

Deep Learning II: Backpropagation

Swati Mishra

Applications of Machine Learning (4AL3)

Fall 2024



ENGINEERING

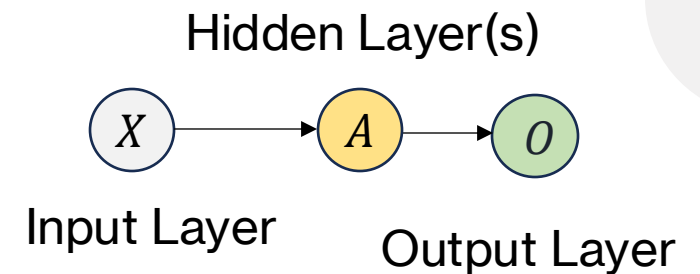
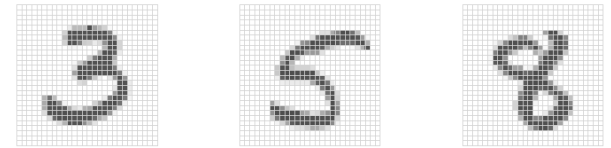
Review

- Neural Network fundamentals
- Architecture: Input, Output, Hidden Units
- Activation functions
- Design considerations for neural architecture

Neural Networks: Architecture

- To design architecture
 - We need a input layer.
 - Size: 784
 - We need output layers.
 - Size: 10 (0-9)
 - We need hidden layers.
 - Size: L1 = 512 , L2 = 512
 - Activation Function: ReLu

0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9



Neural Networks: Architecture

- To design architecture
 - We need a input layer.
 - Size: 784
 - We need output layers.
 - Size: 10 (0-9)
 - We need hidden layers.
 - Size: L1 = 512 , L2 = 512
 - Activation Function: ReLu

```
class NeuralNetwork(nn.Module):
    def __init__(self):
        super().__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(28*28, 512),
            nn.ReLU(),
            nn.Linear(512, 512),
            nn.ReLU(),
            nn.Linear(512, 10),
        )

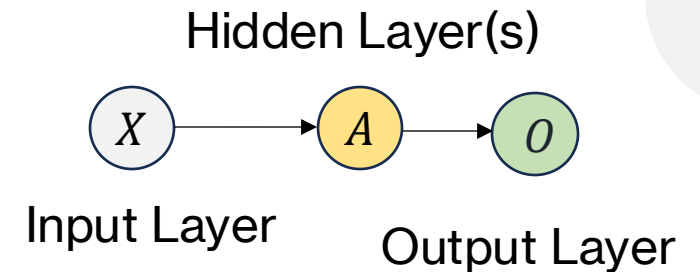
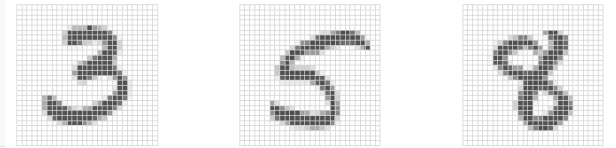
[33] ✓ 0.0s

model = NeuralNetwork()
print(model)

[34] ✓ 0.0s

... NeuralNetwork(
  (flatten): Flatten(start_dim=1, end_dim=-1)
  (linear_relu_stack): Sequential(
    (0): Linear(in_features=784, out_features=512, bias=True)
    (1): ReLU()
    (2): Linear(in_features=512, out_features=512, bias=True)
    (3): ReLU()
    (4): Linear(in_features=512, out_features=10, bias=True)
  )
)
```

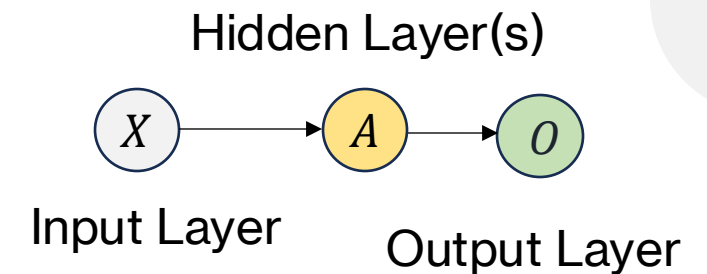
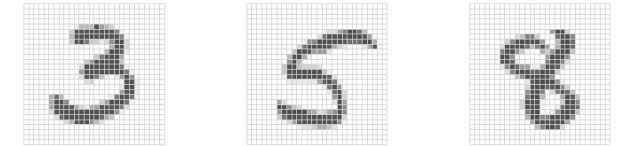
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9



Neural Networks: Training

- To train the model
 - We need a learning algorithm
 - Type:??
 - We need a cost function
 - Type:??
 - We need a training objective
 - Type: ??

0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9



Forward Propagation

- Forward propagation algorithm in neural network

Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: \mathbf{x} , the input to process

Require: \mathbf{y} , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

for $k = 1, \dots, l$ **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

end for

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

Source: Deep Learning Book

Forward Propagation

- Forward propagation algorithm in neural network

Require: Network depth, l

Require: $W^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $b^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: x , the input to process

Require: y , the target output

$$h^{(0)} = x$$

for $k = 1, \dots, l$ **do**

$$a^{(k)} = b^{(k)} + W^{(k)} h^{(k-1)}$$

$$h^{(k)} = f(a^{(k)})$$

end for

$$\hat{y} = h^{(l)}$$

$$J = L(\hat{y}, y) + \lambda \Omega(\theta)$$

```
import numpy as np
w = np.array([[0.9, 0.3, 0.4],
              [0.2, 0.8, 0.2],
              [0.1, 0.5, 0.6]])
x = np.array([[0.9], [0.1], [0.8]])
w.dot(x)
✓ 0.0s

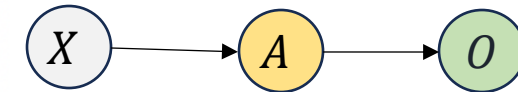
def sigmoid(z):
    return 1/(1 + np.exp(-z))
[6] ✓ 0.0s
```

$$X = \begin{bmatrix} 0.9 \\ 0.1 \\ 0.8 \end{bmatrix}$$

Input units

$$O = \begin{bmatrix} 0.72 \\ 0.70 \\ 0.77 \end{bmatrix}$$

Output units



$$A = \begin{bmatrix} 0.76 \\ 0.60 \\ 0.65 \end{bmatrix}$$

$$X_{hidden} = W_{XA} \cdot X$$

$$X_{output} = W_{AO} \cdot X_{hidden}$$

$$W_{XA} = \begin{bmatrix} 0.9 & 0.3 & 0.4 \\ 0.2 & 0.8 & 0.2 \\ 0.1 & 0.5 & 0.6 \end{bmatrix}$$

$$W_{AO} = \begin{bmatrix} 0.3 & 0.7 & 0.5 \\ 0.6 & 0.5 & 0.2 \\ 0.8 & 0.1 & 0.9 \end{bmatrix}$$

Forward Propagation

- Forward propagation algorithm in neural network

Require: Network depth, l

Require: $W^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $b^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: x , the input to process

Require: y , the target output

$$h^{(0)} = x$$

for $k = 1, \dots, l$ **do**

$$a^{(k)} = b^{(k)} + W^{(k)} h^{(k-1)}$$

$$h^{(k)} = f(a^{(k)})$$

end for

$$\hat{y} = h^{(l)}$$

$$J = L(\hat{y}, y) + \lambda \Omega(\theta)$$

```
import numpy as np
w = np.array([[0.9, 0.3, 0.4],
              [0.2, 0.8, 0.2],
              [0.1, 0.5, 0.6]])
x = np.array([[0.9], [0.1], [0.8]])
w.dot(x)
✓ 0.0s

def sigmoid(z):
    return 1/(1 + np.exp(-z))
[6] ✓ 0.0s
```

$$X = \begin{bmatrix} 0.9 \\ 0.1 \\ 0.8 \end{bmatrix}$$

Input units

$$O = \begin{bmatrix} 0.72 \\ 0.70 \\ 0.77 \end{bmatrix}$$

Output units



$$A = \begin{bmatrix} 0.76 \\ 0.60 \\ 0.65 \end{bmatrix}$$

$$X_{hidden} = W_{XA} \cdot X$$

$$X_{output} = W_{AO} \cdot X_{hidden}$$

$$W_{XA} = \begin{bmatrix} 0.9 & 0.3 & 0.4 \\ 0.2 & 0.8 & 0.2 \\ 0.1 & 0.5 & 0.6 \end{bmatrix}$$

$$W_{AO} = \begin{bmatrix} 0.3 & 0.7 & 0.5 \\ 0.6 & 0.5 & 0.2 \\ 0.8 & 0.1 & 0.9 \end{bmatrix}$$

Neural Network as Graph

- To understand Backpropagation, let us consider the neural network to be a computational graph.
- Each node (or unit or neuron) refers to a variable.
 - Variables can be scalar, matrix, vector, tensor, or something else.

Neural Network as Graph

- To understand Backpropagation, let us consider the neural network to be a computational graph.
- Each node (or unit or neuron) refers to a variable.
 - Variables can be **scalar**, matrix, vector, tensor, or something else.

A scalar is just a single number

Neural Network as Graph

- To understand Backpropagation, let us consider the neural network to be a computational graph.
- Each node (or unit or neuron) refers to a variable.
 - Variables can be scalar, **matrix**, vector, tensor, or something else.

A matrix is a 2-D array of numbers, so each element is identified by two indices instead of just one.

Neural Network as Graph

- To understand Backpropagation, let us consider the neural network to be a computational graph.
- Each node (or unit or neuron) refers to a variable.
 - Variables can be scalar, matrix, **vector**, tensor, or something else.

A vector is an array of numbers arranged in order

Neural Network as Graph

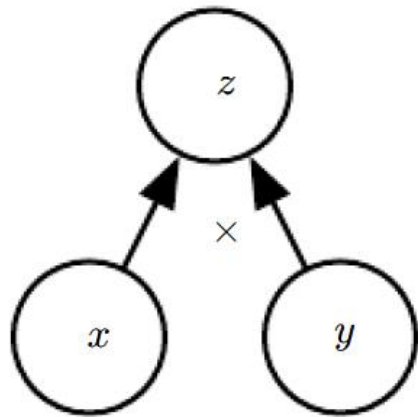
- To understand Backpropagation, let us consider the neural network to be a computational graph.
- Each node (or unit or neuron) refers to a variable.
 - Variables can be scalar, matrix, vector, **tensor**, or something else.

A tensor is an an array with more than two axes. It is essentially an array of numbers arranged on a regular grid with a variable number of axes.

Neural Network as Graph

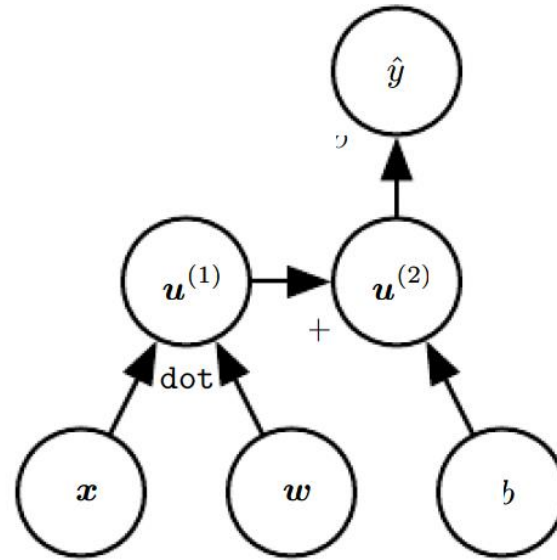
- To understand Backpropagation, let us consider the neural network to be a computational graph.
- Each node (or unit or neuron) refers to a variable.
 - Variables can be scalar, matrix, vector, tensor, or something else.
- An operation is a function of one or more variables
 - Operation returns a single output variable
- If an operation applied to x gives y , then x and y have an edge.

Neural Network as Computational Graph



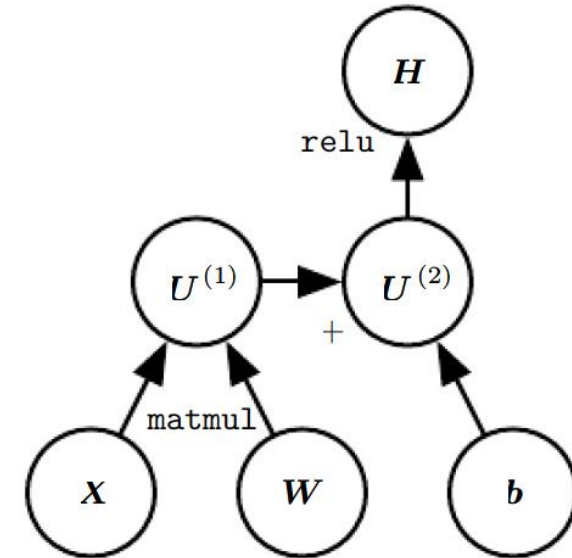
(a)

$$z = xy$$



(b)

$$y' = \sigma(x^T w + b)$$



(c)

$$H = \max\{0, xw + b\}$$

Chain Rule

- To understand Backpropagation, let us consider the neural network to be a computational graph.
- Each node (or unit or neuron) refers to a variable.
 - Variables can be scalar, matrix, vector, tensor, or something else.
- An operation is a function of one or more variables
 - Operation returns a single output variable
- If an operation applied to x gives y , then x and y have an edge.
- Chain rule in calculus: if $y = g(x)$ and $z = f(g(x)) = f(y)$ then, $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

Chain Rule

- Chain rule in calculus: if $y = g(x)$ and $z = f(g(x)) = f(y)$ then, $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$
- Each node (or unit or neuron) refers to a variable.
 - Variables can be scalar, matrix, vector, tensor, or something else
- If we want to calculate the gradient of variable x , we can rewrite the above in a vector notation,

$$\nabla_x z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^\top \nabla_y z, \quad \text{This is Jacobian gradient product}$$

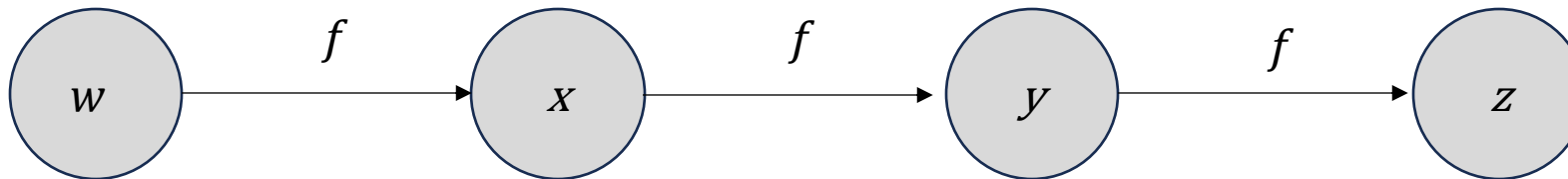
- Backpropagation algorithm consists of performing this product for each operation in the graph.

Chain Rule

- Chain rule in calculus: if $y = g(x)$ and $z = f(g(x)) = f(y)$ then, $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$ $\nabla_x z = \left(\frac{\partial y}{\partial x} \right)^\top \nabla_y z,$

This is Jacobian gradient product

- Consider the computational graph below:



$$x = f(w)$$

$$y = f(x)$$

$$z = f(y)$$

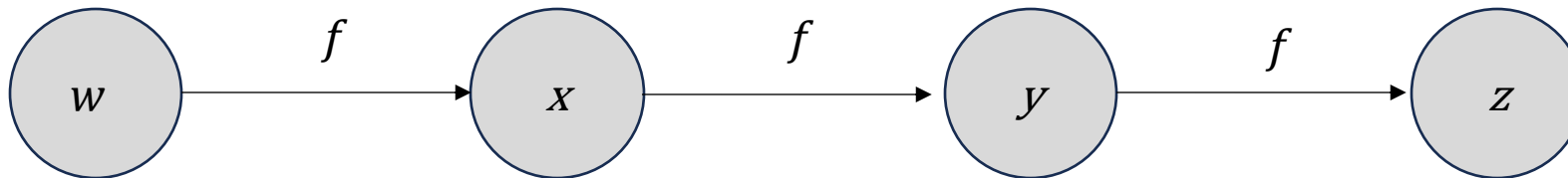
$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} = f'(y) f'(x) f'(w) = f'(f(f(w))) f'(f(w)) f'(w)$$

Chain Rule

- Chain rule in calculus: if $y = g(x)$ and $z = f(g(x)) = f(y)$ then, $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$ $\nabla_x z = \left(\frac{\partial y}{\partial x} \right)^\top \nabla_y z,$

This is Jacobian gradient product

- Consider the computational graph below:



$$x = f(w)$$

$$y = f(x)$$

$$z = f(y)$$

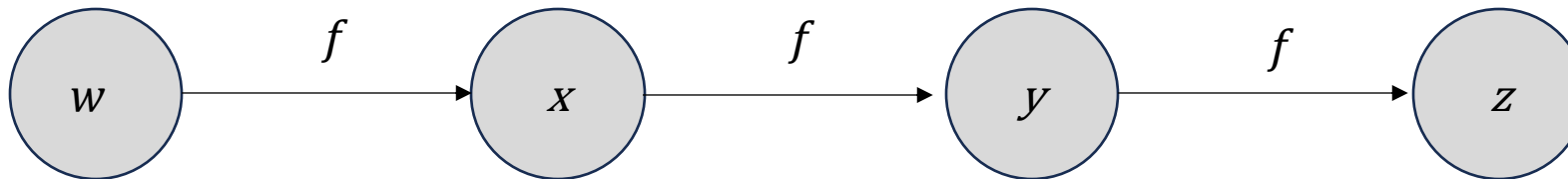
$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} = f'(y) f'(x) f'(w) = f'(f(f(w))) f'(f(w)) f'(w)$$

Chain Rule

- Chain rule in calculus: if $y = g(x)$ and $z = f(g(x)) = f(y)$ then, $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$ $\nabla_x z = \left(\frac{\partial y}{\partial x} \right)^\top \nabla_y z$,

This is Jacobian gradient product

- Consider the computational graph below:



$$x = f(w)$$

$$y = f(x)$$

$$z = f(y)$$

$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} = f'(y)f'(x) f'(w) = f'(f(f(w))) f'(f(w)) f'(w)$$

Backpropagation

The back-propagation algorithm:

- To compute the gradient of some scalar z with respect to one of its ancestors x in the graph
 - We compute the gradient with respect to each parent of z in the graph by multiplying the current gradient by the Jacobian of the operation that produced z .
 - We continue multiplying by Jacobians traveling backwards through the graph in this way until we reach x .
 - For any node that may be reached by going backwards from z through two or more paths, we simply sum the gradients arriving from different paths at that node.

Backpropagation

The back-propagation algorithm:

- To compute the gradient of some scalar z with respect to one of its ancestors x in the graph
 - We compute the gradient with respect to each parent of z in the graph by multiplying the current gradient by the Jacobian of the operation that produced z .
 - We continue multiplying by Jacobians traveling backwards through the graph in this way until we reach x .
 - For any node that may be reached by going backwards from z through two or more paths, we simply sum the gradients arriving from different paths at that node.

Backpropagation

The back-propagation algorithm:

- To compute the gradient of some scalar z with respect to one of its ancestors x in the graph
 - We compute the gradient with respect to each parent of z in the graph by multiplying the current gradient by the Jacobian of the operation that produced z .
 - We continue multiplying by Jacobians traveling backwards through the graph in this way until we reach x .
- For any node that may be reached by going backwards from z through two or more paths, we simply sum the gradients arriving from different paths at that node.

Training with Backpropagation

- Backpropagation Algorithm

Require: \mathbb{T} , the target set of variables whose gradients must be computed.

Require: \mathcal{G} , the computational graph

Require: z , the variable to be differentiated

Let \mathcal{G}' be \mathcal{G} pruned to contain only nodes that are ancestors of z and descendants of nodes in \mathbb{T} .

Initialize `grad_table`, a data structure associating tensors to their gradients

`grad_table[z] \leftarrow 1`

for \mathbf{V} in \mathbb{T} **do**

`build_grad(\mathbf{V} , \mathcal{G} , \mathcal{G}' , grad_table)`

end for

Return `grad_table` restricted to \mathbb{T}

Require: \mathbf{V} , the variable whose gradient should be added to \mathcal{G} and `grad_table`.

Require: \mathcal{G} , the graph to modify.

Require: \mathcal{G}' , the restriction of \mathcal{G} to nodes that participate in the gradient.

Require: `grad_table`, a data structure mapping nodes to their gradients

if \mathbf{V} is in `grad_table` **then**

`Return grad_table[V]`

end if

$i \leftarrow 1$

for \mathbf{C} in `get_consumers(\mathbf{V} , \mathcal{G}')` **do**

`op \leftarrow get_operation(\mathbf{C})`

`D \leftarrow build_grad(\mathbf{C} , \mathcal{G} , \mathcal{G}' , grad_table)`

`$\mathbf{G}^{(i)} \leftarrow \text{op.bprop}(\text{get_inputs}(\mathbf{C}, \mathcal{G}'), \mathbf{V}, \mathbf{D})$`

$i \leftarrow i + 1$

end for

`$\mathbf{G} \leftarrow \sum_i \mathbf{G}^{(i)}$`

`grad_table[V] = \mathbf{G}`

 Insert \mathbf{G} and the operations creating it into \mathcal{G}

Return \mathbf{G}

Source: Deep Learning Book

Training with Backpropagation

- Backpropagation Algorithm

Require: \mathbb{T} , the target set of variables whose gradients must be computed.

Require: \mathcal{G} , the computational graph

Require: z , the variable to be differentiated

Let \mathcal{G}' be \mathcal{G} pruned to contain only nodes that are ancestors of z and descendants of nodes in \mathbb{T} .

Initialize `grad_table`, a data structure associating tensors to their gradients

`grad_table[z] ← 1`

for \mathbf{V} in \mathbb{T} do

`build_grad(V, G, G', grad_table)`

end for

Return `grad_table` restricted to \mathbb{T}

Setting table to store and retrieve gradients

Require: \mathbf{V} , the variable whose gradient should be added to \mathcal{G} and `grad_table`.

Require: \mathcal{G} , the graph to modify.

Require: \mathcal{G}' , the restriction of \mathcal{G} to nodes that participate in the gradient.

Require: `grad_table`, a data structure mapping nodes to their gradients

if \mathbf{V} is in `grad_table` then

 Return `grad_table[V]`

end if

$i \leftarrow 1$ returns the children of \mathbf{V}

for \mathbf{C} in `get_consumers(V, G')` do

`op ← get_operation(C)` returns the operation that computes \mathbf{V}

$\mathbf{D} \leftarrow \text{build_grad}(\mathbf{C}, \mathcal{G}, \mathcal{G}', \text{grad_table})$

$\mathbf{G}^{(i)} \leftarrow \text{op.bprop}(\text{get_inputs}(\mathbf{C}, \mathcal{G}'), \mathbf{V}, \mathbf{D})$

$i \leftarrow i + 1$

end for

$\mathbf{G} \leftarrow \sum_i \mathbf{G}^{(i)}$

`grad_table[V] = G`

Insert \mathbf{G} and the operations creating it into \mathcal{G}

Return \mathbf{G}

Tensor based representation of
Jacobian –based gradient product

$$\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^\top \nabla_{\mathbf{y}} z,$$

get_inputs returns the parents of \mathbf{V}

Source: Deep Learning Book

Training with Backpropagation

- Backward computation of the graph

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

for $k = l, l - 1, \dots, 1$ **do**

Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if f is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

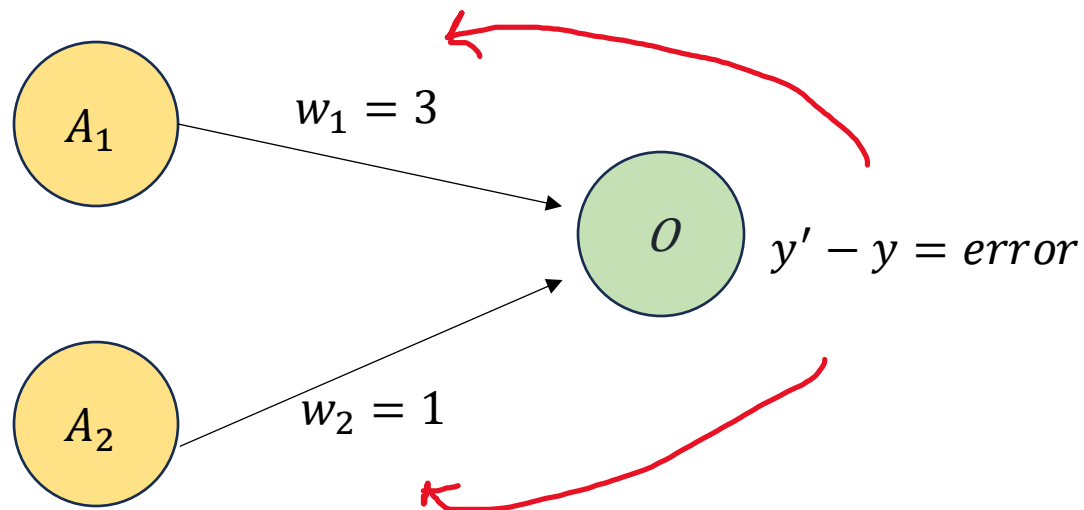
Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

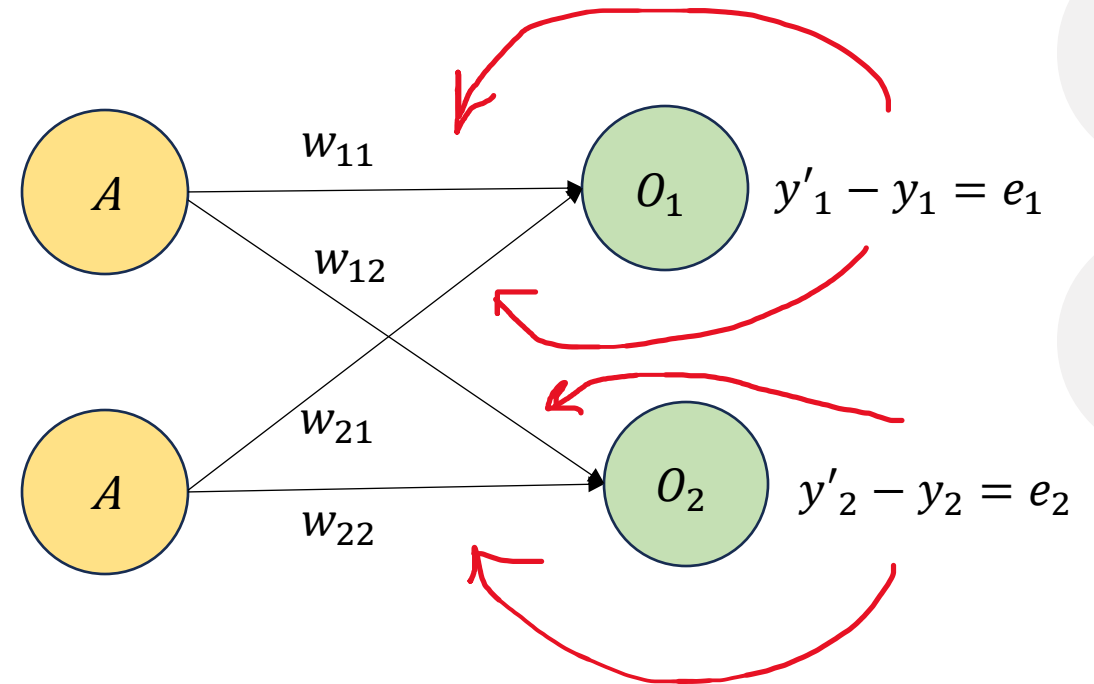
end for

Backpropagation

- Forward propagation takes the signal forward from input to output.
- Backpropagation takes the signal (in this case **error**) from output to input.

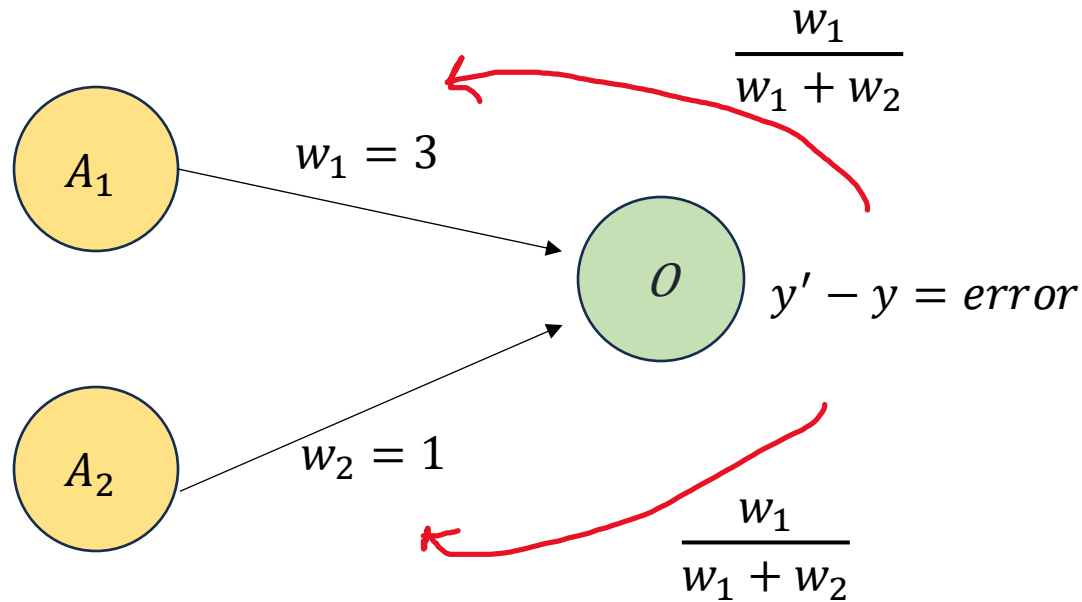


Error not distributed equally

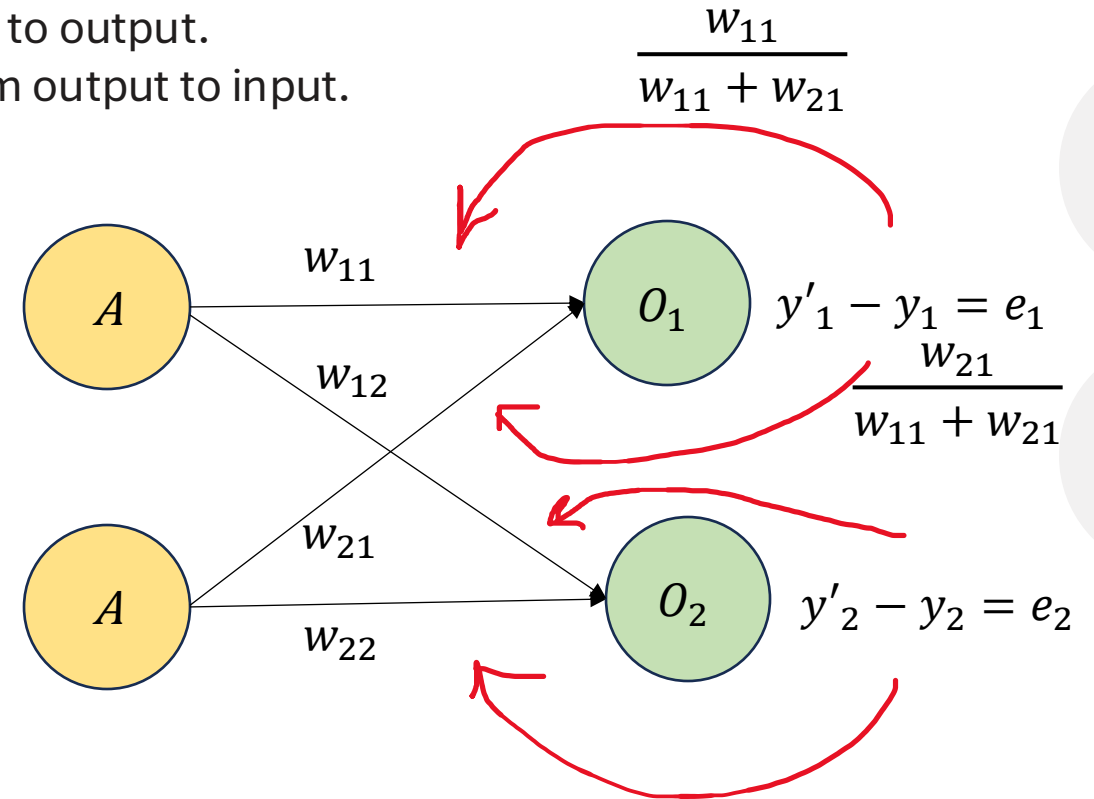


Backpropagation

- Forward propagation takes the signal forward from input to output.
- Backpropagation takes the signal (in this case **error**) from output to input.

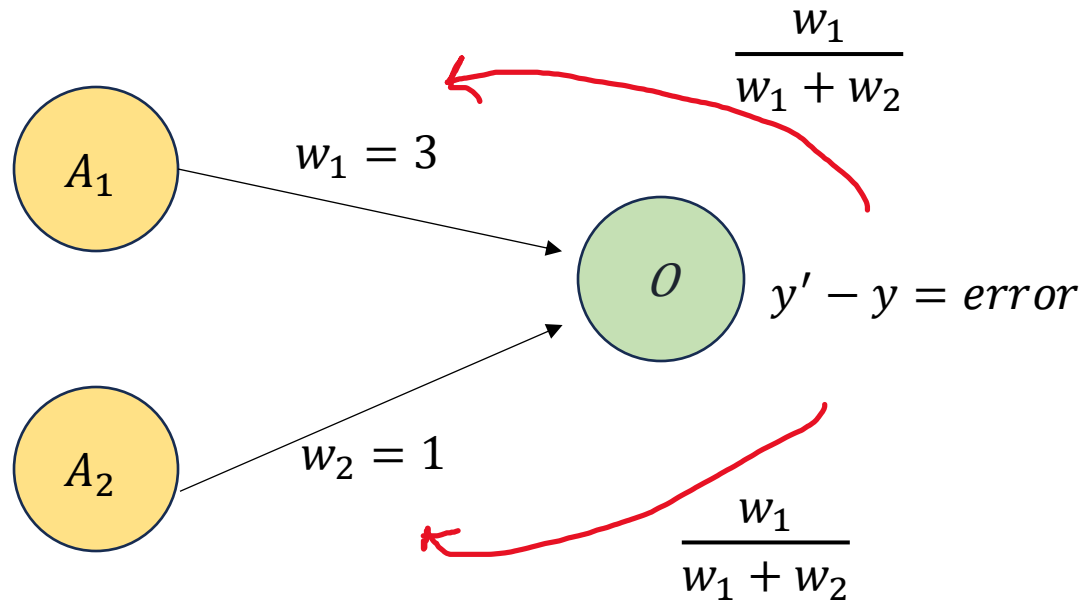


Error is distributed with respect to weights

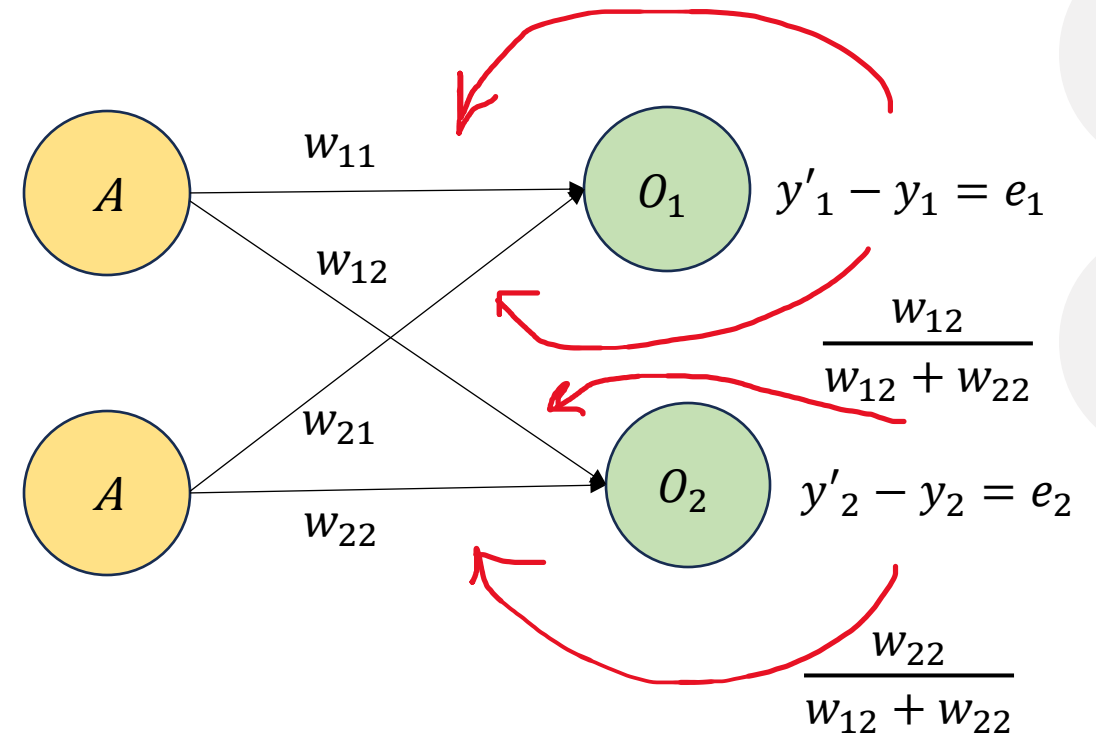


Backpropagation

- Forward propagation takes the signal forward from input to output.
- Backpropagation takes the signal (in this case **error**) from output to input.



Error is distributed with respect to weights



Backpropagation

- Forward propagation takes the signal forward from input to output.
- Backpropagation takes the signal (in this case **error**) from output to input
- What do we do with non-zero error?
 - We update the weights

Backpropagation

- Forward propagation takes the signal forward from input to output.
- Backpropagation takes the signal (in this case **error**) from output to input
- What do we do with non-zero error?
 - We update the weights
- How do we update the weights if error is non-zero?
 - Gradient Descent
- For any node w_{jk} , the new weights associated with it are

$$w'_{jk} = w_{jk} - \alpha \frac{\partial e}{\partial w_{jk}}$$

Backpropagation

- Forward propagation takes the signal forward from input to output.
- Backpropagation takes the signal (in this case **error**) from output to input
- What do we do with non-zero error?
 - We update the weights
- How do we update the weights if error is non-zero?
 - Gradient Descent
- For any node w_{jk} , the new weights associated with it are

$$w'_{jk} = w_{jk} - \alpha \frac{\partial e}{\partial w_{jk}}$$



But gradient of what?

Backpropagation

- Forward propagation takes the signal forward from input to output.
- Backpropagation takes the signal (in this case **error**) from output to input
- What do we do with non-zero error?
 - We update the weights
- How do we update the weights if error is non-zero?
 - Gradient Descent

- For any node w_{jk} , the new weights associated with it are

$$w'_{jk} = w_{jk} - \alpha \frac{\partial e}{\partial w_{jk}}$$

- Backpropagation computes the gradient of units with respect to variables.



But gradient of what?

Backpropagation

- Forward propagation takes the signal forward from input to output.
- Backpropagation takes the signal (in this case **error**) from output to input
- What do we do with non-zero error?
 - We update the weights
- How do we update the weights if error is non-zero?
 - Gradient Descent



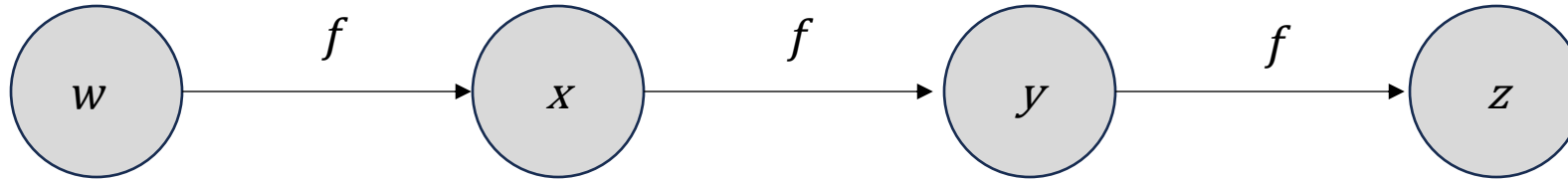
But gradient of what?

- For any node w_{jk} , the new weights associated with it are

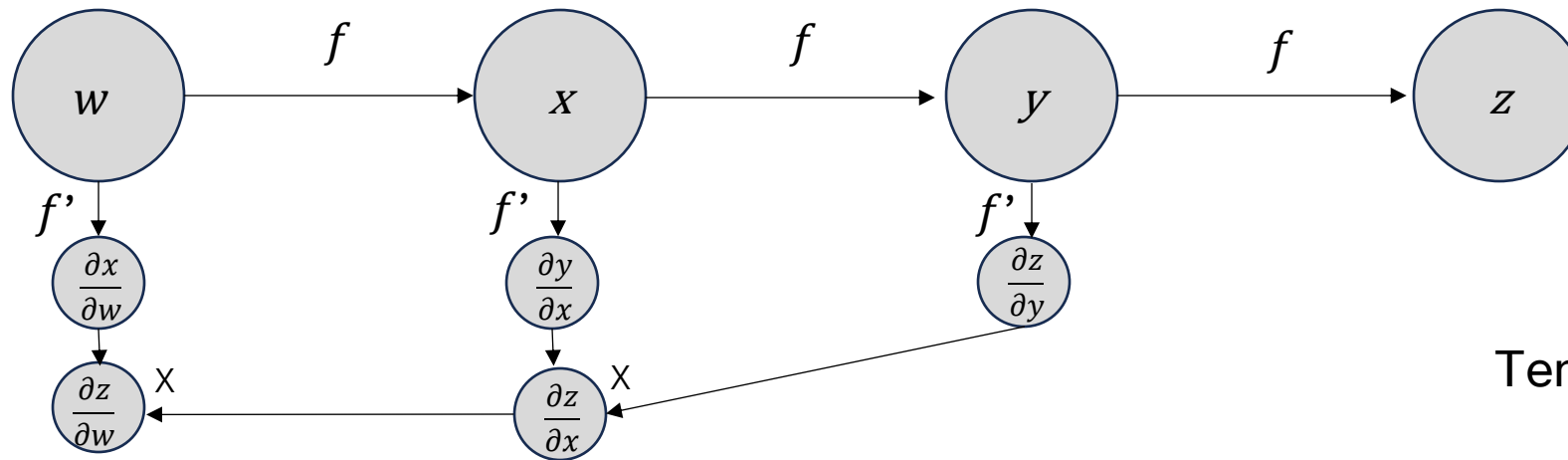
$$w'_{jk} = w_{jk} - \alpha \frac{\partial e}{\partial w_{jk}}$$

- Backpropagation computes the gradient of units with respect to variables.
- The backpropagation starts at the end and recursively applies the chain rule to compute the gradients all the way to the inputs of the network. The gradients can be thought of as flowing backwards through the networks.

Backpropagation



$$\frac{\partial z}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} = f'(y)f'(x) f'(w) = f'(f(f(w))) f'(f(w)) f'(w)$$

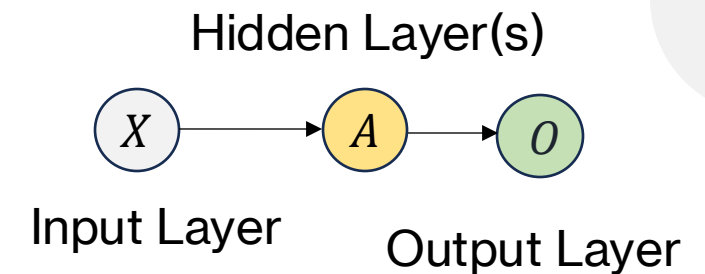
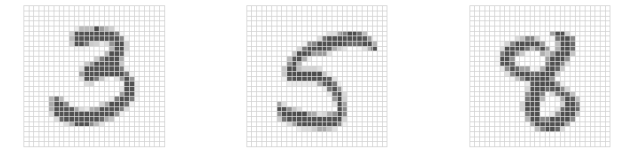


Tensorflow uses this!

Neural Networks: Training

- To train the model
 - We need a learning algorithm
 - Backpropagation
 - Gradient Descent
 - We need a cost function
 - ??
 - We need a training objective
 - ??

0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9

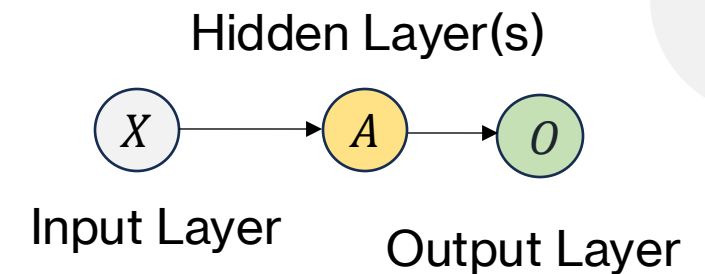
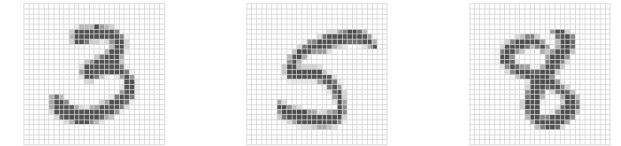


Neural Networks: Training

- To train the model
 - We need a learning algorithm
 - Backpropagation
 - Gradient Descent
 - We need a cost function
 - Cross-entropy loss
 - Mean Absolute Error
 - Mean Squared Error
 - We need a training objective
 - ??

$$\sum_{i=1}^n (y_i - f(x_i))^2.$$

0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9

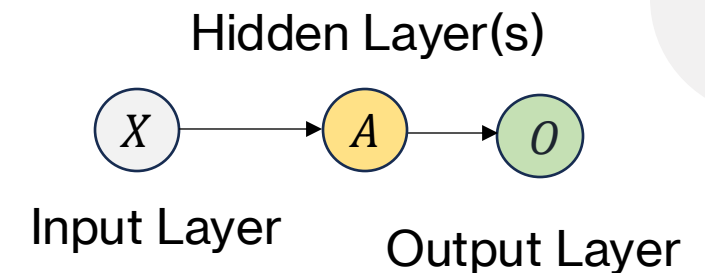
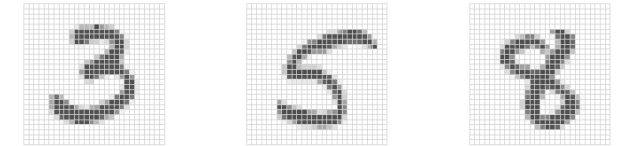


Neural Networks: Training

- To train the model
 - We need a learning algorithm
 - Backpropagation
 - Gradient Descent
 - We need a cost function
 - Cross-entropy loss
 - Mean Absolute Error
 - Mean Squared Error
 - We need a training objective
 - minimize the negative multinomial log-likelihood
 - minimize mean (squared or absolute) error

$$\sum_{i=1}^n (y_i - f(x_i))^2.$$

0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9



Readings

Required Readings:

Introduction to Statistical Learning

- Chapter 10 – Section 10.7 page 427 - 429

Supplemental Readings:

Deep Learning

- Chapter 6 – page 168 - 224

Thank You
