# Outline

- Basic concepts
- Decision tree induction
- Evaluation of classifiers
- Naïve Bayesian classification
- Naïve Bayes for text classification
- Support vector machines
- **Linear regression and gradient descent**
- Neural networks
- K-nearest neighbor
- Ensemble methods
- Summary

# Linear regression

- ## Supervised learning has two main types
  - Classification: discrete predictive/output variable
  - Regression: continuous predictive/output variable

- ## We first study linear regression, i.e., the predictive function $h$ is a linear function.

Many slides are borrowed from Jia-Bin Huang and Matt Gormley

# An example: housing price prediction

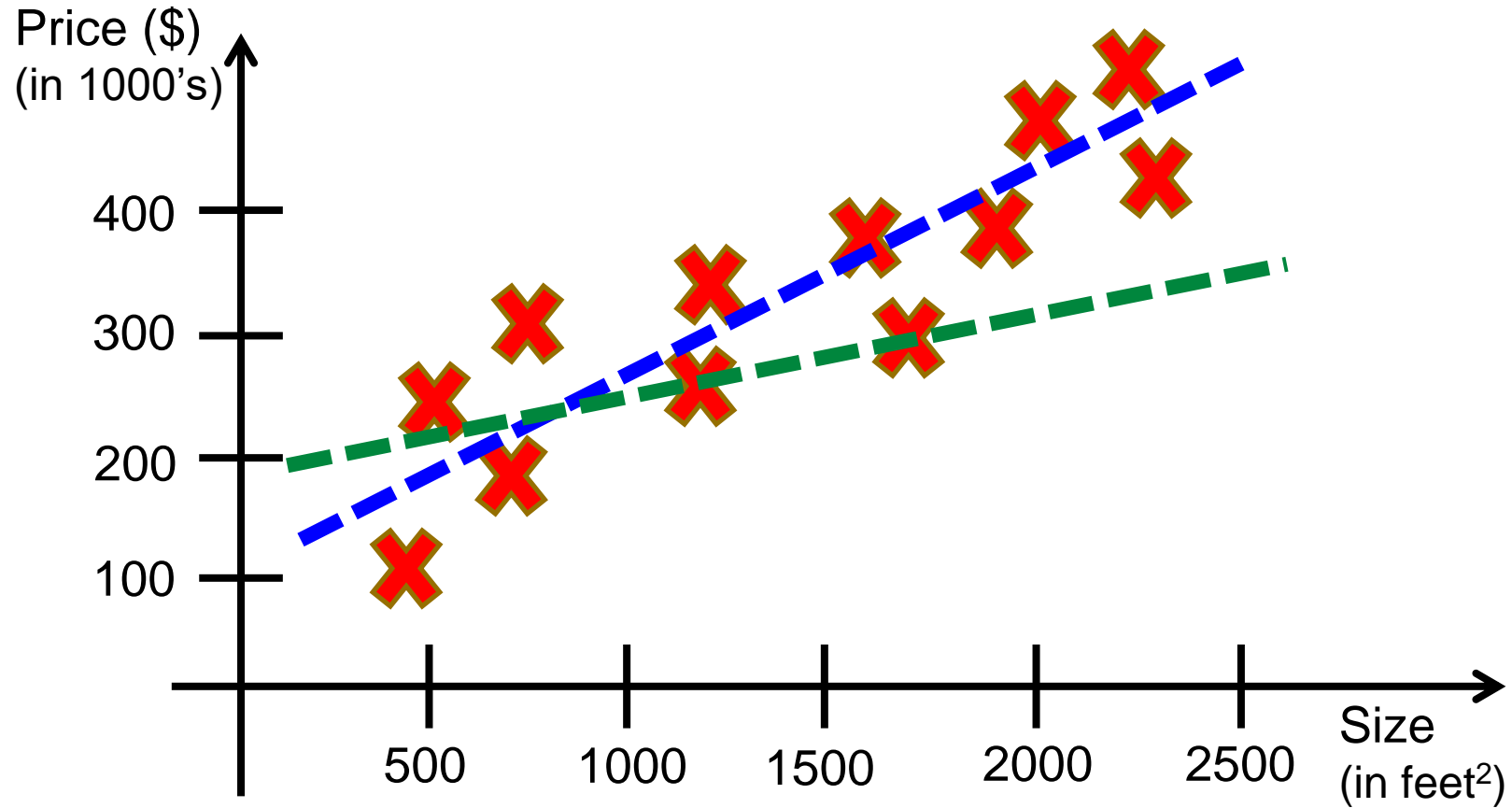- Given the size of a house, predict the price of the house.

- Notation:
  - $n$: Number of training examples
  - $x$: Input variable / feature (Size)
  - $y$: Output variable / target variable (Price)
  - $(x, y)$: One training example in general
  - $(x^i, y^i)$: $i^{th}$ training example

Training data

| Size in feet^2 (x) | Price ($) in 1000's (y) |
|---|---|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| … | … |

# Training data and linear function

# Model representation

- This is a *univariate linear regression problem* as it has only one input variable $x$.

- The linear regression model in this case is as follows

$$y = h_{\boldsymbol{\theta}}(x) = \theta_0 + \theta_1 x$$

  - There are two parameters $\theta_0$ and $\theta_1$.

  - $\boldsymbol{\theta}$ represents the parameter vector, i.e., $(\theta_0, \theta_1)$

- We use the training set to learn this model by optimizing a cost function, also called a *loss function (L)*.
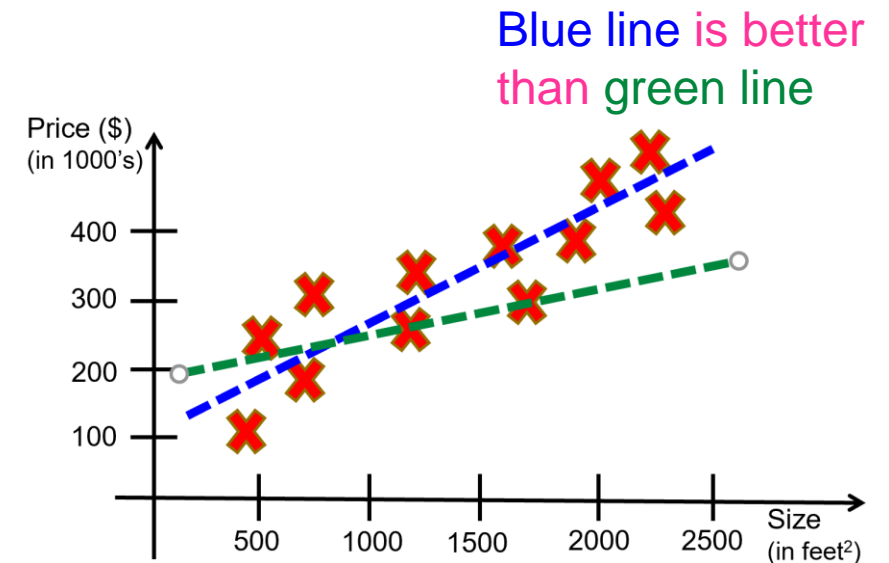
# Loss function

- **Idea:** select $\theta_0$, $\theta_1$ so that $h_{\boldsymbol{\theta}}(x)$ is close to $y$ for the training example $(x, y)$. This is expressed with a loss function.

- Loss function ($L$) used by linear regression:

$$L(\boldsymbol{\theta}) = L(\theta_0, \theta_1) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(x^i) - y^i \right)^2$$

where $h_{\boldsymbol{\theta}}(x^i) = \theta_0 + \theta_1 x^i$

Blue line is better than green line



- **Learning goal:** $\underset{\theta_0,\, \theta_1}{\arg\min}\ L(\theta_0, \theta_1)$

# Solve the minimization problem

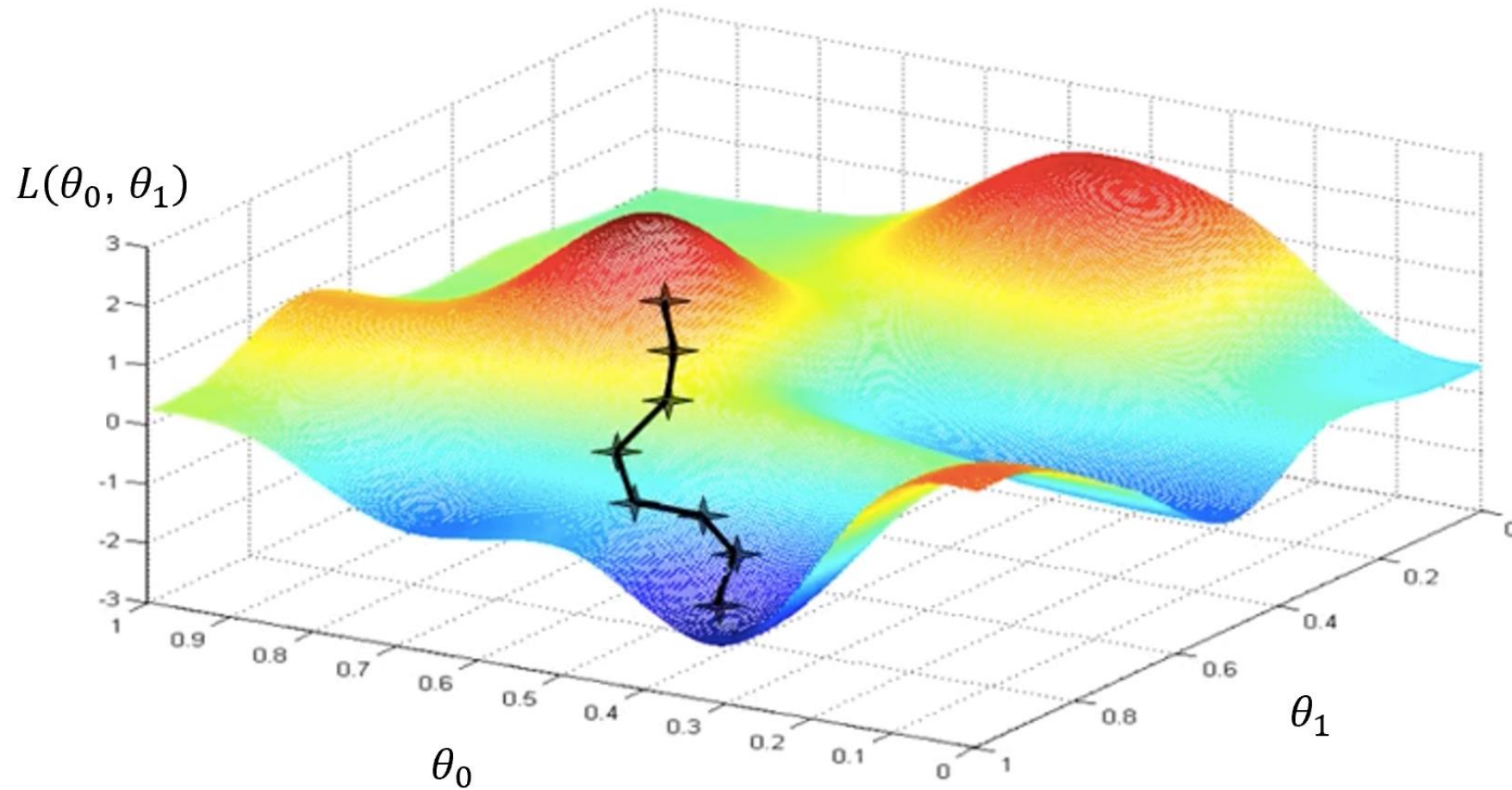- The learning is done using a general technique called
  - gradient descent

# Gradient descent

■ Recall our univariate linear regression problem

❑ Loss function: $L(\theta_0, \theta_1)$

❑ Goal: $\underset{\theta_0, \theta_1}{\operatorname{argmin}} \; L(\theta_0, \theta_1)$

Steps:

■ Start with some initial $\theta_0, \theta_1$

■ Keep changing $\theta_0, \theta_1$ to reduce $L(\theta_0, \theta_1)$
  until we hopefully end up at minimum

# An illustration

# Keep going downhill

Learning rule: $\quad \theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} L(\theta_0, \theta_1)$



$$\frac{\partial}{\partial \theta_1} L(\theta_0, \theta_1) < 0$$

$$\frac{\partial}{\partial \theta_1} L(\theta_0, \theta_1) > 0$$

# Gradient descent algorithm

Repeat until convergence

$$\{$$

$$\theta_j := \theta_j - \alpha \, \frac{\partial}{\partial \theta_j} L(\theta_0, \theta_1) \quad \text{(for } j = 0 \text{ and } j = 1)$$

$$\}$$

- $\alpha$: Learning rate (step size)
- $\frac{\partial}{\partial \theta_j} L(\theta_0, \theta_1)$: derivative (rate of change)

# How to update

**Correct:** **simultaneous update**

- $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} L(\theta_0, \theta_1)$

- $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} L(\theta_0, \theta_1)$

- $\theta_0 := \text{temp0}$
- $\theta_1 := \text{temp1}$

**Incorrect:**

- $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} L(\theta_0, \theta_1)$

- $\theta_0 := \text{temp0}$

- $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} L(\theta_0, \theta_1)$

- $\theta_1 := \text{temp1}$

# Learning rate

Too big learning rate          Small learning rate

# Recall: Loss function and learning goal

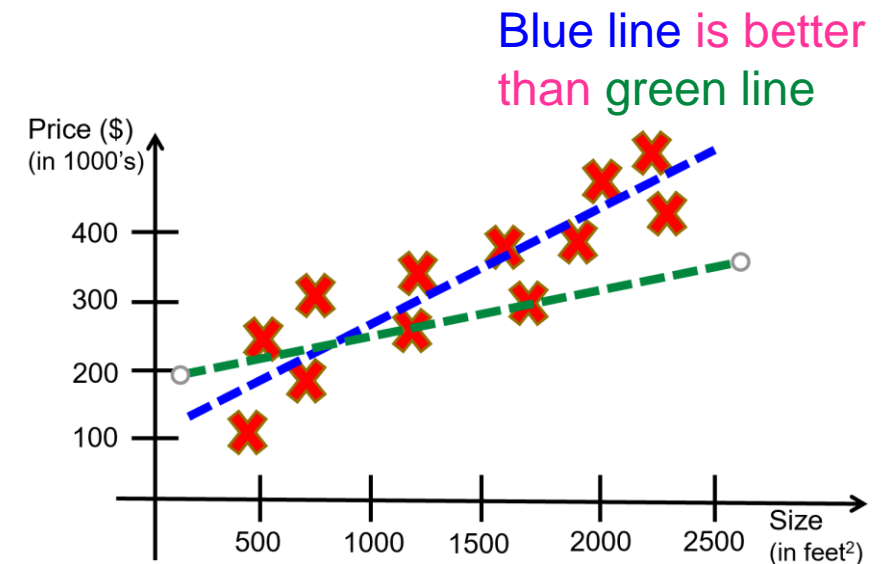- Recall: Loss function (*L*) used by linear regression is:

$$L(\mathbf{\theta}) = L(\theta_0, \theta_1) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\mathbf{\theta}}(x^i) - y^i \right)^2$$

where $h_{\mathbf{\theta}}(x^i) = \theta_0 + \theta_1 x^i$

$h_{\mathbf{\theta}}(x^i)$ is an estimate of $y^i$

- Learning goal:

$$\underset{\theta_0, \theta_1}{\text{argmin}} \; L(\theta_0, \theta_1)$$

Blue line is better than green line

Price ($)
(in 1000's)

400

300

200

100

500   1000   1500   2000   2500

Size
(in feet²)

# Computing partial derivative

- $\dfrac{\partial}{\partial \theta_j} L(\theta_0, \theta_1) = \dfrac{\partial}{\partial \theta_j} \dfrac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(x^i) - y^i \right)^2$

$$= \dfrac{\partial}{\partial \theta_j} \dfrac{1}{2n} \sum_{i=1}^{n} \left( \theta_0 + \theta_1 x^i - y^i \right)^2$$

- $j = 0: \quad \dfrac{\partial}{\partial \theta_0} L(\theta_0, \theta_1) = \dfrac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(x^i) - y^i \right)$

- $j = 1: \quad \dfrac{\partial}{\partial \theta_1} L(\theta_0, \theta_1) = \dfrac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(x^i) - y^i \right) x^i$

# Gradient descent for linear regression

Repeat until convergence
{

$$\theta_0 := \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(x^i) - y^i \right)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(x^i) - y^i \right) x_i$$

}

- Update $\theta_0$ and $\theta_1$ simultaneously

# Batch gradient descent

- Each step or update of gradient descent uses all ($n$) the training examples.

  - Sum over all $n$ training examples for each step – slow

  - It is also memory demanding if the training data is huge.

- In a normal learning process, training needs many steps before convergence.

- The training process that covers all the training examples once is called an epoch.

  - In batch gradient descent, each step is an epoch.

# Stochastic gradient descent (SGD)

- **SGD with one example per step**: In SGD each step uses a single training example. Before each epoch, the data should be shuffled.
  - SGD converges faster when the dataset is large as it causes updates to the parameters more frequently.
    - The loss may fluctuate as only one example is used in each step.

- **SGD with minibatch**: each update/step uses a random ***minibatch*** of ***m*** out of ***n*** examples.
  - It is efficient, more stable, and more likely to jump out of a local minimum

- **Batch Gradient Descent** is more suitable for convex loss functions as it can converge directly to minima.

# Convex and non-convex function

- **A convex function has one minimum.**
  - For all $0 \leq \lambda \leq 1$ and all $x_1, x_2$ in a convex set $X$ (e.g., an interval [a, b]), the following holds

    $$f(\lambda x_1 + (1 - \lambda) x_2) \leq \lambda f(x_1) + (1 - \lambda) f(x_2)$$

- **A non-convex function has local minima (valleys) that are not global minimum.**



Convex

Non-convex

# Multivariate linear regression

- In our previous linear regression problem, we use only one input variable/feature (univariate). In general, the problem can have any number of input variables. Let the number of variables be *k*, $x_1, x_2, \ldots, x_k$.

- Training data: $D = \{\mathbf{x}^i, y^i\}_{i=1}^n$

- Multivariate linear regression model is

$$y = h_{\boldsymbol{\theta}}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_k x_k$$

where $\boldsymbol{\theta}$ is the vector of all $\theta_i$ and $\mathbf{x}$ is the vector of all $x_i$.

# Multivariate linear regression (cont.)

- For convenience of notation, define $x_0 = 1$ ($x_0^j = 1$ for all examples $j$)

- $\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix} \in R^{k+1}$ $\quad\quad$ $\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_k \end{bmatrix} \in R^{k+1}$

- $y = h_{\boldsymbol{\theta}}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_k x_k = \boldsymbol{\theta}^\top \mathbf{x}$

# Univariate and multivariate gradient descent

- Univariate ($k = 1$)

  Repeat until convergence {

  $$\theta_0 := \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(x^i) - y^i \right)$$

  $$\theta_1 := \theta_1 - \alpha \frac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(x^i) - y^i \right) x^i$$

  }

- Multivariate ($k > 1$)

  Repeat until convergence {

  $$\theta_j := \theta_j - \alpha \frac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}(\mathbf{x}^i) - y^i \right) x_j^i$$

  }

  Simultaneously update
  $\theta_j$, for $j = 0, 1, \cdots, k$

# Outline

- Basic concepts
- Decision tree induction
- Evaluation of classifiers
- Naïve Bayesian classification
- Naïve Bayes for text classification
- Support vector machines
- Linear regression and gradient descent
- **Neural networks**
- K-nearest neighbor
- Ensemble methods
- Summary

# Some example successes of neural networks

# Resurgence of neural networks

- Origin: Algorithms that try to mimic the brain (1943).
- Was very widely used in 80s and early 90s; popularity diminished in late 90s.
- Recent resurgence: State-of-the-art results in many applications.
- It works especially well for computer vision and natural language processing (including speech recognition).
  - It has revolutionized the two fields in recent years.
  - It has spread to almost every machine learning area and application in practice.

# A single neuron in the brain

# The first neural network (McCulloch & Pitts, 1943)



In 1943 American neurophysiologist and cybernetician of the University of Illinois at Chicago Warren McCulloch and self-taught logician and cognitive psychologist Walter Pitts published "A Logical Calculus of the ideas Imminent in Nervous Activity," describing the "McCulloch - Pitts neuron," "the first mathematical model of a neural network.

Building on ideas in Alan Turing's "On Computable Numbers", McCulloch and Pitts's paper provided a way to describe brain functions in abstract terms, and showed that simple elements connected in a neural network can have immense computational power. The paper

# Simple model of a neuron (McCulloch & Pitts, 1943)



Bias Weight

$a_0 = 1$

$w_{0,j}$

$a_j = \sigma(in_j)$

$in_j$

$\sigma$

$S$

$w_{i,j}$

$a_i$

$a_j$

Input Links

Input Function

Activation Function

Output

Output Links

- Inputs $a_i$ come from the output of node $i$ to this node $j$ (or from "outside")
- Each input link has a **weight** $w_{i,j}$
- There is an additional fixed input $a_0$ (bias) with weight $w_{0,j}$
- The total input is $in_j = \Sigma_i\, w_{i,j}\, a_i$
- The output is $a_j = \sigma(in_j) = \sigma(\Sigma_i\, w_{i,j}\, a_i) = \sigma(\mathbf{w.a})$

# Logistic regression in a figure



"Bias unit"

$x_0$

$w_0$

$x_1$

$w_1$

$h_{\mathbf{W}}(\mathbf{x}) = \dfrac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}$

$x_2$

$w_2$

$y$

Output

$x_3$

$w_3$

$\Sigma \quad \mathbf{w}^\top \mathbf{x} = z$

input

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \qquad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

"Weights"
"Parameters"

$$y = h_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x}),$$

$$\text{where } \sigma(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid (logistic)
function

# An artificial neuron: a logistic unit

- ## A neuron is a logistic unit
  - $\sigma(\mathbf{w}^\top\mathbf{x})$ is called activation function.
  - Activation function does not have to be sigmoid.

- ## A neural network is a composition of many logistic units organized in layers.
  - It can also be seen as a logistic regression model with one or more hidden layers.



output

$h_{\mathbf{w}}(\mathbf{x})$

Layer 1    layer 2 (hidden)    layer 3

# Neural network: an example



$a_i^{(j)}$ = "activation" of unit $i$ in layer $j$

$\mathbf{W}^{(j)}$ = matrix of weights controlling function mapping from layer $j$ to layer $j+1$
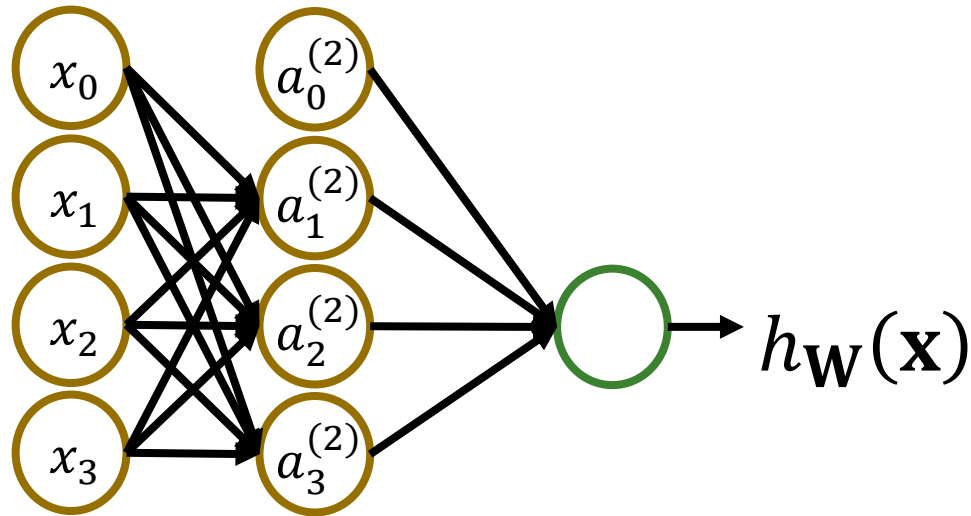
$s_j$ units in layer $j$

$s_{j+1}$ units in layer $j+1$

$$a_1^{(2)} = \sigma\left(\mathbf{W}_{10}^{(1)}x_0 + \mathbf{W}_{11}^{(1)}x_1 + \mathbf{W}_{12}^{(1)}x_2 + \mathbf{W}_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = \sigma\left(\mathbf{W}_{20}^{(1)}x_0 + \mathbf{W}_{21}^{(1)}x_1 + \mathbf{W}_{22}^{(1)}x_2 + \mathbf{W}_{23}^{(1)}x_3\right)$$

$$a_3^{(2)} = \sigma\left(\mathbf{W}_{30}^{(1)}x_0 + \mathbf{W}_{31}^{(1)}x_1 + \mathbf{W}_{32}^{(1)}x_2 + \mathbf{W}_{33}^{(1)}x_3\right)$$

$$h_{\mathbf{W}}(x) = \sigma\left(\mathbf{W}_{10}^{(2)}a_0^{(2)} + \mathbf{W}_{11}^{(2)}a_1^{(2)} + \mathbf{W}_{12}^{(2)}a_2^{(2)} + \mathbf{W}_{13}^{(2)}a_3^{(2)}\right)$$

# Neural network: an example

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\mathbf{z}^{(2)} = \begin{bmatrix} \mathbf{z}_1^{(2)} \\ \mathbf{z}_2^{(2)} \\ \mathbf{z}_3^{(2)} \end{bmatrix}$$

$$a_1^{(2)} = \sigma\left(\mathbf{W}_{10}^{(1)} x_0 + \mathbf{W}_{11}^{(1)} x_1 + \mathbf{W}_{12}^{(1)} x_2 + \mathbf{W}_{13}^{(1)} x_3\right) = \sigma(z_1^{(2)})$$

$$a_2^{(2)} = \sigma\left(\mathbf{W}_{20}^{(1)} x_0 + \mathbf{W}_{21}^{(1)} x_1 + \mathbf{W}_{22}^{(1)} x_2 + \mathbf{W}_{23}^{(1)} x_3\right) = \sigma(z_2^{(2)})$$

$$a_3^{(2)} = \sigma\left(\mathbf{W}_{30}^{(1)} x_0 + \mathbf{W}_{31}^{(1)} x_1 + \mathbf{W}_{32}^{(1)} x_2 + \mathbf{W}_{33}^{(1)} x_3\right) = \sigma(z_3^{(2)})$$

$$h_{\mathbf{W}}(x) = \sigma\left(\mathbf{W}_{10}^{(2)} a_0^{(2)} + \mathbf{W}_{11}^{(2)} a_1^{(2)} + \mathbf{W}_{12}^{(2)} a_2^{(2)} + \mathbf{W}_{13}^{(2)} a_3^{(2)}\right) = \sigma(z^{(3)})$$

# Neural network: an example



"Pre-activation"

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \qquad \mathbf{z}^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$a_1^{(2)} = \sigma(z_1^{(2)})$$
$$a_2^{(2)} = \sigma(z_2^{(2)})$$
$$a_3^{(2)} = \sigma(z_3^{(2)})$$
$$h_{\mathbf{W}}(\mathbf{x}) = \sigma(z^{(3)})$$

$$\mathbf{z}^{(2)} = \mathbf{W}^{(1)}\mathbf{x} = \mathbf{W}^{(1)}\mathbf{a}^{(1)} \quad // \text{ layer 1: } \mathbf{a}^{(1)} = \mathbf{x}$$
$$\mathbf{a}^{(2)} = \sigma(\mathbf{z}^{(2)})$$

Add $a_0^{(2)} = 1$
$$\mathbf{z}^{(3)} = \mathbf{W}^{(2)}\mathbf{a}^{(2)}$$
$$h_{\mathbf{W}}(\mathbf{x}) = \mathbf{a}^{(3)} = \sigma(\mathbf{z}^{(3)})$$

# Neural network learning its own features

- Other machine learning models directly use the input features to build models.

- But a neural network can learn higher level features that consider the interactions of the input features.

# More layers

# More layers give different levels of abstraction

- We don't know the "right" levels of abstraction

- So let the model figure it out!

- **Face Recognition:**
  - Deep network can build up increasingly higher levels of abstraction
  - Lines, parts, regions

Feature representation

3rd layer
"Objects"

2nd layer
"Object parts"

1st layer
"Edges"

Pixels

Example from
Honglak Lee
(NIPS 2010)

# Multiple classes

- With multiple classes in a classification problem, we will need multiple output units, one output unit per class.



Layer 1    Layer 2    Layer 3    Layer 4
Input layer                      Output layer

# Activation function

- So far, we've assumed that the activation function is always the sigmoid/logistic function. In fact, it is not widely used any more.

"Bias unit"

$x_0$
$x_1$
$x_2$
$\ldots$
$x_n$

$w_0$
$w_1$
$w_2$
$w_n$

$\Sigma$

$\mathbf{z} = \mathbf{w}^\top \mathbf{x}$
$= \sum_{i=1}^{n} w_i x_i$

$a = \sigma(z) = \dfrac{1}{1 + e^{-z}}$

$\sigma(z)$

$z = \mathbf{w}^\top \mathbf{x}$

# Two more activation functions, Tanh and ReLu

## Sigmoid Function



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

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

## Hyperbolic Tangent



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

$$\sigma'(z) = 1 - \sigma(z)^2$$

## Rectified Linear Unit (ReLU)



$$\sigma(z) = \max(0, z)$$

$$\sigma'(z) = \begin{cases} 1, & z > 0 \\ 0, & otherwise \end{cases}$$

# An example: recognizing hand-written digits

- Each hand-written digit is a 28x28 = 784 image

- We want to build a neural network to recognize 10 digits: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9

28

28

Subsequent slides are based on 4 videos starting from the following:
https://www.youtube.com/watch?v=aircAruvnKk

# A network for recognizing of hand-written digits

- **This is the simplest network, called Multilayer perceptron (MLP)**
- **One input layer**
- **Two hidden layers**
  - 2 is an arbitrary choice
  - Each has 16 neurons or units
- **One output layer with 10 units for the 10 digits**
- **All units are fully connected.**



Input layer | 2 hidden layers (16 neurons each) | output layer

# Each neuron is a function, computing an activation value based on all its inputs



- These 784 neurons form the first layer.

- The value held in each output neuron basically tells how likely the input image is each digit.

- Activations of one layer determine the activations of the next layer



$28 \times 28 = 784$

0.58   "Activation"

# Intuitive idea of layers

- The first layer just the gray scale value of each pixel in the image.

- The second layer may capture some low-level features, e.g., edges of different orientations.

- The third layer may capture some high-level features such as loops, strokes, and lines.

- The final layer tells which combination of the subcomponents corresponds to each digit.

# Let us look at a particular neuron



- How does it pick up a small patten?

- For the value of this neuron, we compute

    $$w_1 a_1 + w_2 a_2 + w_3 a_3 + \ldots + w_n a_n + b$$

    - Which may be any value. In this case, we want the values between 0 and 1, we use squash function sigmoid ($\sigma$)



- $\sigma(w_1 a_1 + w_2 a_2 + w_3 a_3 + \ldots + w_n a_n + b)$

Weights

$w_1: 2.07$
$w_2: 2.31$
$w_3: 3.64$
$w_4: 1.87$
$w_5: -1.51$
$w_6: -0.43$
$w_7: 2.01$
$w_8: 1.07$

$784$

# How many parameters?

- Each neuron in one layer is connected with every neuron in the next layer (fully connected).
- We have
  - Number of parameters (or weights): 784x16 + 16x16 + 16x10
  - Number of biases: 16 + 16 + 10
- Total number of parameters:  13,002
  - These all can be tuned and changed.
- Learning: find suitable values for all these parameters to solve the problem at hand, e.g., classifying hand-written digits.

# This network is a function with 13,002 parameters



Sigmoid

$$a_0^{(1)} = \sigma\left( w_{0,0}\, a_0^{(0)} + w_{0,1}\, a_1^{(0)} + \cdots + w_{0,n}\, a_n^{(0)} + b_0 \right)$$

Bias

$$\sigma\left( \begin{bmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix} \right)$$

$$\mathbf{a}^{(1)} = \sigma\left( \mathbf{W}\mathbf{a}^{(0)} + \mathbf{b} \right)$$

# Learning

- ■ **Use a lot of training examples**
  - ❑ Images of handwritten digits with the correct labels (what numbers they correspond to)
- ■ **to adjust those 13,002 weights and biases to improve the performance on training data.**
  - ❑ Hopefully, the resulting network also generalizes to test data.
- ■ **An algorithm is needed:** <span style="color:red">backpropagation</span>

# Training is an optimization problem.

- Trying to find a minima for a cost function $C(x)$

- At the beginning, we just give those weights and biases some random values.

- The cost function basically shows how bad the prediction is.

# We start with a random initialization

- Input 3 gets nonsense results at the output layer.
- Use cost function to measure the difference.

# Square loss (cost) function

- We take the squared difference of what the system gives and what is correct.

Cost of 3

$$0.1863 \leftarrow (0.43 - 0.00)^2 +$$
$$0.0809 \leftarrow (0.28 - 0.00)^2 +$$
$$0.0357 \leftarrow (0.19 - 0.00)^2 +$$
$$0.0138 \leftarrow (0.88 - 1.00)^2 +$$
$$3.37 \quad 0.5242 \leftarrow (0.72 - 0.00)^2 +$$
$$0.0001 \leftarrow (0.01 - 0.00)^2 +$$
$$0.4079 \leftarrow (0.64 - 0.00)^2 +$$
$$0.7388 \leftarrow (0.86 - 0.00)^2 +$$
$$0.9817 \leftarrow (0.99 - 0.00)^2 +$$
$$0.3998 \leftarrow (0.63 - 0.00)^2$$

What's the "cost" of this difference?

Utter trash

# Cost will be small if the classification is correct.

# Cost average over all training data

- The average cost gives an idea how good the network is in classification.

- Training algorithm basically changes all 13002 those weights and biases to get better cost.

  - How to do that?

Here we only show only one training example
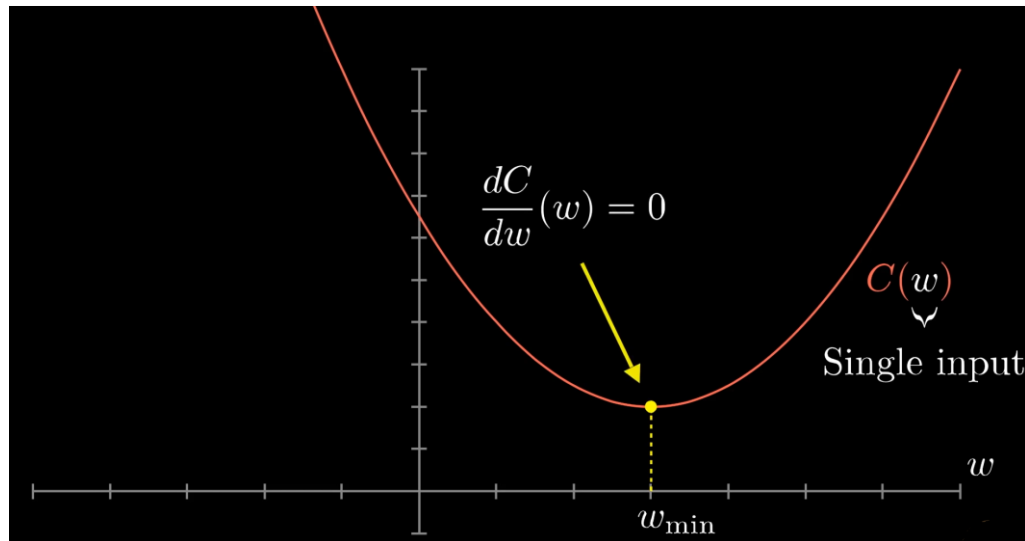


Average cost of all training data...

$$\text{Cost of } \boxed{3} \begin{cases} (0.67 - 0.00)^2 + \\ (0.08 - 0.00)^2 + \\ (0.27 - 0.00)^2 + \\ (0.51 - 1.00)^2 + \\ (0.72 - 0.00)^2 + \\ (0.19 - 0.00)^2 + \\ (0.50 - 0.00)^2 + \\ (0.02 - 0.00)^2 + \\ (0.50 - 0.00)^2 + \\ (0.78 - 0.00)^2 \end{cases}$$

What's the "cost" of this difference?

Utter trash

# How do we optimize? Let us consider only one weight $w$ first

## For a simple function



$$\frac{dC}{dw}(w) = 0$$

$C(w)$

Single input

$w$

$w_{\min}$

## For a complex function



Sometimes infeasible

$$\frac{dC}{dw}(w) = 0$$

$C(w)$

$w$

$w_{\min}$

- **Very difficult for our cost function with 13,002 variables.**
  - We need <span style="color:red">gradient decent</span>.

# How is gradient descent used

- Let us put all the 13,002 weights and biases in a single vector and all the negative gradients of them into another vector.

- We can nudge or change the weights and biases to reduce the cost and to minimize it.

- The algorithm doing this is backpropagation.

$$\vec{\mathbf{W}} = \begin{bmatrix} 2.25 \\ -1.57 \\ 1.98 \\ \vdots \\ -1.16 \\ 3.82 \\ 1.21 \end{bmatrix} \qquad -\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.18 \\ 0.45 \\ -0.51 \\ \vdots \\ 0.40 \\ -0.32 \\ 0.82 \end{bmatrix}$$

13,002 weights and biases

How to nudge all weights and biases

# Meaning of those gradient numbers

- ## We can see
  - ❑ what weight should increase and
  - ❑ what should decrease
  - ❑ what change means a lot

$$\vec{\mathbf{W}} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_{13,000} \\ w_{13,001} \\ w_{13,002} \end{bmatrix}$$

$$-\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -1.25 \\ \vdots \\ 0.78 \\ -0.37 \\ 0.16 \end{bmatrix}$$

$w_0$  should increase somewhat
$w_1$  should increase a little
$w_2$  should decrease a lot

$w_{13,000}$  should increase a lot
$w_{13,001}$  should decrease somewhat
$w_{13,002}$  should increase a little
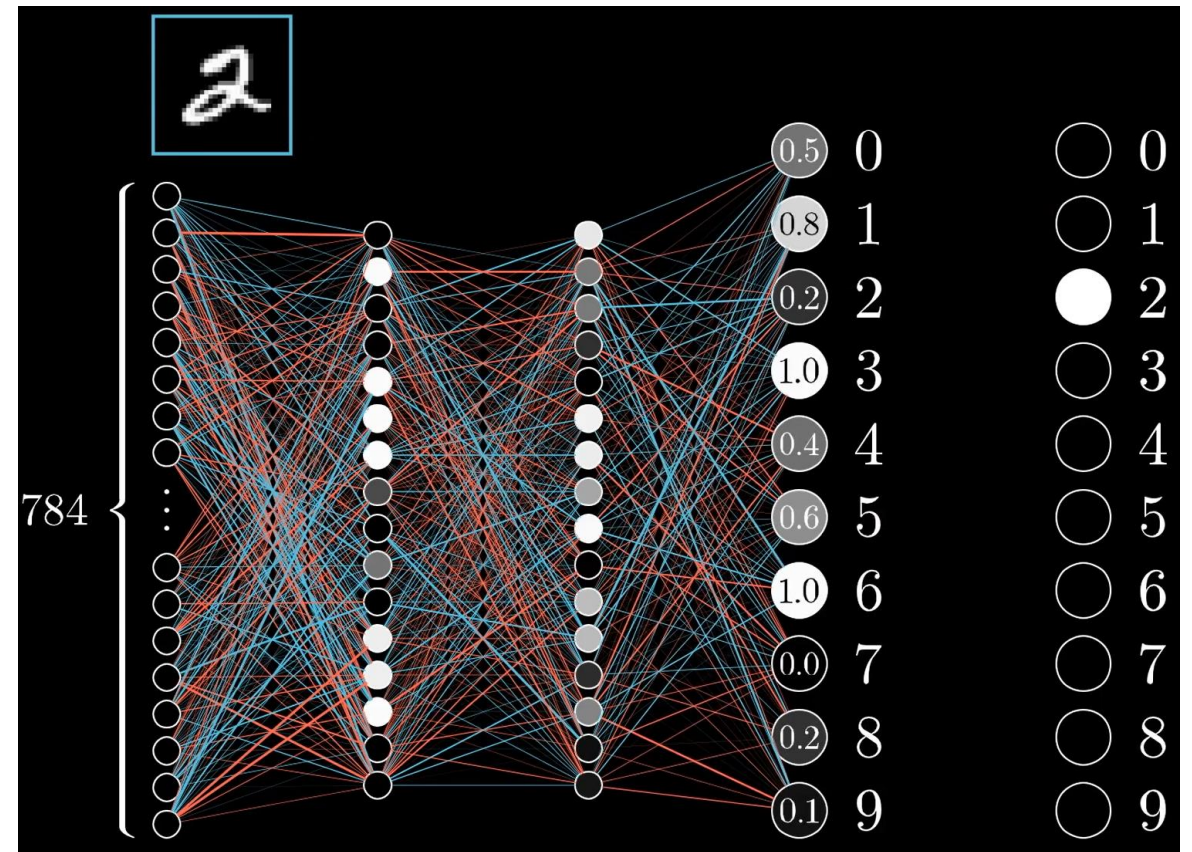
# Backpropagation

- The backpropagation algorithm was originally introduced in the 1970s,

- but its importance wasn't fully appreciated until a <u>famous 1986 paper</u> by <u>David Rumelhart</u>, <u>Geoffrey Hinton</u>, and <u>Ronald Williams</u>.

- That paper describes several neural networks where backpropagation works far faster than earlier approaches,
  - making it possible to use neural nets to solve problems which had previously been insoluble.

- Today, the backpropagation algorithm is the workhorse of learning in neural networks.

# Training: backpropagation algorithm

- **Step 1: initialize the weights and biases.**
  - ❑ Weights in the network are initialized to random numbers from interval [-1,1]
  - ❑ Each unit has a BIAS associated with it
  - ❑ Biases are similarly initialized to random numbers from the interval [-1,1]

- **Step 2: feed the training sample**

- **Step 3: propagate the inputs forward; we compute the net input and output of each unit in the hidden and output layers.**

- **Step 4: back-propagate the error.**

- **Step 5: update weights and biases to reflect the propagated errors.**

- **Step 6: terminating conditions.**

# Intuition of backpropagation

- Since in each step the cost is over all training examples, let us focus on a single example.

- The network isn't well trained, the output activations are pretty random for the input image of 2.

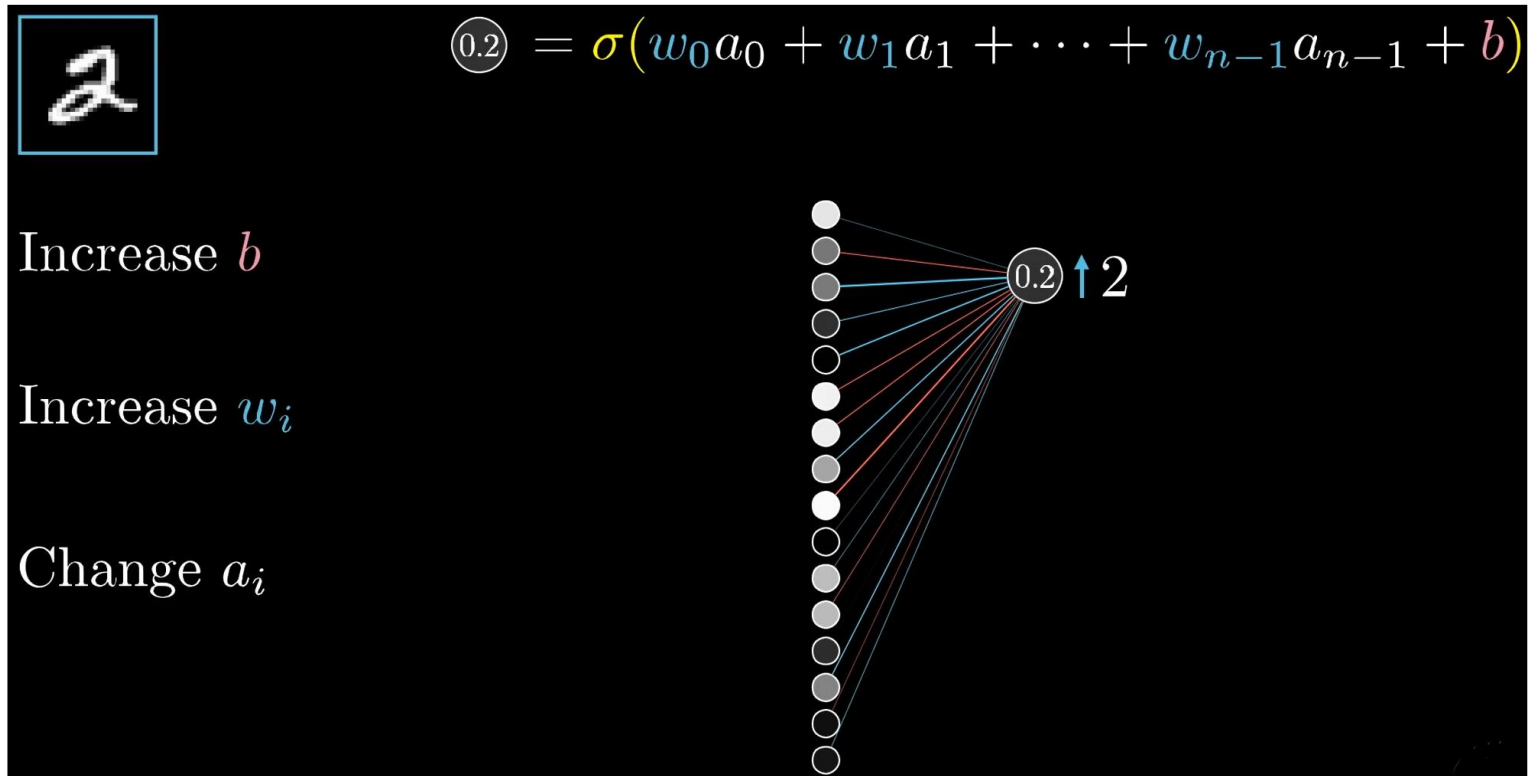- So we need to adjust those weights and biases.

# Intuition of backpropagation (cont.)

- We know which activation should go up and which should go down.
- In this case, the target value for 2 should 1.0 and the others should be 0.0.
- We should nudge activation value for the number '2' up & the rest down.
  - For 7, 8, 9, the values are small.
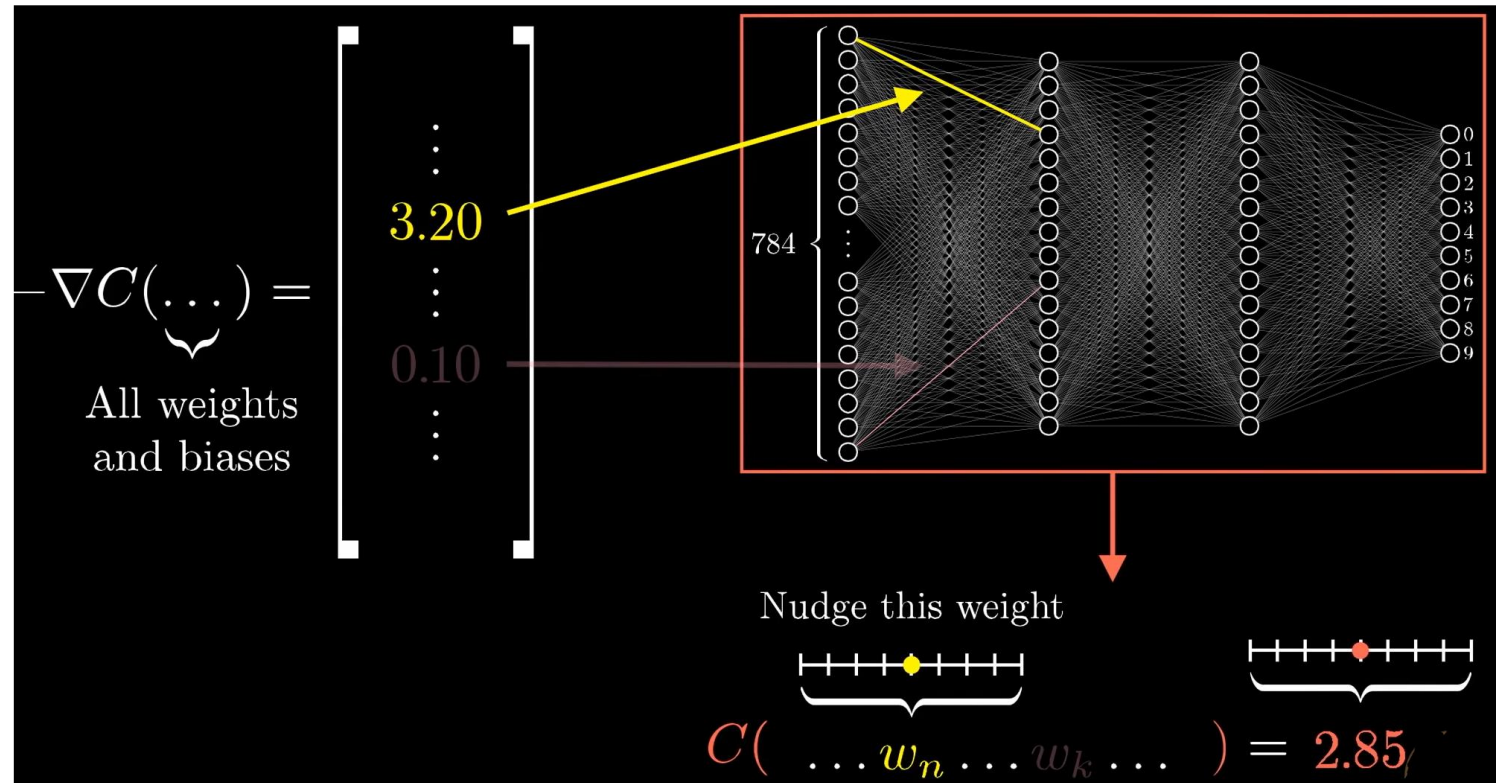  - The size of each nudge should be in proportion to its target value

# Let us look at neuron for 2 only

- **We can nudge weights, the bias and activations.**

  

  $$0.2 = \sigma(w_0 a_0 + w_1 a_1 + \cdots + w_{n-1} a_{n-1} + b)$$

  - Note that we cannot change activations,

  - but only the weights and biases of the previous layers, which affect the activations
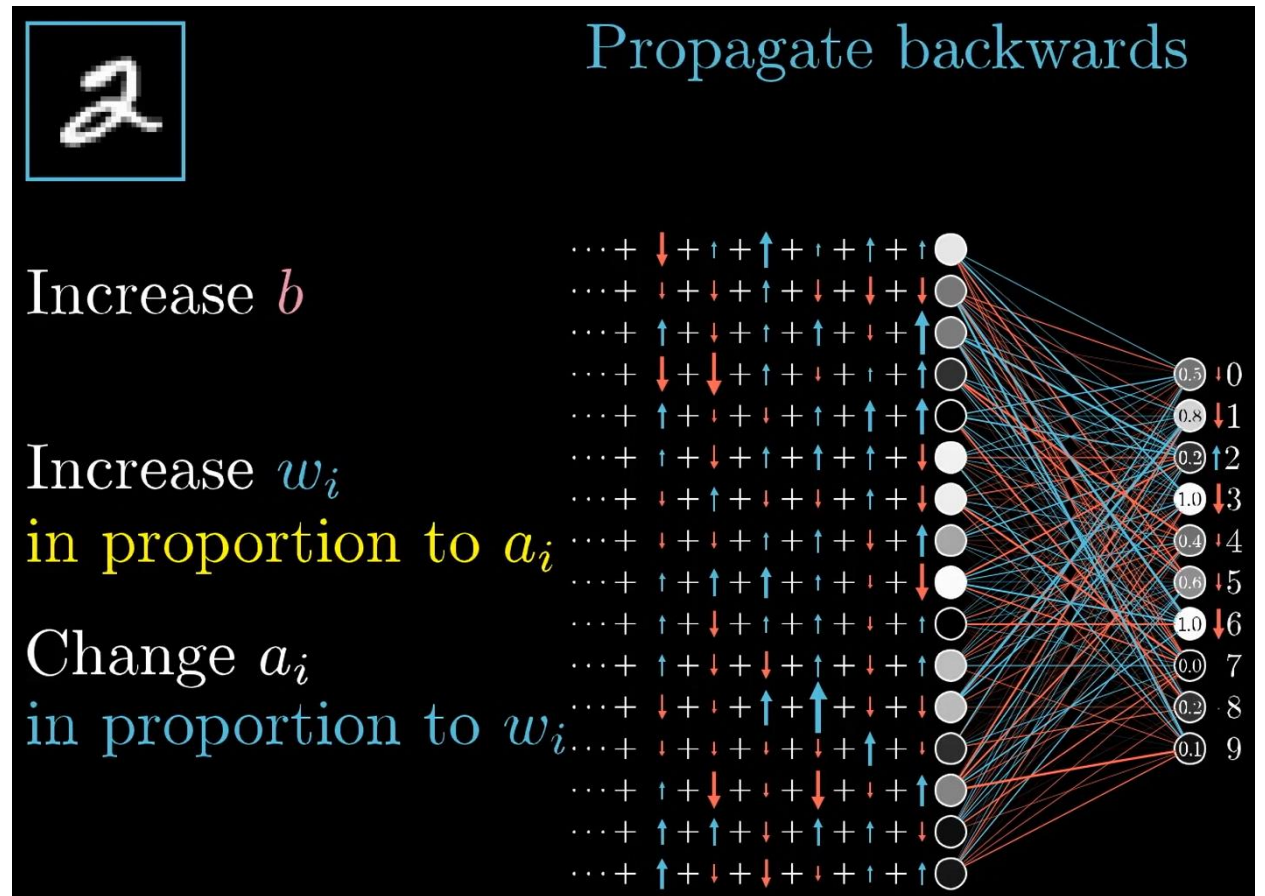
  Increase $b$

  Increase $w_i$

  Change $a_i$

# The effect of gradient

- The gradients tell us which weight or bias should be nudged up or which down,

- but which nudge will give us the best effect "best bang for the buck".

# Considering all output neurons

- We have only considered the output neuron for 2.

- We also need to consider all the output neurons and how they should be nudged and their effect on the second last layer.

# The idea of backpropagation

- Finally, we sum up all the effects to get what should happen to the second to the last layer.

- Then we can recursively apply the same process to the previous layer and so on.

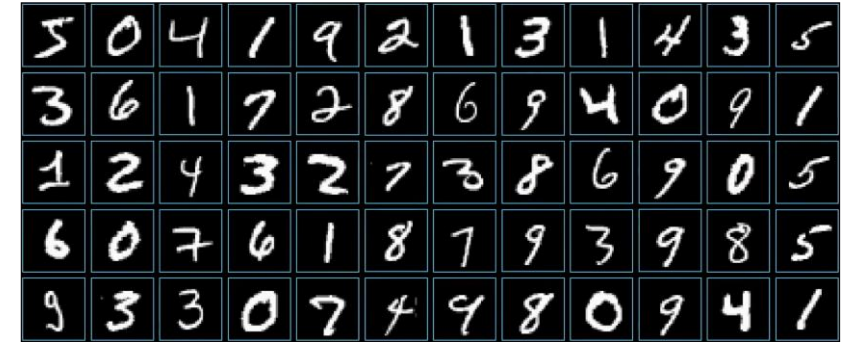  - So that their weights and biases can be adjusted.

# Considering all training examples

- **So far, we have only looked at one training example of 2.**
  - ❑ We can get how much change should be applied to each weight and bias.
  - ❑ But we need to average over all training data to get their desired changes

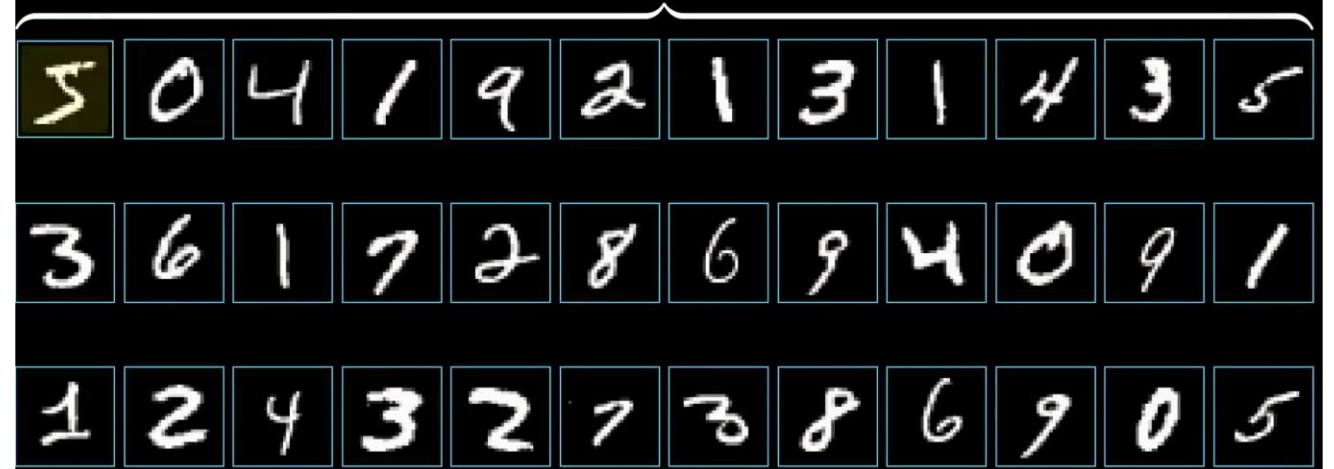| | ![2] | ![5] | ![0] | ![4] | ![1] | ![9] | | Average over all training data |
|---|---|---|---|---|---|---|---|---|
| $w_0$ | $-0.08$ | $+0.02$ | $-0.02$ | $+0.11$ | $-0.05$ | $-0.14$ | $\cdots$ | $\rightarrow$ $-0.08$ |
| $w_1$ | $-0.11$ | $+0.11$ | $+0.07$ | $+0.02$ | $+0.09$ | $+0.05$ | $\cdots$ | $\rightarrow$ $+0.12$ |
| $w_2$ | $-0.07$ | $-0.04$ | $-0.01$ | $+0.02$ | $+0.13$ | $-0.15$ | $\cdots$ | $\rightarrow$ $-0.06$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| $w_{13,001}$ | $+0.13$ | $+0.08$ | $-0.06$ | $-0.09$ | $-0.02$ | $+0.04$ | $\cdots$ | $\rightarrow$ $+0.04$ |

# Stochastic gradient descent



- It takes too long to go though all the training data and all those computations to calculate each nudge/change.

- In practice, we use stochastic gradient descent.

  - We shuffle the data & divide them in minibatches and
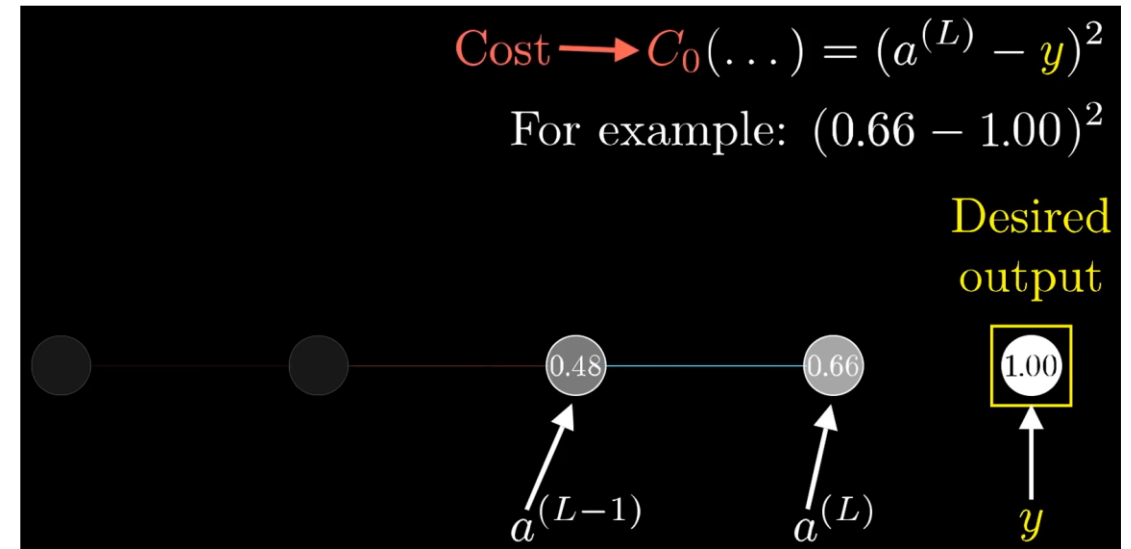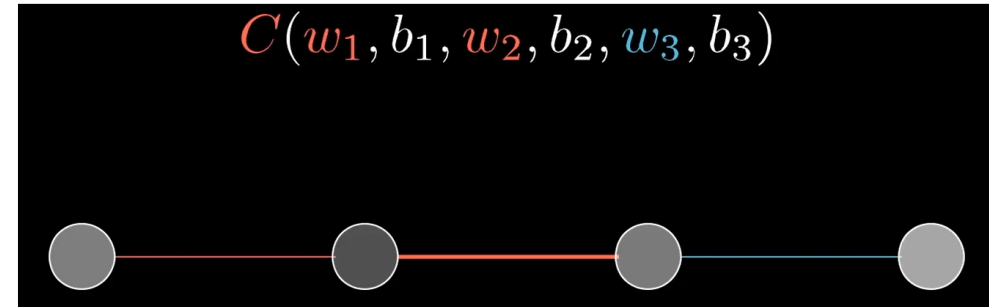  - work on each minibatch in each step.
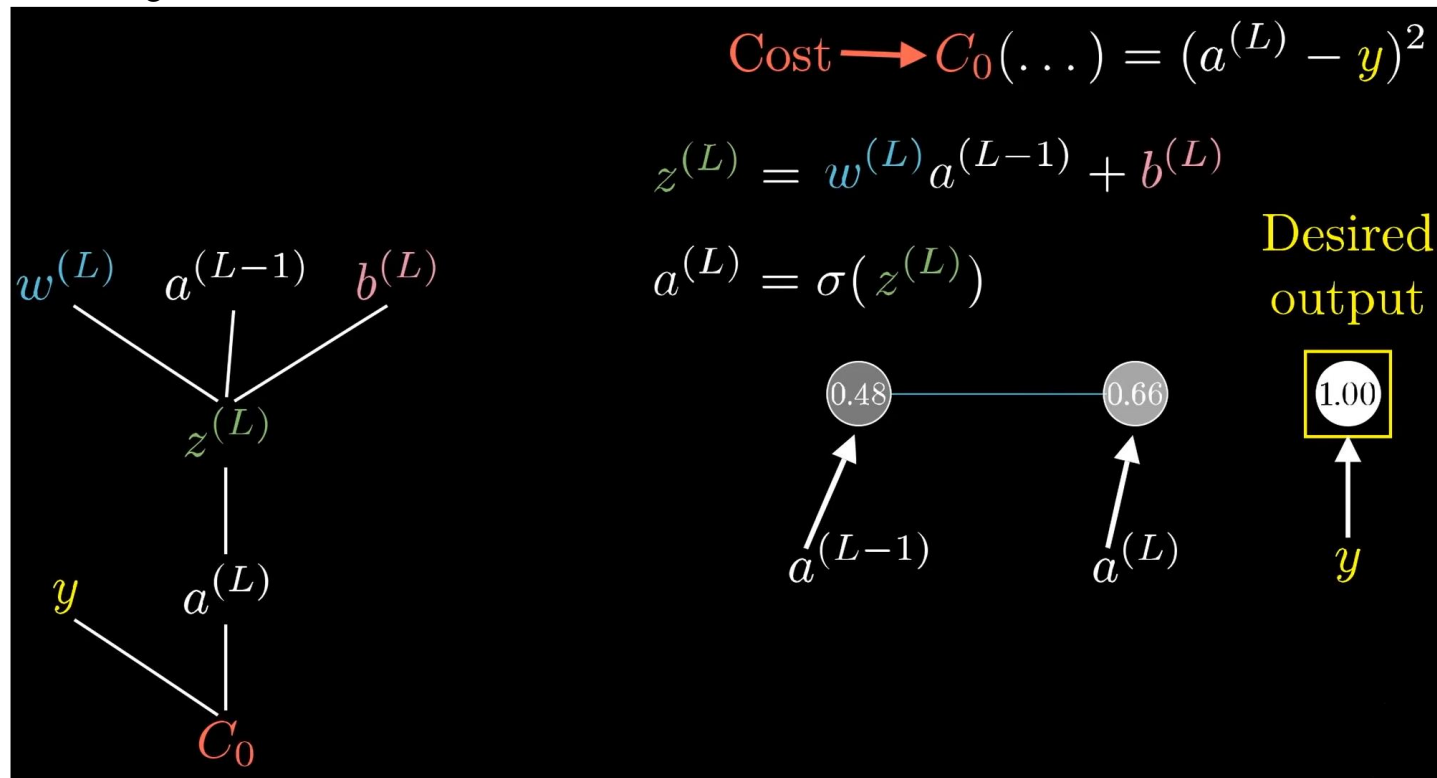
Computing based on minibatches

# Math of backpropagation

- We start with a very simple case:
  - one neuron in each layer
- Further, we will focus on the last two layers.
  - For a training example with class y, the last neuron is for the class (i.e., 1.00)
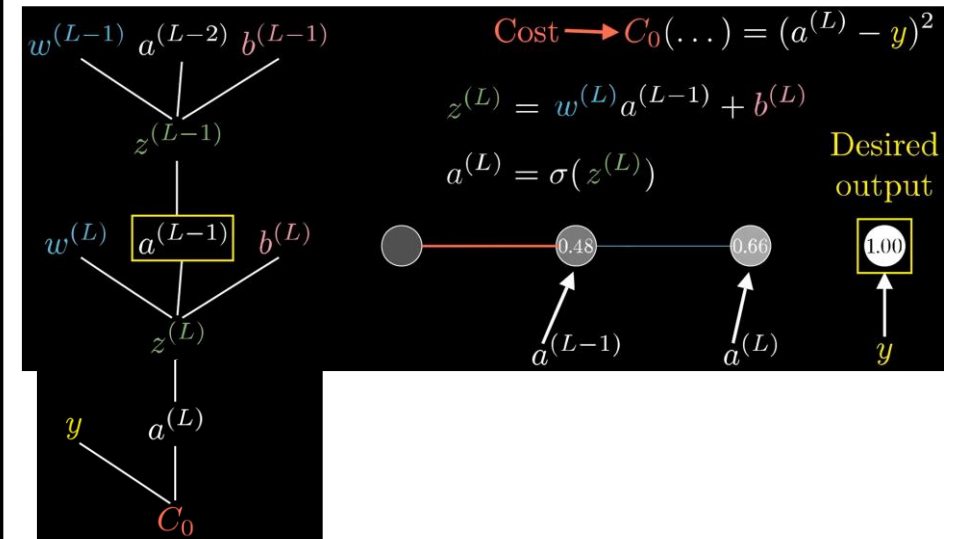- We work on one training example first.

$$C(w_1, b_1, w_2, b_2, w_3, b_3)$$



$$\text{Cost} \longrightarrow C_0(\dots) = (a^{(L)} - y)^2$$

For example: $(0.66 - 1.00)^2$

Desired output

$$a^{(L-1)} \qquad a^{(L)} \qquad y$$

# Model the two layers

- Let us see the flow structure for 2 layers. $C_0$ is the cost of one training example

- Note that we can go to the next level too, but we will not focus on that



$$\text{Cost} \longrightarrow C_0(\ldots) = (a^{(L)} - y)^2$$

$$z^{(L)} = w^{(L)}a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$

# How sensitive cost is to a small change in weight?

- Each term is just a numerical value with a number line.
- To get the sensitivity, we take partial derivatives
  - Chain rule

$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C0}{\partial a^{(L)}}$$

**Chain rule**

$$C_0(\ldots) = (a^{(L)} - y)^2$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$

Desired output

# Compute all derivatives

$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C0}{\partial a^{(L)}}$$

$$\frac{\partial C0}{\partial a^{(L)}} = 2(a^{(L)} - y)$$

$$\frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)})$$

$$\frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)}$$

$$C_0 = (a^{(L)} - y)^2$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$



0.48 —— 0.66    1.00

$a^{(L-1)}$     $a^{(L)}$     $y$

# Consider all training examples

- We have only considered one example and its cost $C_0$.

- To consider all training examples, we average the gradients

$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}} = a^{(L-1)} \sigma'(z^{(L)}) 2(a^{(L)} - y)$$
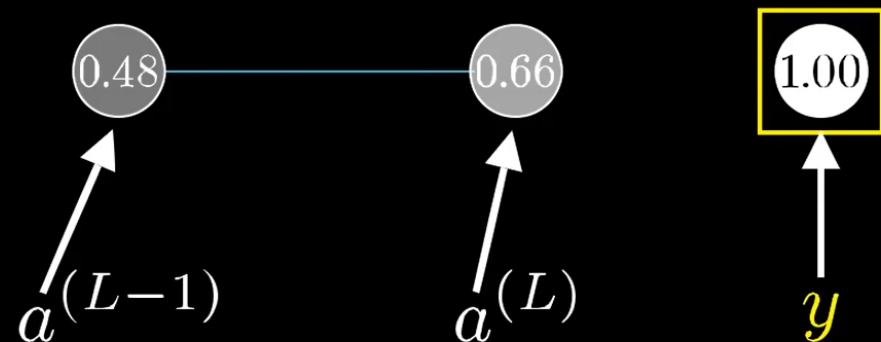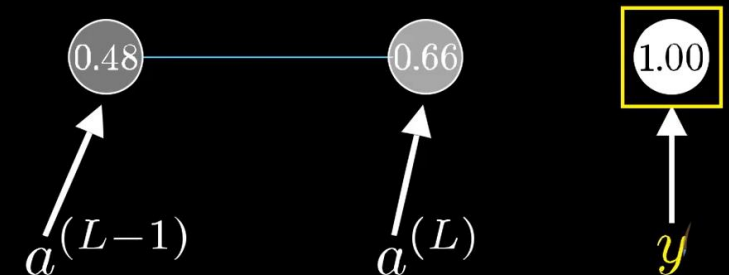
Average of all training examples

$$\frac{\partial C}{\partial w^{(L)}} = \frac{1}{n} \overbrace{\sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^{(L)}}}$$

Derivative of full cost function

$$C_0 = (a^{(L)} - y)^2$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$

0.48 —— 0.66    1.00

$a^{(L-1)}$    $a^{(L)}$    $y$

# Take derivative of the bias



$$\frac{\partial C_0}{\partial b^{(L)}} = \frac{\partial z^{(L)}}{\partial b^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C0}{\partial a^{(L)}} = 1\sigma'(z^{(L)})2(a^{(L)} - y)$$

$$C_0 = (a^{(L)} - y)^2$$

$$z^{(L)} = w^{(L)}a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$

$w^{(L)}$  $a^{(L-1)}$  $b^{(L)}$

$z^{(L)}$

$y$  $a^{(L)}$

$C_0$

0.48 — 0.66    1.00
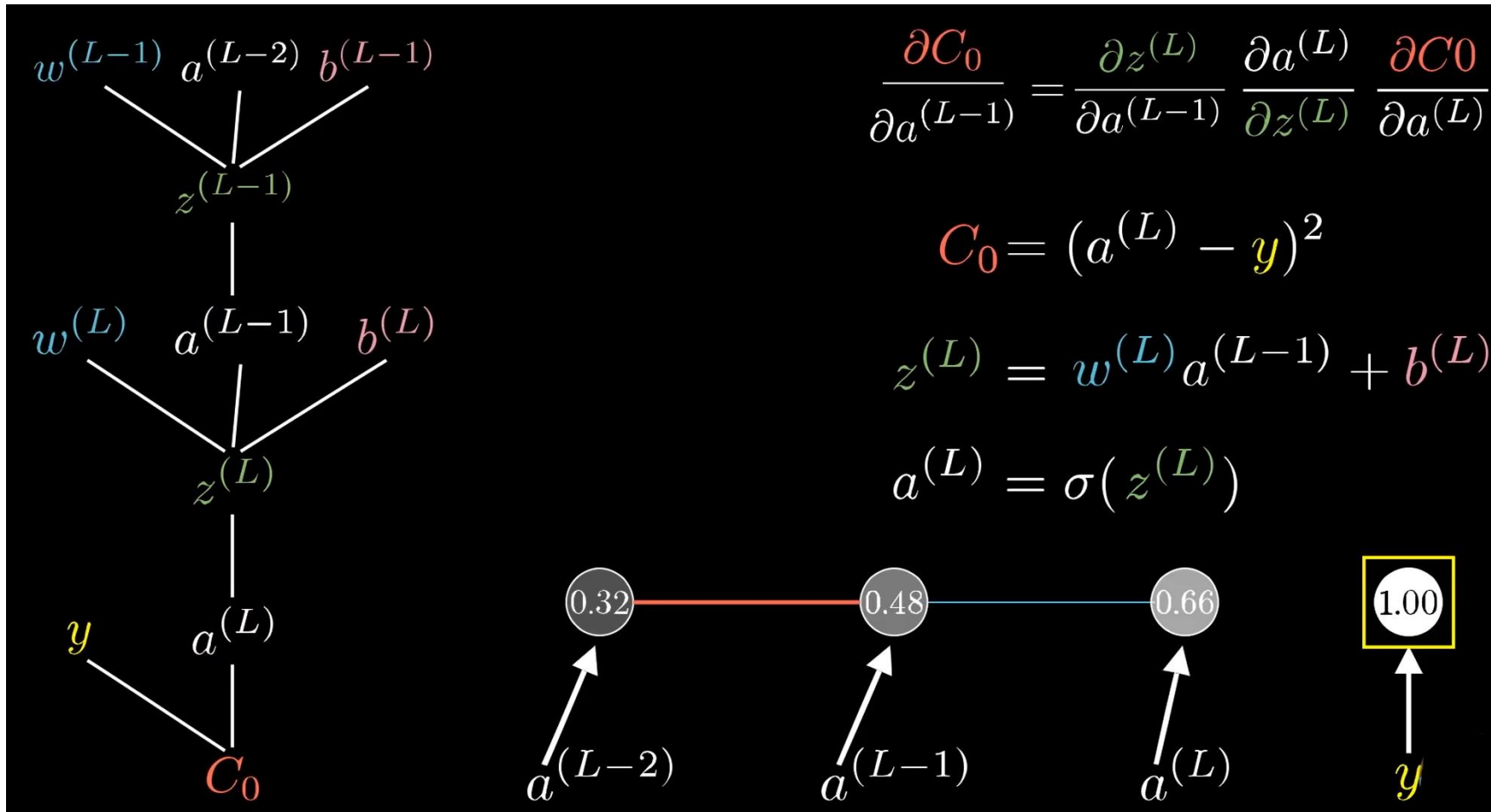
$a^{(L-1)}$    $a^{(L)}$    $y$

# Take derivative of the activation (propagate back)

$$\frac{\partial C_0}{\partial a^{(L-1)}} = \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C0}{\partial a^{(L)}} = \quad w^{(L)} \sigma'(z^{(L)}) 2(a^{(L)} - y)$$

$$C_0 = (a^{(L)} - y)^2$$

$$w^{(L)} \quad a^{(L-1)} \quad b^{(L)}$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

$$z^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$

$$y \quad a^{(L)}$$

| 0.48 | — | 0.66 | | 1.00 |

$$C_0$$

$$a^{(L-1)} \quad a^{(L)} \quad y$$

# Iterating the same chain rule idea backward to the previous layer and so on

$$\frac{\partial C_0}{\partial a^{(L-1)}} = \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C0}{\partial a^{(L)}}$$

$$C_0 = (a^{(L)} - y)^2$$

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$

# General case: more neurons in each layer

- **Need more indices and everything else is basically the same.**

$$z_j^{(L)} = w_{j0}^{(L)} a_0^{(L-1)} + w_{j1}^{(L)} a_1^{(L-1)} + w_{j2}^{(L)} a_2^{(L-1)} + b_j^{(L)}$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$

$$C_0 = \sum_{j=0}^{n_L - 1} (a_j^{(L)} - y_j)^2$$

(0.10)

$$w_{jk}^{(L)}$$

(0.83) $\leftarrow a_j^{(L)}$

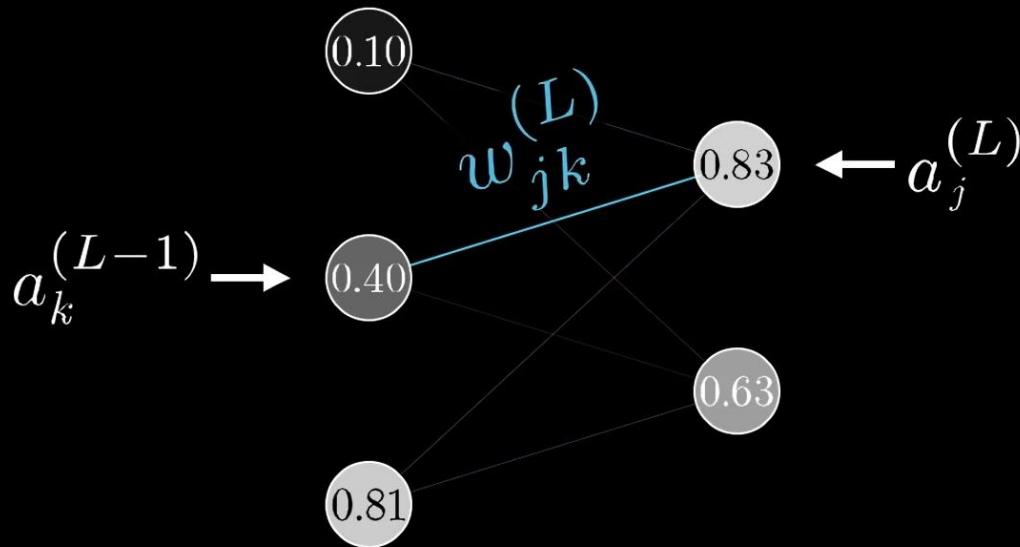$a_k^{(L-1)} \rightarrow$ (0.40)

(0.63)

(0.81)

# Derivatives on weights and biases are the same

$$\frac{\partial C_0}{\partial w_{jk}^{(L)}} = \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}}$$

$$z_j^{(L)} = \cdots + w_{jk}^{(L)} a_k^{(L-1)} + \cdots$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$

$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

$w_{jk}^{(L)}$

$a_j^{(L)}$

$a_k^{(L-1)}$

0.10

0.83

0.40

0.63

0.81

# Derivative on the activation **changes**

- Since the neuron $(a_k^{(L-1)})$ influences the cost function through multiple different paths (2 in this case).

$$\frac{\partial C_0}{\partial a_k^{(L-1)}} = \underbrace{\sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}}}_{\text{Sum over layer L}}$$

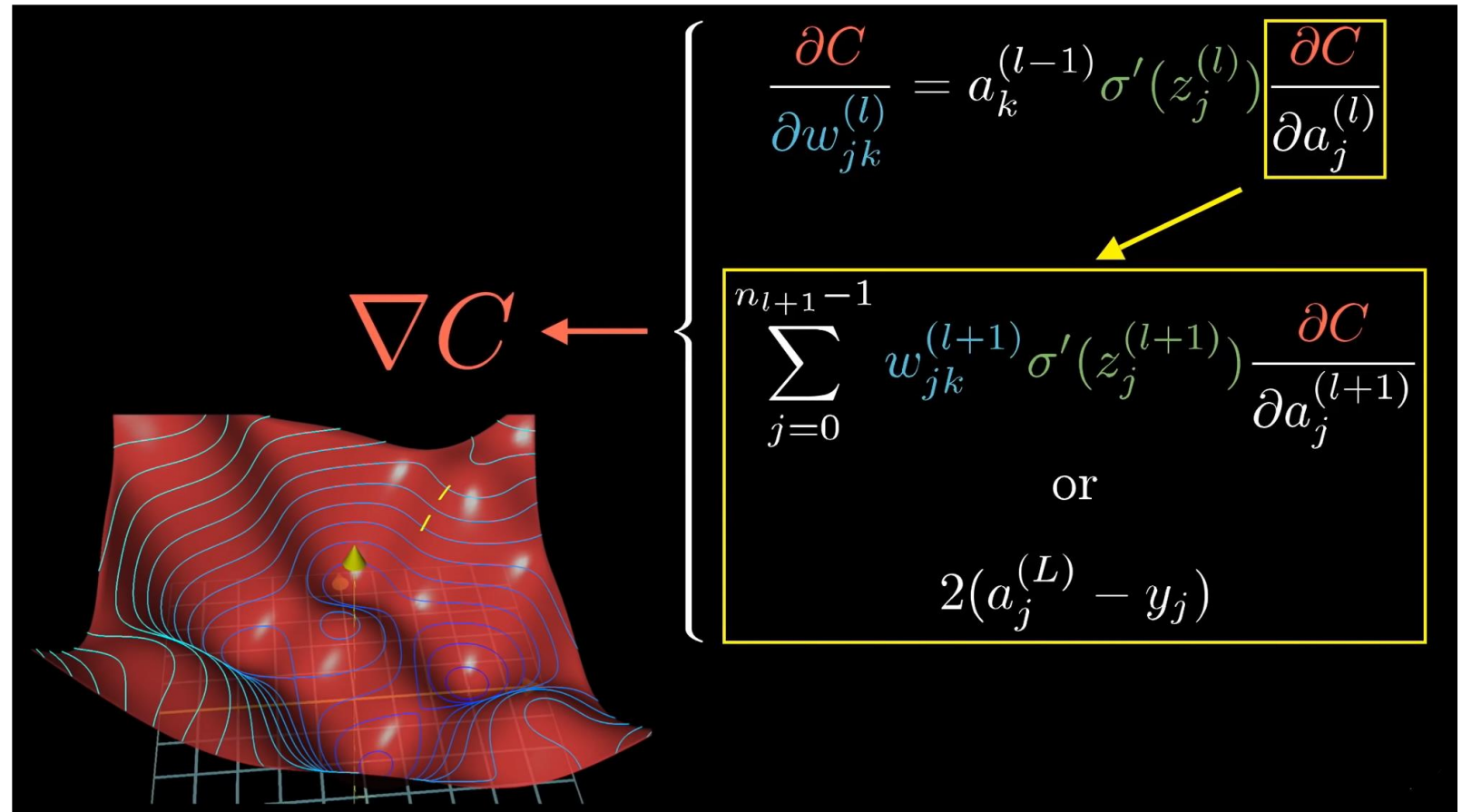$$z_j^{(L)} = \cdots + w_{jk}^{(L)} a_k^{(L-1)} + \cdots$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$

$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

# With all the gradients, we apply gradient descent

- The expression "or" means that at the last layer (which is different from other layers), we take the derivative on the cost.

- Note the typo:

  *l -> L*

$$\nabla C \longleftarrow \begin{cases} \dfrac{\partial C}{\partial w_{jk}^{(l)}} = a_k^{(l-1)} \sigma'(z_j^{(l)}) \boxed{\dfrac{\partial C}{\partial a_j^{(l)}}} \\ \\ \boxed{\begin{aligned} &\sum_{j=0}^{n_{l+1}-1} w_{jk}^{(l+1)} \sigma'(z_j^{(l+1)}) \dfrac{\partial C}{\partial a_j^{(l+1)}} \\ &\qquad\qquad \text{or} \\ &\qquad 2(a_j^{(L)} - y_j) \end{aligned}} \end{cases}$$

# Watch these YouTube videos about neural network and backpropagation

- https://www.youtube.com/watch?v=aircAruvnKk
  - There are 4 videos introducing neural networks and backpropagation. Most of our slides are based on these videos.
- https://www.youtube.com/watch?v=IN2XmBhILt4https://www.youtube.com/watch?v=iyn2zdALii8
- https://www.youtube.com/watch?v=GKZoOHXGcLo
- A playlist:
  - https://www.youtube.com/watch?v=CqOfi41LfDw&list=PLblh5JKOoLUIxGDQs4LFFD--41Vzf-ME1