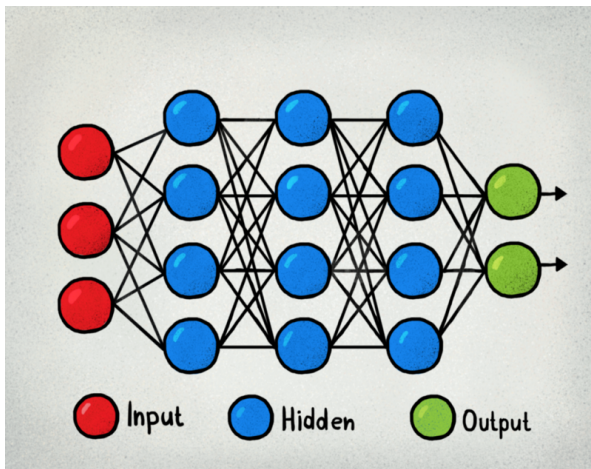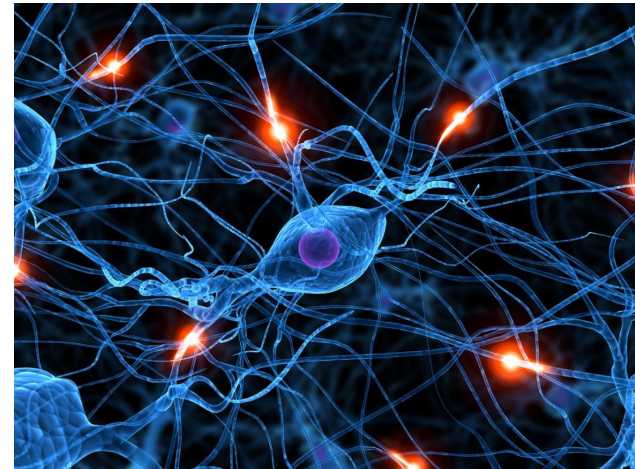# CSE 3521: Neural Networks

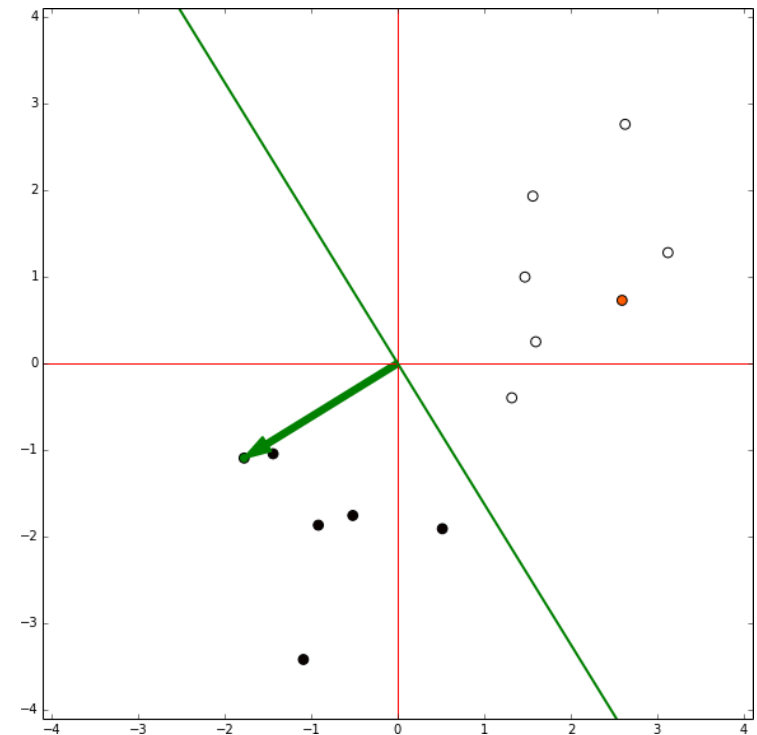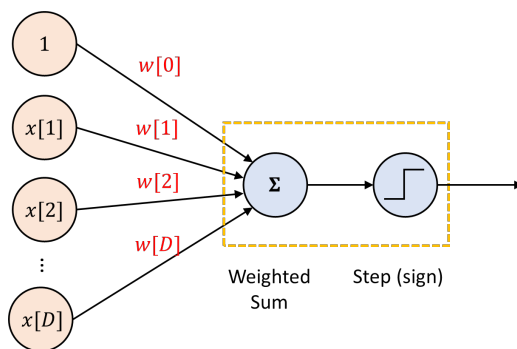

OR



THE OHIO STATE UNIVERSITY

# Today

- Learning "deeper" networks beyond one-layer perceptrons
  - Losses and gradients
  - Back-propagation
  - Stochastic gradient descent

- Training particulars
  - Regularization or weight decay

- Discriminative vs. generative models

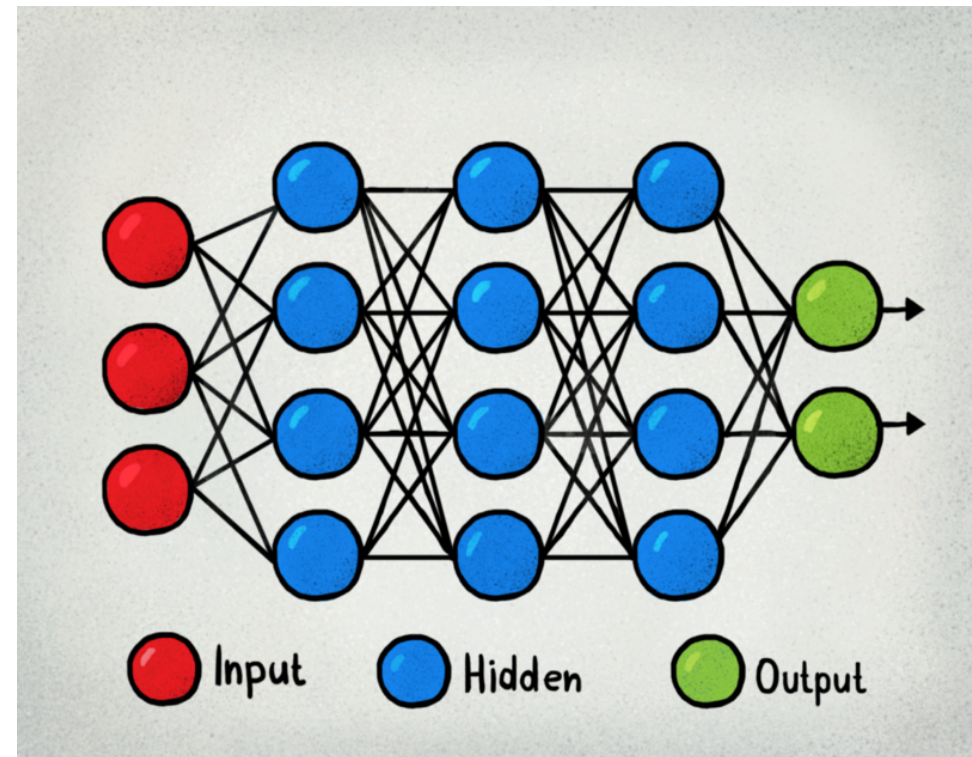# Perceptron Algorithm

- Let $x \in \mathbb{R}^D$ and $w \in \mathbb{R}^D$ ($b$ is merged into $w$), $y \in \{-1,1\}$

- Initialize weight vector $w = 0$
- Loop for T iterations
  - Loop for all training examples $x_n$ (random order!)
  - Predict $\hat{y}_n = \text{sign}(w^T x_n)$
  - If $\hat{y}_n \neq y_n$
    - Update: $w \leftarrow w + \eta(y_n x_n)$
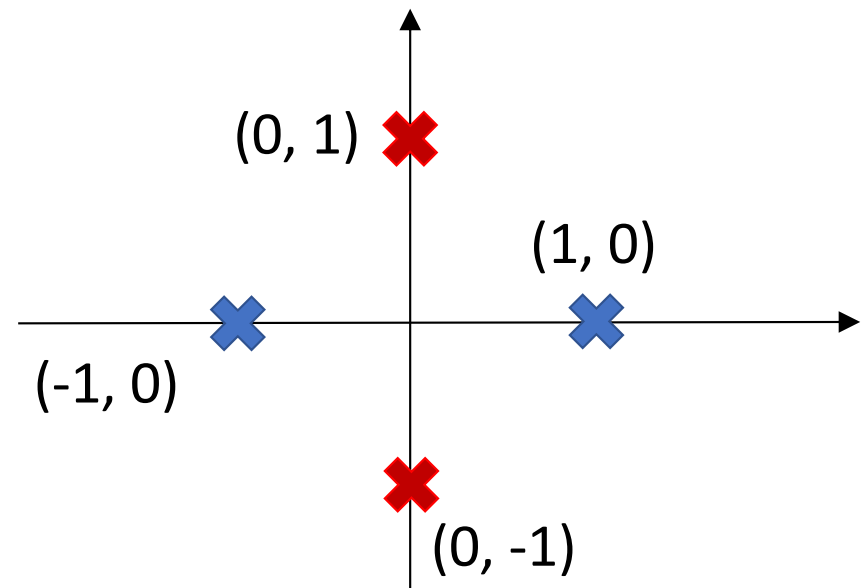
# Multi-layer perceptron

- How to learn?

- Perceptron algorithm?
  - Only look at one instance a time
  - Can update only a single perceptron
  - Not converge for not linear separable data

# Why multiple layers?

$$\text{sign}\left(\sum_0^2 w[d]x[d]\right)$$

1

$x[1]$

$x[2]$

$w[0]$

$w[1]$

$w[2]$

(0, 1)

(1, 0)

(-1, 0)

(0, -1)

Impossible to classify the 4 points correctly!

# Why multiple layers?

$$\hat{y} = \text{sign}\left(\sum_0^2 v[d]z[d]\right)$$

# Why multiple layers?



Let
$w[1,0] = -0.5$
$w[1,1] = 1$
$w[1,2] = 0$

$z[1] = -1$    $z[1] = 1$

(0, 1)

(1, 0)

(-1, 0)

(0, -1)

$$z[d'] = \text{sign}\left(\sum_0^2 w[d',d]x[d]\right)$$

# Why multiple layers?

1

$x[1]$

$w[2,0]$

$w[2,1]$

$x[2]$

$z[2]$

$w[2,2]$

Let
$w[2,0] = -0.5$
$w[2,1] = -1$
$w[2,2] = 0$

$$z[d'] = \text{sign}\left(\sum_0^2 w[d',d]x[d]\right)$$

(0, 1)

(1, 0)

(-1, 0)

(0, -1)

$z[2] = 1$     $z[2] = -1$

# Why multiple layers?

$$\hat{y} = \text{sign}\left(\sum_0^2 v[d]z[d]\right)$$



$z[1] = -1$    $z[1] = 1$

$v[0]$
$v[1]$
$v[2]$

Let
$v[0] = +0.5$
$v[1] = 1$
$v[2] = 1$

$(0, 1)$
$(1, 0)$
$(-1, 0)$
$(0, -1)$

$z[2] = 1$    $z[2] = -1$

# Why multiple layers?

$$\hat{y} = \text{sign}\left(\sum_0^2 v[d]z[d]\right)$$

Let
$v[0] = +0.5$
$v[1] = 1$
$v[2] = 1$

(0, 1)
(1, 0)
(-1, 0)
(0, -1)

$\hat{y} = 1$
$\hat{y} = 1$
$\hat{y} = -1$

1
1
$x[1]$
$z[1]$
$x[2]$
$z[2]$
$\hat{y}$
$v[0]$
$v[1]$
$v[2]$

# Today

- Learning "deeper" networks beyond one-layer perceptrons
  - Losses and gradients
  - Back-propagation
  - Stochastic gradient descent


- Training particulars
  - Regularization or weight decay


- Discriminative vs. generative models

# Losses: $\hat{y}$ vs. $y$



$$\hat{y} = \text{sign}\left(\sum_{0}^{2} v[d]z[d]\right)$$

# Temporally take out "1" for simplicity



$$\hat{y} = \text{sign}\left( \sum_{1}^{2} v[d]z[d] \right)$$

# Re-written with linear algebra



$$\hat{y} = \text{sign}\left(\sum_1^2 v[d]z[d]\right) = \text{sign}(\boldsymbol{v}^T\boldsymbol{z}) = \text{sign}(\boldsymbol{v}^T\text{sign}(\boldsymbol{Wx}))$$

$\begin{bmatrix} z[1] \\ z[2] \end{bmatrix} = \text{sign}(\begin{bmatrix} w[1,1] & w[1,2] \\ w[2,1] & w[2,2] \end{bmatrix}\begin{bmatrix} x[1] \\ x[2] \end{bmatrix})$

# Using sigmoid for inner layers



$$\hat{y} = \text{sign}\left(\sum_{1}^{2} v[d]z[d]\right) = \text{sign}(\boldsymbol{v}^T \boldsymbol{z}) = \text{sign}(\boldsymbol{v}^T \rho(\boldsymbol{Wx}))$$

$$\begin{bmatrix} z[1] \\ z[2] \end{bmatrix} = \rho\left(\begin{bmatrix} w[1,1] & w[1,2] \\ w[2,1] & w[2,2] \end{bmatrix}\begin{bmatrix} x[1] \\ x[2] \end{bmatrix}\right)$$

# What are the parameters?

$$\hat{y} = \text{sign}\left(\sum_{1}^{2} v[d]z[d]\right) = \text{sign}(\boldsymbol{v}^T\boldsymbol{z}) = \text{sign}(\boldsymbol{v}^T \rho(\boldsymbol{W}\boldsymbol{x}))$$

# What are the parameters?

$$\hat{y} = \text{sign}\left(\sum_{1}^{2} v[d]z[d]\right) = \text{sign}(\boldsymbol{v}^T \boldsymbol{z}) = \text{sign}(\boldsymbol{v}^T \rho(\boldsymbol{W}\boldsymbol{x}))$$

Answer: $\boldsymbol{v}, \boldsymbol{W}$

# Losses and gradients for one data instance



$$\hat{y} = \text{sign}\left(\sum_1^2 v[d]z[d]\right) = \text{sign}(\boldsymbol{v}^T\boldsymbol{z}) = \text{sign}(\boldsymbol{v}^T\rho(\boldsymbol{W}\boldsymbol{x}))$$

$$l(y, \hat{y}) \approx l(y, \boldsymbol{v}^T\rho(\boldsymbol{W}\boldsymbol{x}))$$

For example, logistic loss (binary entropy loss):
$$l(y, \boldsymbol{v}^T\rho(\boldsymbol{W}\boldsymbol{x})) = -y \times \log\rho(\boldsymbol{v}^T\rho(\boldsymbol{W}\boldsymbol{x})) - (1-y) \times \log(1 - \rho(\boldsymbol{v}^T\rho(\boldsymbol{W}\boldsymbol{x})))$$

# Losses and gradients for one data instance



$$l(y, \hat{y}) \approx l(y, \boldsymbol{v}^T \rho(\boldsymbol{W} \boldsymbol{x}))$$

$$\nabla_{\boldsymbol{v}} l(y, \hat{y}) = \nabla_{\boldsymbol{v}} l\big(y, \boldsymbol{v}^T \rho(\boldsymbol{W} \boldsymbol{x})\big) = \nabla_{\boldsymbol{v}} l(y, \boldsymbol{v}^T \boldsymbol{z})$$

$$\nabla_{\boldsymbol{W}} l(y, \hat{y})?$$

# Losses and gradients for one data instance



$$l(y, \hat{y}) \approx l(y, \boldsymbol{v}^T \rho(\boldsymbol{Wx}))$$

$$\nabla_{\boldsymbol{v}} l(y, \hat{y}) = \nabla_{\boldsymbol{v}} \, l\big(y, \boldsymbol{v}^T \rho(\boldsymbol{Wx})\big) = \nabla_{\boldsymbol{v}} \, l(y, \boldsymbol{v}^T \boldsymbol{z})$$

$\nabla_{\boldsymbol{W}} l(y, \hat{y})$ by chain rules: $\quad \nabla_{\boldsymbol{z}} l(y, \hat{y}) = \begin{bmatrix} \frac{\partial l}{\partial z[1]} \\ \frac{\partial l}{\partial z[2]} \end{bmatrix}, \quad \nabla_{\boldsymbol{W}} l(y, \hat{y}) = \sum_{d=1}^{2} \frac{\partial l}{\partial z[d]} \times \nabla_{\boldsymbol{W}} z[d]$

# Review chain rules

- If $f(x) = A(B(C(x, w)))$: assuming $x, w$ and all functions output a scalar
- $\dfrac{\partial f}{\partial x} = \dfrac{\partial f}{\partial A} \times \dfrac{\partial A}{\partial B} \times \dfrac{\partial B}{\partial C} \times \dfrac{\partial C}{\partial x}$
- $\dfrac{\partial f}{\partial w} = \dfrac{\partial f}{\partial A} \times \dfrac{\partial A}{\partial B} \times \dfrac{\partial B}{\partial C} \times \dfrac{\partial C}{\partial w}$
- Where
  - $A = A(B(C(x, w)))$
  - $B = B\big(C(x, w)\big)$
  - $C = C(x, w)$

# Backpropagation

$$x \Rightarrow z^{(1)} \Rightarrow z^{(2)} \Rightarrow z^{(3)} \Rightarrow \hat{y} \Rightarrow y$$

$$\frac{\partial l}{\partial z^{(3)}} = \frac{\partial l}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z^{(3)}}$$

$$\frac{\partial l}{\partial \hat{y}}$$

$$\frac{\partial l}{\partial W^{(4)}} = \frac{\partial l}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial W^{(4)}}$$



Input   Hidden   Output

$$W^{(1)} \quad W^{(2)} \quad W^{(3)} \quad W^{(4)}$$

# Backpropagation

$$x \implies z^{(1)} \implies z^{(2)} \implies z^{(3)} \implies \hat{y} \implies y$$

$$\frac{\partial l}{\partial z^{(2)}} = \frac{\partial l}{\partial z^{(3)}} \times \frac{\partial z^{(3)}}{\partial z^{(2)}} \qquad \frac{\partial l}{\partial z^{(3)}} \qquad \frac{\partial l}{\partial \hat{y}}$$

$$\frac{\partial l}{\partial W^{(3)}} = \frac{\partial l}{\partial z^{(3)}} \times \frac{\partial z^{(3)}}{\partial W^{(3)}} \qquad \frac{\partial l}{\partial W^{(4)}}$$



Input  Hidden  Output

$$W^{(1)} \qquad W^{(2)} \qquad W^{(3)} \qquad W^{(4)}$$

# Backpropagation

$$x \implies z^{(1)} \implies z^{(2)} \implies z^{(3)} \implies \hat{y} \implies y$$

$$\frac{\partial l}{\partial z^{(1)}} = \frac{\partial l}{\partial z^{(2)}} \times \frac{\partial z^{(2)}}{\partial z^{(1)}} \qquad \frac{\partial l}{\partial z^{(2)}} \qquad \frac{\partial l}{\partial z^{(3)}} \qquad \frac{\partial l}{\partial \hat{y}}$$

$$\frac{\partial l}{\partial W^{(2)}} = \frac{\partial l}{\partial z^{(2)}} \times \frac{\partial z^{(2)}}{\partial W^{(2)}} \qquad \frac{\partial l}{\partial W^{(3)}} \qquad \frac{\partial l}{\partial W^{(4)}}$$



Input    Hidden    Output

$$W^{(1)} \qquad W^{(2)} \qquad W^{(3)} \qquad W^{(4)}$$

# Backpropagation

$$x \Rightarrow z^{(1)} \Rightarrow z^{(2)} \Rightarrow z^{(3)} \Rightarrow \hat{y} \Rightarrow y$$

$$\frac{\partial l}{\partial z^{(1)}} \Leftarrow \frac{\partial l}{\partial z^{(2)}} \Leftarrow \frac{\partial l}{\partial z^{(3)}} \Leftarrow \frac{\partial l}{\partial \hat{y}} \Leftarrow$$

$$\frac{\partial l}{\partial W^{(1)}} = \frac{\partial l}{\partial z^{(1)}} \times \frac{\partial z^{(1)}}{\partial W^{(1)}} \qquad \frac{\partial l}{\partial W^{(2)}} \qquad \frac{\partial l}{\partial W^{(3)}} \qquad \frac{\partial l}{\partial W^{(4)}}$$
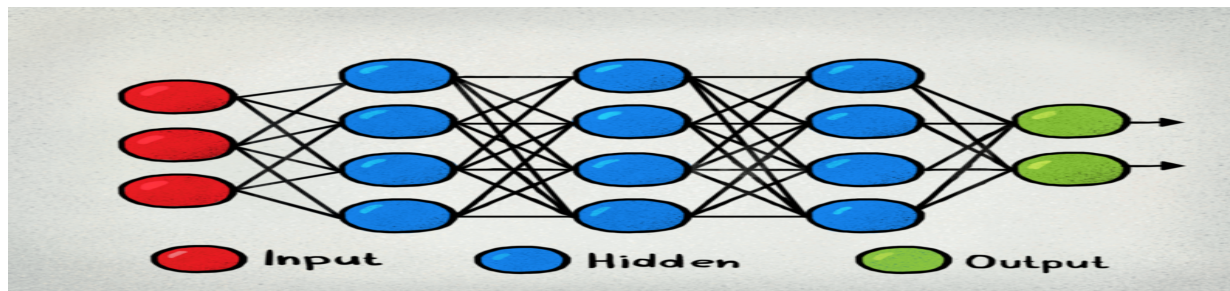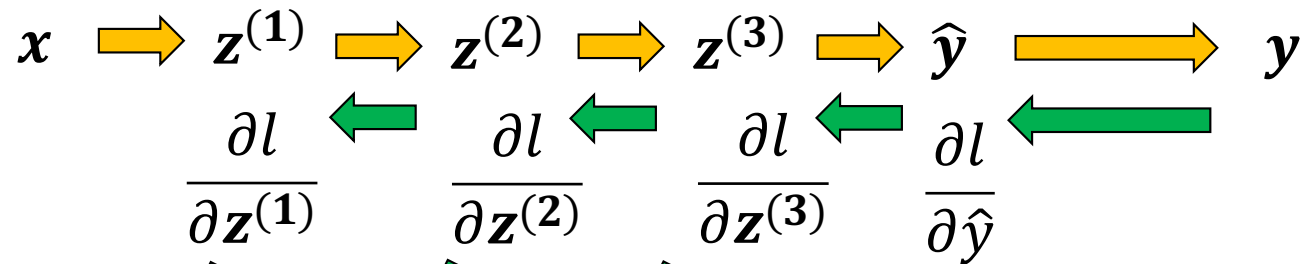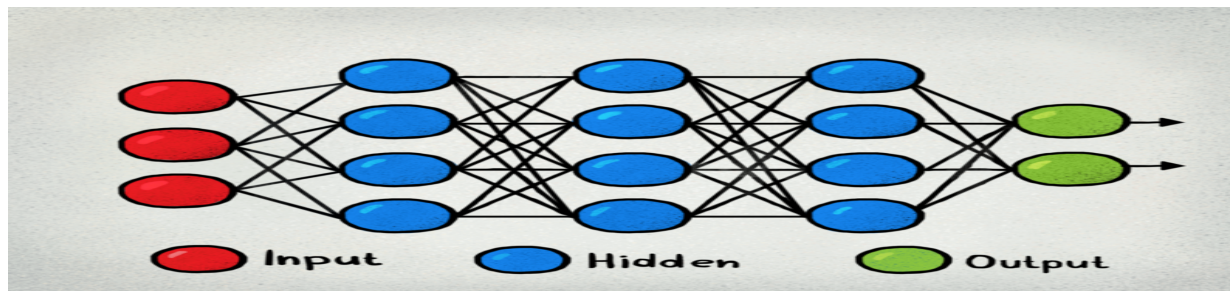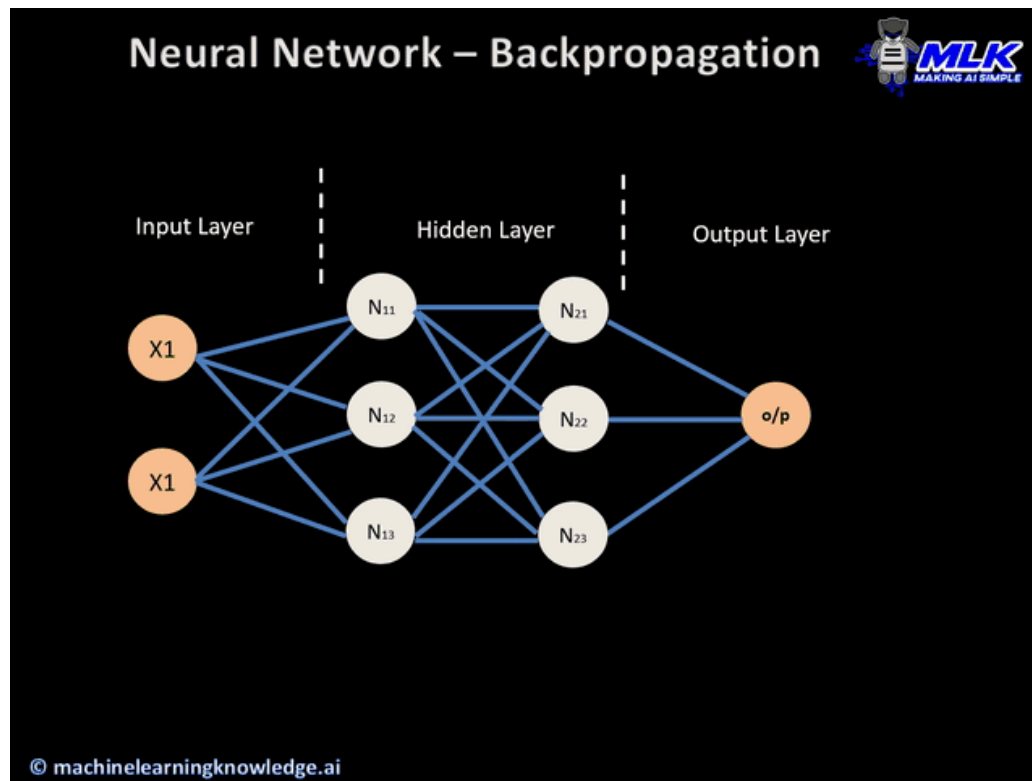


$$W^{(1)} \qquad W^{(2)} \qquad W^{(3)} \qquad W^{(4)}$$

# Illustration

# Backpropagation

- When it is not a scalar:

  ○ $\dfrac{\partial l}{\partial \boldsymbol{z}^{(2)}} = \dfrac{\partial l}{\partial \boldsymbol{z}^{(3)}} \times \dfrac{\partial \boldsymbol{z}^{(3)}}{\partial \boldsymbol{z}^{(2)}}$ $\quad\rightarrow \dfrac{\partial l}{\partial \boldsymbol{z}^{(2)}} = \sum_{d=1}^{D} \dfrac{\partial l}{\partial \boldsymbol{z}^{(3)}[\text{d}]} \times \dfrac{\partial \boldsymbol{z}^{(3)}[\text{d}]}{\partial \boldsymbol{z}^{(2)}}$

  ○ $\dfrac{\partial l}{\partial \boldsymbol{W}^{(3)}} = \dfrac{\partial l}{\partial \boldsymbol{z}^{(3)}} \times \dfrac{\partial \boldsymbol{z}^{(3)}}{\partial \boldsymbol{W}^{(3)}}$ $\quad\rightarrow \dfrac{\partial l}{\partial \boldsymbol{W}^{(3)}} = \sum_{d=1}^{D} \dfrac{\partial l}{\partial \boldsymbol{z}^{(3)}[\text{d}]} \times \dfrac{\partial \boldsymbol{z}^{(3)}[\text{d}]}{\partial \boldsymbol{W}^{(3)}}$

# Training a neural network: gradient descent

- Let $x \in \mathbb{R}^D$; $y$ as true label, and $\{(x_n, y_n)\}$ as the training data
- $\boldsymbol{\theta}$ as all parameters, $\hat{y} = f_{\boldsymbol{\theta}}(x)$ as the neural network's prediction

- Initialize $\boldsymbol{\theta}$ with [with some specific methods]

- Loop for T "epochs"
  - $\nabla_{\boldsymbol{\theta}} L = \frac{1}{N} \sum_{n=1}^{N} \nabla_{\boldsymbol{\theta}} l\,(y_n, \hat{y}_n)$
  - $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \times \nabla_{\boldsymbol{\theta}} L$

# Stochastic gradient descent

- Let $x \in \mathbb{R}^D$; $y$ as true label, and $\{(x_n, y_n)\}$ as the training data
- $\boldsymbol{\theta}$ as all parameters, $\hat{y} = f_{\boldsymbol{\theta}}(x)$ as the neural network's prediction

- Initialize $\boldsymbol{\theta}$ with [with some specific methods]

- Loop for T "epochs"
  - Loop for all training examples $x_n$ (random order!)
  - $\nabla_{\boldsymbol{\theta}} L = \nabla_{\boldsymbol{\theta}} l(y_n, \hat{y}_n)$
  - $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \times \nabla_{\boldsymbol{\theta}} L$

# "Mini-batch" Stochastic gradient descent

- Let $x \in \mathbb{R}^D$; $y$ as true label, and $\{(x_n, y_n)\}$ as the training data
- $\boldsymbol{\theta}$ as all parameters, $\hat{y} = f_{\boldsymbol{\theta}}(x)$ as the neural network's prediction

- Initialize $\boldsymbol{\theta}$ with [with some specific methods]

- Loop for T "epochs"
  - Loop for "B sampled examples from $\{(x_n, y_n)\}$" (called batch) without replacement (random order!)
  - $\nabla_{\boldsymbol{\theta}} L = \frac{1}{N} \sum_{b=1}^{B} \nabla_{\boldsymbol{\theta}} l(y_b, \hat{y}_b)$
  - $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \times \nabla_{\boldsymbol{\theta}} L$