# Optimization and Backpropagation
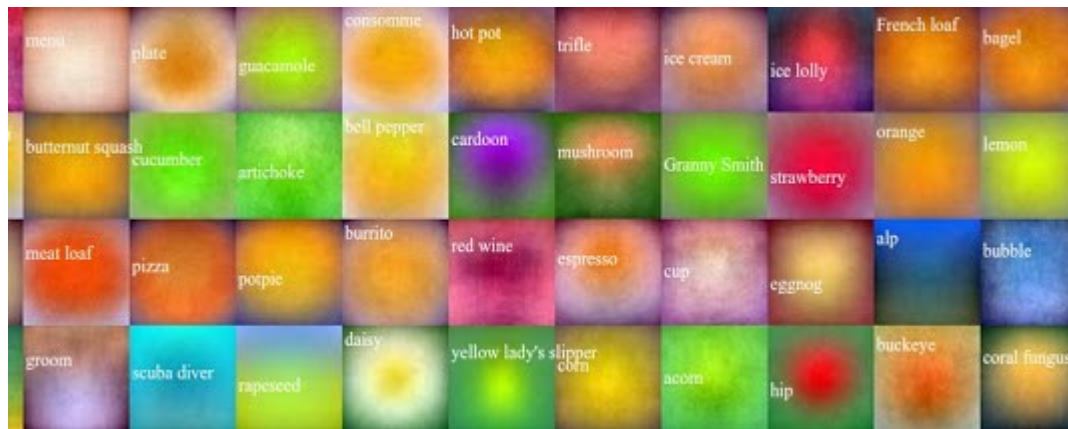
# Lecture 3 Recap

# Neural Network

- Linear score function $f = Wx$



On CIFAR-10
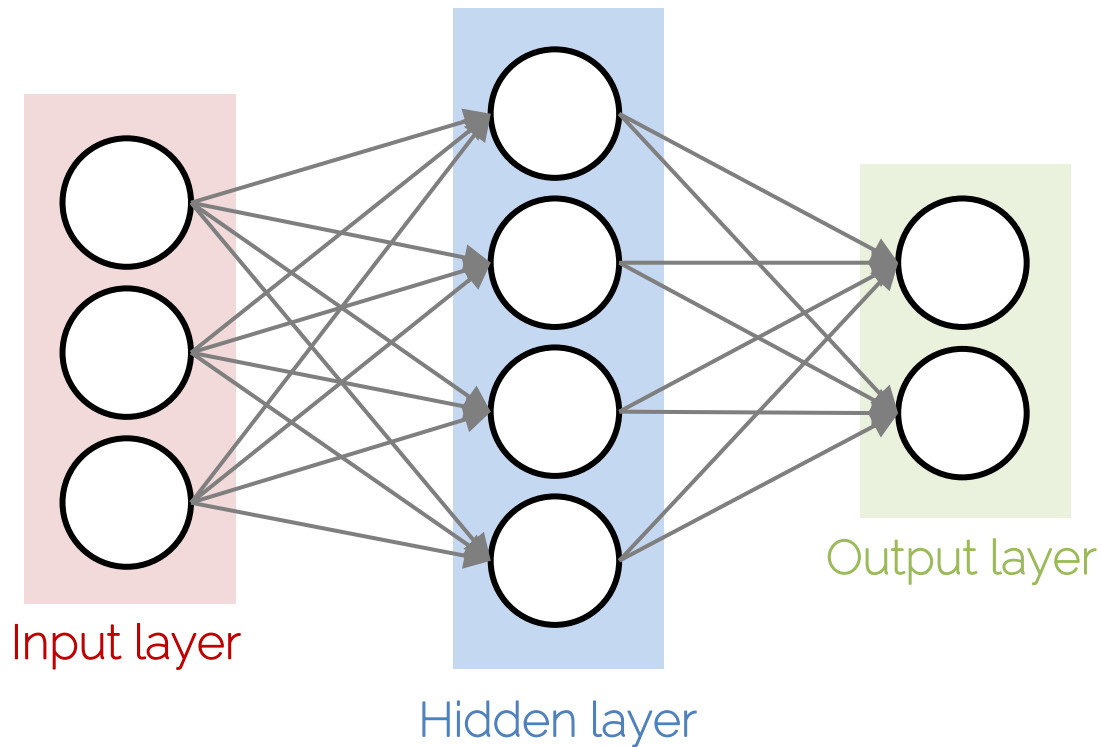


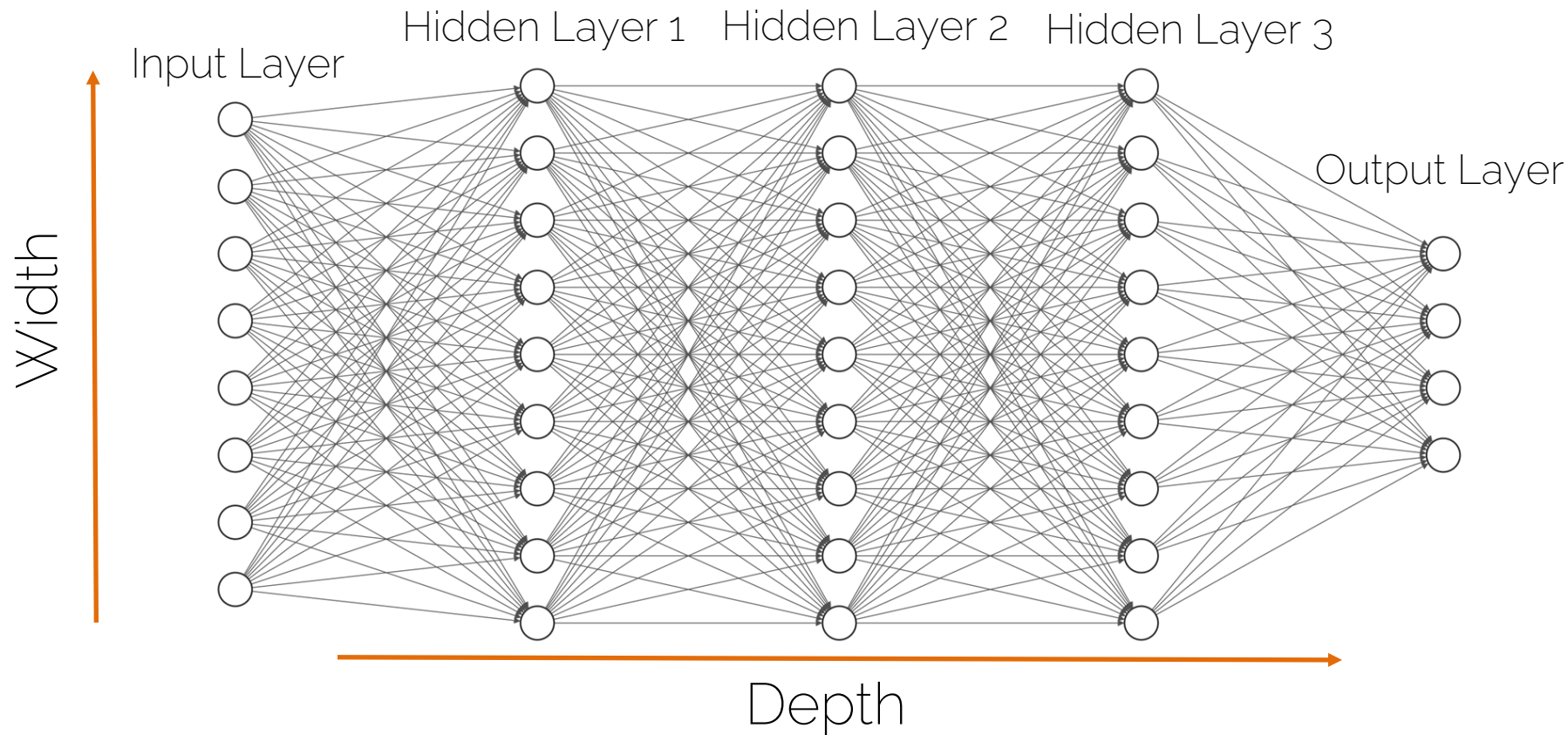On ImageNet

Credit:
Li/Karpathy/Johnson

# Neural Network

- Linear score function $f = Wx$

- Neural network is a nesting of 'functions'
  - 2-layers: $f = W_2 \max(0, W_1 x)$
  - 3-layers: $f = W_3 \max(0, W_2 \max(0, W_1 x))$
  - 4-layers: $f = W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x)))$
  - 5-layers: $f = W_5 \sigma(W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x))))$
  - ... up to hundreds of layers

# Neural Network



Input layer

Hidden layer

Output layer

Credit: Li/Karpathy/Johnson

# Neural Network



Input Layer

Hidden Layer 1    Hidden Layer 2    Hidden Layer 3

Output Layer

Width

Depth

# Activation Functions

Sigmoid: $\sigma(x) = \dfrac{1}{(1+e^{-x})}$

Leaky ReLU: $\max(0.1x, x)$

tanh: $\tanh(x)$

Parametric ReLU: $\max(\alpha x, x)$

Maxout $\max(w_1^T x + b_1, w_2^T x + b_2)$

ReLU: $\max(0, x)$

ELU $f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$

# Loss Functions

- Measure the goodness of the predictions (or equivalently, the network's performance)

- Regression loss
  - L1 loss $\boldsymbol{L}(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_{i}^{n} \|y_i - \widehat{y}_i\|_1$
  - MSE loss $\boldsymbol{L}(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_{i}^{n} \|y_i - \widehat{y}_i\|_2^2$

- Classification loss (for multi-class classification)
  - Cross Entropy loss $E(y, \widehat{y}; \theta) = -\sum_{i=1}^{n} \sum_{k=1}^{K} (y_{ik} \cdot \log \widehat{y}_{ik})$

# Computational Graphs

- Neural network is a computational graph

  – It has compute nodes

  – It has edges that connect nodes

  – It is directional

  – It is organized in 'layers'

# Backprop

# The Importance of Gradients

- Our optimization schemes are based on computing gradients

$$\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

- One can compute gradients analytically but what if our function is too complex?

- Break down gradient computation

Backpropagation

Done by many people before, but often credited to Rumelhart 1986

# Backprop: Forward Pass

- $f(x, y, z) = (x + y) \cdot z$

Initialization $x = 1, y = -3, z = 4$

# Backprop: Backward Pass

$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$



$$d = x + y \qquad \frac{\partial d}{\partial x} = 1, \frac{\partial d}{\partial y} = 1$$

$$f = d \cdot z \qquad \frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$ ?

# Backprop: Backward Pass

$f(x, y, z) = (x + y) \cdot z$

with $x = 1, y = -3, z = 4$

$d = x + y$ $\qquad \frac{\partial d}{\partial x} = 1, \frac{\partial d}{\partial y} = 1$

$f = d \cdot z$ $\qquad \frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$ ?

# Backprop: Backward Pass

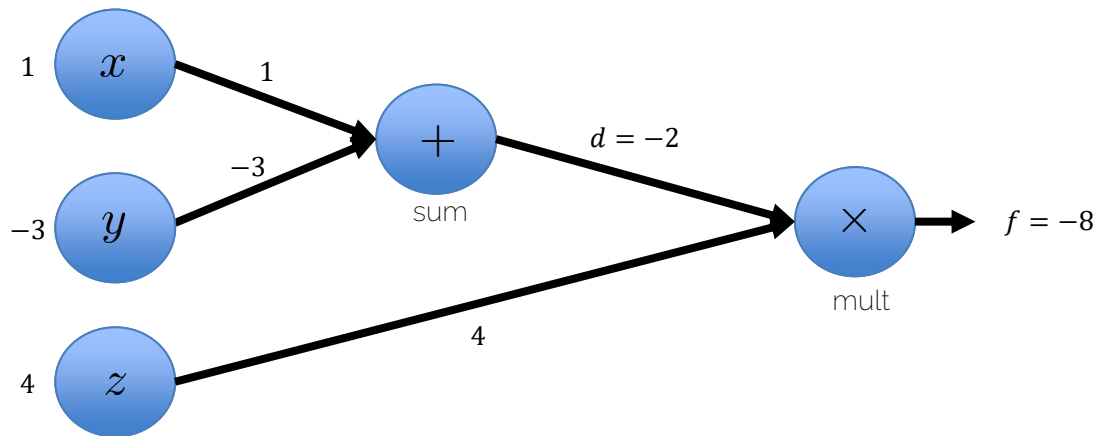$f(x, y, z) = (x + y) \cdot z$

with $x = 1, y = -3, z = 4$

$d = x + y$     $\dfrac{\partial d}{\partial x} = 1, \dfrac{\partial d}{\partial y} = 1$

$f = d \cdot z$     $\dfrac{\partial f}{\partial d} = z, \boxed{\dfrac{\partial f}{\partial z} = d}$

What is $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$ ?



1   $x$

$-3$   $y$

$-3$

1

$+$   sum

$d = -2$

$\times$   mult

$f = -8$

1

$4$   $z$

$-2$

4

$-2$

$\boxed{\dfrac{\partial f}{\partial z}}$

# Backprop: Backward Pass

$$f(x, y, z) = (x + y) \cdot z$$

with $x = 1, y = -3, z = 4$



$d = x + y$      $\dfrac{\partial d}{\partial x} = 1, \dfrac{\partial d}{\partial y} = 1$

$f = d \cdot z$      $\boxed{\dfrac{\partial f}{\partial d} = z} \dfrac{\partial f}{\partial z} = d$

What is $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$ ?
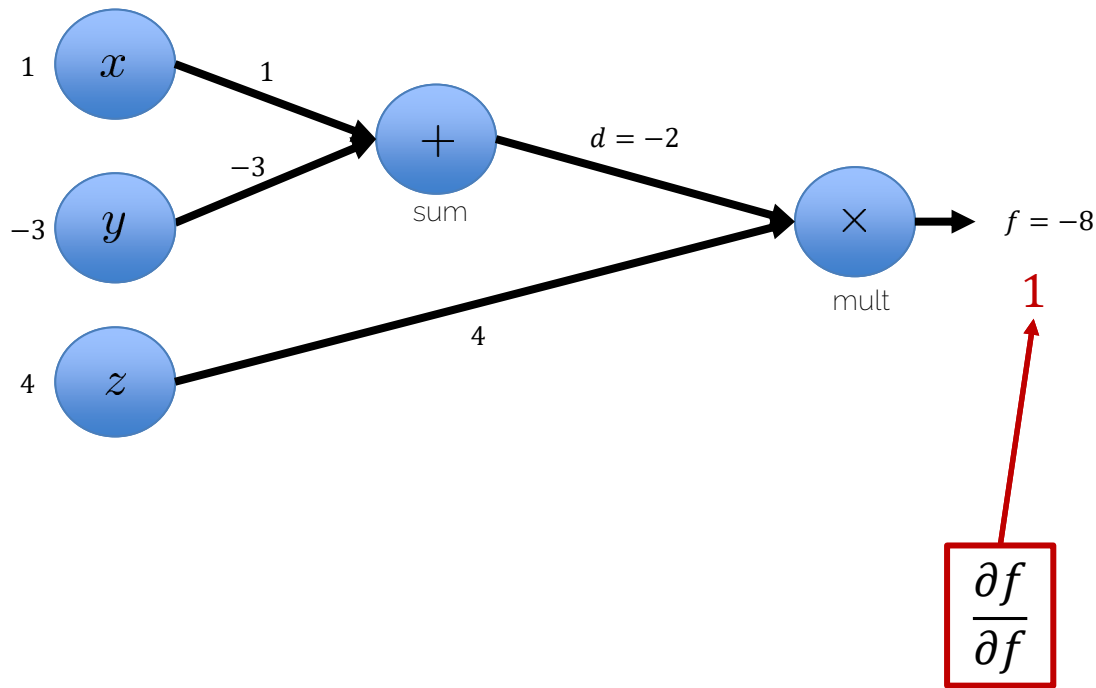
# Backprop: Backward Pass

$$f(x, y, z) = (x + y) \cdot z$$
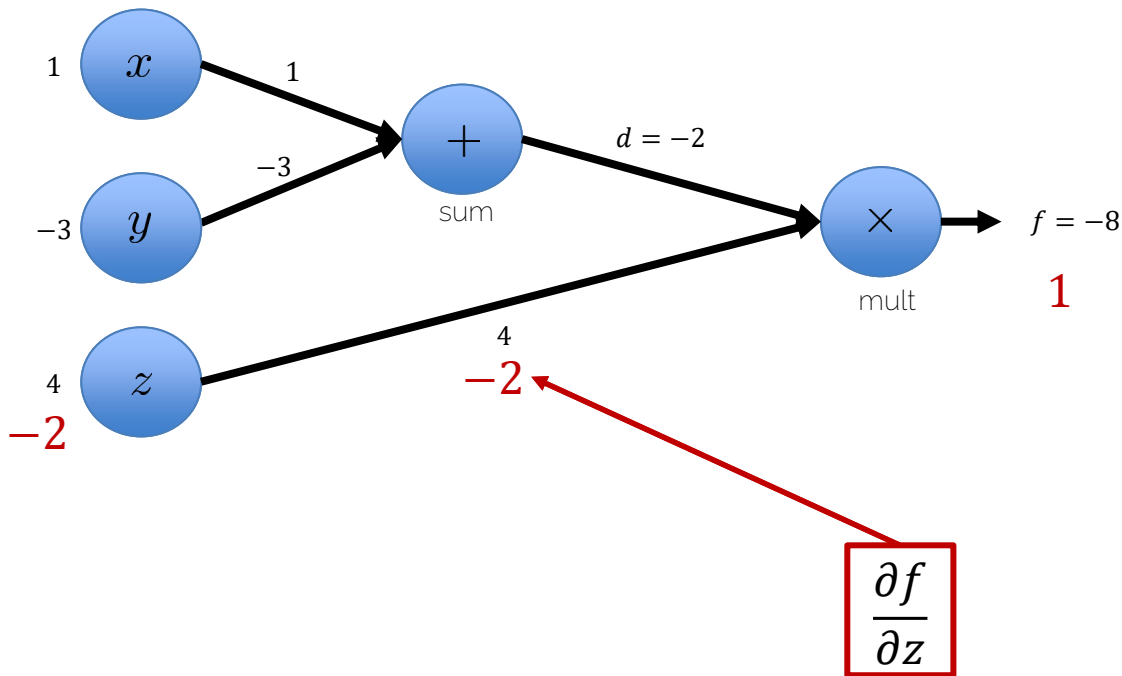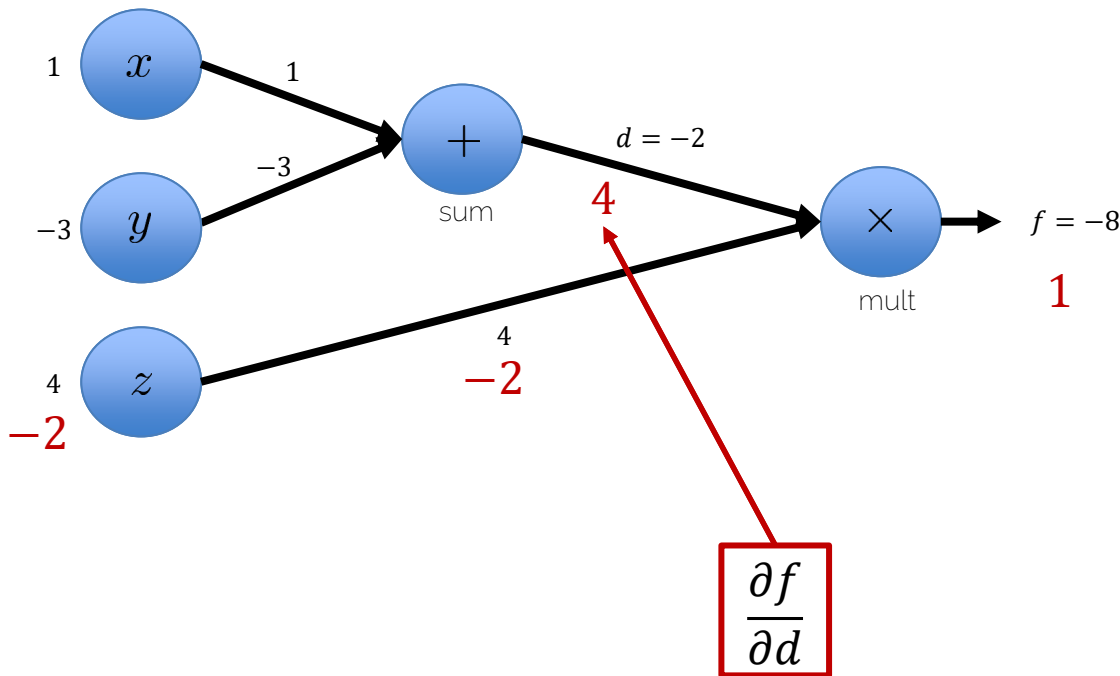
with $x = 1, y = -3, z = 4$



$$d = x + y$$

$$\frac{\partial d}{\partial x} = 1, \boxed{\frac{\partial d}{\partial y} = 1}$$

$$f = d \cdot z$$

$$\frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$ ?

Chain Rule:

$$\boxed{\frac{\partial f}{\partial y} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial y}}$$

$$\rightarrow \frac{\partial f}{\partial y} = 4 \cdot 1 = 4$$

# Backprop: Backward Pass

$$f(x, y, z) = (x + y) \cdot z$$
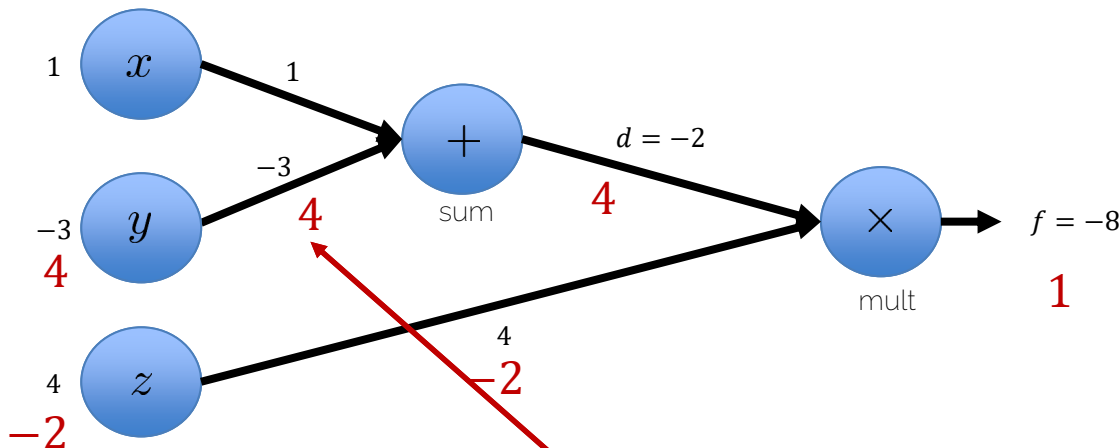
with $x = 1, y = -3, z = 4$

$d = x + y$

$$\boxed{\frac{\partial d}{\partial x} = 1}, \frac{\partial d}{\partial y} = 1$$

$f = d \cdot z$

$$\frac{\partial f}{\partial d} = z, \frac{\partial f}{\partial z} = d$$

What is $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$ ?



Chain Rule:

$$\boxed{\frac{\partial f}{\partial y} = \frac{\partial f}{\partial d} \cdot \frac{\partial d}{\partial y}}$$

$$\boxed{\frac{\partial f}{\partial x}}$$

$$\rightarrow \frac{\partial f}{\partial x} = 4 \cdot 1 = 4$$

# Compute Graphs ->  Neural Networks

- $x_k$ input variables
- $w_{l,m,n}$ network weights (note 3 indices)
  - $l$ which layer
  - $m$ which neuron in layer
  - $n$ which weight in neuron
- $\hat{y}_i$ computed output ($i$ output dim; $n_{out}$)
- $y_i$ ground truth targets
- $L$ loss function

# Compute Graphs -> Neural Networks



Input layer    Output layer

$x_0$

$x_1$

$\hat{y}_0$

$y_0$

e.g., class label/ regression target

$x_0$

$x_1$

$* w_0$

$* w_1$

$+$

$-y_0$

$x*x$

Loss/ cost

Input    Weights (unknowns!)    L2 Loss function

# Compute Graphs -> Neural Networks

Input layer   Output layer

$x_0$

$x_1$

$\hat{y}_0$   $y_0$

e.g., class label/
regression target

$x_0$

$x_1$

$* w_0$

$* w_1$

$+$

$\max(0, x)$

$-y_0$

$x * x$

Loss/
cost

Input

Weights
(unknowns!)

ReLU Activation
(not arguing this
is the right choice here)

L2 Loss
function

We want to compute gradients w.r.t. all weights $\boldsymbol{W}$

# Compute Graphs -> Neural Networks

Input layer

Output layer



$x_0$

$\hat{y}_0$  ↔  $y_0$

$x_1$

$\hat{y}_1$  →  $y_1$

$\hat{y}_2$  →  $y_2$

$* \, w_{0,0}$

$* \, w_{0,1}$

$+$  →  $-y_0$  →  $x^*x$  →  Loss/cost

$x_0$

$* \, w_{1,0}$

$* \, w_{1,1}$

$+$  →  $-y_0$  →  $x^*x$  →  Loss/cost

$x_1$

$* \, w_{2,0}$

$* \, w_{2,1}$

$+$  →  $-y_0$  →  $x^*x$  →  Loss/cost

We want to compute gradients w.r.t. all weights $\boldsymbol{W}$

# Compute Graphs -> Neural Networks

Input layer          Output layer



$$\hat{y}_i = A(b_i + \sum_k x_k w_{i,k})$$

Activation   bias
function

Goal: We want to compute gradients of the loss function $L$ w.r.t. all weights $\boldsymbol{W}$

$$L = \sum_i L_i$$

$L$: sum over loss per sample, e.g. L2 loss $\longrightarrow$ simply sum up squares:
$$L_i = (\hat{y}_i - y_i)^2$$

$\longrightarrow$ use chain rule to compute partials

$$\frac{\partial L}{\partial w_{i,k}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w_{i,k}}$$

We want to compute gradients w.r.t. all weights $\boldsymbol{W}$ AND all biases $\boldsymbol{b}$

# NNs as Computational Graphs

- We can express any kind of functions in a computational graph, e.g. $f(\boldsymbol{w}, \boldsymbol{x}) = \dfrac{1}{1 + e^{-(b + w_0 x_0 + w_1 x_1)}}$



Sigmoid function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

# NNs as Computational Graphs

- $f(\boldsymbol{w}, \boldsymbol{x}) = \dfrac{1}{1+e^{-(b+w_0 x_0 + w_1 x_1)}}$

# NNs as Computational Graphs

- $f(\boldsymbol{w}, \boldsymbol{x}) = \dfrac{1}{1 + e^{-(b + w_0 x_0 + w_1 x_1)}}$



$$g(x) = \frac{1}{x} \quad \Rightarrow \frac{\partial g}{\partial x} = -\frac{1}{x^2}$$

$$g_\alpha(x) = \alpha + x \Rightarrow \frac{\partial g}{\partial x} = 1$$

$$g(x) = e^x \quad \Rightarrow \frac{\partial g}{\partial x} = e^x$$

$$g_\alpha(x) = \alpha x \quad \Rightarrow \frac{\partial g}{\partial x} = \alpha$$

$$1 \cdot -\frac{1}{1.37^2} = -0.53$$

# NNs as Computational Graphs

- $f(\boldsymbol{w}, \boldsymbol{x}) = \dfrac{1}{1+e^{-(b+w_0 x_0 + w_1 x_1)}}$

$$g(x) = \frac{1}{x} \quad \Rightarrow \frac{\partial g}{\partial x} = -\frac{1}{x^2}$$

$$g_\alpha(x) = \alpha + x \Rightarrow \frac{\partial g}{\partial x} = 1$$

$$g(x) = e^x \quad \Rightarrow \frac{\partial g}{\partial x} = e^x$$

$$g_\alpha(x) = \alpha x \quad \Rightarrow \frac{\partial g}{\partial x} = \alpha$$

$-0.53 \cdot 1 = -0.53$

# NNs as Computational Graphs

- $f(\boldsymbol{w}, \boldsymbol{x}) = \dfrac{1}{1 + e^{-(b + w_0 x_0 + w_1 x_1)}}$

$$g(x) = \frac{1}{x} \quad\Rightarrow\quad \frac{\partial g}{\partial x} = -\frac{1}{x^2}$$

$$g_\alpha(x) = \alpha + x \Rightarrow \frac{\partial g}{\partial x} = 1$$

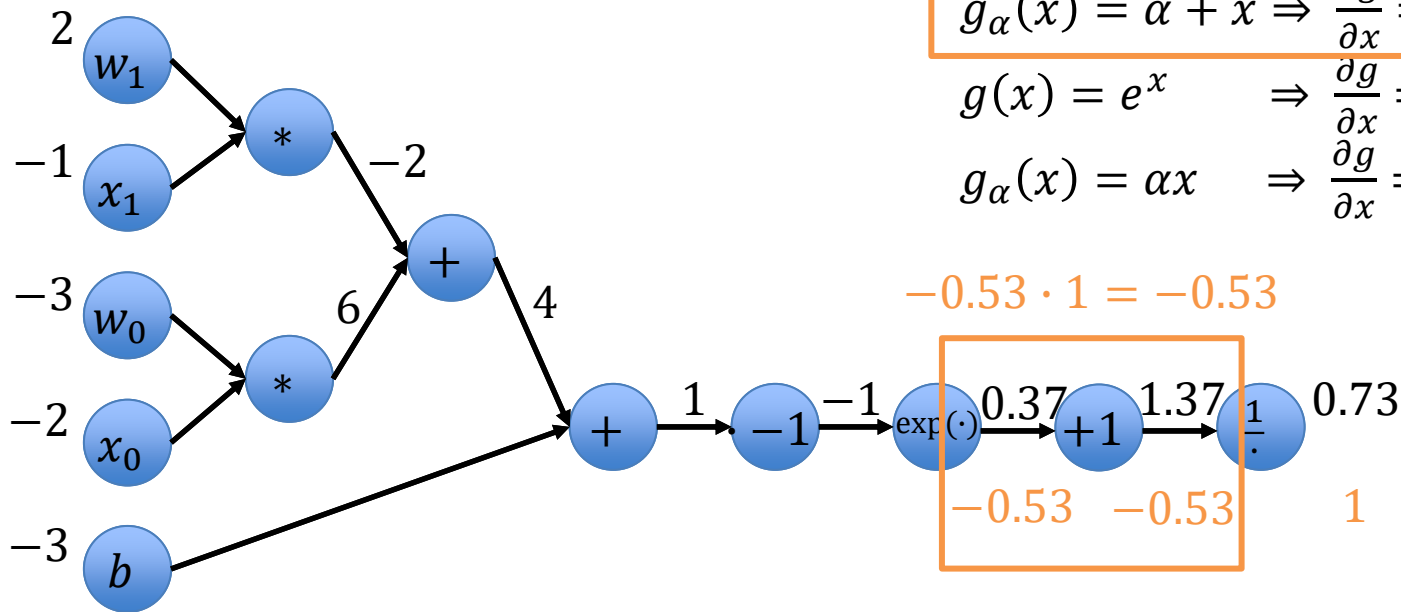$$g(x) = e^x \quad\Rightarrow\quad \frac{\partial g}{\partial x} = e^x$$

$$g_\alpha(x) = \alpha x \quad\Rightarrow\quad \frac{\partial g}{\partial x} = \alpha$$



$-0.53 \cdot \mathrm{e}^{-1} = -0.2$

# NNs as Computational Graphs

- $f(\boldsymbol{w}, \boldsymbol{x}) = \dfrac{1}{1+e^{-(b+w_0 x_0 + w_1 x_1)}}$

$g(x) = \dfrac{1}{x} \qquad \Rightarrow \dfrac{\partial g}{\partial x} = -\dfrac{1}{x^2}$

$g_\alpha(x) = \alpha + x \Rightarrow \dfrac{\partial g}{\partial x} = 1$

$g(x) = e^x \qquad \Rightarrow \dfrac{\partial g}{\partial x} = e^x$

$g_\alpha(x) = \alpha x \qquad \Rightarrow \dfrac{\partial g}{\partial x} = \alpha$



$-0.2 \cdot -1 = 0.2$

# NNs as Computational Graphs

- $f(\boldsymbol{w}, \boldsymbol{x}) = \dfrac{1}{1+e^{-(b+w_0 x_0 + w_1 x_1)}}$

$$g(x) = \frac{1}{x} \qquad \Rightarrow \frac{\partial g}{\partial x} = -\frac{1}{x^2}$$
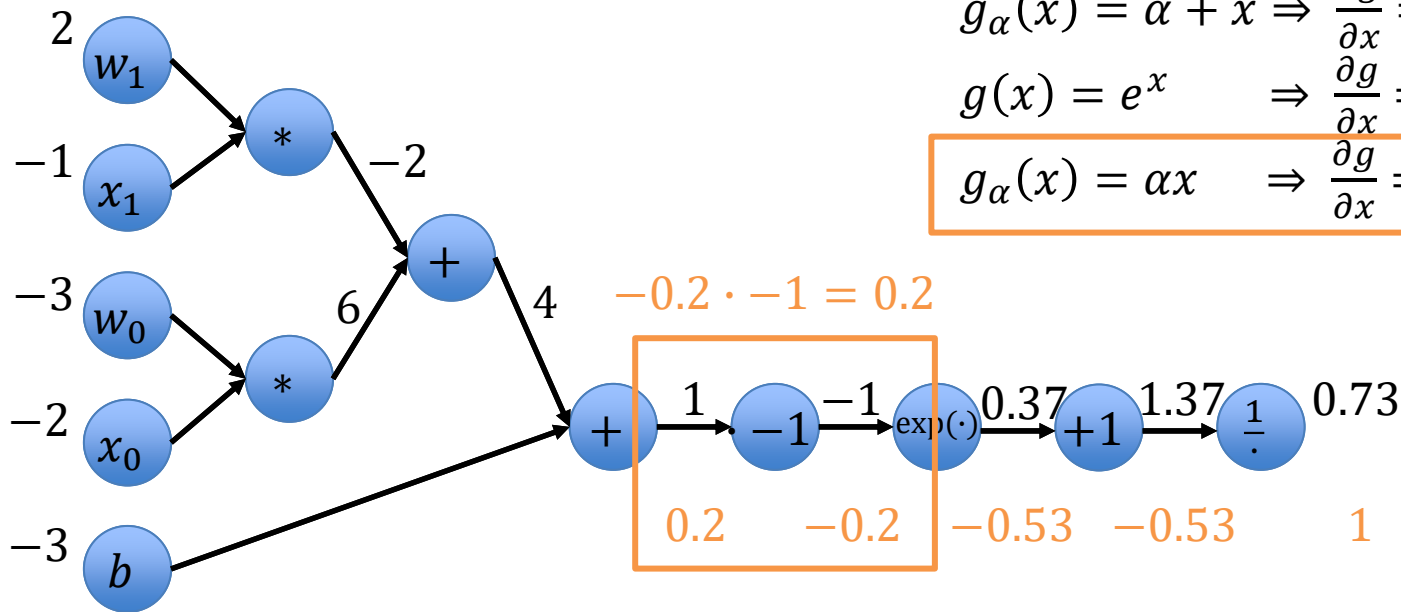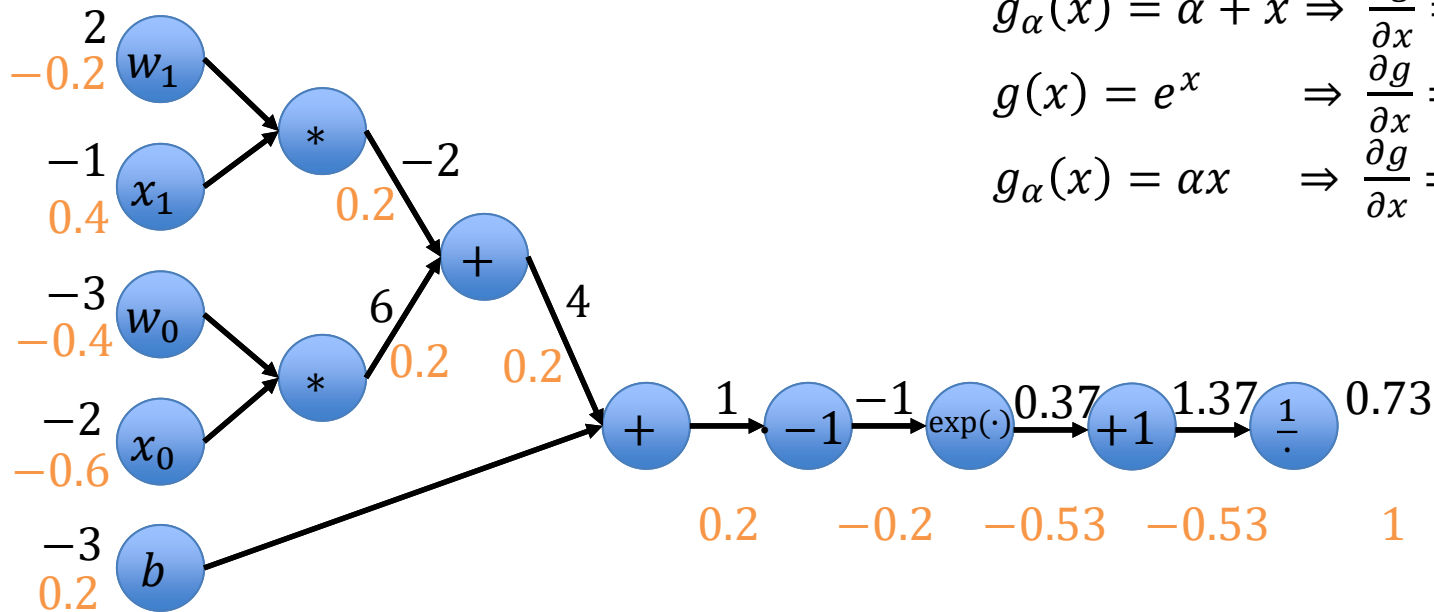
$$g_\alpha(x) = \alpha + x \Rightarrow \frac{\partial g}{\partial x} = 1$$

$$g(x) = e^x \qquad \Rightarrow \frac{\partial g}{\partial x} = e^x$$

$$g_\alpha(x) = \alpha x \quad \Rightarrow \frac{\partial g}{\partial x} = \alpha$$

# Gradient Descent

# Gradient Descent

$$x^* = \arg\min f(x)$$

first order optimization algorithm

Initialization

★ Optimum

# Gradient Descent

- From derivative to gradient

$$\frac{\mathrm{d}f(x)}{\mathrm{d}x} \longrightarrow \nabla_x f(x)$$

Direction of greatest increase of the function

check it

- Gradient steps in direction of negative gradient

$\nabla_x f(x)$

$x$

$$x' = x - \alpha \nabla_x f(x)$$

Learning rate

# Gradient Descent for Neural Networks

# Gradient Descent for Neural Networks

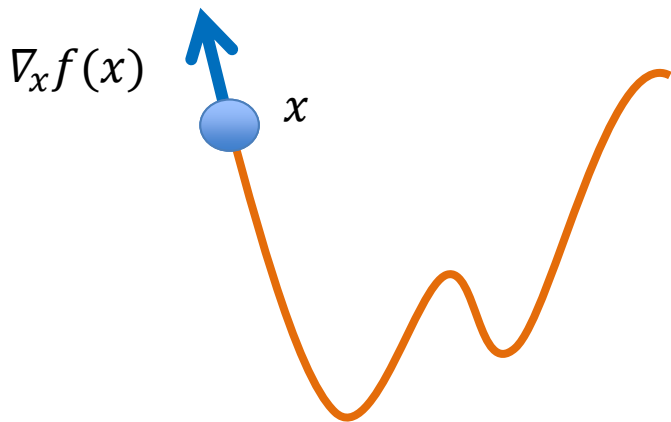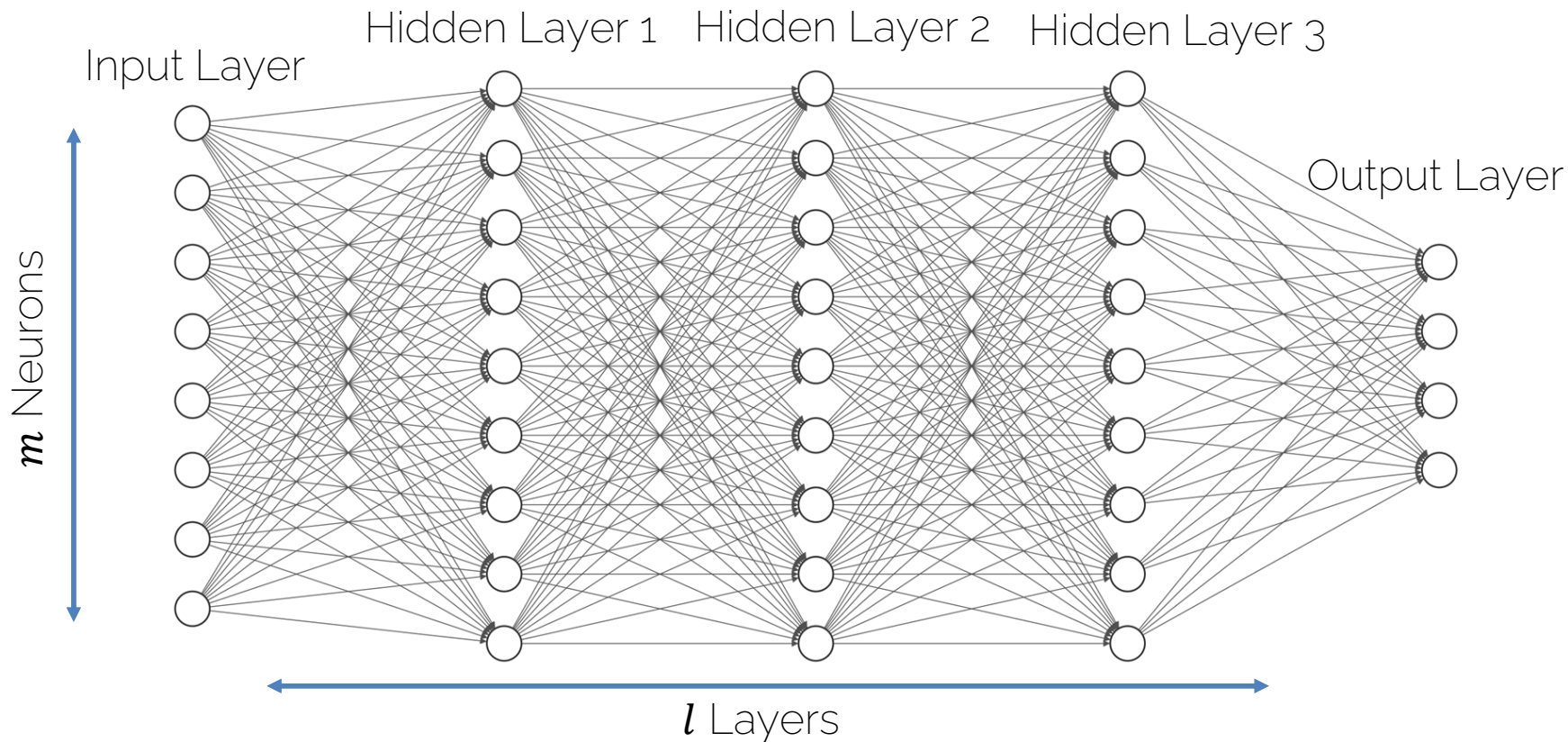For a given training pair $\{\boldsymbol{x}, \boldsymbol{y}\}$, we want to update all weights, i.e., we need to compute the derivatives w.r.t. to all weights:

$$\nabla_{\boldsymbol{W}} f_{\{\boldsymbol{x}, \boldsymbol{y}\}}(\boldsymbol{W}) = \begin{bmatrix} \dfrac{\partial f}{\partial w_{0,0,0}} \\ \dots \\ \dots \\ \dfrac{\partial f}{\partial w_{l,m,n}} \end{bmatrix}$$

Gradient step:

$$\boldsymbol{W'} = \boldsymbol{W} - \alpha \nabla_{\boldsymbol{W}} f_{\{\boldsymbol{x}, \boldsymbol{y}\}}(\boldsymbol{W})$$



$m$ Neurons

Input Layer    Hidden Layer 1    Hidden Layer 2    Hidden Layer 3

Output Layer

$l$ Layers

# NNs can Become Quite Complex...

- These graphs can be huge!



Convolution
AvgPool
MaxPool
Concat
Dropout
Fully connected
Softmax

[Szegedy et al.,CVPR'15] Going Deeper with Convolutions

# The Flow of the Gradients

- Many many many many of these nodes form a neural network

NEURONS

- Each one has its own work to do

FORWARD AND BACKWARD PASS

# The Flow of the Gradients



$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$

$x$

Activations

$\frac{\partial z}{\partial x}$

$\frac{\partial z}{\partial y}$

$f$

$z = f(x, y)$

$\frac{\partial L}{\partial z}$

$y$

Activation function

# Gradient Descent for Neural Networks

Loss function
$$L_i = (\hat{y}_i - y_i)^2$$



input layer

hidden layer

output layer

$$\hat{y}_i = A(b_{1,i} + \sum_j h_j w_{1,i,j})$$

$$h_j = A(b_{0,j} + \sum_k x_k w_{0,j,k})$$

Just simple:
$$A(x) = \max(0, x)$$

# Gradient Descent for Neural Networks



$$h_j = A(b_{0,j} + \sum_k x_k w_{0,j,k})$$

$$\hat{y}_i = A(b_{1,i} + \sum_j h_j w_{1,i,j})$$

$$L_i = (\hat{y}_i - y_i)^2$$

Backpropagation

Just go through layer by layer

$$\frac{\partial L}{\partial w_{1,i,j}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w_{1,i,j}}$$

$$\frac{\partial L_i}{\partial \hat{y}_i} = 2(\hat{y}_i - y_i)$$

sunu tam anlamadim backpropta eger 0 dan kucuk bir deger cikarsa guncelleme yapmayacak miyiz

$$\frac{\partial \hat{y}_i}{\partial w_{1,i,j}} = h_j \quad \text{if} > 0, \text{else } 0$$

$$\frac{\partial L}{\partial w_{0,j,k}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial h_j} \cdot \frac{\partial h_j}{\partial w_{0,j,k}}$$

...

# Gradient Descent for Neural Networks



input layer

hidden layer

output layer

$$h_j = A(b_{0,j} + \sum_k x_k w_{0,j,k})$$

$$\hat{y}_i = A(b_{1,i} + \sum_j h_j w_{1,i,j})$$

$$L_i = (\hat{y}_i - y_i)^2$$

How many unknown weights?

- Output layer: $2 \cdot 4 + 2$

- Hidden Layer: $4 \cdot 3 + 4$

  #neurons · #input channels + #biases

Note that some activations have also weights

# Derivatives of Cross Entropy Loss



input layer

hidden layer

output layer

Gradients of weights of last layer:

$$\frac{\partial L}{\partial w_{ji}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial s_i} \cdot \frac{\partial s_i}{\partial w_{ji}}$$

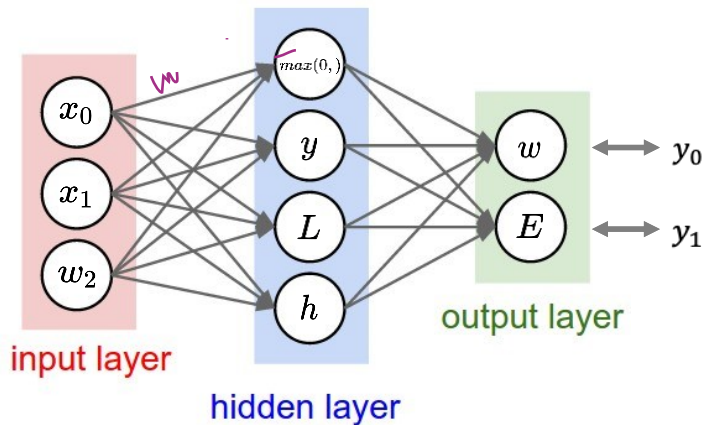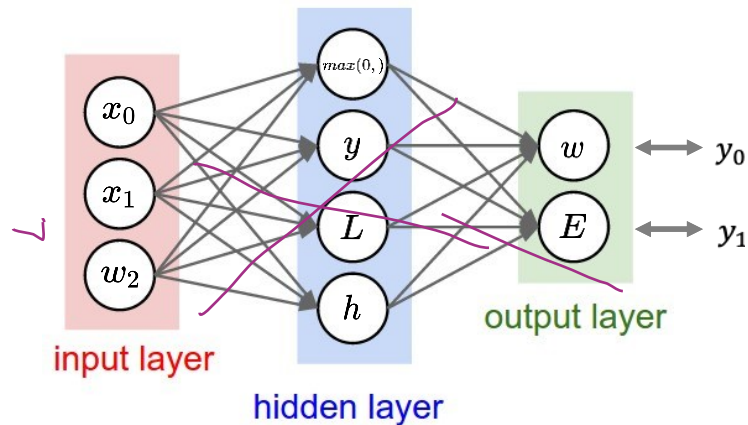$$\frac{\partial L}{\partial \hat{y}_i} = \frac{-y_i}{\hat{y}_i} + \frac{1 - y_i}{1 - \hat{y}_i} = \frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)},$$

$$\frac{\partial \hat{y}_i}{\partial s_i} = \hat{y}_i(1 - \hat{y}_i),$$

$$\frac{\partial s_i}{\partial w_{ji}} = h_j$$

Binary Cross Entropy loss

$$L = -\sum_{i=1}^{n_{out}} (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$

$$\hat{y}_i = \frac{1}{1 + e^{-s_i}} \qquad s_i = \sum_j h_j w_{ji}$$

output          scores

$$\Rightarrow \frac{\partial L}{\partial w_{ji}} = (\hat{y}_i - y_i) h_j, \quad \frac{\partial L}{\partial s_i} = \hat{y}_i - y_i$$

# Derivatives of Cross Entropy Loss

Gradients of weights of first layer:

$$\boxed{\frac{\partial L}{\partial h_j}} = \sum_{i=1}^{n_{out}} \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial s_j} \frac{\partial s_j}{\partial h_j} = \sum_{i=1}^{n_{out}} \frac{\partial L}{\partial \hat{y}_i} \hat{y}_i (1 - \hat{y}_i) w_{ji} = \sum_{i=1}^{n_{out}} (\hat{y}_i - y_i) w_{ji}$$
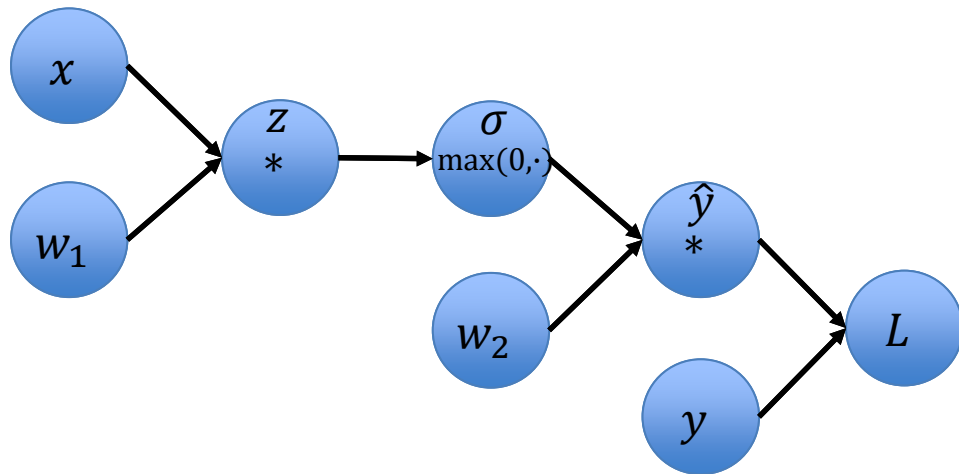
$$\boxed{\frac{\partial L}{\partial s_j^1}} = \boxed{\sum_{i=1}^{n_{out}} \frac{\partial L}{\partial s_i} \frac{\partial s_i}{\partial h_j} \frac{\partial h_j}{\partial s_j^1}} = \boxed{\sum_{i=1}^{n_{out}} (\hat{y}_i - y_i) w_{ji}} \left( h_j (1 - h_j) \right)$$

$$\frac{\partial L}{\partial w_{kj}^1} = \boxed{\sum_{i=1}^{n_{out}} \frac{\partial L}{\partial s_j^1} \frac{\partial s_j^1}{\partial w_{kj}^1}} = \boxed{\sum_{i=1}^{n_{out}} (\hat{y}_i - y_i) w_{ji} \left( h_j (1 - h_j) \right) x_k}$$

# Back to Compute Graphs & NNs

- Inputs $\boldsymbol{x}$ and targets $\boldsymbol{y}$
- Two-layer NN for regression with ReLU activation
- Function we want to optimize:

$$\sum_{i=1}^{n} \|w_2 \max(0, w_1 x_i) - y_i\|_2^2$$
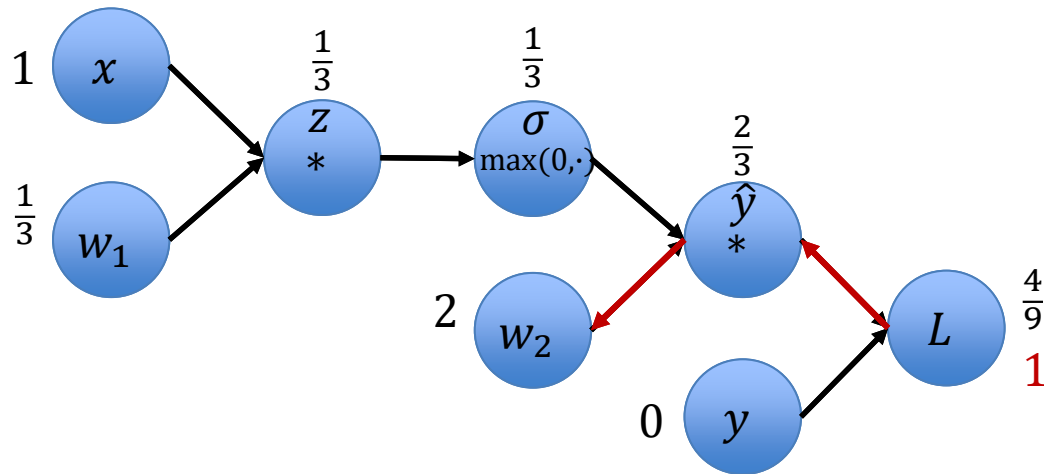
# Gradient Descent for Neural Networks

Initialize $x = 1, \; y = 0,$
$$w_1 = \tfrac{1}{3}, w_2 = 2$$

$$L(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_i^n \|\hat{y}_i - y_i\|_2^2$$

In our case $n, d = 1$:

$$L = (\hat{y} - y)^2 \quad \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\hat{y} = w_2 \cdot \sigma \qquad \Rightarrow \frac{\partial \hat{y}}{\partial w_2} = \sigma$$

1    $x$

$\tfrac{1}{3}$    $w_1$

$\tfrac{1}{3}$   $z$   $*$

$\tfrac{1}{3}$   $\sigma$   $\max(0,\cdot)$

$\tfrac{2}{3}$   $\hat{y}$   $*$

2    $w_2$

$L$   $\tfrac{4}{9}$

   1

0    $y$

Backpropagation
$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$
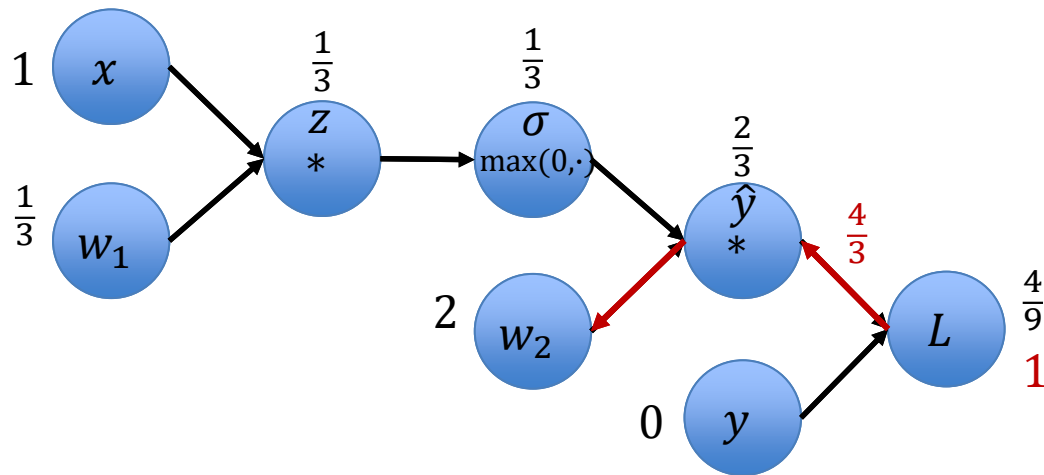
# Gradient Descent for Neural Networks

Initialize $x = 1$, $y = 0$,
$w_1 = \frac{1}{3}, w_2 = 2$

$L(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_i^n \|\widehat{y}_i - y_i\|_2^2$

In our case $n, d = 1$:

$L = (\widehat{y} - y)^2 \quad \Rightarrow \boxed{\dfrac{\partial L}{\partial \widehat{y}} = 2(\widehat{y} - y)}$

$\widehat{y} = w_2 \cdot \sigma \qquad \Rightarrow \dfrac{\partial \widehat{y}}{\partial w_2} = \sigma$

Backpropagation

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \widehat{y}} \cdot \frac{\partial \widehat{y}}{\partial w_2}$$

$$2 \cdot \frac{2}{3}$$
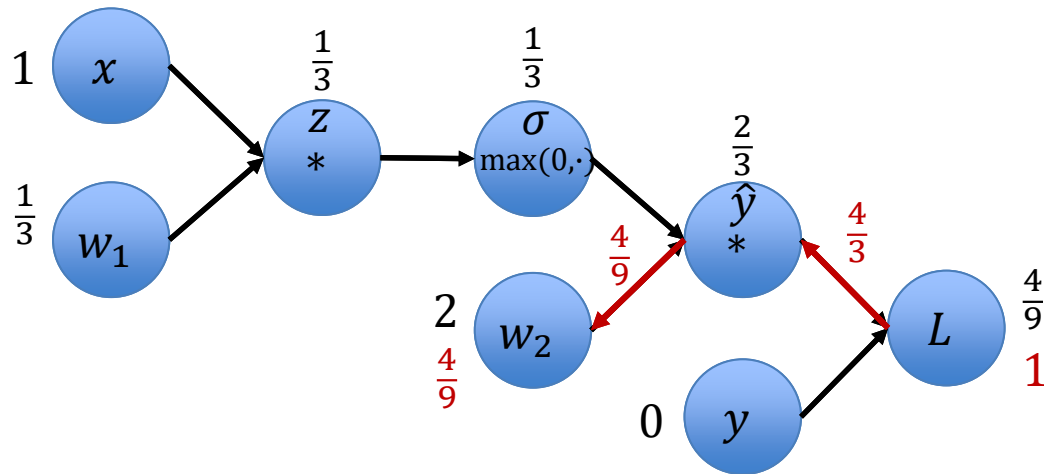
# Gradient Descent for Neural Networks

Initialize $x = 1$, $y = 0$,
$$w_1 = \frac{1}{3}, w_2 = 2$$

$$L(\boldsymbol{y}, \hat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_i^n ||\hat{y}_i - y_i||_2^2$$

In our case $n, d = 1$:

$$L = (\hat{y} - y)^2 \quad \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\hat{y} = w_2 \cdot \sigma \qquad \Rightarrow \boxed{\frac{\partial \hat{y}}{\partial w_2} = \sigma}$$

1   $x$

$\frac{1}{3}$   $w_1$

$z$ $*$   $\frac{1}{3}$

$\sigma$ $\max(0,\cdot)$   $\frac{1}{3}$

$\hat{y}$ $*$   $\frac{2}{3}$

$\frac{4}{9}$

2   $w_2$   $\frac{4}{9}$

$\frac{4}{3}$

$L$   $\frac{4}{9}$   1

0   $y$

Backpropagation
$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

$$2 \cdot \frac{2}{3} \cdot \frac{1}{3}$$

# Gradient Descent for Neural Networks
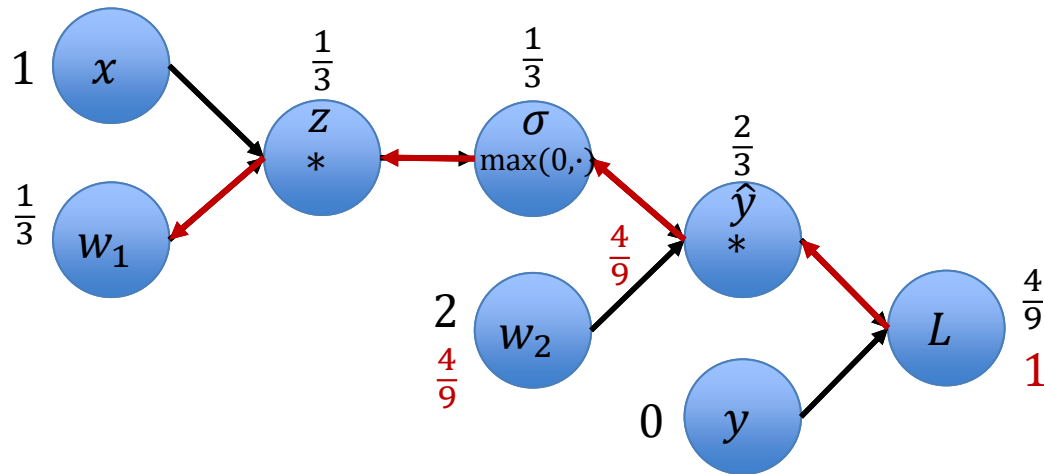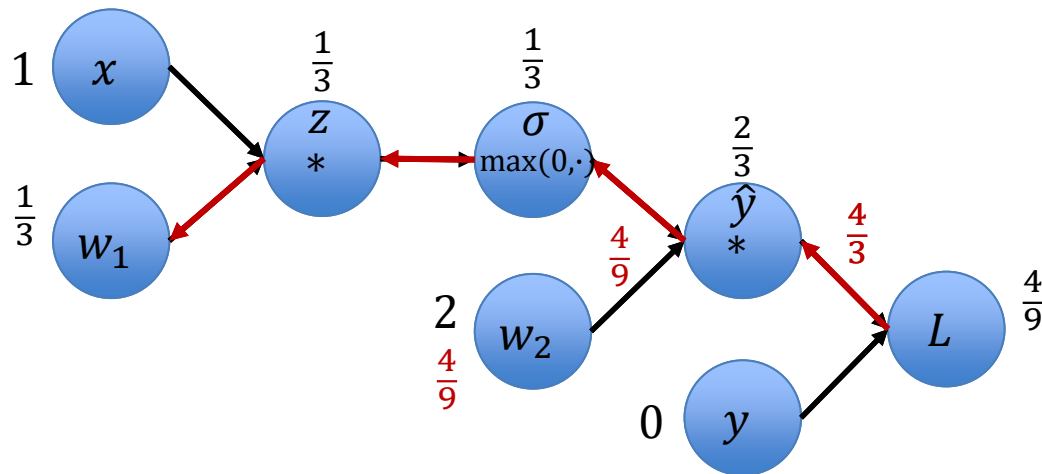
Initialize $x = 1$, $y = 0$,
$$w_1 = \frac{1}{3}, w_2 = 2$$

In our case $n, d = 1$:

$L = (\hat{y} - y)^2 \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$

$\hat{y} = w_2 \cdot \sigma \quad \Rightarrow \frac{\partial \hat{y}}{\partial \sigma} = w_2$

$\sigma = \max(0, z) \Rightarrow \frac{\partial \sigma}{\partial z} = \begin{cases} 1 \text{ if } x > 0 \\ \quad 0 \text{ else} \end{cases}$

$z = x \cdot w_1 \quad \Rightarrow \frac{\partial z}{\partial w_1} = x$



$1 \quad x$

$\frac{1}{3} \quad w_1$

$\frac{1}{3} \quad z \quad *$

$\frac{1}{3} \quad \sigma \quad \max(0, \cdot)$

$\frac{2}{3} \quad \hat{y} \quad *$

$\frac{4}{9}$

$2 \quad w_2$
$\frac{4}{9}$

$0 \quad y$

$L \quad \frac{4}{9}$

$1$

Backpropagation
$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

# Gradient Descent for Neural Networks

Initialize $x = 1, \ y = 0,$
$\qquad w_1 = \frac{1}{3}, w_2 = 2$

In our case $n, d = 1$:

$L = (\hat{y} - y)^2 \quad \Rightarrow \boxed{\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)}$

$\hat{y} = w_2 \cdot \sigma \qquad \Rightarrow \frac{\partial y}{\partial \sigma} = w_2$

$\sigma = \max(0, z) \Rightarrow \frac{\partial \sigma}{\partial z} = \begin{cases} 1 \text{ if } x > 0 \\ \quad 0 \text{ else} \end{cases}$

$z = x \cdot w_1 \qquad \Rightarrow \frac{\partial z}{\partial w_1} = x$

1  $x$

$\frac{1}{3}$  $w_1$

$\frac{1}{3}$  $z$ $*$

$\frac{1}{3}$  $\sigma$ $\max(0, \cdot)$

$\frac{4}{9}$

$\frac{2}{3}$  $\hat{y}$ $*$

2  $w_2$

$\frac{4}{9}$

$\frac{4}{3}$

$L$  $\frac{4}{9}$

0  $y$

Backpropagation
$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$2 \cdot \frac{2}{3}$

# Gradient Descent for Neural Networks

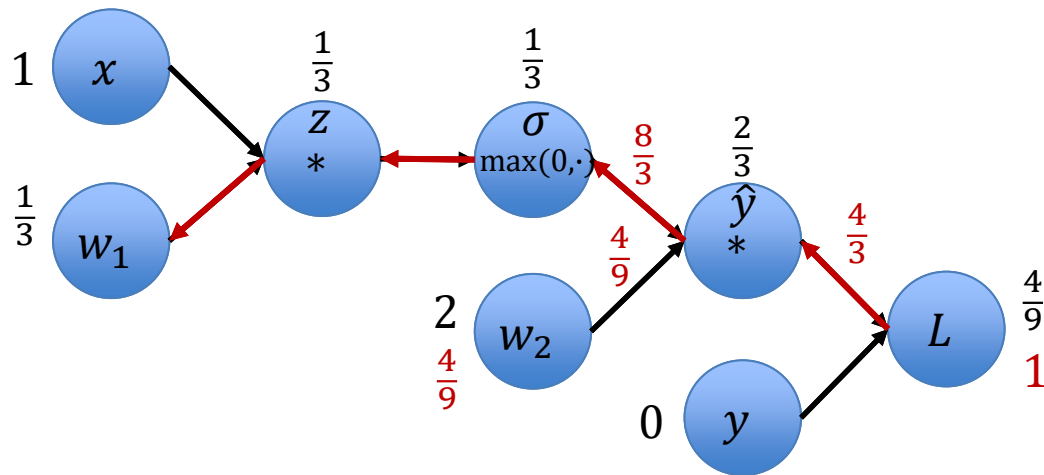Initialize $x = 1, \ y = 0,$
$$w_1 = \tfrac{1}{3}, w_2 = 2$$

In our case $n, d = 1$:

$L = (\hat{y} - y)^2 \quad \Rightarrow \dfrac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$

$\hat{y} = w_2 \cdot \sigma \quad \Rightarrow \boxed{\dfrac{\partial \hat{y}}{\partial \sigma} = w_2}$

$\sigma = \max(0, z) \Rightarrow \dfrac{\partial \sigma}{\partial z} = \begin{cases} 1 \text{ if } x > 0 \\ \ 0 \text{ else} \end{cases}$

$z = x \cdot w_1 \quad \Rightarrow \dfrac{\partial z}{\partial w_1} = x$

Backpropagation

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$2 \cdot \tfrac{2}{3} \cdot 2$$

# Gradient Descent for Neural Networks
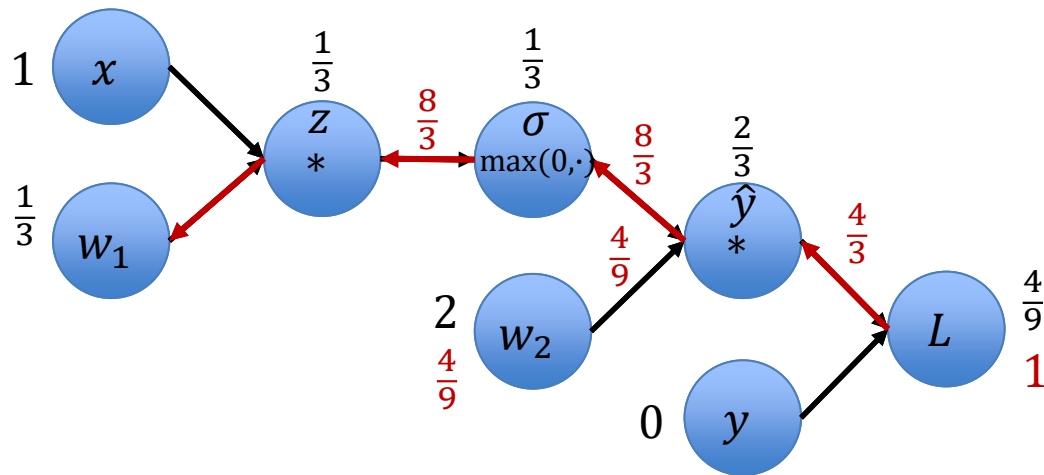
Initialize $x = 1$, $y = 0$,
$$w_1 = \frac{1}{3}, w_2 = 2$$

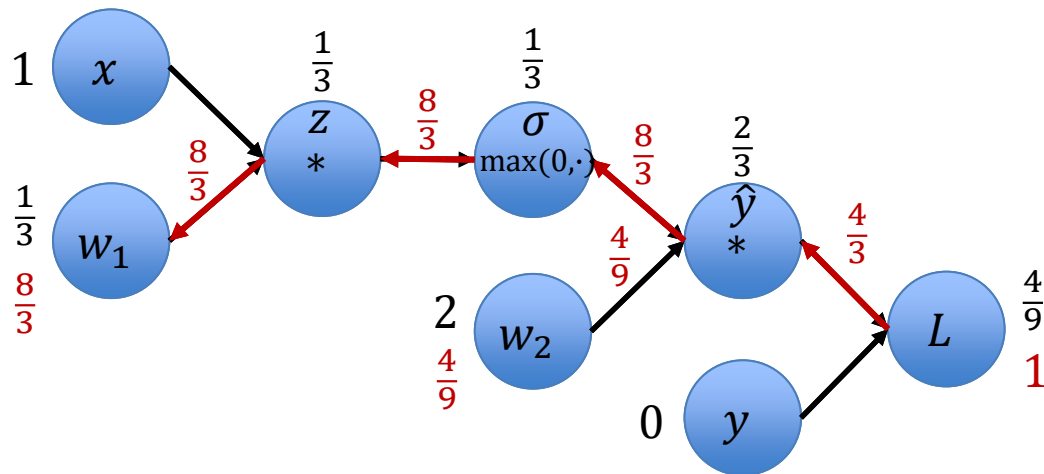In our case $n, d = 1$:

$L = (\hat{y} - y)^2 \quad \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$

$\hat{y} = w_2 \cdot \sigma \quad \Rightarrow \frac{\partial \hat{y}}{\partial \sigma} = w_2$

$\sigma = \max(0, z) \Rightarrow \boxed{\frac{\partial \sigma}{\partial z} = \begin{cases} 1 \text{ if } x > 0 \\ \quad 0 \text{ else} \end{cases}}$

$z = x \cdot w_1 \quad \Rightarrow \frac{\partial z}{\partial w_1} = x$

Backpropagation
$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$
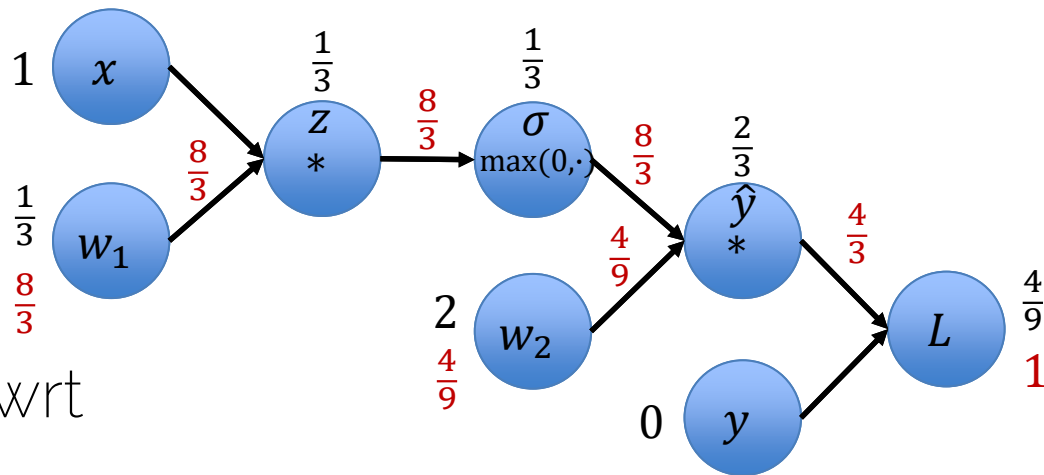
$$2 \cdot \frac{2}{3} \cdot 2 \cdot 1$$

# Gradient Descent for Neural Networks

Initialize $x = 1, \ y = 0,$
$$w_1 = \frac{1}{3}, w_2 = 2$$

In our case $n, d = 1$:

$L = (\hat{y} - y)^2 \quad \Rightarrow \frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)$

$\hat{y} = w_2 \cdot \sigma \qquad \Rightarrow \frac{\partial \hat{y}}{\partial \sigma} = w_2$

$\sigma = \max(0, z) \Rightarrow \frac{\partial \sigma}{\partial z} = \begin{cases} 1 \text{ if } x > 0 \\ \phantom{1} 0 \text{ else} \end{cases}$

$z = x \cdot w_1 \qquad \Rightarrow \boxed{\dfrac{\partial z}{\partial w_1} = x}$



Backpropagation
$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$2 \cdot \frac{2}{3} \ \cdot 2 \ \cdot 1 \quad \cdot 1$$

# Gradient Descent for Neural Networks

- Function we want to optimize:

$$f(x, \boldsymbol{w}) = \sum_{i=1}^{n} \|w_2 \max(0, w_1 x_i) - y_i\|_2^2$$



- Computed gradients wrt to weights $\boldsymbol{w_1}$ and $\boldsymbol{w_2}$

- Now: update the weights

$$\boldsymbol{w'} = \boldsymbol{w} - \alpha \cdot \nabla_{\boldsymbol{w}} f = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} - \alpha \cdot \begin{pmatrix} \nabla_{w_1} f \\ \nabla_{w_2} f \end{pmatrix}$$

$$= \begin{pmatrix} \frac{1}{3} \\ 2 \end{pmatrix} - \alpha \cdot \begin{pmatrix} \frac{8}{3} \\ \frac{4}{9} \end{pmatrix}$$

But: how to choose a good learning rate $\boldsymbol{\alpha}$ ?

# Gradient Descent

- How to pick good learning rate?

  look other models learning rates adsljldas at see which is converges best

- How to compute gradient for single training pair?

- How to compute gradient for large training set?
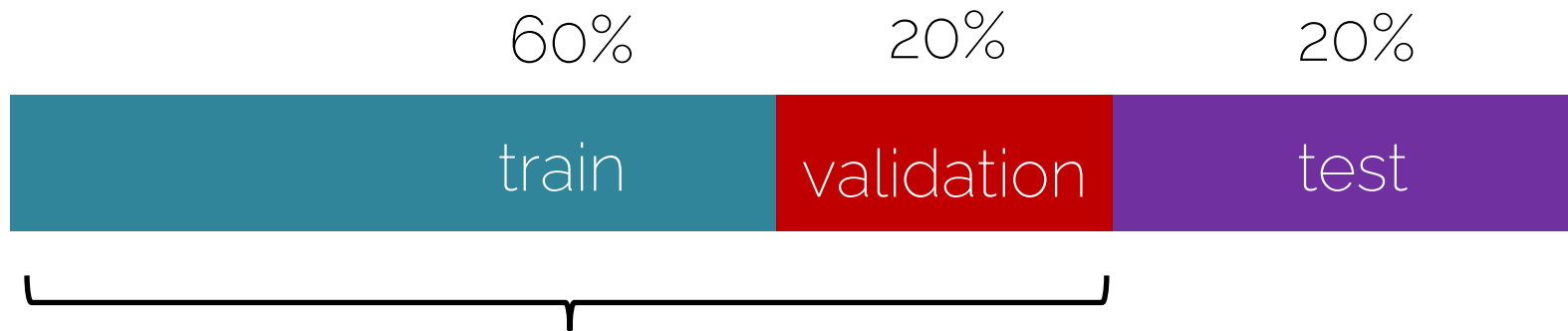
  divide data to batches and compute batches gradient

- How to speed things up? More to see in next lectures...

# Regularization
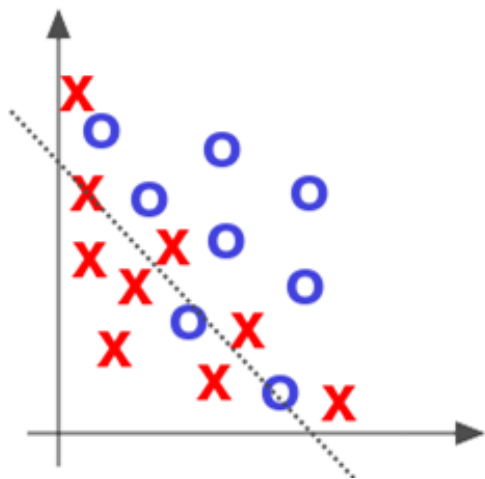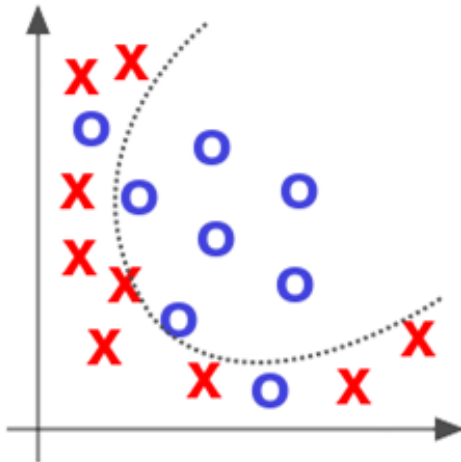
# Recap: Basic Recipe for ML

- Split your data

| 60% | 20% | 20% |
|---|---|---|
| train | validation | test |

Find your hyperparameters

Other splits are also possible (e.g., 80%/10%/10%)

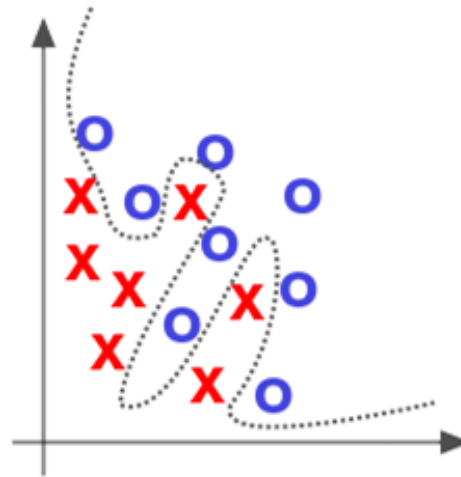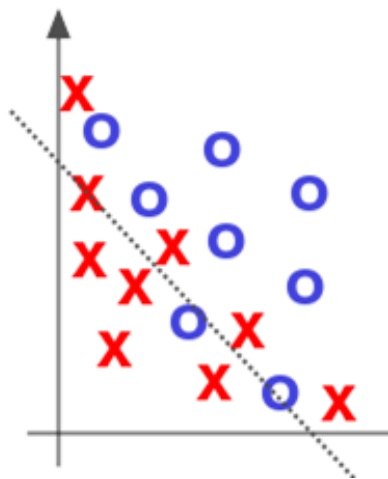# Over- and Underfitting



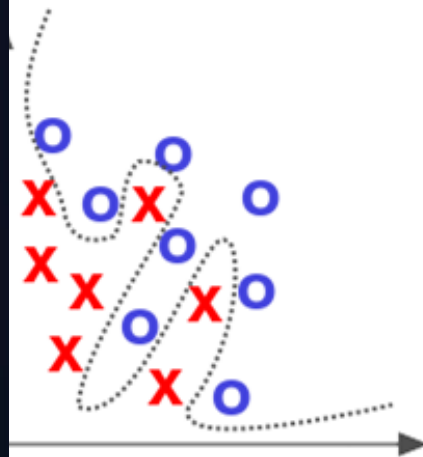Underfitted        Appropriate        Overfitted

Source: Deep Learning by Adam Gibson, Josh Patterson, O'Reily Media Inc., 2017

# Over- and Underfitting



Underfitted

Overfitted

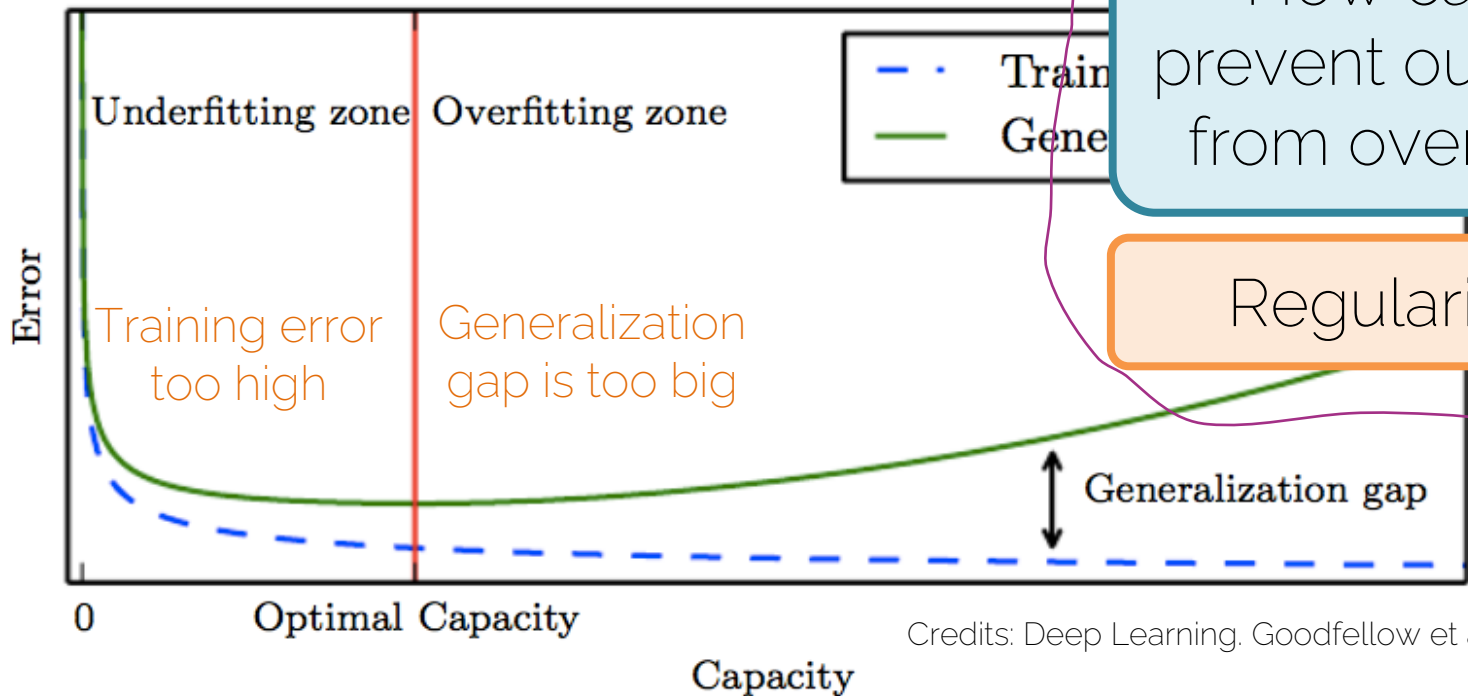ML Engineers looking at their classification model running on the test set.

@debo

Labels

...eily Media Inc., 2017

# Training a Neural Network

- Training/ Validation curve



**Underfitting zone** | **Overfitting zone**

Training error too high

Generalization gap is too big

Error

Generalization gap

0    Optimal Capacity

Capacity

Tra...

Gene...

How can we prevent our model from overfitting?

Regularization

Credits: Deep Learning. Goodfellow et al.

# Regularization

- Loss function $L(\boldsymbol{y}, \widehat{\boldsymbol{y}}, \boldsymbol{\theta}) = \sum_{i=1}^{n}(\hat{y}_i - y_i)^2$ [ ]

- Regularization techniques
  - L2 regularization
  - L1 regularization
  - Max norm regularization
  - Dropout
  - Early stopping
  - ...

Add **regularization term** to loss function

# Regularization

- Loss function $L(\boldsymbol{y}, \hat{\boldsymbol{y}}, \boldsymbol{\theta}) = \sum_{i=1}^{n}(\hat{y}_i - y_i)^2 \boxed{+ \lambda R(\boldsymbol{\theta})}$

  Model weights

- Regularization techniques
  - L2 regularization
  - L1 regularization

    Add **regularization term** to loss function

  - Max norm regularization
  - Dropout
  - Early stopping
  - ...

    More details later

# Regularization: Example

if we dont make regularization, model can ignore some features like in the first example,

- Input: 3 features $\boldsymbol{x} = [1, 2, 1]$

- Two linear classifiers that give the same result:

- $\theta_1 = [0, 0.75, 0]$     ⟶     Ignores 2 features

- $\theta_2 = [0.25, 0.5, 0.25]$  ⟶  Takes information from all features

# Regularization: Example

- Loss $L(\boldsymbol{y}, \widehat{\boldsymbol{y}}, \boldsymbol{\theta}) = \sum_{i=1}^{n}(x_i \theta_{ji} - y_i)^2 \boxed{+ \lambda R(\boldsymbol{\theta})}$

- L2 regularization $R(\boldsymbol{\theta}) = \sum_{i=1}^{n} \theta_i^2$

$\theta_1 \longrightarrow 0 + 0.75^2 + 0 = 0.5625$

so if weights are high in some spesific term, that weight would be high and reg cost high
We would like to take information from every features

$\theta_2 \longrightarrow 0.25^2 + 0.5^2 + 0.25^2 = \boxed{0.375}$   Minimization

$x = [1, 2, 1], \theta_1 = [0, 0.75, 0], \theta_2 = [0.25, 0.5, 0.25]$

# Regularization: Example

- Loss $L(\boldsymbol{y}, \widehat{\boldsymbol{y}}, \boldsymbol{\theta}) = \sum_{i=1}^{n}(x_i\theta_{ji} - y_i)^2 \boxed{+ \lambda R(\boldsymbol{\theta})}$

- L1 regularization $R(\boldsymbol{\theta}) = \sum_{i=1}^{n}|\theta_i|$

vice versa in here . Maybe you may want these result in image processing .

$\theta_1 \longrightarrow 0 + 0.75 + 0 = \boxed{0.75}$ Minimization

$\theta_2 \longrightarrow 0.25 + 0.5 + 0.25 = 1$

$x = [1, 2, 1], \theta_1 = [0, 0.75, 0], \theta_2 = [0.25, 0.5, 0.25]$

# Regularization: Example

- Input: 3 features $\boldsymbol{x} = [1, 2, 1]$

- Two linear classifiers that give the same result:

- $\theta_1 = [0, 0.75, 0]$ $\longrightarrow$ Ignores 2 features

- $\theta_2 = [0.25, 0.5, 0.25]$ $\longrightarrow$ Takes information from all features

# Regularization: Example

- Input: 3 features $\boldsymbol{x} = [1, 2, 1]$

- Two linear classifiers that give the same result:

- $\theta_1 = [0, 0.75, 0]$ $\longrightarrow$ L1 regularization enforces **sparsity**

- $\theta_2 = [0.25, 0.5, 0.25]$ $\longrightarrow$ L2 regularization enforces that the weights have **similar values**

# Regularization: Effect

- Dog classifier takes different inputs

As i said that can be good for image classfacition

Furry

Has two eyes

Has a tail

Has paws

Has two ears

L1 regularization will focus all the attention to a few key features

# Regularization: Effect

- Dog classifier takes different inputs

Furry

Has two eyes

Has a tail

Has paws

Has two ears

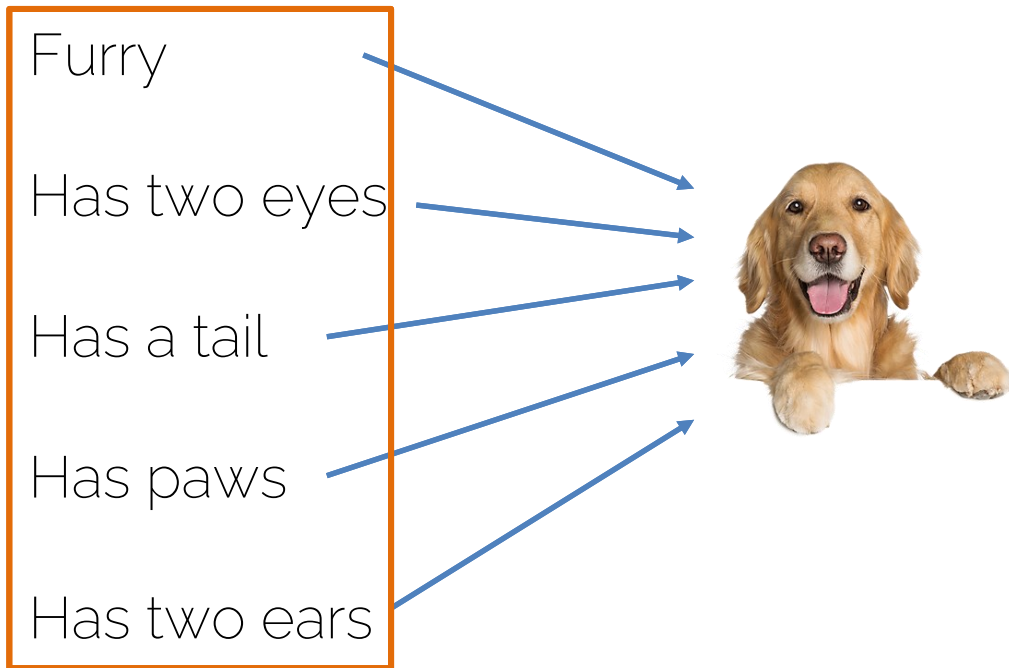**L2 regularization** will take all information into account to make decisions
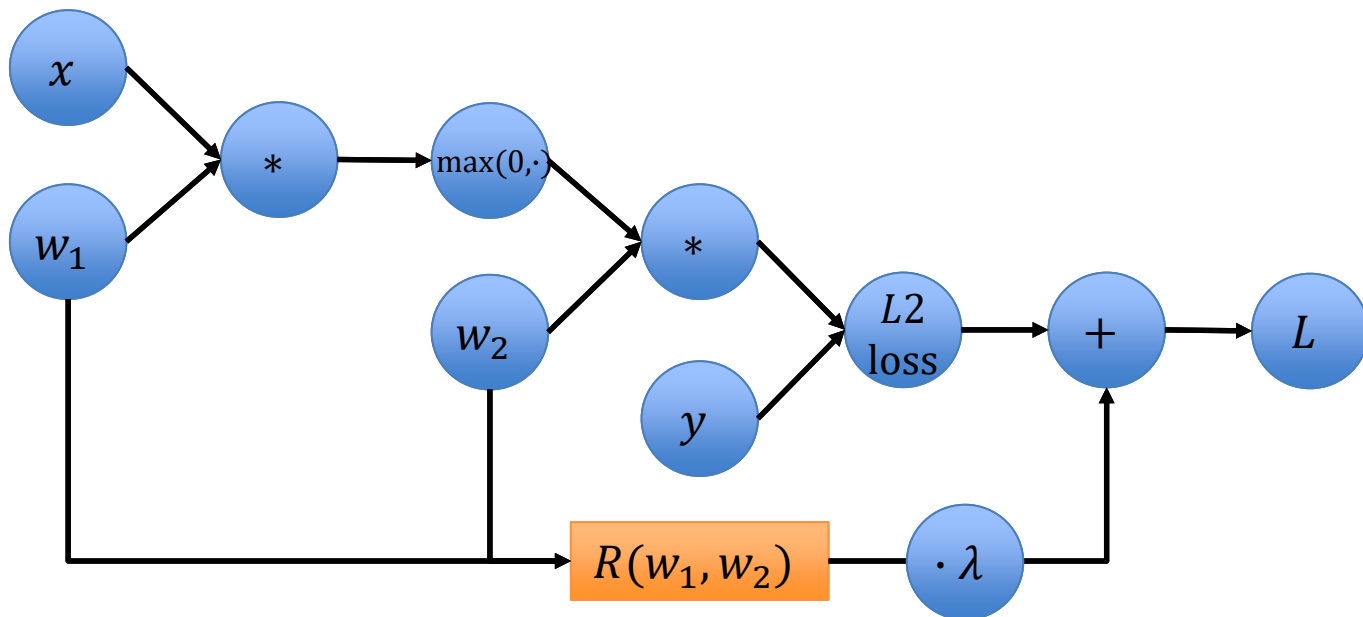
# Regularization for Neural Networks



Combining nodes:
Network output + L2-loss +
regularization

$$\sum_{i=1}^{n} \|w_2 \max(0, w_1 x_i) - y_i\|_2^2 + \lambda R(w_1, w_2)$$
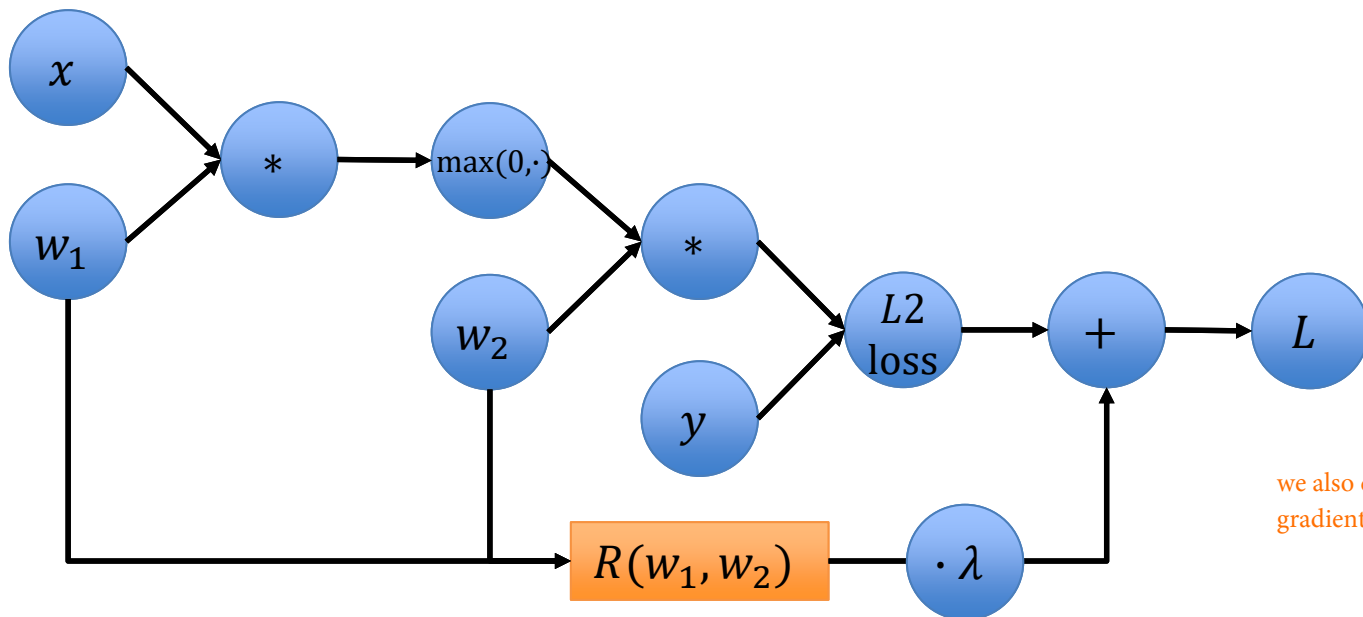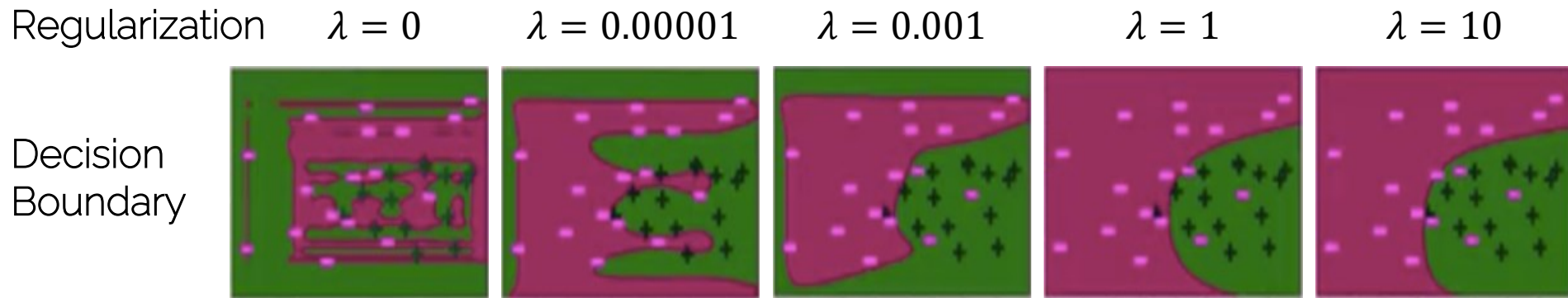
# Regularization for Neural Networks



Combining nodes:
Network output + L2-loss +
regularization

$$\sum_{i=1}^{n} \|w_2 \max(0, w_1 x_i) - y_i\|_2^2 + \lambda \left\| \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \right\|_2^2$$

# Regularization for Neural Networks



we also calculate backward gradients for regularization too

Combining nodes:
Network output + L2-loss + regularization

$$\sum_{i=1}^{n} \|w_2 \max(0, w_1 x_i) - y_i\|_2^2 + \lambda(w_1^2 + w_2^2)$$

# Regularization

Regularization     $\lambda = 0$      $\lambda = 0.00001$      $\lambda = 0.001$      $\lambda = 1$      $\lambda = 10$

Decision
Boundary



Credit: University of Washington

What happens to the training error? increase
also making training more difficult
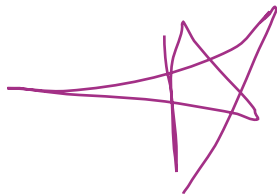
What is the goal of regularization?

# Regularization

- Any strategy that aims to

Lower validation error

Increasing training error

# Next Lecture

- This week:
  - Check exercises!
  - Check piazza / post questions ☺

- Next lecture
  - Optimization of Neural Networks
  - In particular, introduction to SGD (our main method!)

See you next week ☺

# Further Reading

- Backpropagation
  - Chapter 6.5 (6.5.1 - 6.5.3) in
    http://www.deeplearningbook.org/contents/mlp.html
  - Chapter 5.3 in Bishop, Pattern Recognition and Machine Learning
  - http://cs231n.github.io/optimization-2/

- Regularization
  - Chapter 7.1 (esp. 7.1.1 & 7.1.2)
    http://www.deeplearningbook.org/contents/regularization.html
  - Chapter 5.5 in Bishop, Pattern Recognition and Machine Learning