{l,i}} = z_{l,i} - t_{n,i}$$

# Updating Weights of Output Units

- If  $U_{l,i}$  is an output unit, then we have computed all the terms we need for  $\frac{\partial E_n}{\partial w_{l,i,j}}$ .

$$\frac{\partial E_n}{\partial w_{l,i,j}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial w_{l,i,j}}$$

$$\frac{\partial E_n}{\partial w_{l,i,j}} = (z_{l,i} - t_{n,i}) * z_{l,i} * (1 - z_{l,i}) * z_{l-1,j}$$

- So, if  $U_{l,i}$  is an output unit, we update  $w_{l,i,j}$  as:

$$w_{l,i,j} = w_{l,i,j} - \eta (z_{l,i} - t_{n,i}) * z_{l,i} * (1 - z_{l,i}) * z_{l-1,j}$$

# Updating Weights of Output Units

- Similarly, if  $U_{l,i}$  is an output unit, we can compute  $\frac{\partial E_n}{\partial b_{l,i}}$ .

$$\frac{\partial E_n}{\partial b_{l,i}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial b_{l,i}}$$
$$\frac{\partial E_n}{\partial b_{l,i}} = (z_{l,i} - t_{n,i}) * z_{l,i} * (1 - z_{l,i}) * 1$$

- So, if  $U_{l,i}$  is an output unit, we update  $b_{l,i}$  using:

$$b_{l,i} = b_{l,i} - \eta (z_{l,i} - t_{n,i}) * z_{l,i} * (1 - z_{l,i})$$

# Computing $\frac{\partial E_{n,l}}{\partial z_{l,i}}$ , Case 2: If Unit $U_{l,i}$ Is a Hidden Unit

- We want to compute  $\frac{\partial E_{n,l}}{\partial z_{l,i}}$  when  $U_{l,i}$  is a hidden unit.
- We use the chain rule, to relate  $E_{n,l}$  to  $E_{n,l+1}$ .

$$\frac{\partial E_{n,l}}{\partial z_{l,i}} = \sum_{k=1}^{J_{l+1}} \left( \frac{\partial E_{n,l+1}}{\partial z_{l+1,k}} * \frac{\partial z_{l+1,k}}{\partial a_{l+1,k}} * \frac{\partial a_{l+1,k}}{\partial z_{l,i}} \right)$$

- We need to compute these three terms.



Computing  $\frac{\partial E_{n,l}}{\partial z_{l,i}}$ , Case 2:

If Unit  $U_{l,i}$  Is a Hidden Unit

$$\frac{\partial E_{n,l}}{\partial z_{l,i}} = \sum_{k=1}^{J_{l+1}} \left( \frac{\partial E_{n,l+1}}{\partial z_{l+1,k}} * \frac{\partial z_{l+1,k}}{\partial a_{l+1,k}} * \frac{\partial a_{l+1,k}}{\partial z_{l,i}} \right)$$

$$\begin{aligned} \frac{\partial a_{l+1,k}}{\partial z_{l,i}} &= \frac{\partial \left( b_{l+1,i} + \sum_{k=1}^{J_{l+1}} \left( w_{l+1,i,k} * z_{l,k} \left( a_{l,k}(\mathbf{x}_n, \mathbf{b}, \mathbf{w}) \right) \right) \right)}{\partial z_j} \\ &= w_{l+1,k,i} \end{aligned}$$

# Computing $\frac{\partial E_{n,l}}{\partial z_{l,i}}$ , Case 2: If Unit $U_{l,i}$ Is a Hidden Unit

$$\frac{\partial E_{n,l}}{\partial z_{l,i}} = \sum_{k=1}^{J_{l+1}} \left( \frac{\partial E_{n,l+1}}{\partial z_{l+1,k}} * \frac{\partial z_{l+1,k}}{\partial a_{l+1,k}} * \frac{\partial a_{l+1,k}}{\partial z_{l,i}} \right)$$

$$\begin{aligned} \frac{\partial z_{l+1,k}}{\partial a_{l+1,k}} &= \frac{\partial \left( \sigma(a_{l+1,k}) \right)}{\partial a_{l+1,k}} = \sigma(a_{l+1,k}) * \left( 1 - \sigma(a_{l+1,k}) \right) \\ &= z_{l,i} * (1 - z_{l,i}) \end{aligned}$$

Derivative of the  
sigmoid function

# Computing $\frac{\partial E_{n,l}}{\partial z_{l,i}}$ , Case 2: If Unit $U_{l,i}$ Is a Hidden Unit

$$\frac{\partial E_{n,l}}{\partial z_{l,i}} = \sum_{k=1}^{J_{l+1}} \left( \frac{\partial E_{n,l+1}}{\partial z_{l+1,k}} * \frac{\partial z_{l+1,k}}{\partial a_{l+1,k}} * \frac{\partial a_{l+1,k}}{\partial z_{l,i}} \right)$$

- We can plug in our results for  $\frac{\partial z_{l+1,k}}{\partial a_{l+1,k}}$  and  $\frac{\partial a_{l+1,k}}{\partial z_{l,i}}$ .
- So, the formula becomes:

$$\frac{\partial E_{n,l}}{\partial z_{l,i}} = \sum_{k=1}^{J_{l+1}} \left( \frac{\partial E_{n,l+1}}{\partial z_{l+1,k}} * z_{l+1,k} * (1 - z_{l+1,k}) * w_{l+1,k,i} \right)$$

# Computing $\frac{\partial E_{n,l}}{\partial z_{l,i}}$ , Case 2:

## If Unit $U_{l,i}$ Is a Hidden Unit

$$\frac{\partial E_{n,l}}{\partial z_{l,i}} = \sum_{k=1}^{J_{l+1}} \left( \frac{\partial E_{n,l+1}}{\partial z_{l+1,k}} * z_{l+1,k} (1 - z_{l+1,k}) * w_{l+1,k,i} \right)$$

- Notice that  $\frac{\partial E_{n,l}}{\partial z_{l,i}}$  is defined using  $\frac{\partial E_{n,l+1}}{\partial z_{l+1,k}}$ .
  - This is a **recursive** definition. To compute the values for layer  $l$ , we use the values from the **next** layer (i.e., layer  $l + 1$ ).
  - This is why the whole algorithm is called **backpropagation**.
    - We propagate computations from the output layer backwards towards the input layer.

# Computing $\frac{\partial E_n}{\partial w_{l,i,j}}$ for Hidden Units

- From the previous slides, we have these formulas:

$$\begin{aligned}\frac{\partial E_n}{\partial w_{l,i,j}} &= \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial w_{l,i,j}} \\ &= \frac{\partial E_{n,l}}{\partial z_{l,i}} * (z_{l,i} - t_{n,i}) * z_{l,i} * (1 - z_{l,i}) * z_{l-1,j}\end{aligned}$$

$$\frac{\partial E_{n,l}}{\partial z_{l,i}} = \sum_{k=1}^{J_{l+1}} \left( \frac{\partial E_{n,l+1}}{\partial z_{l+1,k}} * z_{l+1,k} * (1 - z_{l+1,k}) * w_{l+1,k,i} \right)$$

- We can combine these formulas, to compute  $\frac{\partial E_n}{\partial w_{l,i,j}}$  for any weight of any hidden unit.

# Computing $\frac{\partial E_n}{\partial b_{l,i}}$ for Hidden Units

- The formula for  $\frac{\partial E_n}{\partial b_{l,i}}$  is similar, we just replace  $\frac{\partial a_{l,i}}{\partial w_{l,i,j}}$  with  $\frac{\partial a_{l,i}}{\partial b_{l,i}}$ .

$$\frac{\partial E_n}{\partial b_{l,i}} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}} * \frac{\partial a_{l,i}}{\partial b_{l,i}}$$

$$= \frac{\partial E_{n,l}}{\partial z_{l,i}} * (z_{l,i} - t_{n,i}) * z_{l,i} * (1 - z_{l,i}) * 1$$

$$\frac{\partial E_{n,l}}{\partial z_{l,i}} = \sum_{k=1}^{J_{l+1}} \left( \frac{\partial E_{n,l+1}}{\partial z_{l+1,k}} * z_{l+1,k} * (1 - z_{l+1,k}) * w_{l+1,k,i} \right)$$

# Simplifying Notation

- The previous formulas are sufficient and will work, but look complicated.
- We can simplify the formulas considerably, by defining:

$$\delta_{l,i} = \frac{\partial E_{n,l}}{\partial z_{l,i}} * \frac{\partial z_{l,i}}{\partial a_{l,i}}$$

- Then, if we combine calculations we already did:

- If  $U_{l,i}$  is an output unit, then:

$$\delta_{l,i} = (z_{l,i} - t_{n,i}) * z_{l,i} * (1 - z_{l,i})$$

- If  $U_{l,i}$  is a hidden unit, then:

$$\delta_{l,i} = \left( \sum_{k=1}^{J_{l+1}} (\delta_{l+1,k} * w_{l+1,k,i}) \right) * z_{l,i} * (1 - z_{l,i})$$

# Final Backpropagation Formulas

- Using the definition of  $\delta_{l,i}$  from the previous slide, we finally get very simple formulas:

$$\frac{\partial E_n}{\partial w_{l,i,j}} = \delta_{l,i} * z_{l-1,j}$$

$$\frac{\partial E_n}{\partial b_{l,i}} = \delta_{l,i}$$

- Therefore, given a training input  $x_n$ , and given a positive learning rate  $\eta$ , weights  $w_{l,i,j}$  and  $b_{l,i}$  are updated as follows:

$$w_{l,i,j} = w_{l,i,j} - \eta * \delta_{l,i} * z_{l-1,j}$$

$$b_{l,i} = b_{l,i} - \eta * \delta_{l,i}$$



# Backpropagation for One Object

## Step 1: Initialize Input Layer

- We will now see how to apply backpropagation, step by step, in pseudocode style, for a single training object.
- **NOTE: IN THE PSEUDOCODE, ARRAY INDICES START AT 1, NOT 0.**
- First, given a training example  $x_n$ , and its target output  $t_n$ , we must initialize the input units:

// 2D array  $z$  will store, for every unit  $U_{l,i}$ , its output

- `double ** Z = new double*[L]` //  $L$  is the number of layers
- `Z[1] = new double[D]` //  $D$  is the dimensionality of  $x_n$

// Update the input layer, set inputs equal to  $x_n$ .

- For  $I = 1$  to  $D$ :
  - `z[1][I] =  $x_{n,i}$`  //  $x_{n,i}$  is the  $i$ -th dimension of training input  $x_n$ .

# Backpropagation for One Object

## Step 2: Compute Outputs

// we create a 2D array  $a$ , which will store, for every  
// unit  $U_{l,i}$ , the weighted sum of the inputs of  $U_{l,i}$ .

- `double ** a = new double*[L]`

// Update the rest of the layers:

- For  $l = 2$  to  $L$ : //  $L$  is the number of layers
  - `a[l]= new double[Jl]` //  $J_l$  is the number of units in layer  $l$
  - `z[l]= new double[Jl]`
  - For each unit  $U_{l,i}$  in layer  $l$ :
    - $a[l][i] = b_{l,i} + \sum_{j=1}^{J_{l-1}} (w_{l,i,j} z[l-1][j])$  // weighted sum
    - $z[l][i] = \sigma(a[l][i]) = \frac{1}{1 + e^{-a[l][i]}}$  // output of unit  $U_{l,i}$

# Backpropagation for One Object

## Step 3: Compute New $\delta$ Values

// array  $\delta$  will store, for every unit  $U_{l,i}$ , value  $\delta_{l,i}$ .

- `double **  $\delta$  = new double*[L]`
- `$\delta[L]$  = new double*[K] // K is the number of classes`
- For each output unit  $U_{L,i}$ :
  - $\delta[L][i] = (z[L][i] - t_{n,i}) * z[L][i] * (1 - z[L][i])$
- For  $l = L - 1$  to 2: **// MUST be in decreasing order of  $l$** 
  - `$\delta[l]$  = new double[Jl] // Jl is the number of units in layer  $l$`
  - For each unit  $U_{l,i}$  in layer  $l$ :
    - $\delta[l][i] = \left( \sum_{k=1}^{J_{l+1}} (\delta[l+1][k] * w_{l+1,k,i}) \right) * z[l][i] * (1 - z[l][i])$

# Backpropagation for One Object

## Step 4: Update Weights

- For  $l = 2$  to  $L$ : *// Order does not matter here, we can go // from 2 to  $L$  or from  $L$  to 2.*
  - For  $i = 1$  to  $J_l$ :
    - $b_{l,i} = b_{l,i} - \eta * \delta[l][i]$
    - For  $j = 1$  to  $J_{l-1}$ :
      - $w_{l,i,j} = w_{l,i,j} - \eta * \delta[l][i] * z[l-1][j]$

**IMPORTANT:** Do Step 3 before Step 4. Do NOT do steps 3 and 4 as a single loop.

- All  $\delta$  values must be computed using the old values of weights.
- Then, all weights must be updated using the new  $\delta$  values .

# Backpropagation Summary

- Inputs:
  - N D-dimensional training objects  $\mathbf{x}_1, \dots, \mathbf{x}_N$ .
  - The associated target values  $\mathbf{t}_1, \dots, \mathbf{t}_N$ , which are K-dimensional vectors.
- 1. Initialize weights  $b_{l,i}$  and  $w_{l,i,j}$  to small random numbers.
  - For example, set each  $b_{l,i}$  and  $w_{l,i,j}$  to a value between -0.1 and 0.1.
- 2.  $\text{last\_error} = E(\mathbf{b}, \mathbf{w})$  // sum over all training examples
- 3. For  $n = 1$  to  $N$ :
  - Given  $\mathbf{x}_n$ , update weights  $b_{l,i}$  and  $w_{l,i,j}$  as described in the previous slides.
- 4.  $\text{err} = E(\mathbf{b}, \mathbf{w})$  // sum over all training examples
- 5. If  $|\text{err} - \text{last\_error}| < \text{threshold}$ , **exit**. // threshold can be 0.00001.
- 6. Else:  $\text{last\_error} = \text{err}$ , go to step 3.

# Classification with Neural Networks

- Suppose we have  $K$  classes  $C_1, \dots, C_K$ , where  $K > 2$ .
- Each class  $C_k$  corresponds to an output unit  $U_{L,k}$ .
- Given a test pattern  $\mathbf{x}$  to classify:
  - Compute outputs for all units of the network, working from the input layer towards the output layer.
- Find the output unit  $U_{L,k}$  with the highest output  $z_{L,k}$ .
- Return class  $C_k$ .

# Structure of Neural Networks

- Backpropagation describes how to learn weights.
- However, it does not describe how to learn the structure:
  - How many layers?
  - How many units at each layer?
- These are parameters that we have to choose somehow.
- A good way to choose such parameters is by using a validation set, containing examples and their class labels.
  - The validation set should be separate (disjoint) from the training set.

# Structure of Neural Networks

- To choose the best structure for a neural network using a validation set, we try many different parameters (number of layers, number of units per layer).
- For each choice of parameters:
  - We train several neural networks using backpropagation.
  - We measure how well each neural network classifies the validation examples.
  - Why not train just one neural network?



# Structure of Neural Networks

- To choose the best structure for a neural network using a validation set, we try many different parameters (number of layers, number of units per layer).
- For each choice of parameters:
  - We train several neural networks using backpropagation.
  - We measure how well each neural network classifies the validation examples.
  - Why not train just one neural network?
  - Each network is randomly initialized, so after backpropagation it can end up being different from the other networks.
- At the end, we select the neural network that did best on the validation set.