

Information Engineering 2

Machine Learning Q&A

Prof. Dr. Kurt Stockinger

Semesterplan

SW	Datum	Vorlesungsthema	Praktikum
1	23.02.2022	Data Warehousing Einführung	Praktikum 1: KNIME Tutorial
2	02.03.2022	Dimensionale Datenmodellierung 1	Praktikum 1: KNIME Tutorial (Vertiefung)
3	09.03.2022	Dimensionale Datenmodellierung 2	Praktikum 2: Datenmodellierung
4	16.03.2022	Datenqualität und Data Matching	Praktikum 3: Star-Schema, Bonus: Praktikum 4: Slowly Changing Dimensions
5	23.03.2022	Big Data Einführung	DWH Projekt - Teil 1
6	30.03.2022	Spark - Data Frames	DWH Projekt - Teil 2 (Abgabe: 4.4.2022 23:59:59)
7	06.04.2022	Data Storage: Hadoop Distributed File System & Parquet	Praktikum 1: Data Frames
8	13.04.2022	Query Optimization	Praktikum 2: Data Storage
9	20.04.2022	Spark Best Practices & Applications	Praktikum 3: Query Optimization & Performance Analysis
10	27.04.2022	Machine Learning mit Spark 1	Praktikum 3: Query Optimization & Performance Analysis (Vertiefung)
11	04.05.2022	Machine Learning mit Spark 2 + Q&A	Praktikum 4: Machine Learning (Regression)
12	11.05.2022	NoSQL Systems	Big Data Projekt - Teil 1
13	18.05.2022	Keine Vorlesung (Arbeit am Projekt)	Big Data Projekt - Teil 2
14	25.05.2022	Keine Vorlesung (Arbeit am Projekt)	Big Data Projekt - Teil 3 (Abgabe: 30.5.2022 23:59:59)

Online Courses on Introduction to Neural Networks

- Basics of neural networks (20 min):
<https://www.youtube.com/watch?v=aircAruvnKk>
- How neural networks learn (21 min):
<https://www.youtube.com/watch?v=IHZwWFHWa-w>

Transformers (15 min):

<https://youtu.be/4Bdc55j80I8>

Overview

- Simple introduction into **perceptrons and neural networks**
- Highlights of the two videos:
 - **Basics** of neural networks
 - How neural networks **learn**

Input based on Yale Lecture

Andrej Karpathy's Hacker's guide to Neural Networks:
<http://karpathy.github.io/neuralsnets/>

Andrej Karpathy's lecture notes:
<http://cs231n.github.io/>

Geoffrey E. Hinton, Yann LeCun, and Yoshua Bengio (video; NIPS 2015):
<http://research.microsoft.com/apps/video/default.aspx?id=259574>

Michael Nielsen's Neural Networks and Deep Learning:
<http://neuralnetworksanddeeplearning.com/>

Simple Decision

- Assume we want to go to a chocolate festival.
- Three variables go into our decision:
 - Is the weather good?
 - Does your friend want to go with you?
 - Is it near public transportation?

Simple Decision #2

- Assume that answers to these questions are **binary variables** x_i
 - 0: no
 - 1: yes
- Questions and **answers**:
 - Is the weather good? x_1
 - Does your friend want to go with you? x_2
 - Is it near public transportation? X_3
- How do we describe a **decision statement**?

Simple Decision #3

- Use **weights w_i** for importance of features:

$$x_1 \cdot w_1 + x_2 \cdot w_2 + x_3 \cdot w_3 \geq \text{threshold}$$

- Assumption:
 - If statement is true, we attend festival (otherwise not)
- Example:
 - If we hated weather and were more neural about friends and public transport, we would pick the **weights 6, 2, 2**
 - Assume **threshold 5**: we would only go, if the weather was great
 - Threshold 3?

Simple Decision #4

- We can define binary variable y
- Defines if we go to festival or not

$$y = \begin{cases} 0, & x_1 \cdot w_1 + x_2 \cdot w_2 + x_3 \cdot w_3 < \text{threshold} \\ 1, & x_1 \cdot w_1 + x_2 \cdot w_2 + x_3 \cdot w_3 \geq \text{threshold} \end{cases}$$

Simple Decision #5

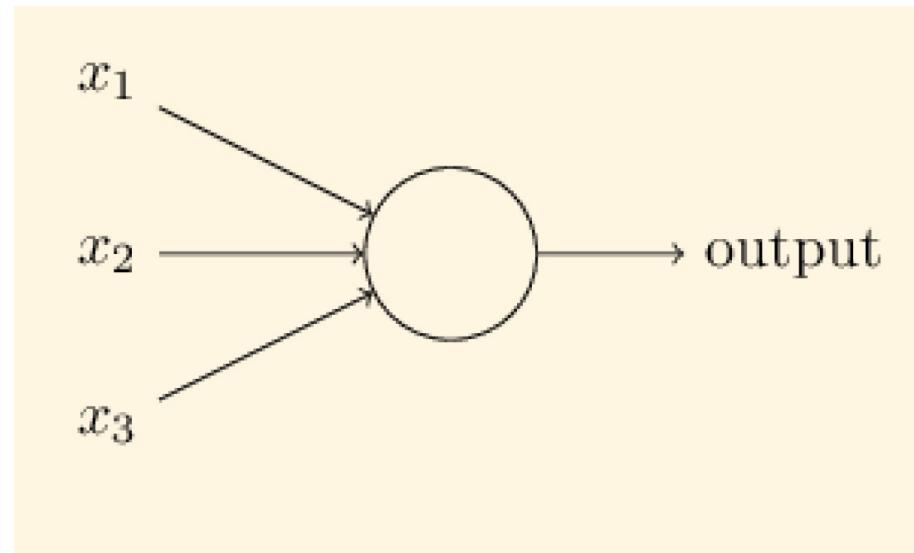
- We can rewrite this as **dot product** of vector of **binary inputs x** and vector of **weights w**
- Reformulate threshold as negative **bias b**

$$y = \begin{cases} 0, & x \cdot w + b < 0 \\ 1, & x \cdot w + b \geq 0 \end{cases}$$

- Goal: Find **separate hyperplane** (like with logistic regression)

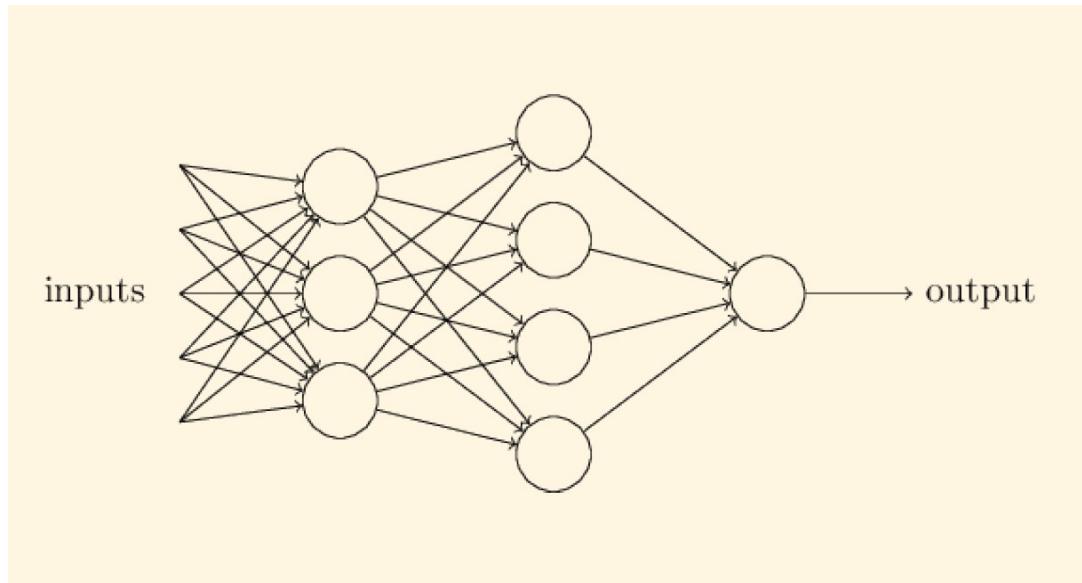
A Perceptron

- Graphical representation of algorithm as a **perceptron**
 - 3 binary inputs
 - 1 binary output
 - Use **weighting scheme** described previously

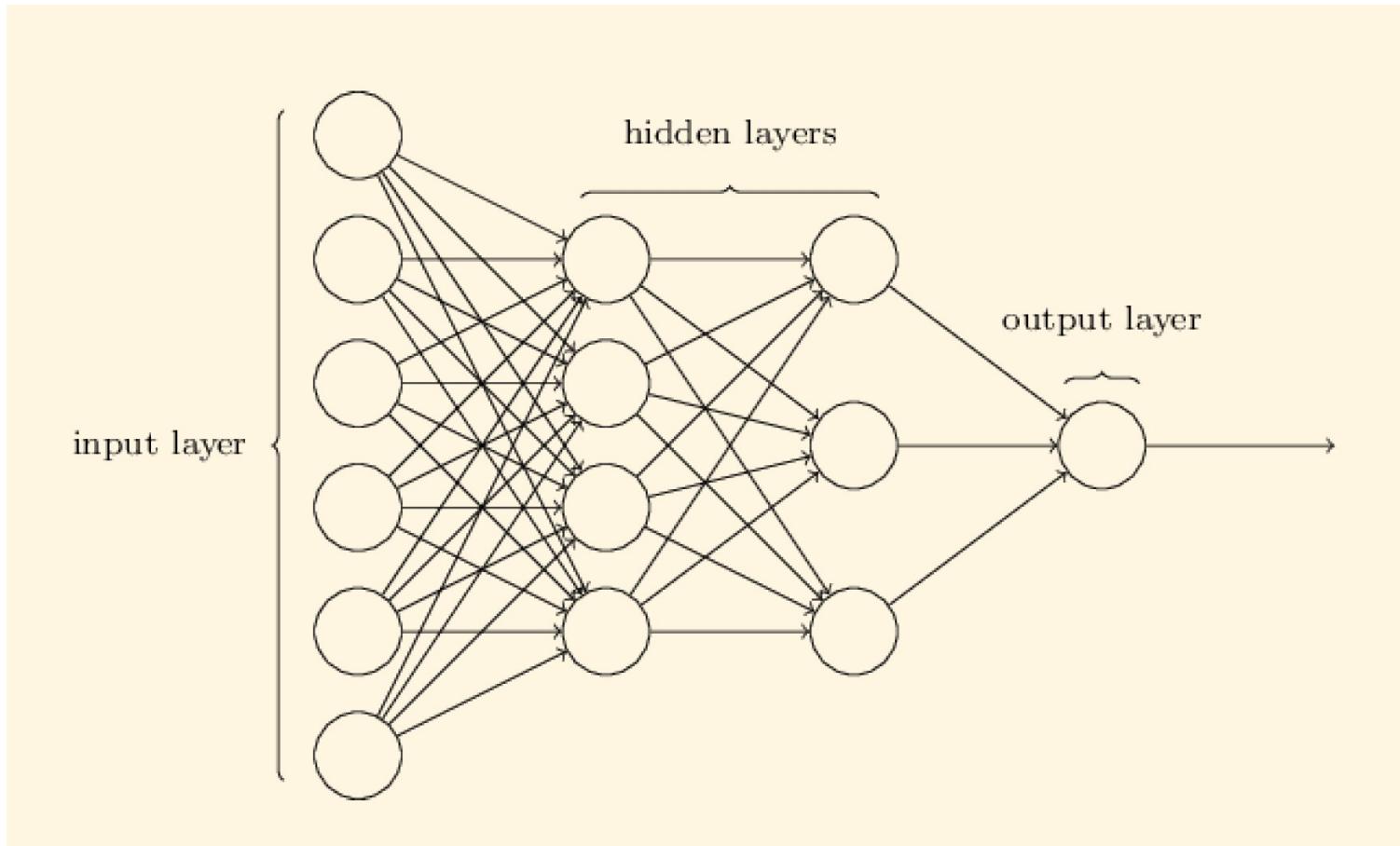


Neural Network = Network of Perceptrons

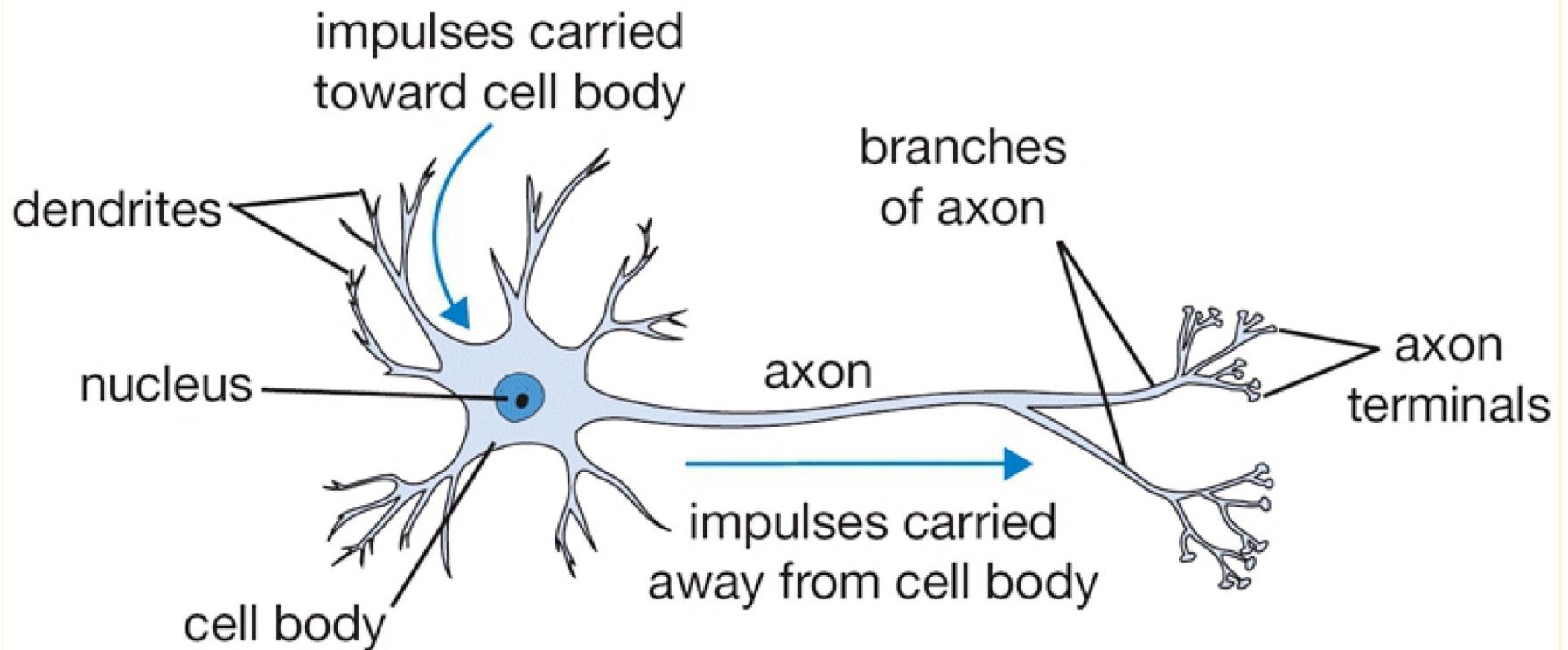
- Output of one perceptron is used as input for other perceptrons



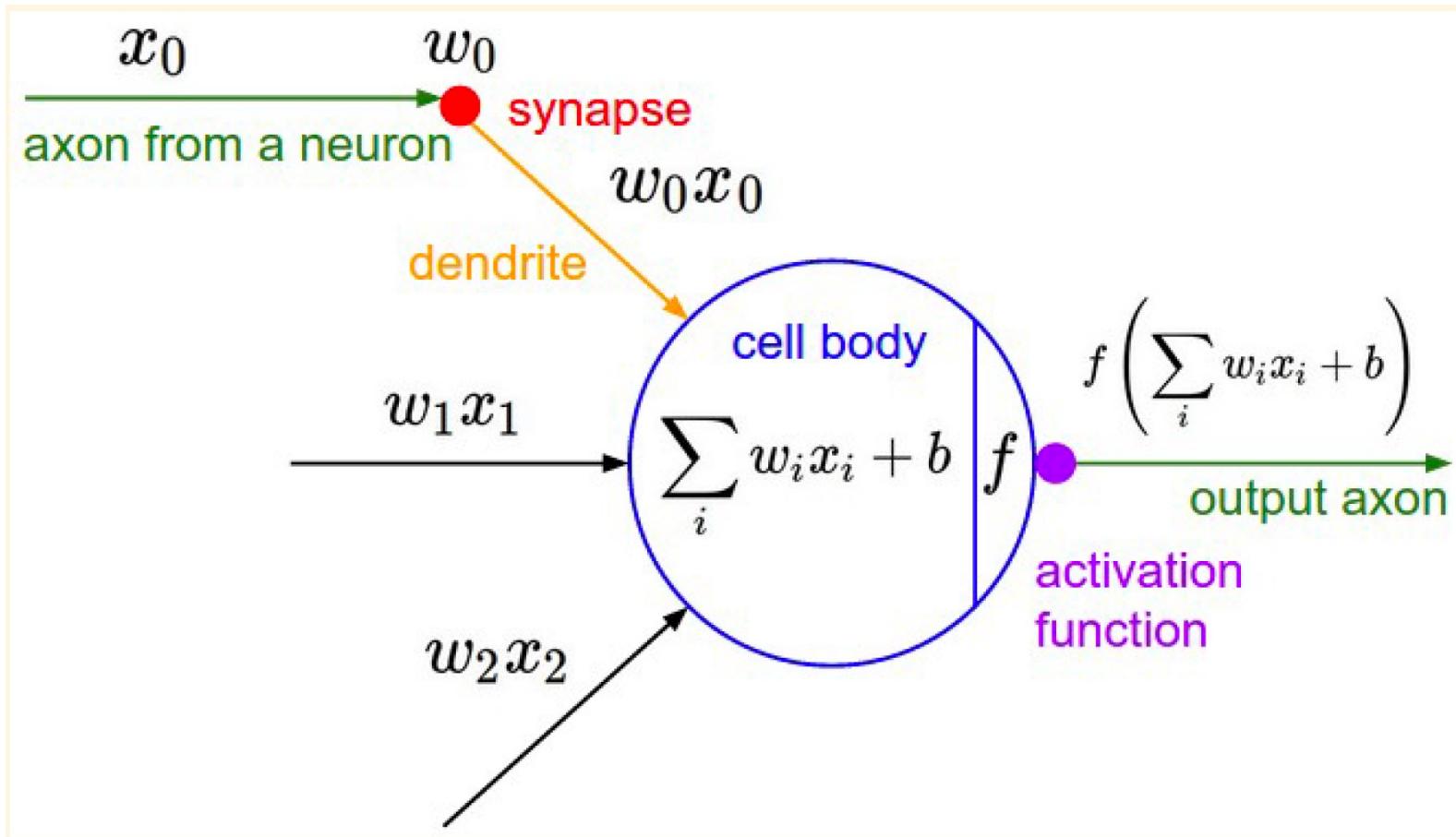
Neural Network with Hidden Layers



Neural Network in Nature



Activation Function for Neurons to Fire

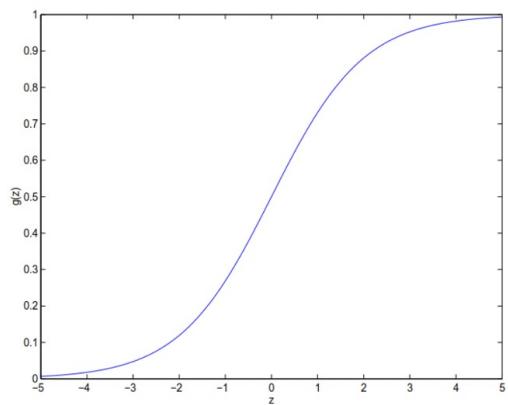


Disadvantage of Simple Approach

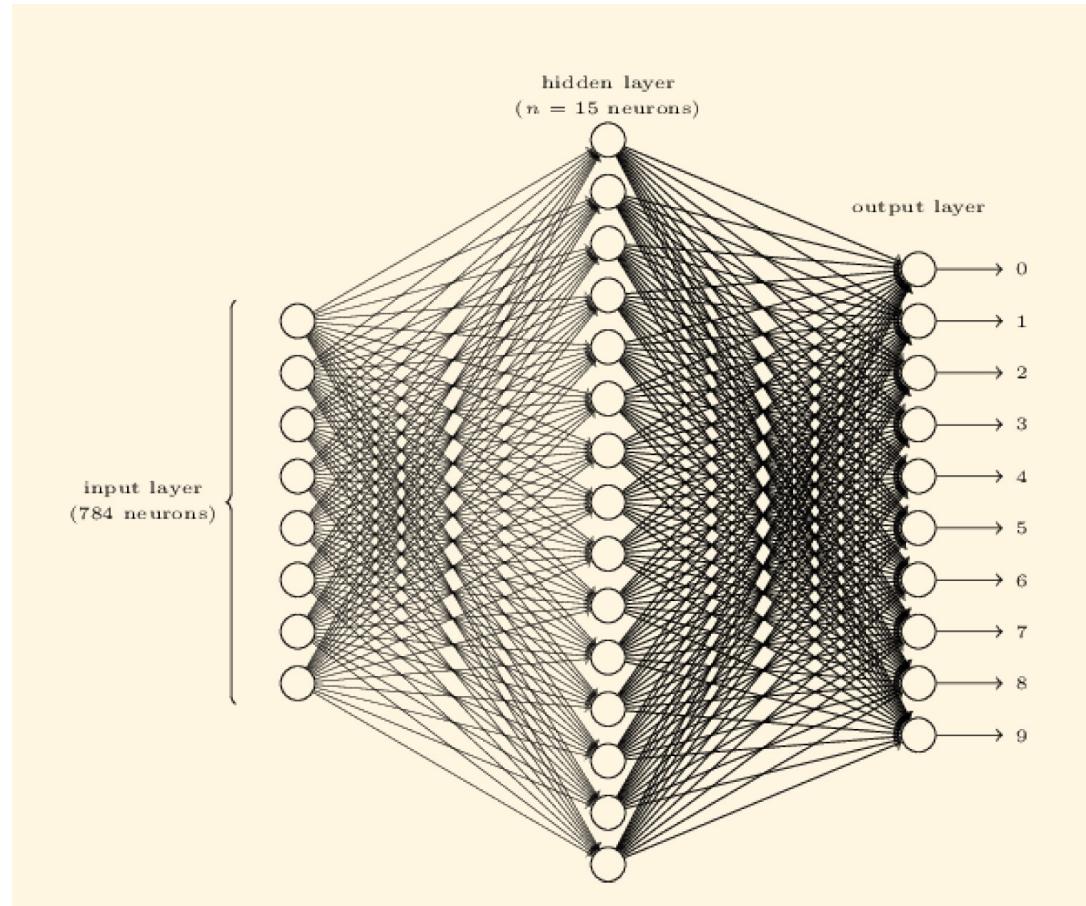
- Each node (or neuron) has only two possible states: 0 or 1
→ Small change in the input values can cause a large change of the output
- Better solution:
 - Output states should be continuous values, e.g. between 0 and 1
- How we do get such an output?

Use Sigmoid Function as Activation Function

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



More Complex Neural Network for Detecting Numbers



Classification Problem

- Goal:
 - Given an input (number), in which class does it fall?
 - E.g. numbers between 0 and 9
- Approach:
 - Find output **neurons with largest output**
 - We need to **learn weights and bias terms**
- Note:
 - If we only have a network with **one node**, we can directly use **logistic regression**

Learning Approach

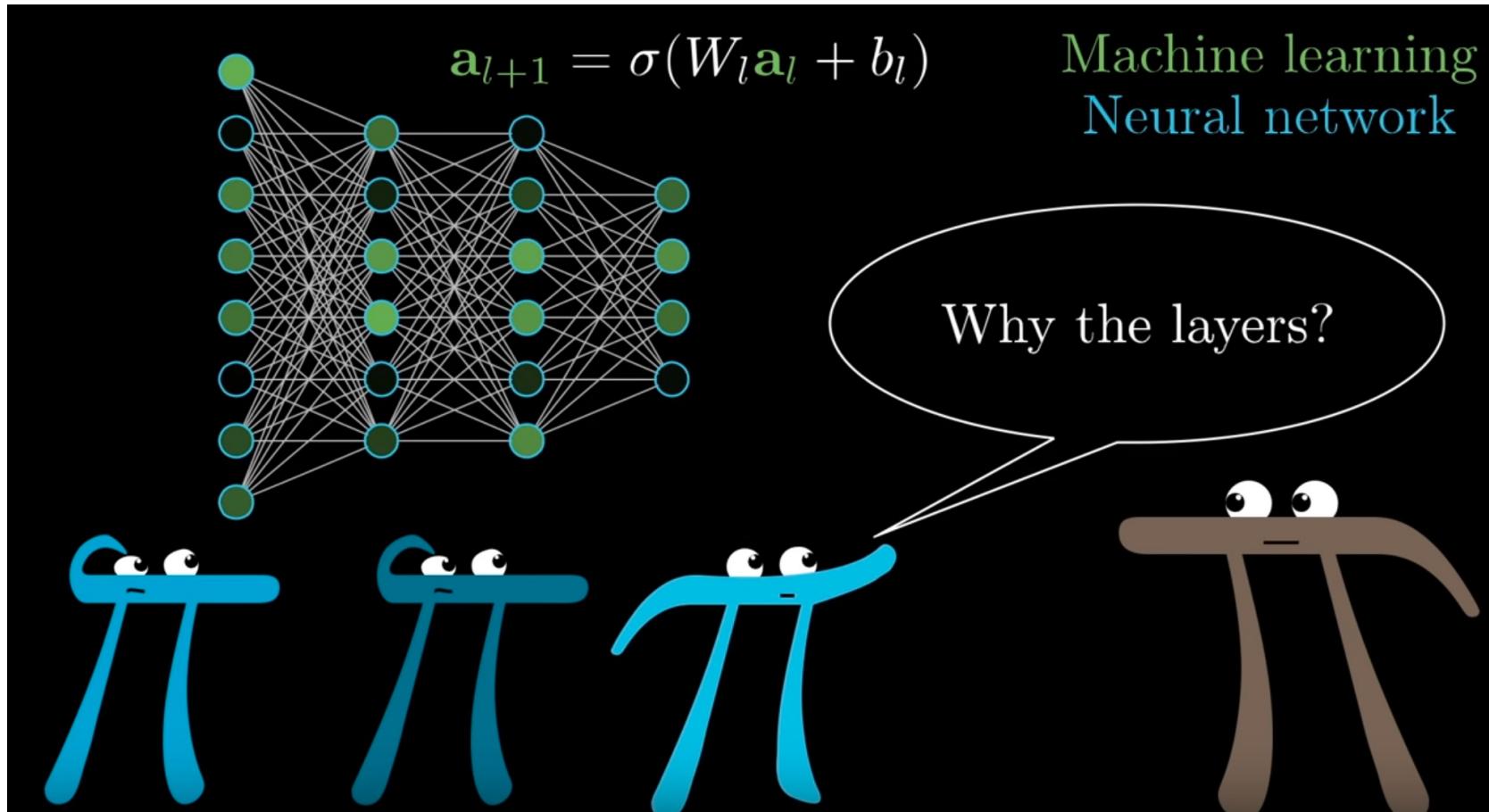
- We need a **cost function (loss function)** that measures how well the network predicts outputs on the test set
- Goal:
 - Find a **set of weights and biases that minimizes the cost**
- E.g. Use **squared error loss**:

$$C(w, b) = \frac{1}{2n} \sum_i (y_i - \hat{y}(x_i))^2$$

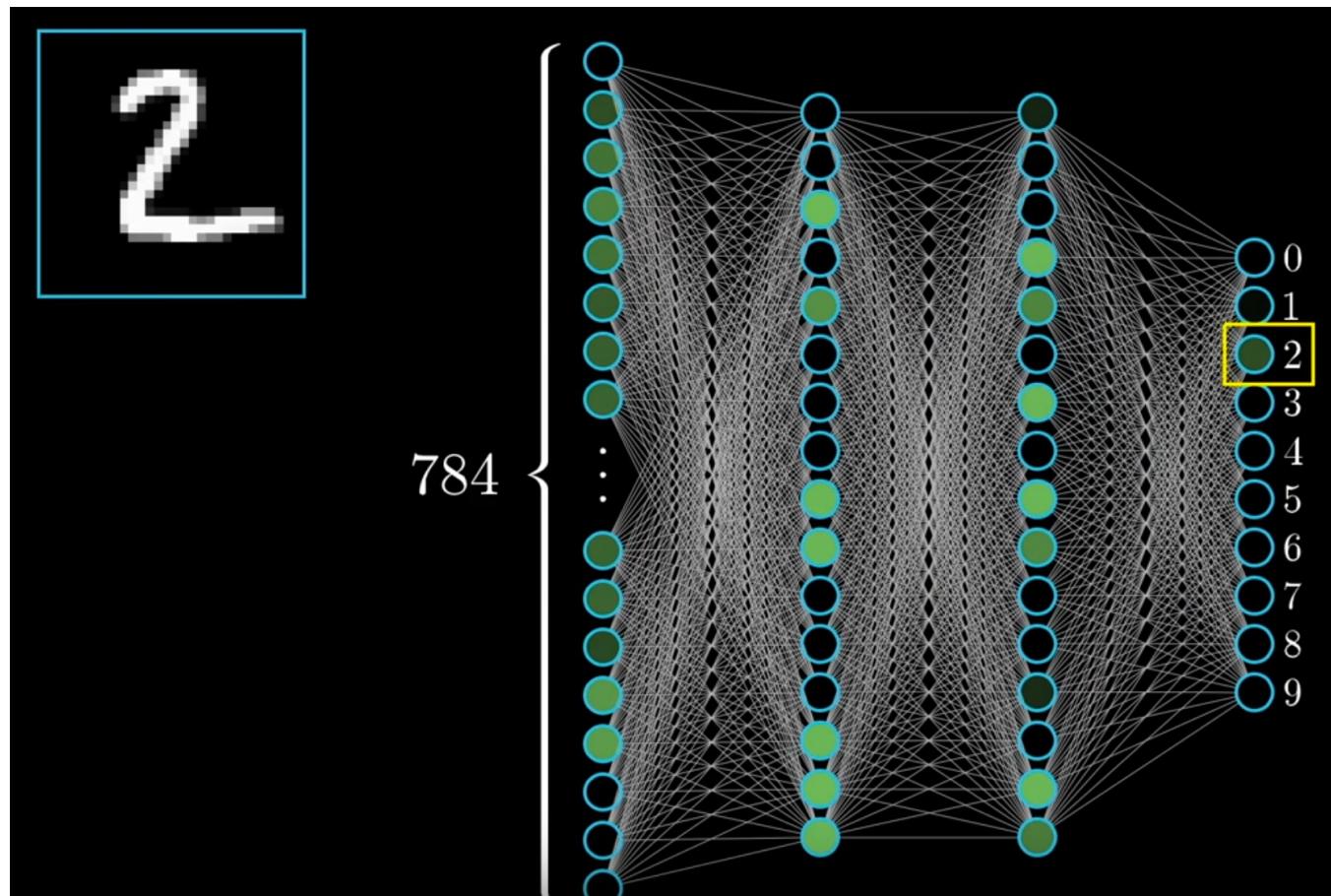
Overview

- Simple introduction into **perceptrons and neural networks**
- **Highlights of the two videos:**
 - **Basics** of neural networks
 - How neural networks **learn**

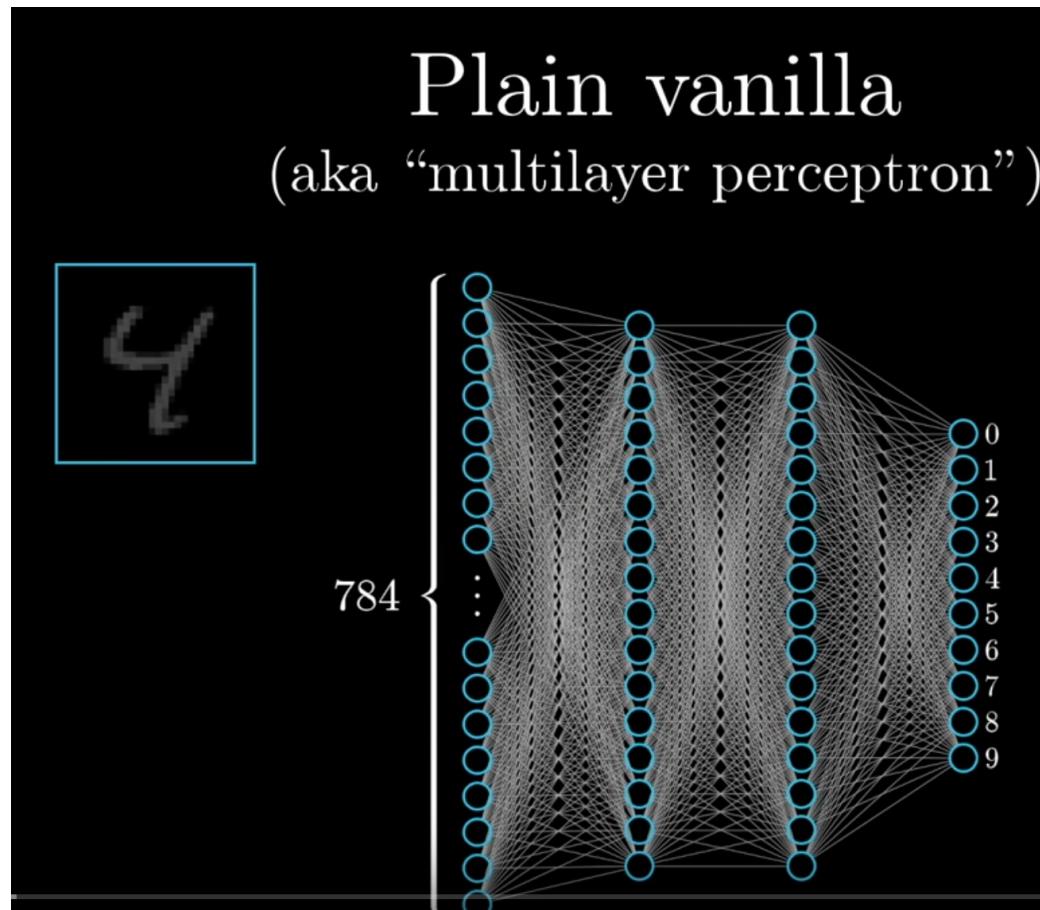
Highlights of the Video



Input 2 → Output 2



Multilayer Perceptron



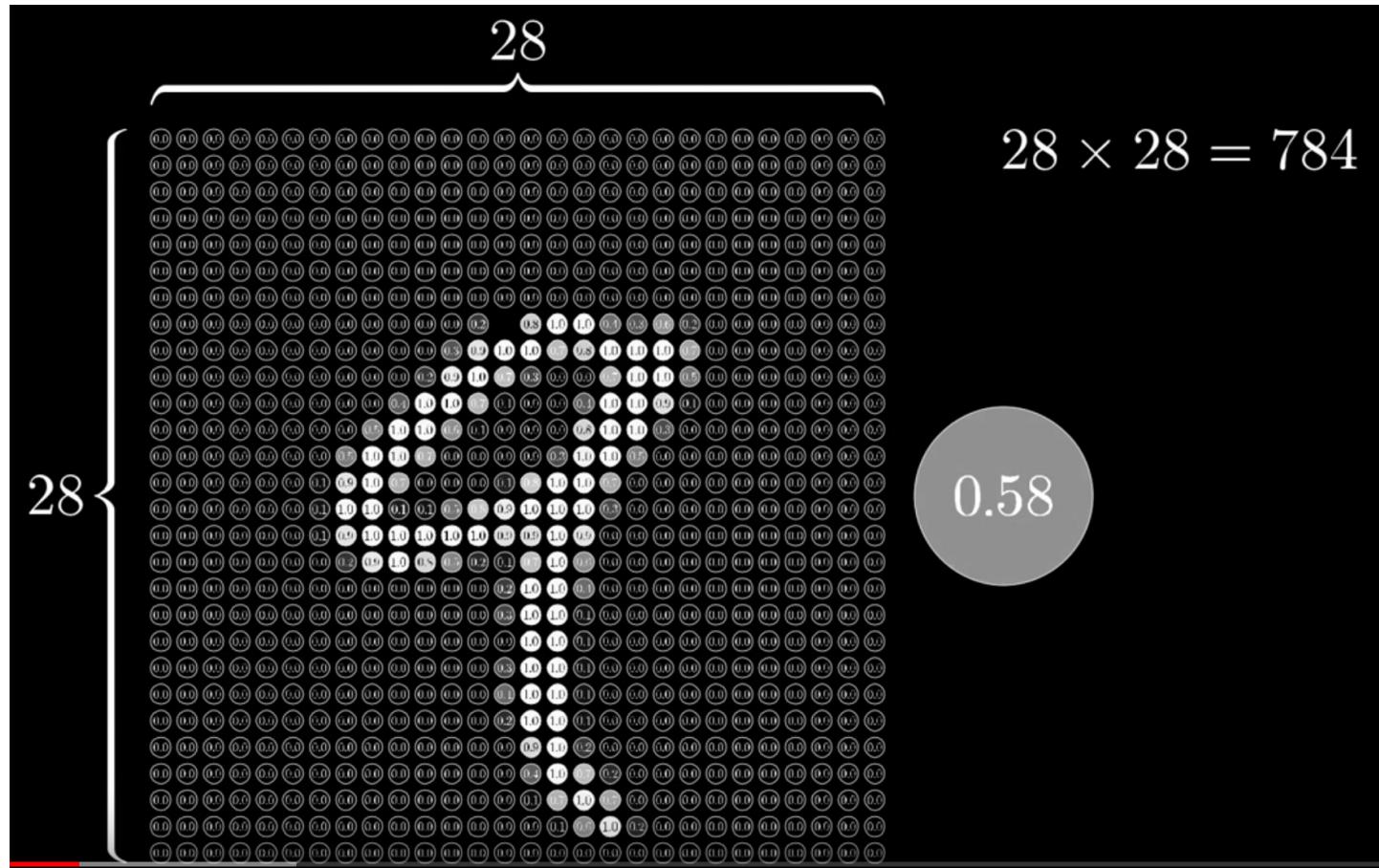
Neuron holds a number = gray scale of pixel



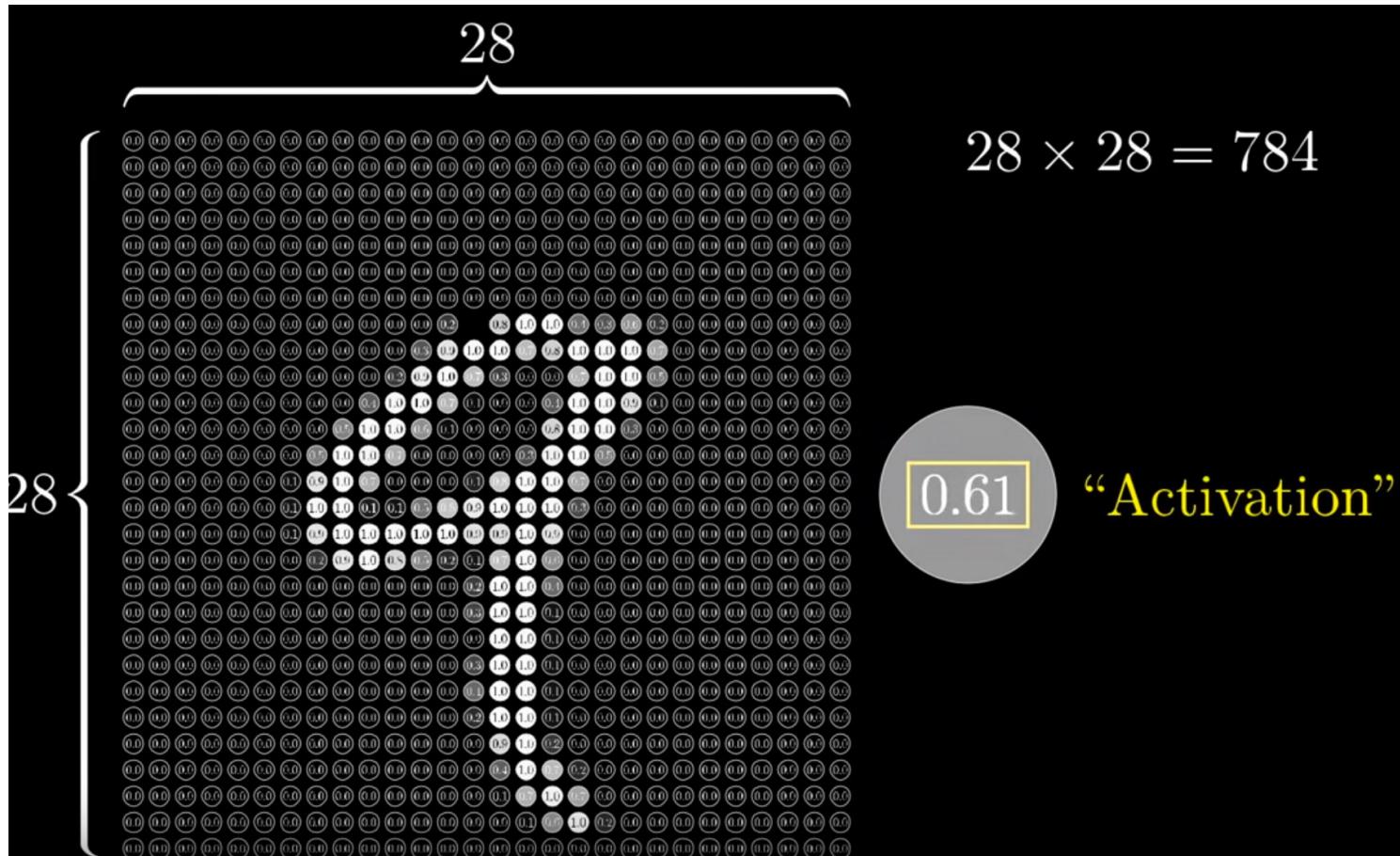
0.8

Neuron → Thing that holds a number

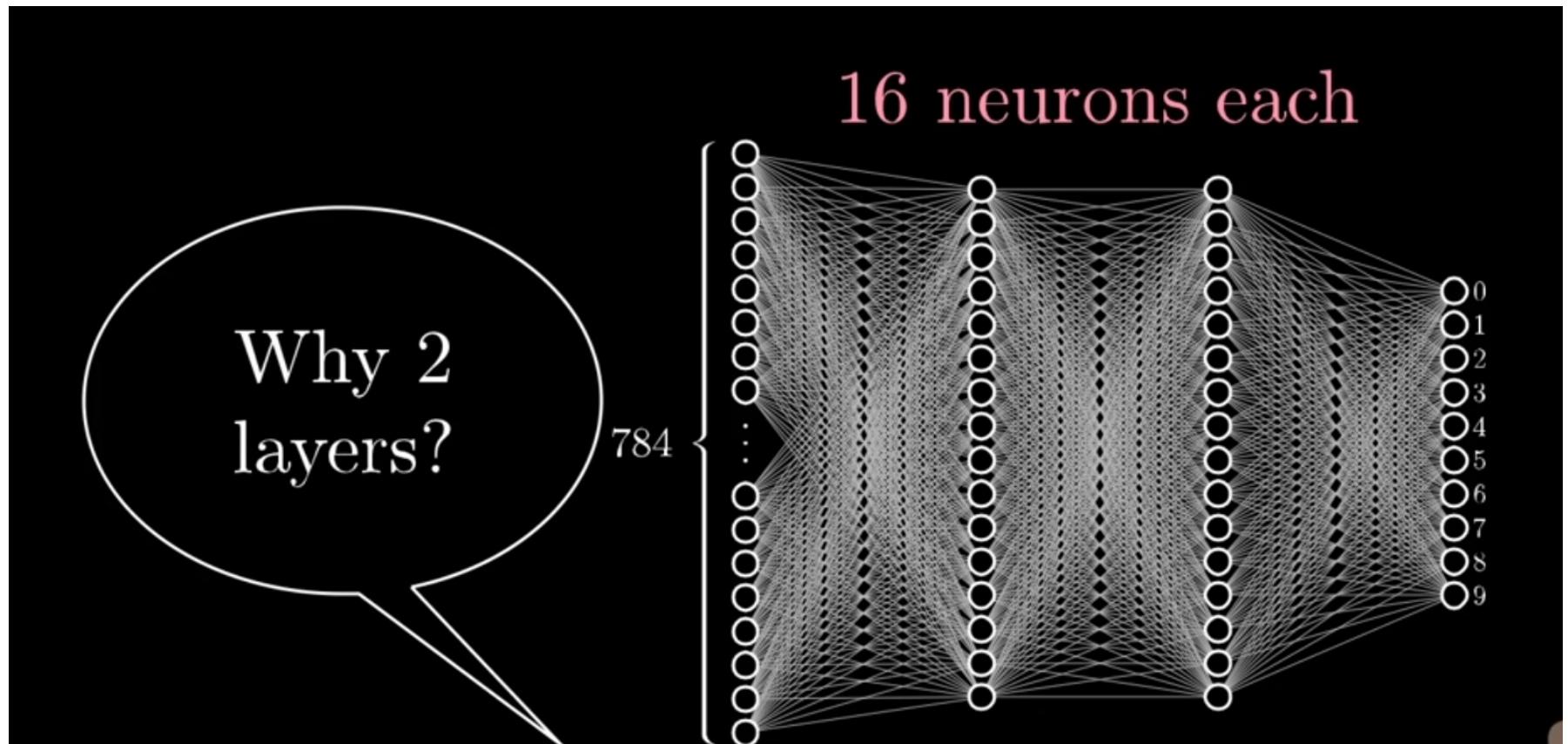
784 neurons to represent 28x28 pixel image



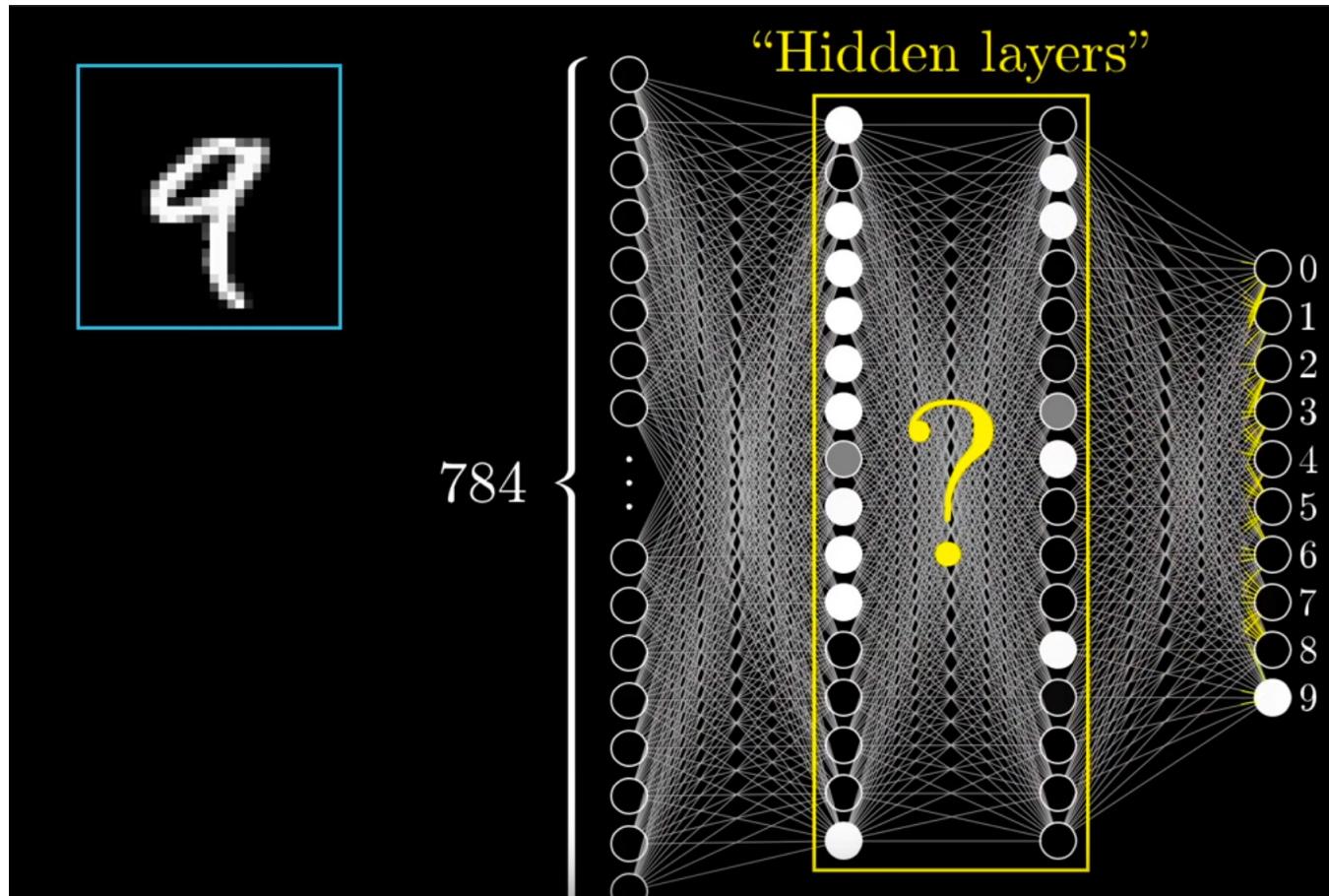
Activation of a neuron



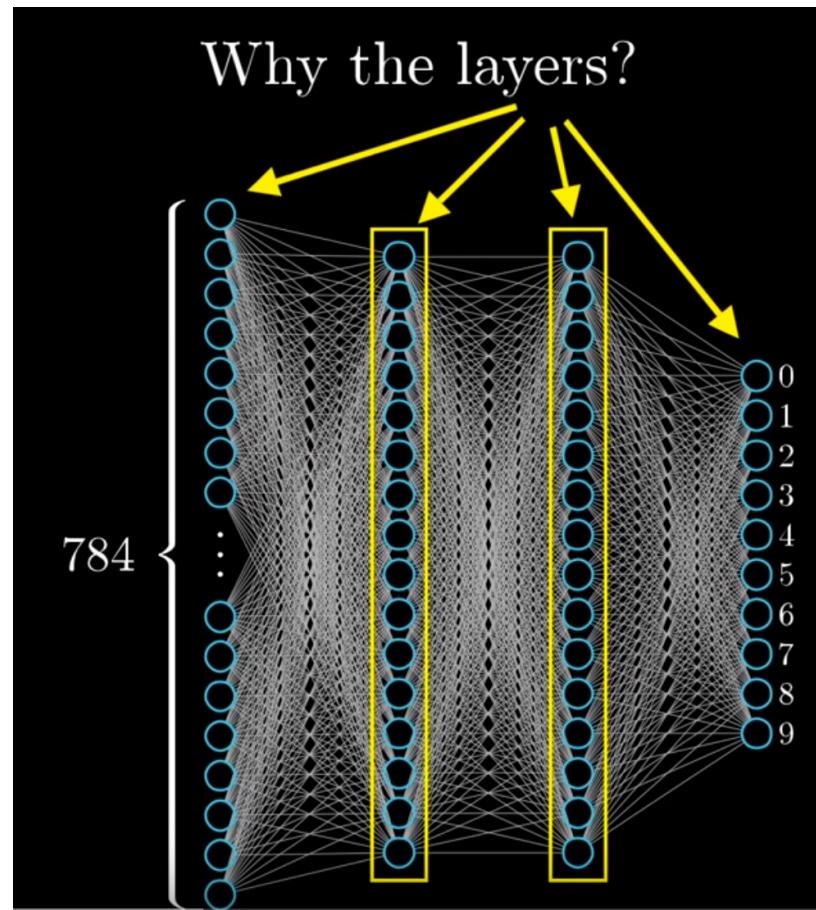
2 hidden layers with 16 neurons each



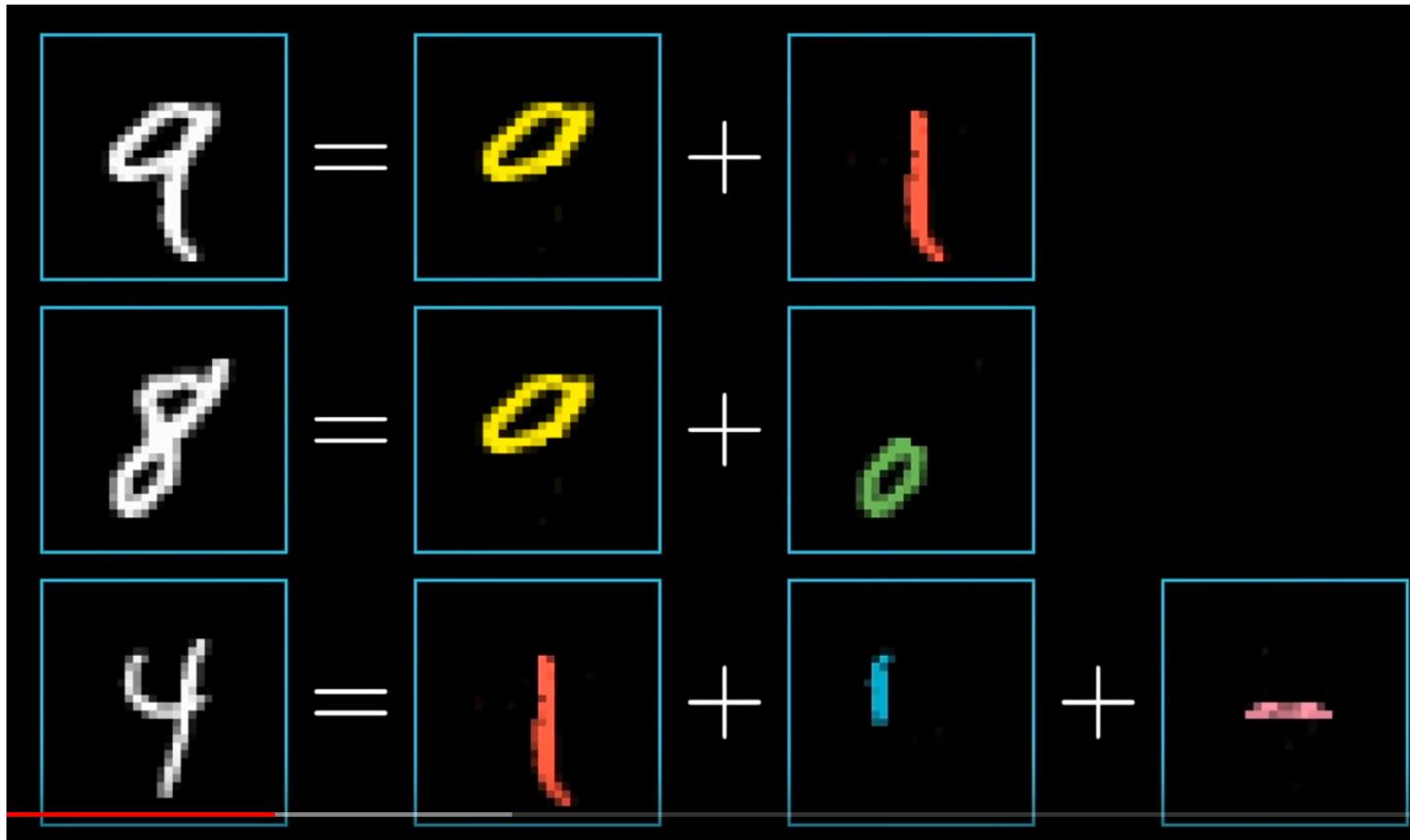
Activation of different layers of trained network



Why hidden layers?

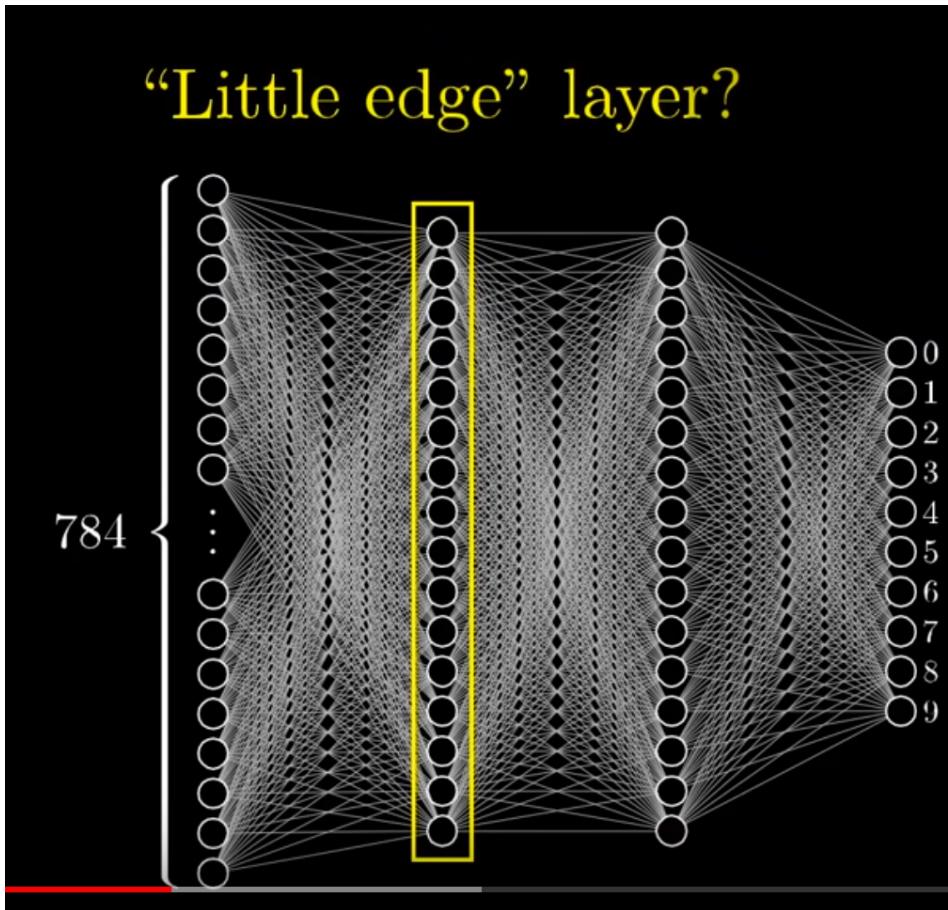


Hidden layers should detect different sub components

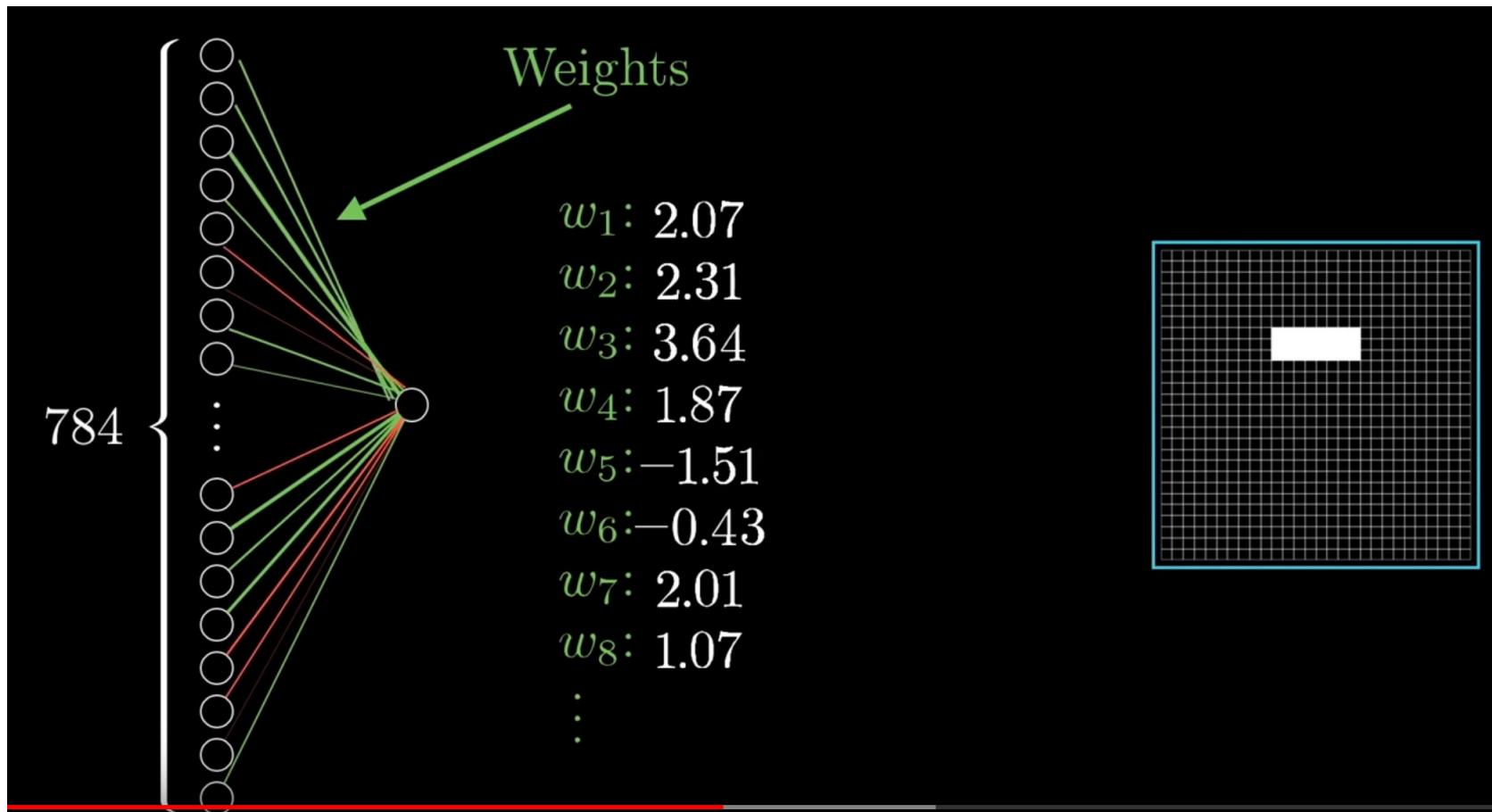


2nd layer detects edges?

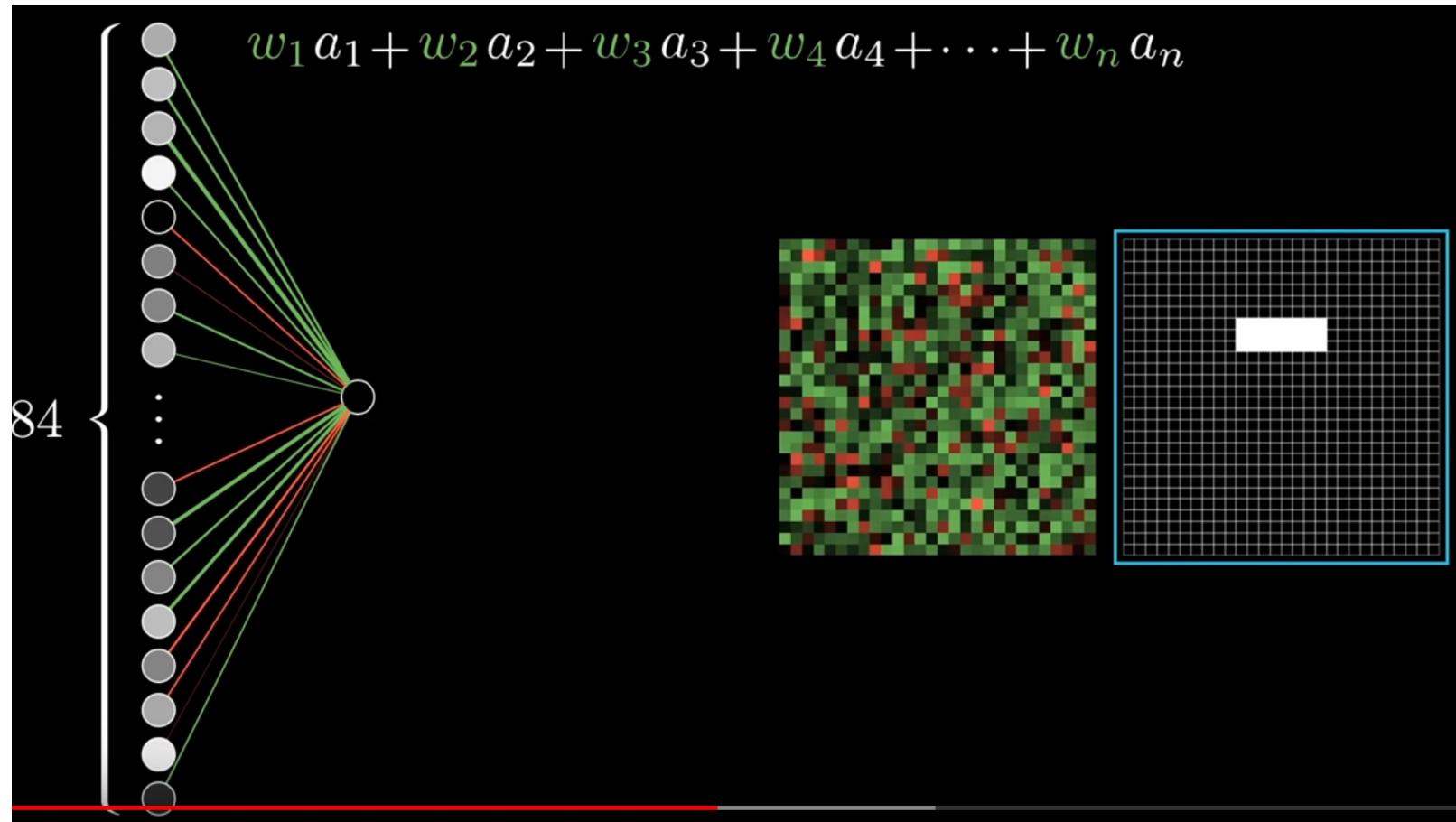
“Little edge” layer?



Weights for the connections between neurons

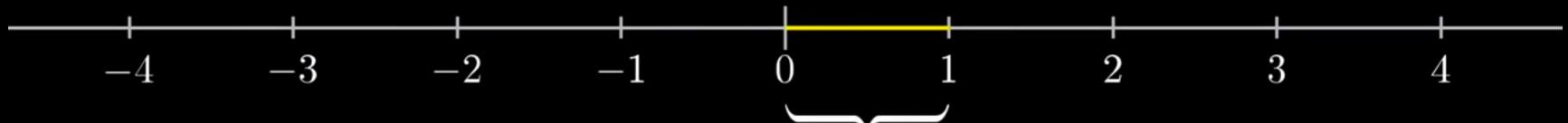


Weighted sum over all weights



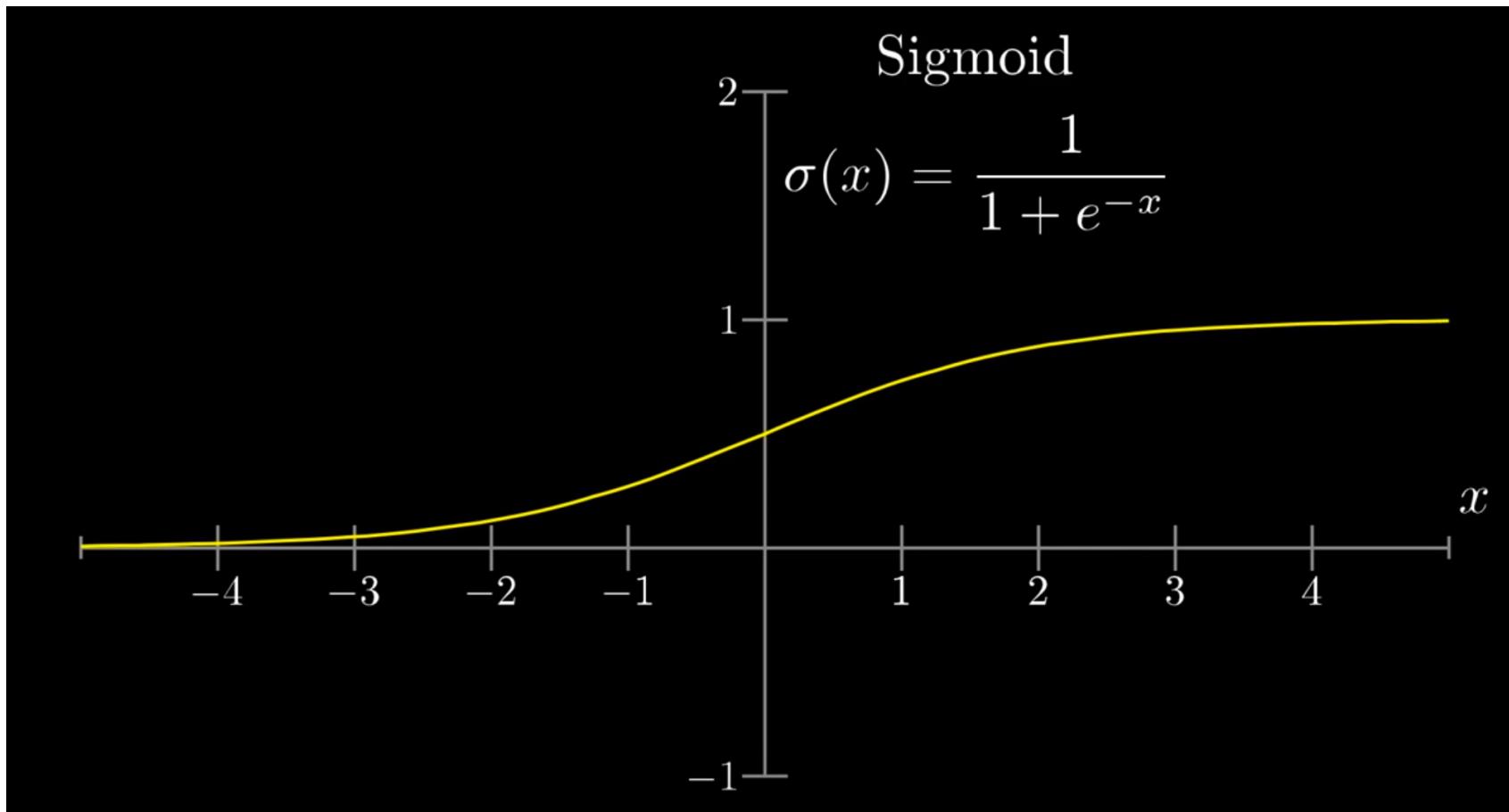
Activation between 0 and 1

$$w_1 a_1 + w_2 a_2 + w_3 a_3 + w_4 a_4 + \cdots + w_n a_n$$

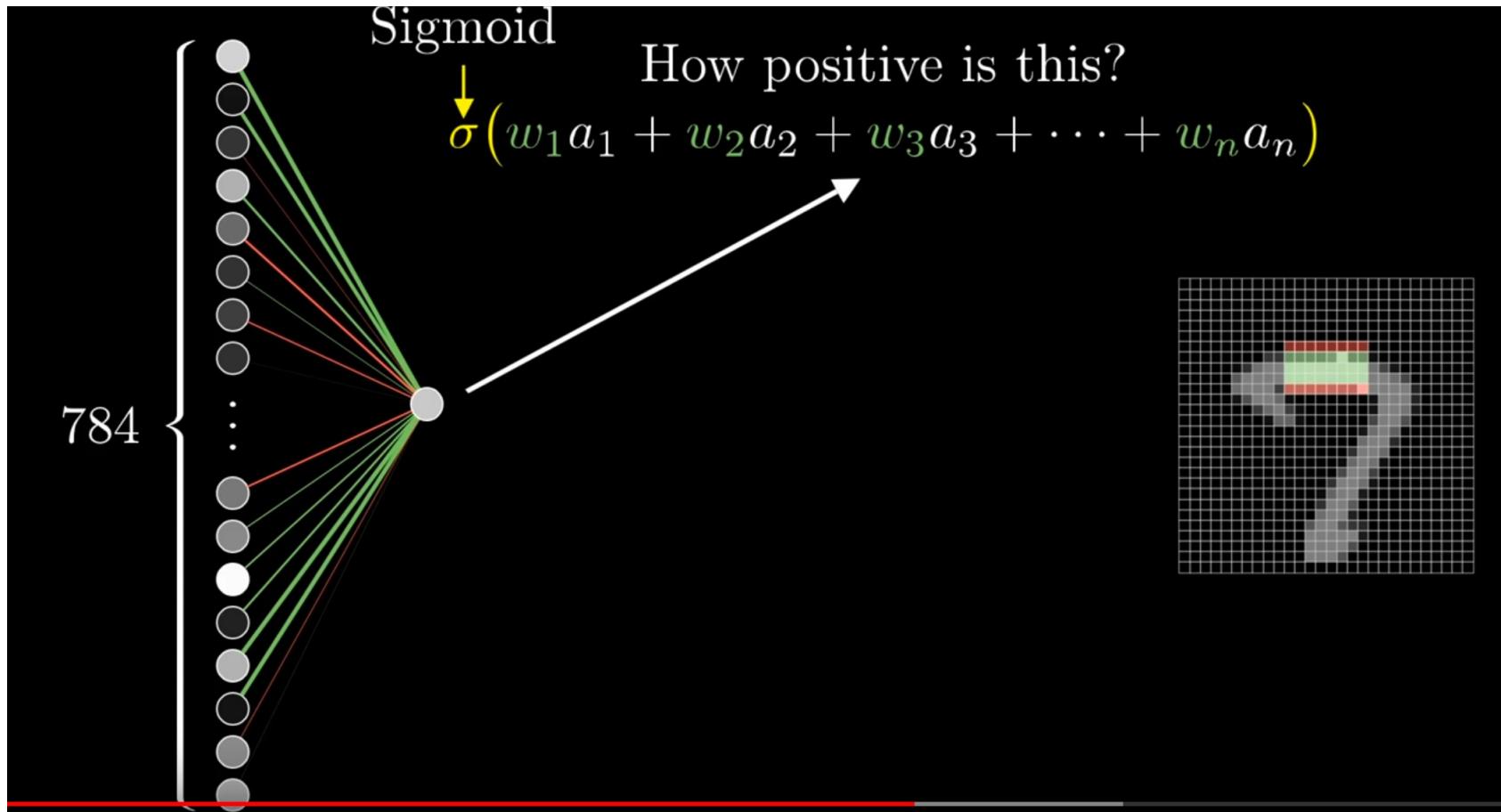


Activations should be in this range

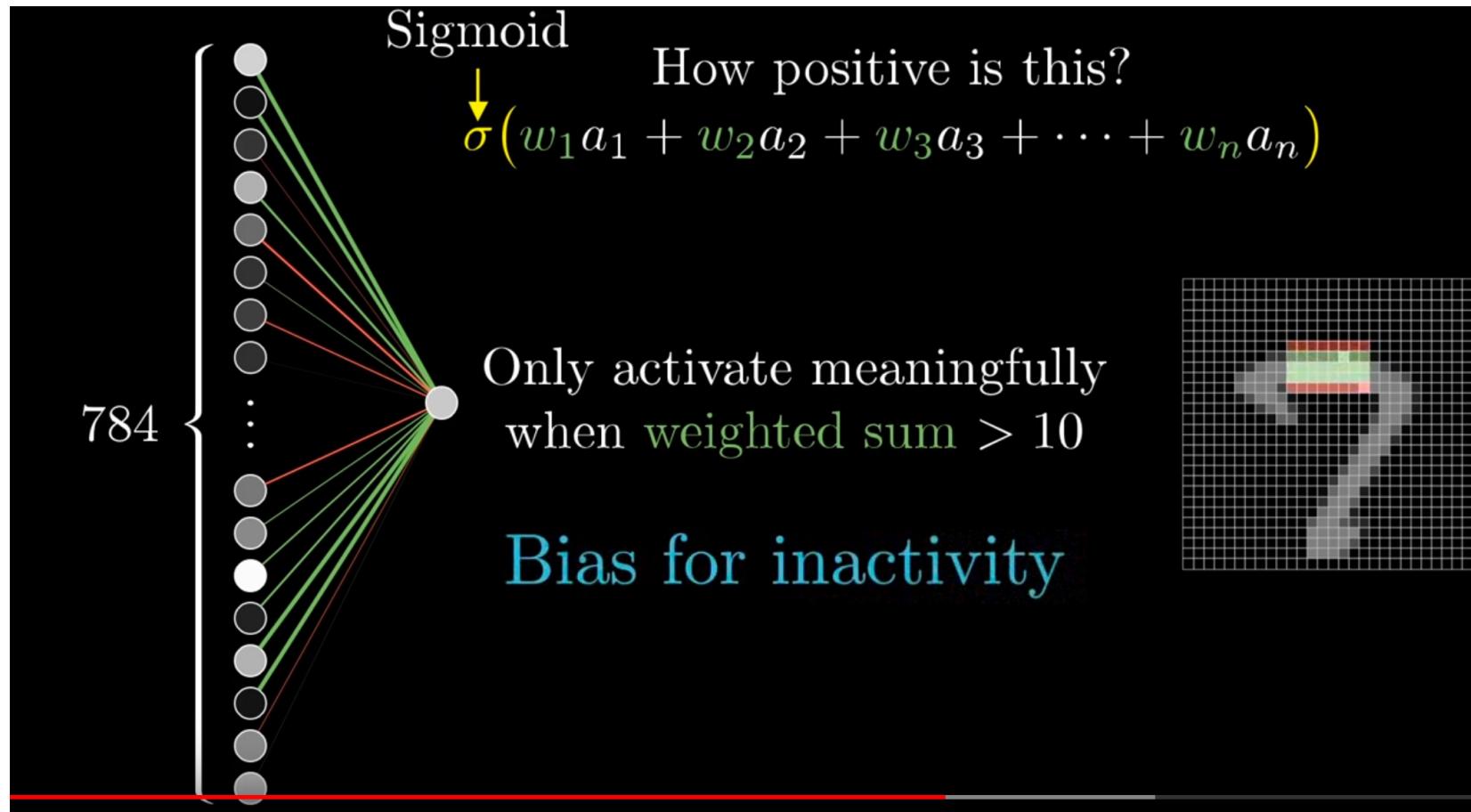
Use Sigmoid function (again)



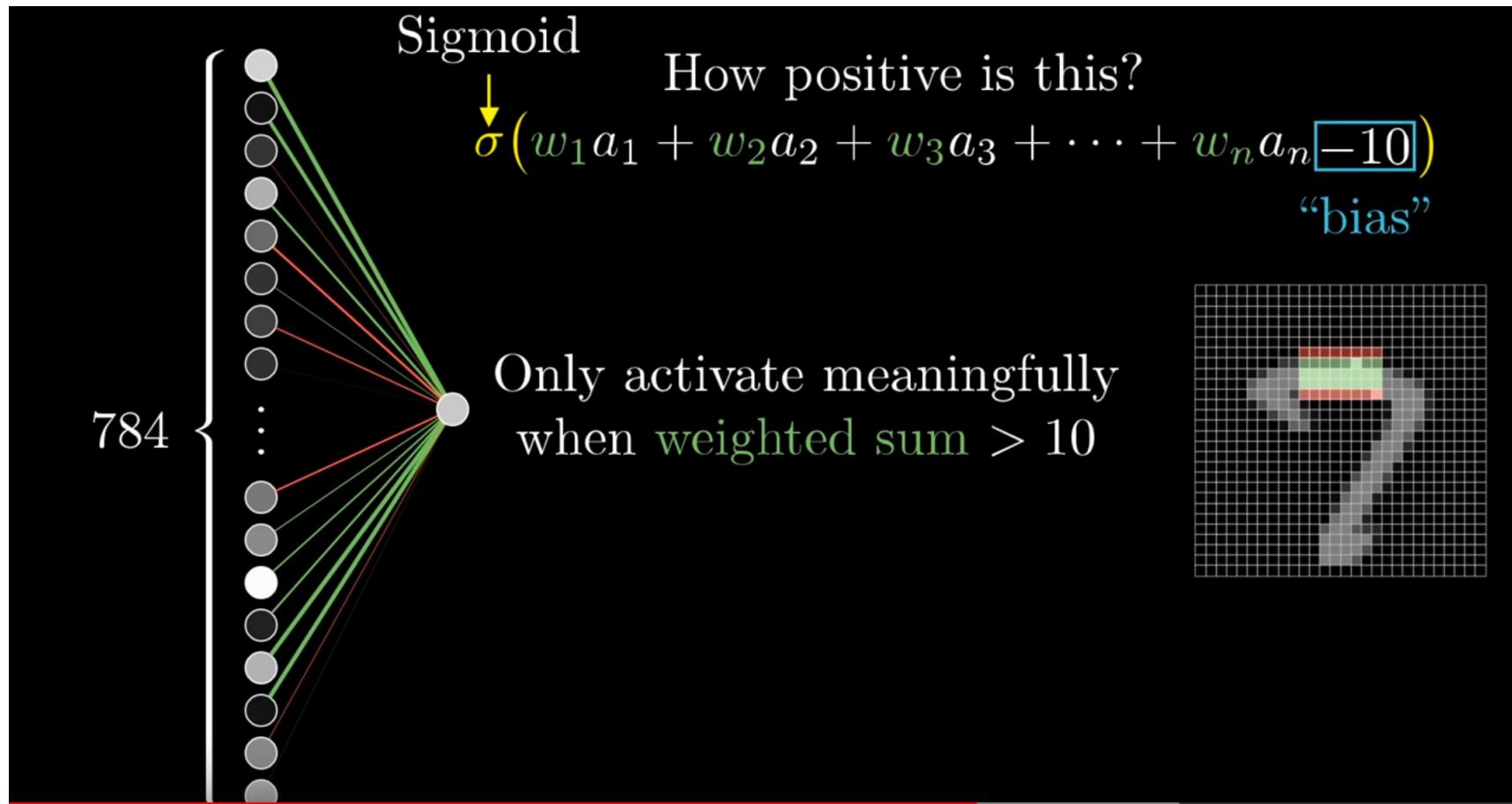
Activation = how positive is the weighted sum



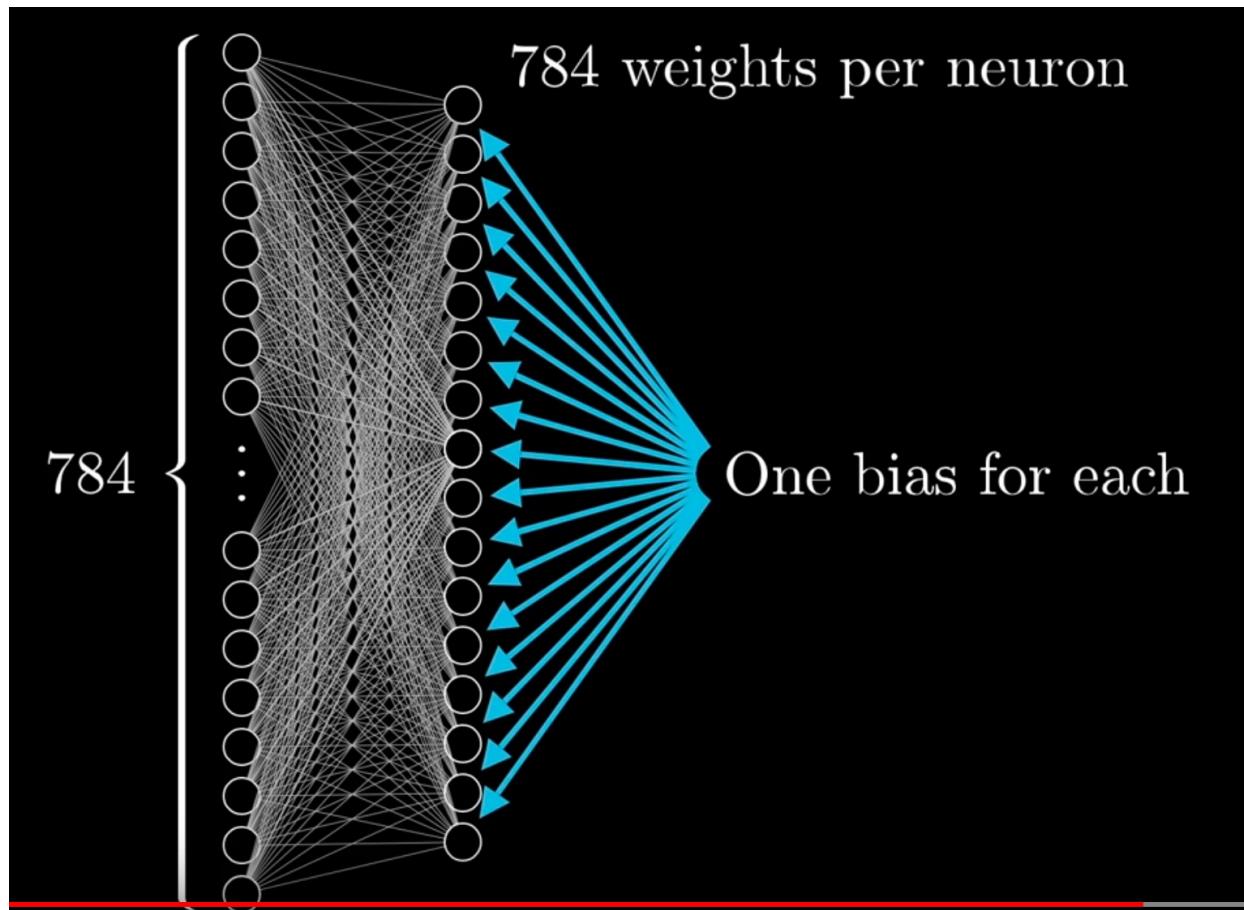
Add bias for inactivity if weighted sum <= 10



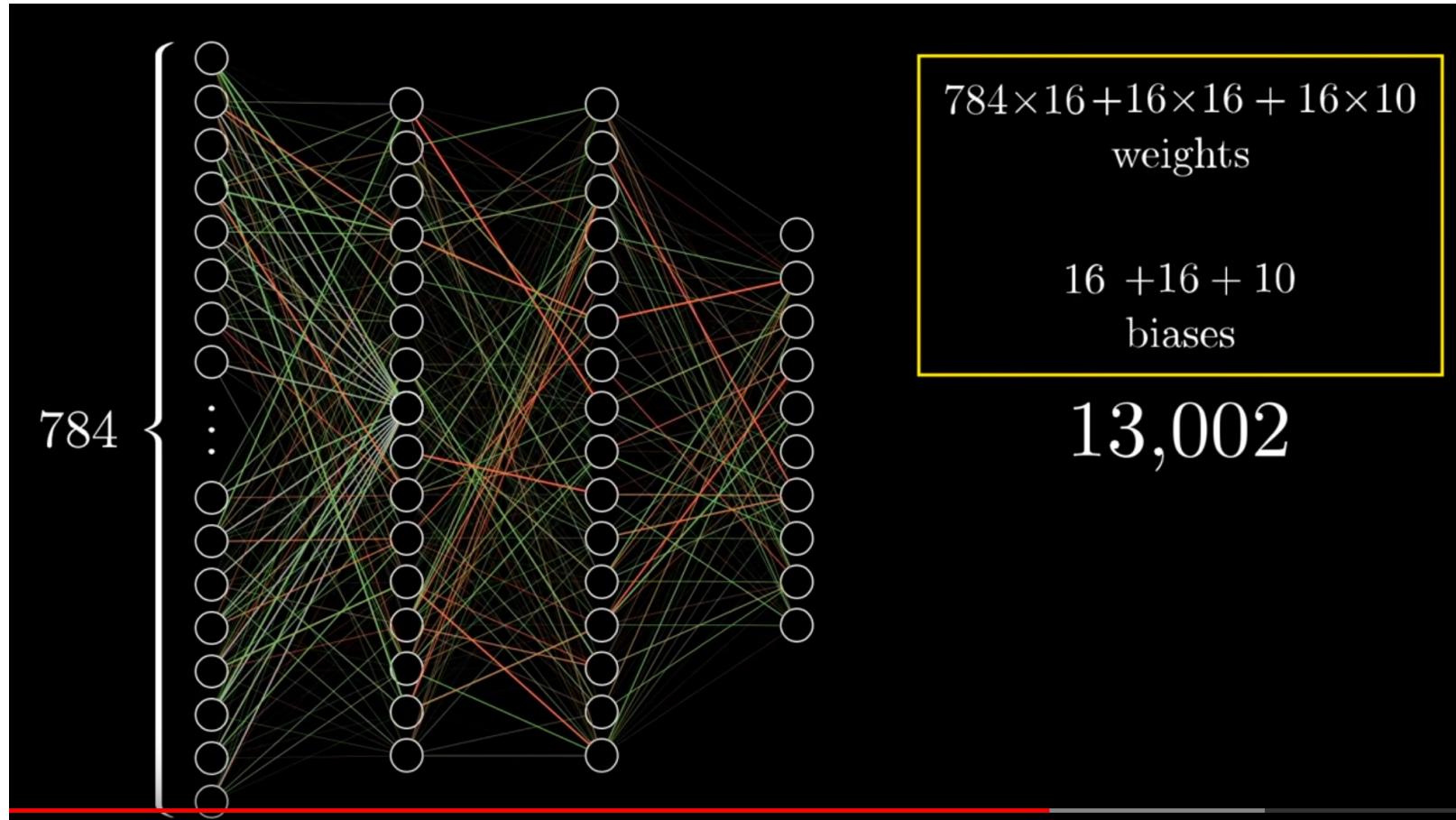
Add bias before applying sigmoid function



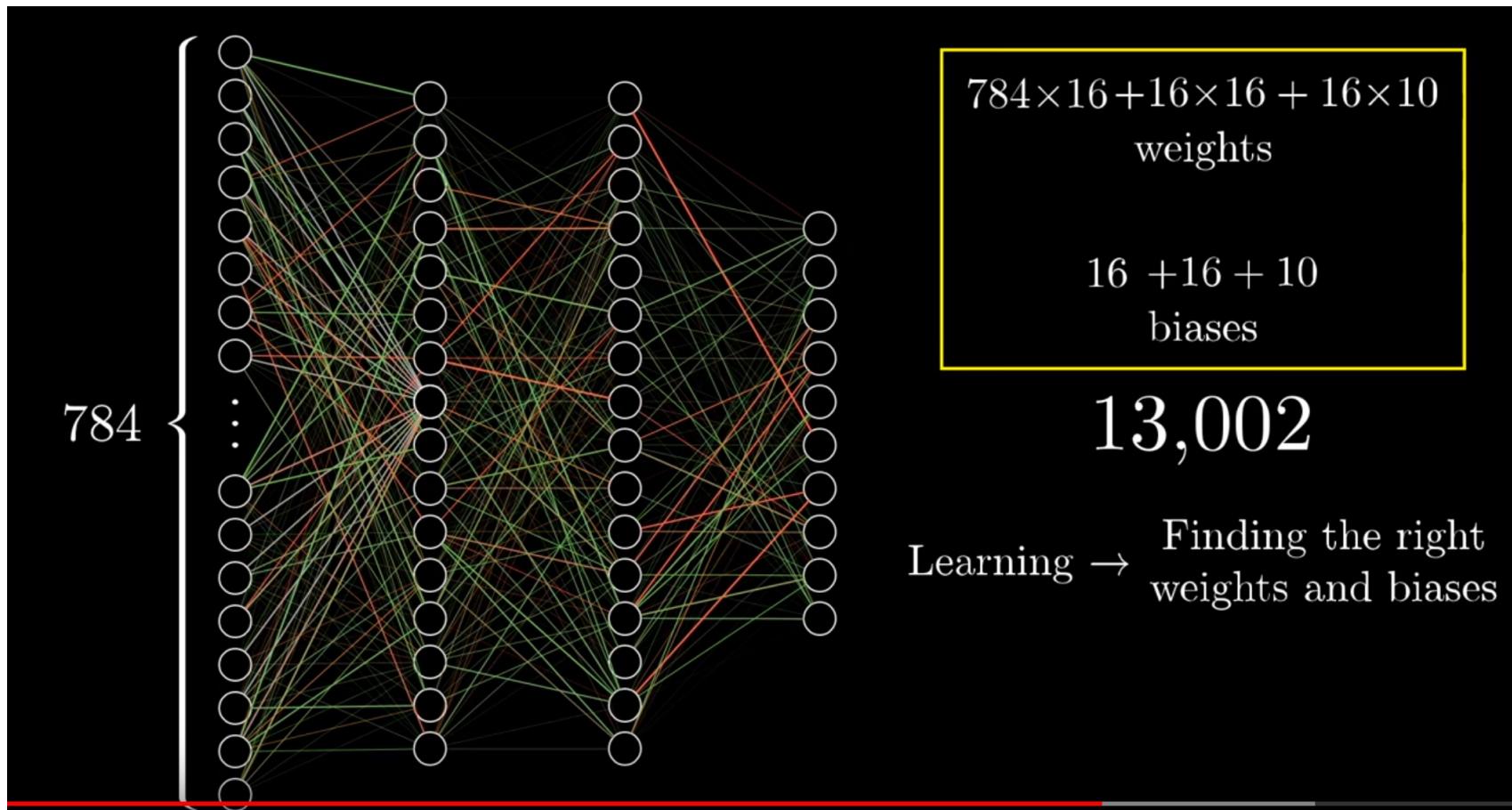
Many weights and biases



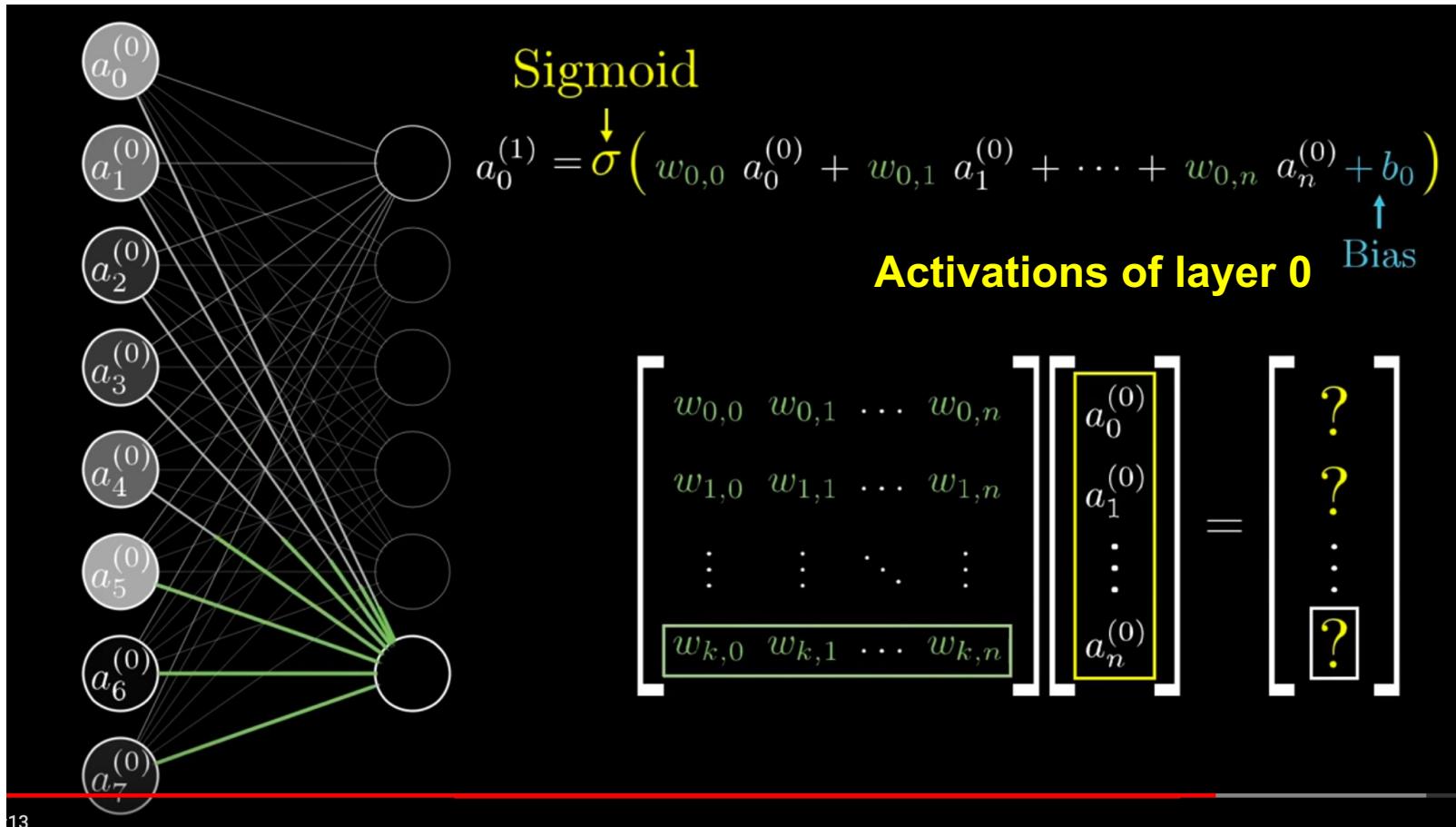
13,002 weights & biases



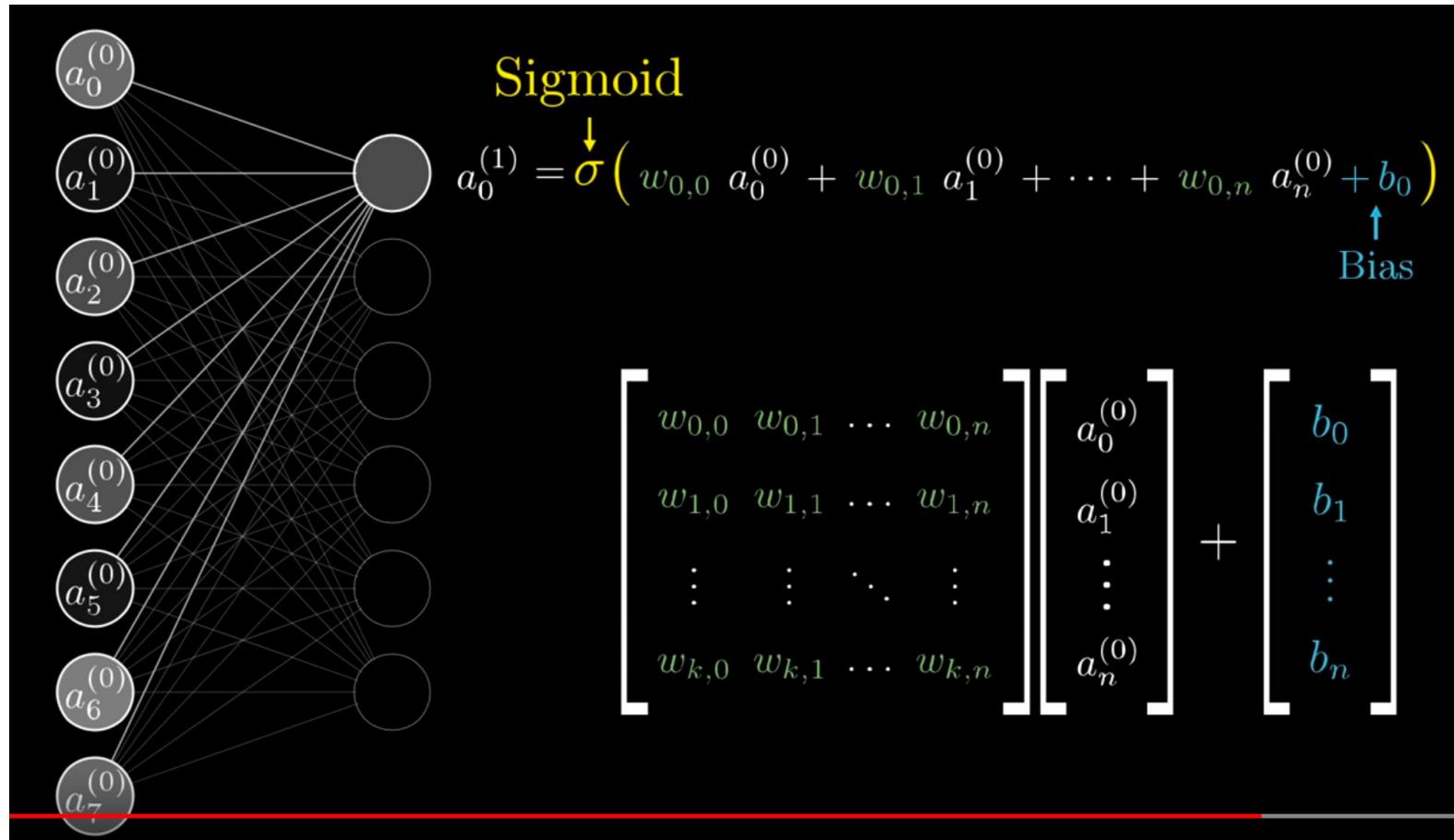
Learning = finding the right weights and biases



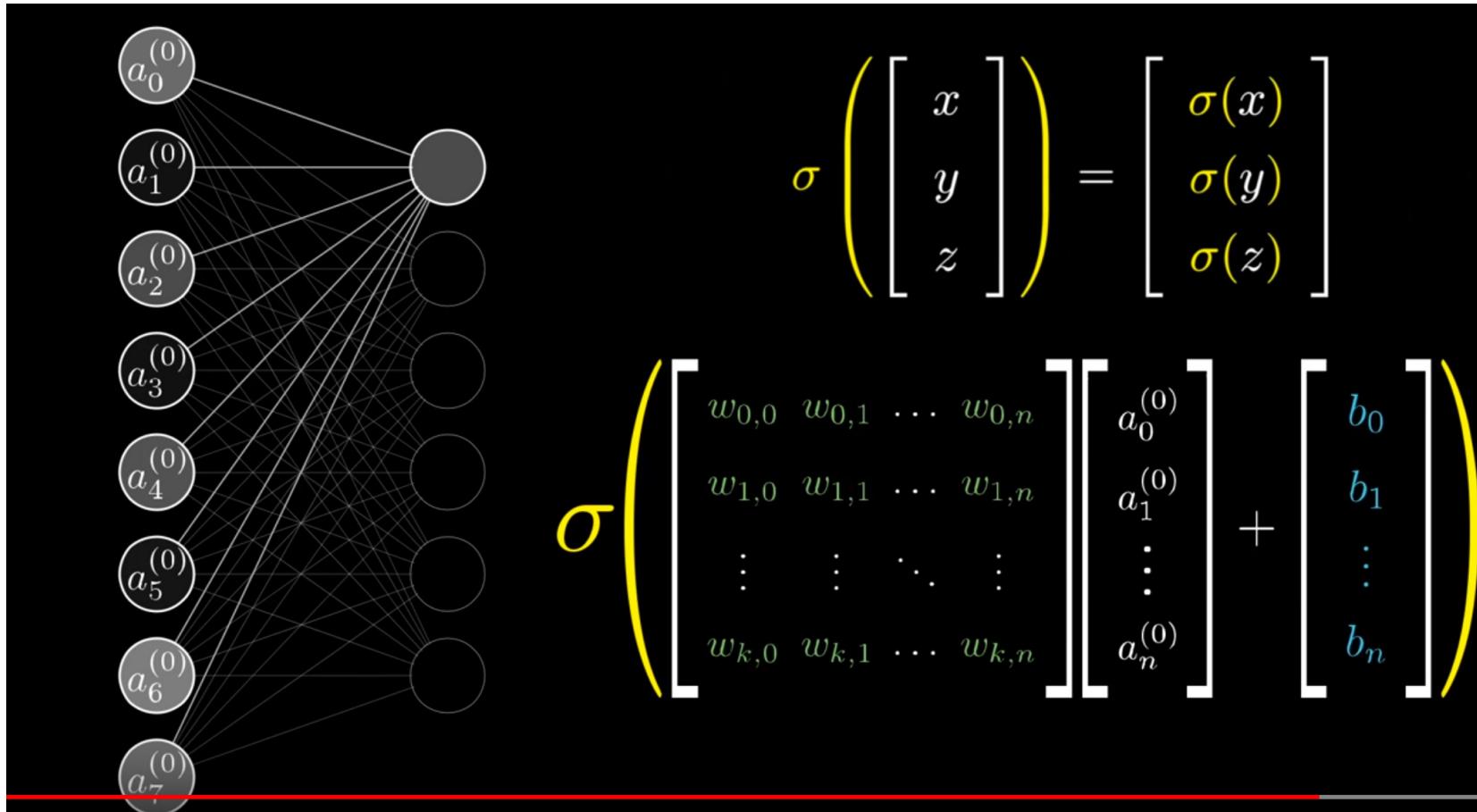
Dot product = Matrix-vector-product



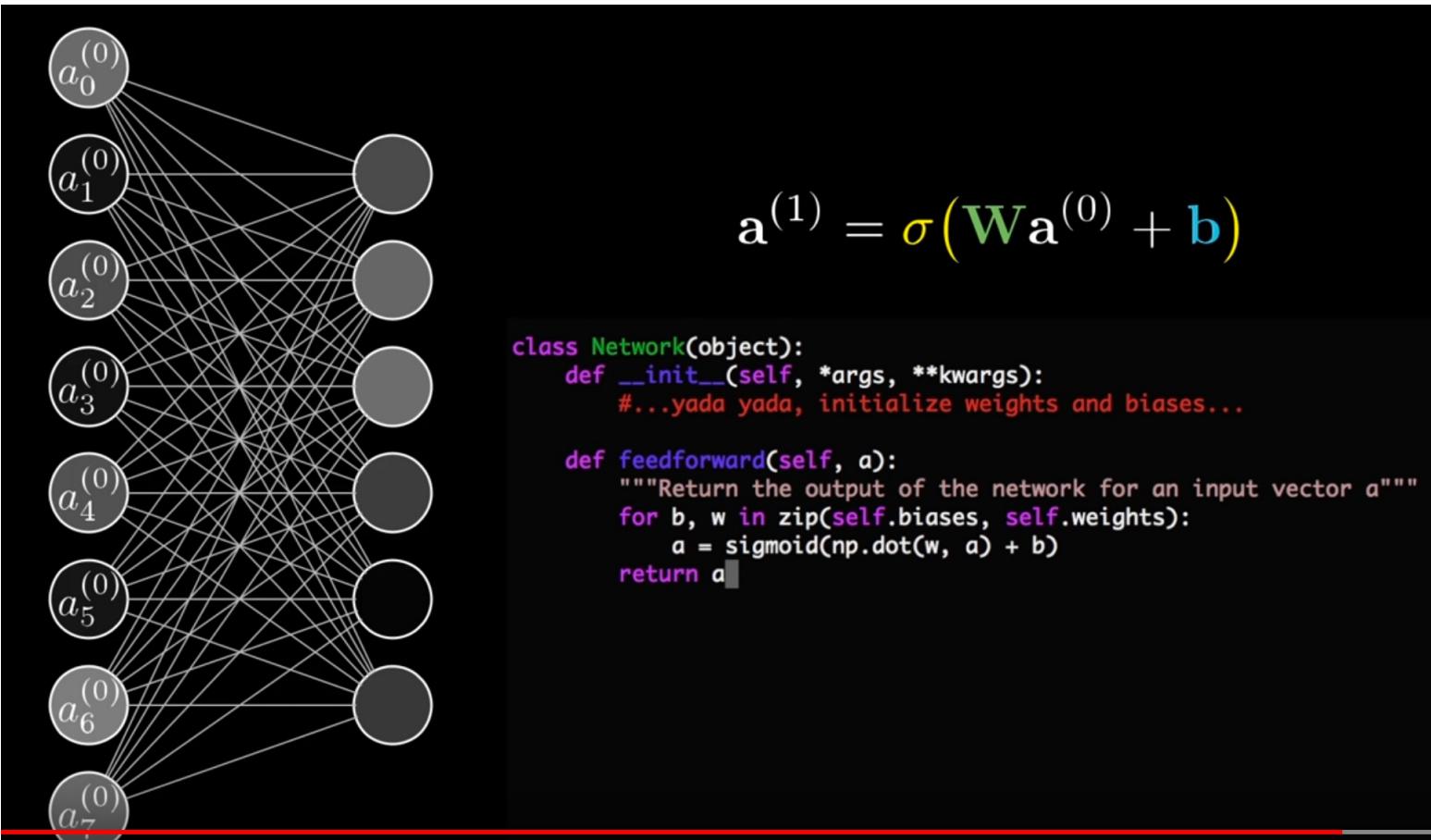
Add bias vector



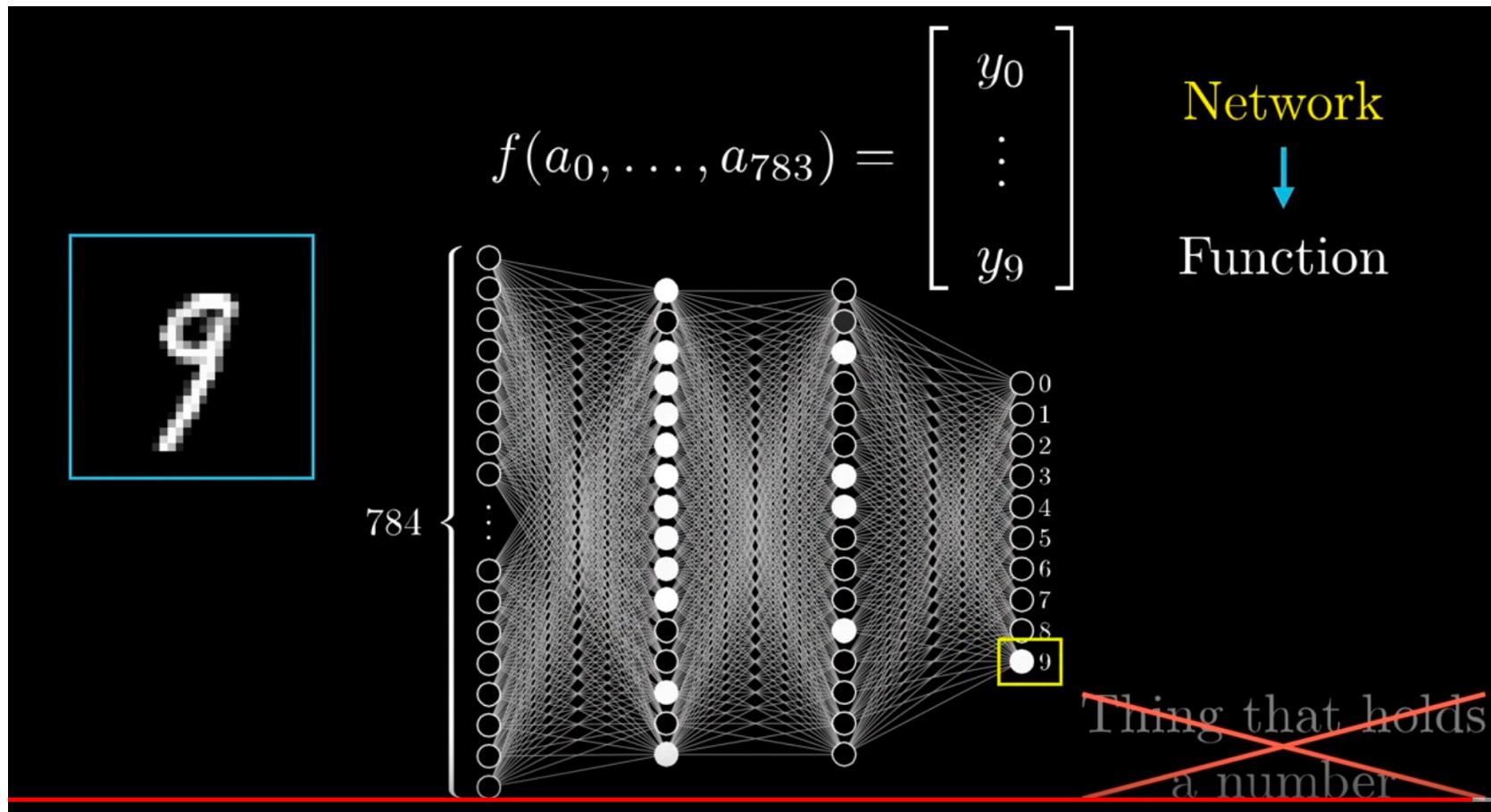
Apply sigmoid function



Code



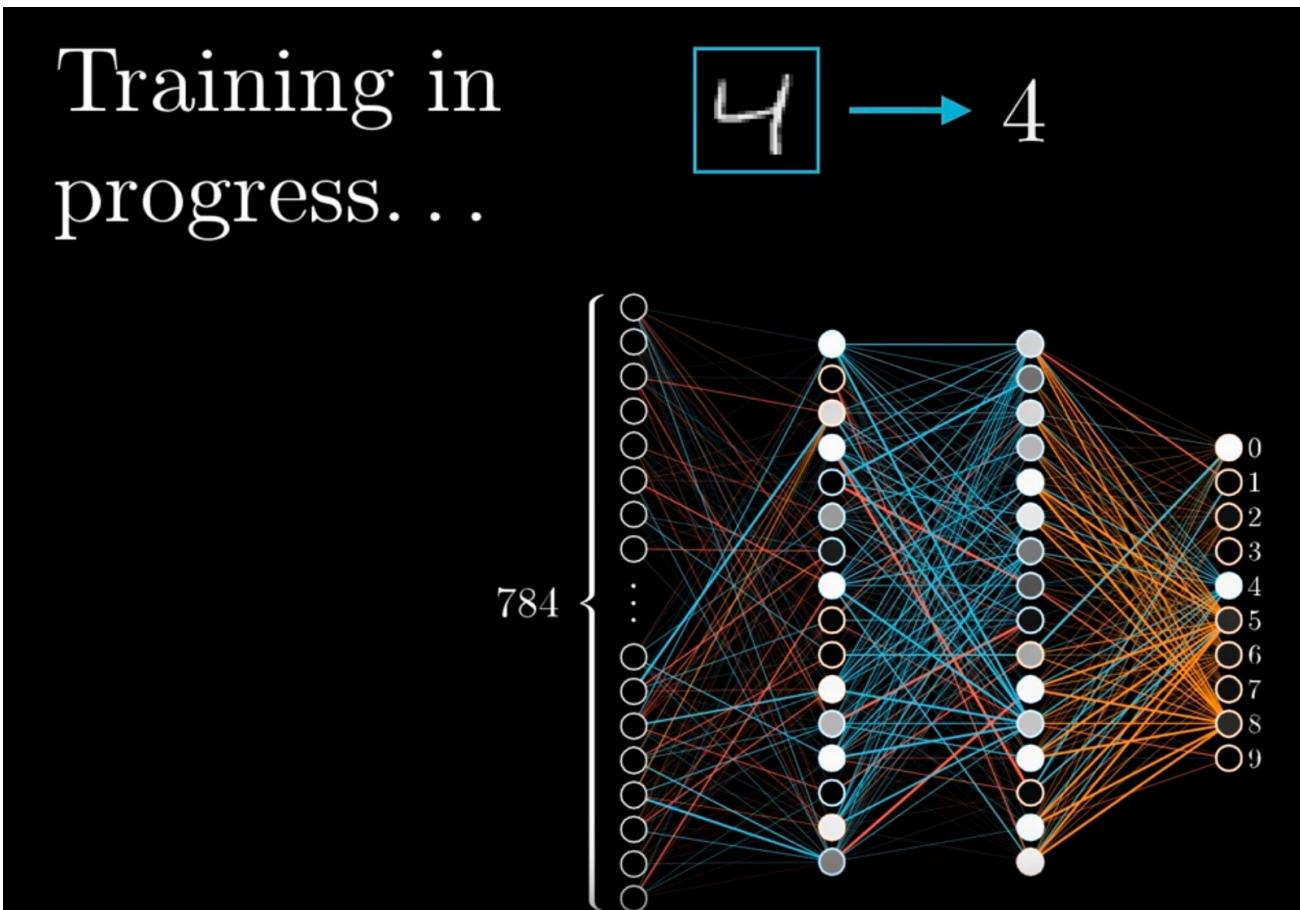
Neuron + network of neurons = function



Overview

- Simple introduction into **perceptrons and neural networks**
- Highlights of the two videos:
 - **Basics** of neural networks
 - How neural networks **learn**

Training with Label “4”



Training and testing

Train on
these

5	→ 5
0	→ 0
4	→ 4
1	→ 1
9	→ 9
2	→ 2
1	→ 1
3	→ 3
1	→ 1
4	→ 4

3	→ 3
5	→ 5
3	→ 3
6	→ 6
1	→ 1
7	→ 7
2	→ 2
8	→ 8
6	→ 6
9	→ 9

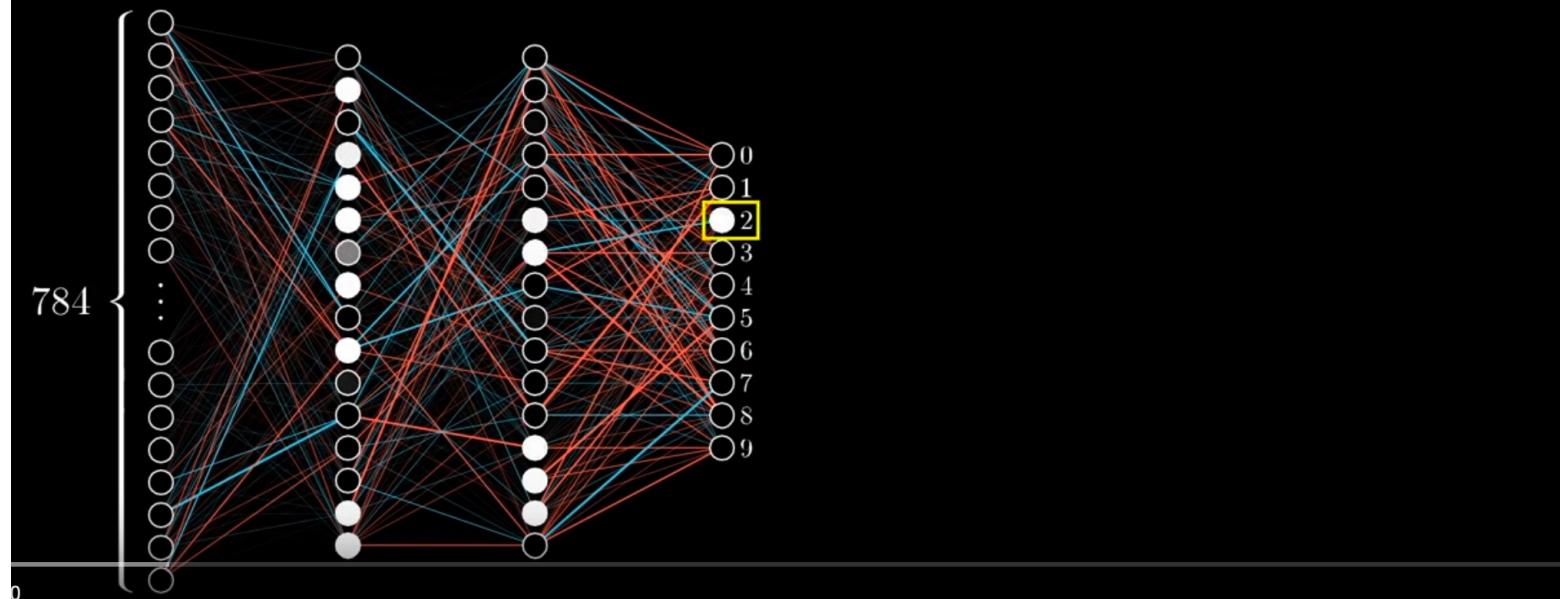
Test on
these

Test evaluation

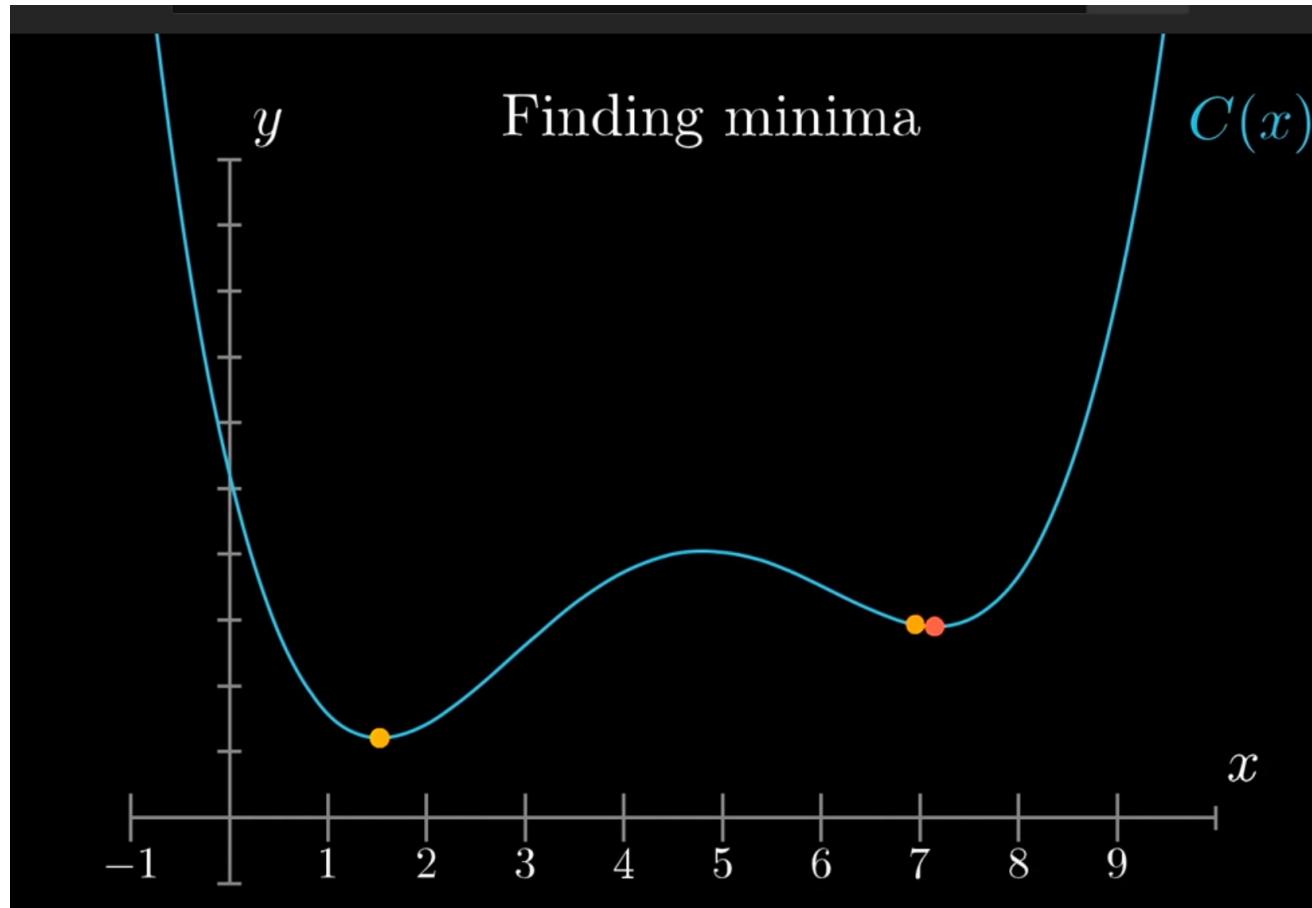
Testing data



$$\frac{\text{Number correct}}{\text{total}} = \frac{47}{48} = 0.979$$

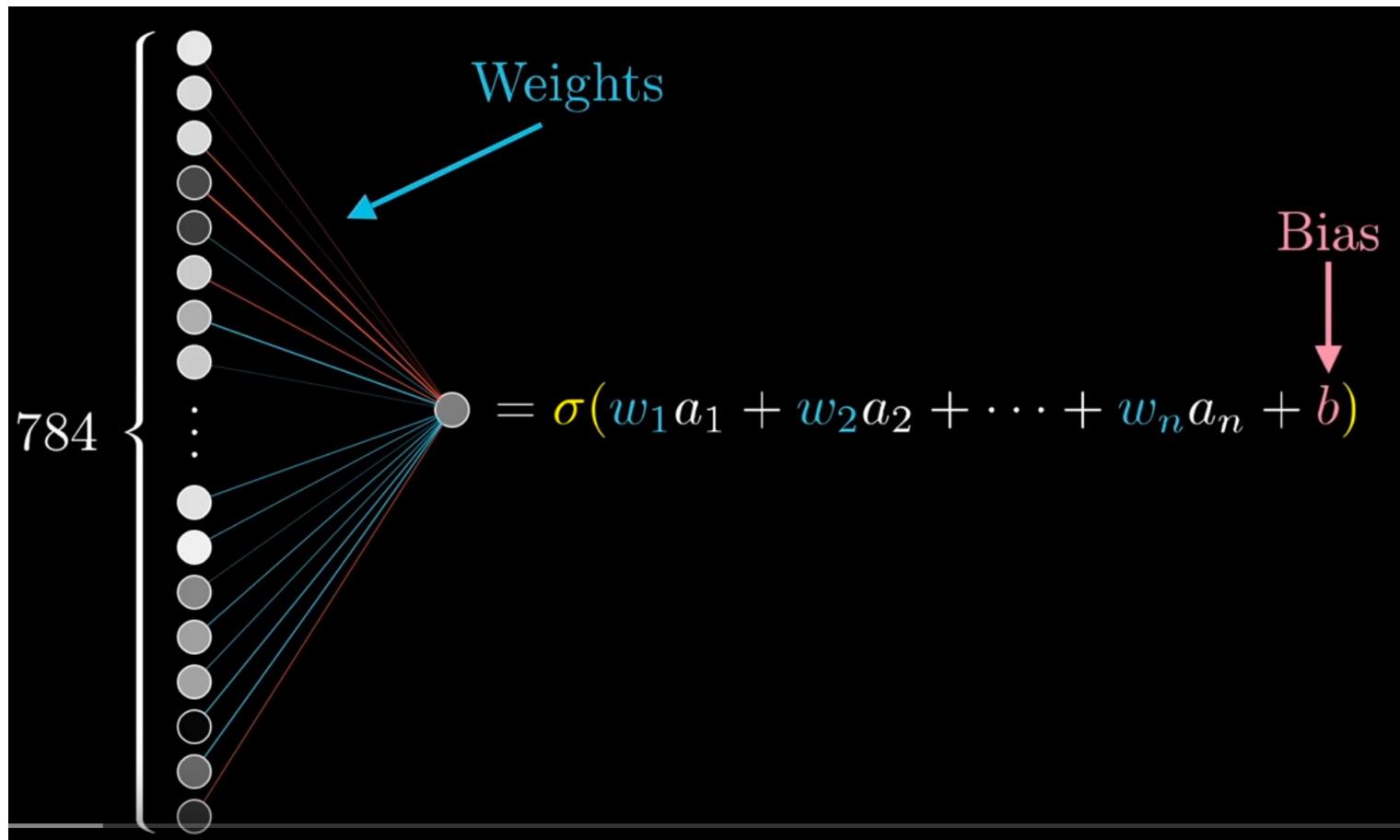


Learning = finding minimum of certain function

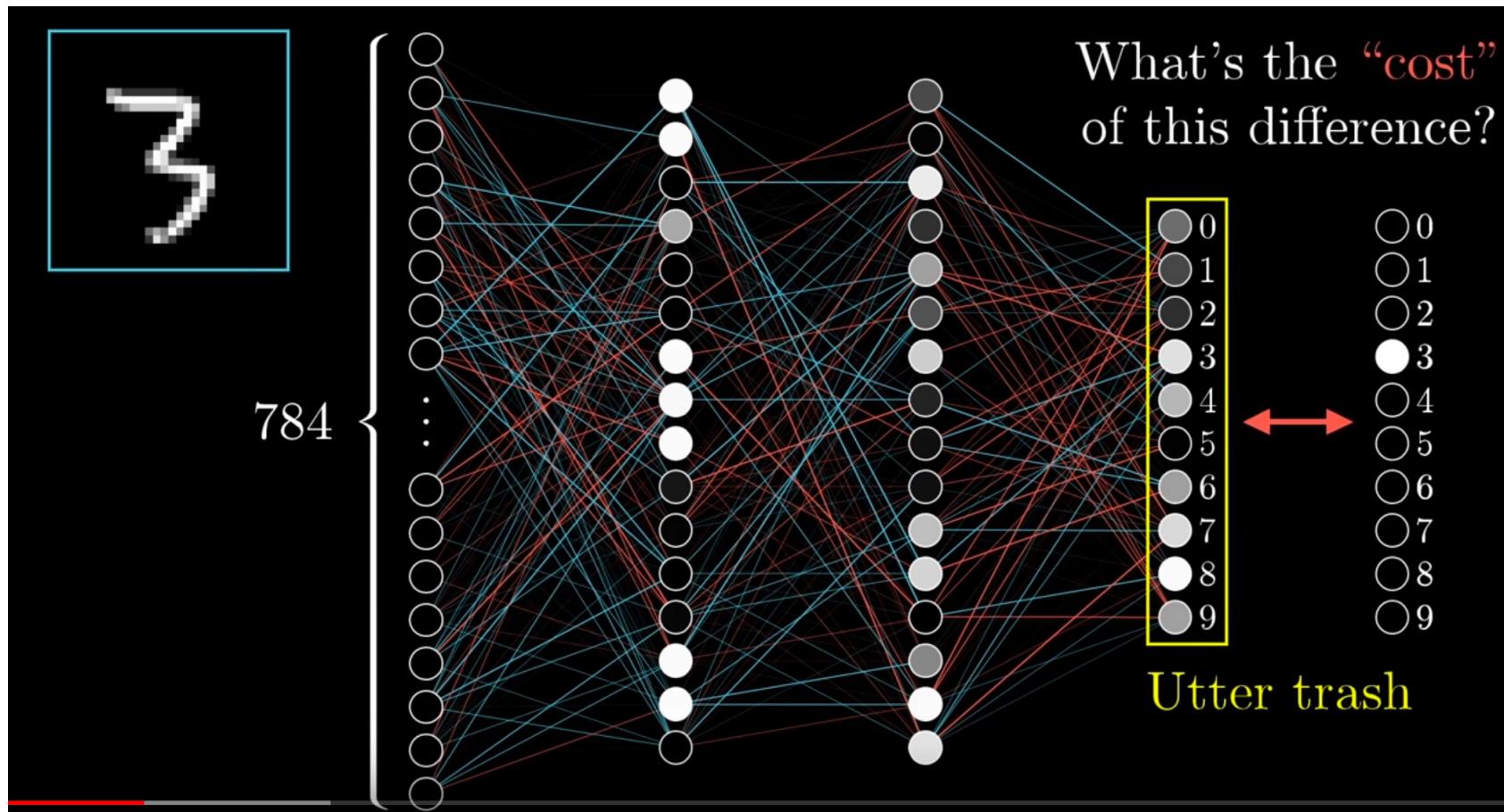


Weight = strength of connection

Bias = is neuron active or not



Random initialization of weights and biases



Calculate difference (cost) between actual value and expected value

Cost of **3**

$$3.37 \left\{ \begin{array}{l} 0.1863 \leftarrow (0.43 - 0.00)^2 + \\ 0.0809 \leftarrow (0.28 - 0.00)^2 + \\ 0.0357 \leftarrow (0.19 - 0.00)^2 + \\ 0.0138 \leftarrow (0.88 - 1.00)^2 + \\ 0.5242 \leftarrow (0.72 - 0.00)^2 + \\ 0.0001 \leftarrow (0.01 - 0.00)^2 + \\ 0.4079 \leftarrow (0.64 - 0.00)^2 + \\ 0.7388 \leftarrow (0.86 - 0.00)^2 + \\ 0.9817 \leftarrow (0.99 - 0.00)^2 + \\ 0.3998 \leftarrow (0.63 - 0.00)^2 \end{array} \right.$$

What's the “cost”
of this difference?

<input type="radio"/>	0
<input type="radio"/>	1
<input type="radio"/>	2
<input checked="" type="radio"/>	3
<input type="radio"/>	4
<input type="radio"/>	5
<input type="radio"/>	6
<input type="radio"/>	7
<input type="radio"/>	8
<input type="radio"/>	9



Utter trash

Calculate average costs of all training examples

Average cost of
all training data...

Cost of 

$$\left\{ \begin{array}{l} (0.01 - 0.00)^2 + \\ (0.11 - 0.00)^2 + \\ (0.02 - 0.00)^2 + \\ (0.93 - 0.00)^2 + \\ (0.77 - 0.00)^2 + \\ (0.98 - 0.00)^2 + \\ (0.20 - 0.00)^2 + \\ (0.90 - 0.00)^2 + \\ (0.94 - 1.00)^2 + \\ (0.03 - 0.00)^2 \end{array} \right.$$

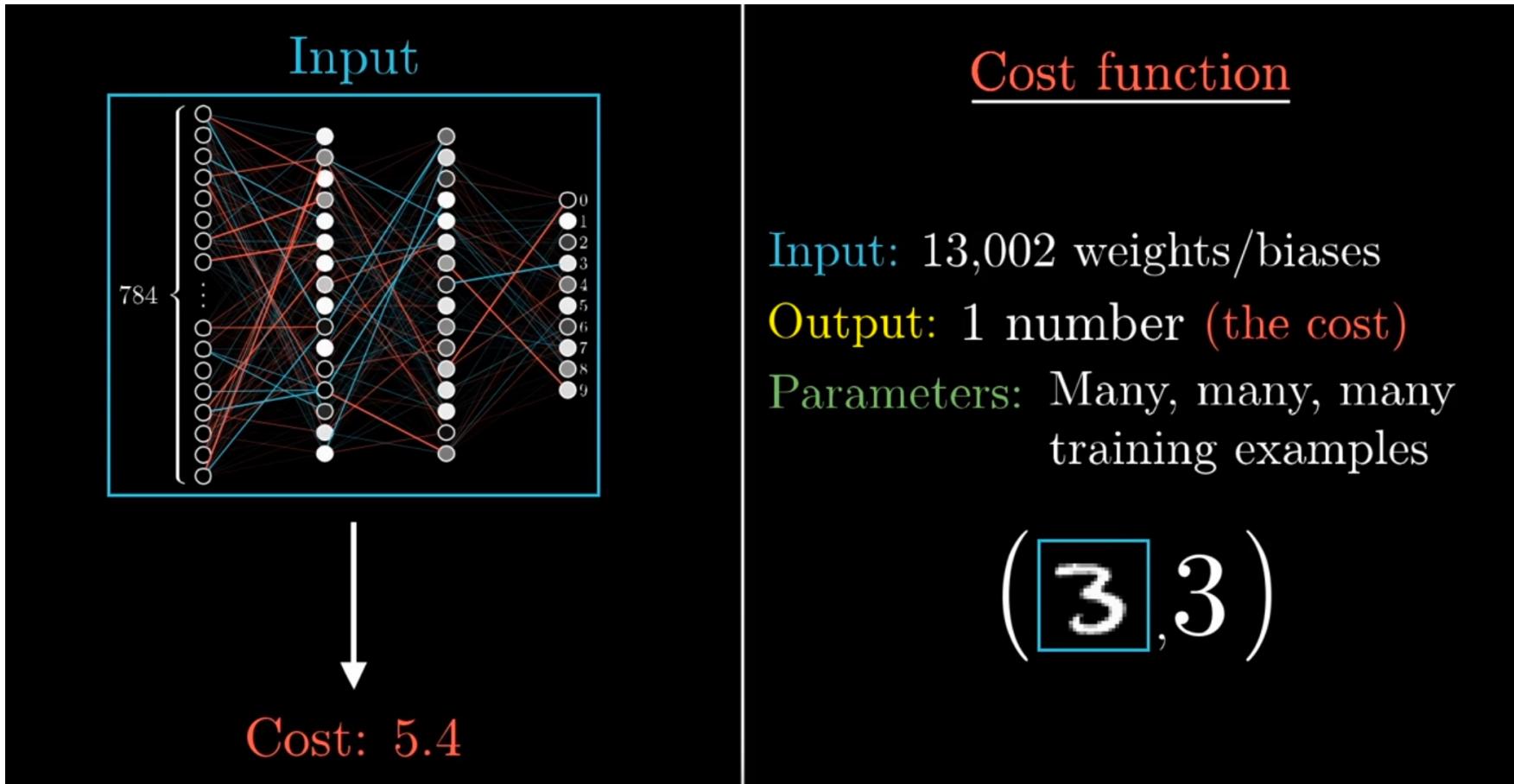
What's the “cost”
of this difference?

<input type="radio"/> 0	○ 0
<input checked="" type="radio"/> 1	○ 1
<input type="radio"/> 2	○ 2
<input checked="" type="radio"/> 3	○ 3
<input checked="" type="radio"/> 4	○ 4
<input checked="" type="radio"/> 5	○ 5
<input checked="" type="radio"/> 6	○ 6
<input checked="" type="radio"/> 7	○ 7
<input checked="" type="radio"/> 8	● 8
<input type="radio"/> 9	○ 9



Utter trash

Cost function

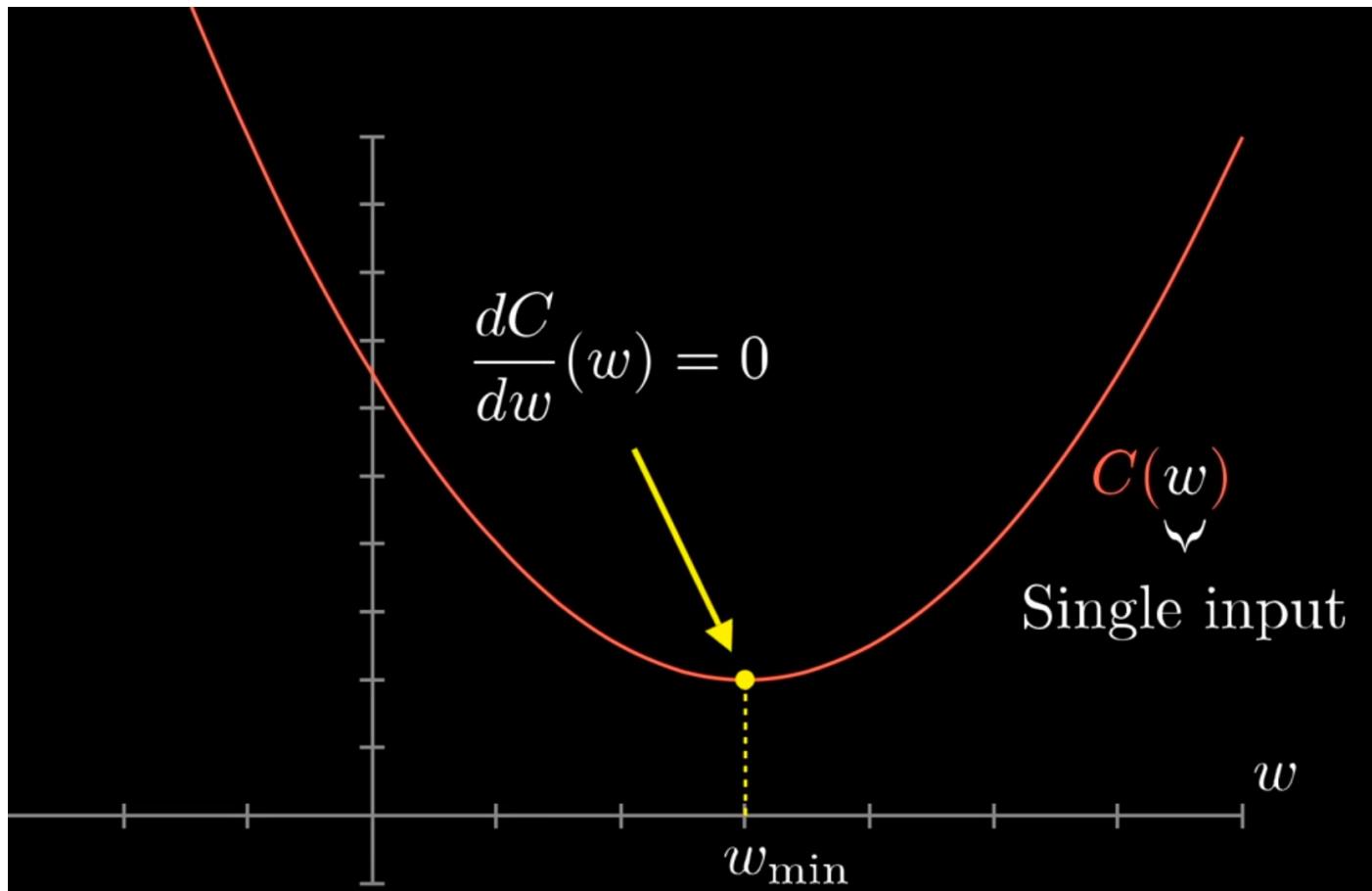


Learning

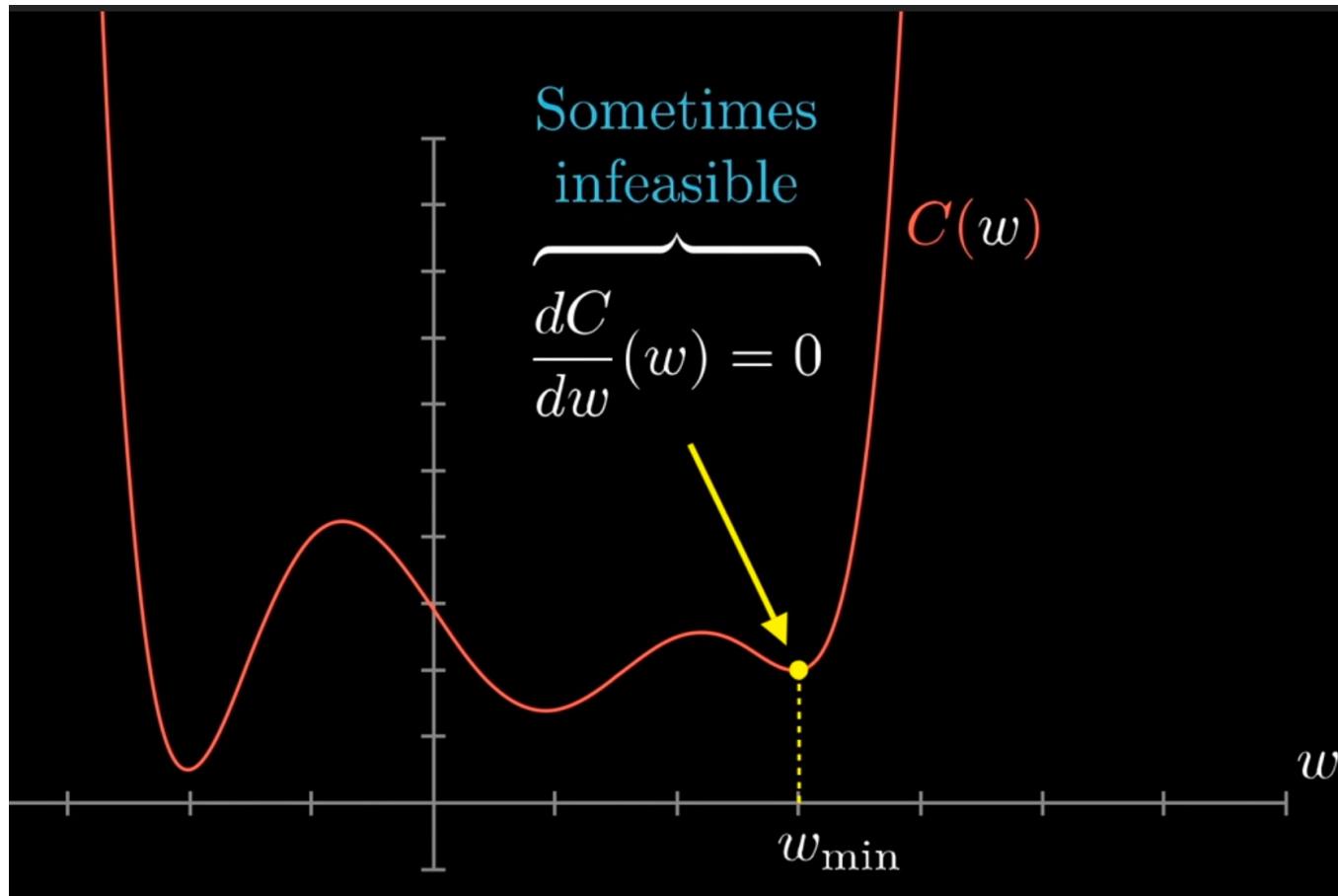


- How to change weights and biases such that cost is minimized?

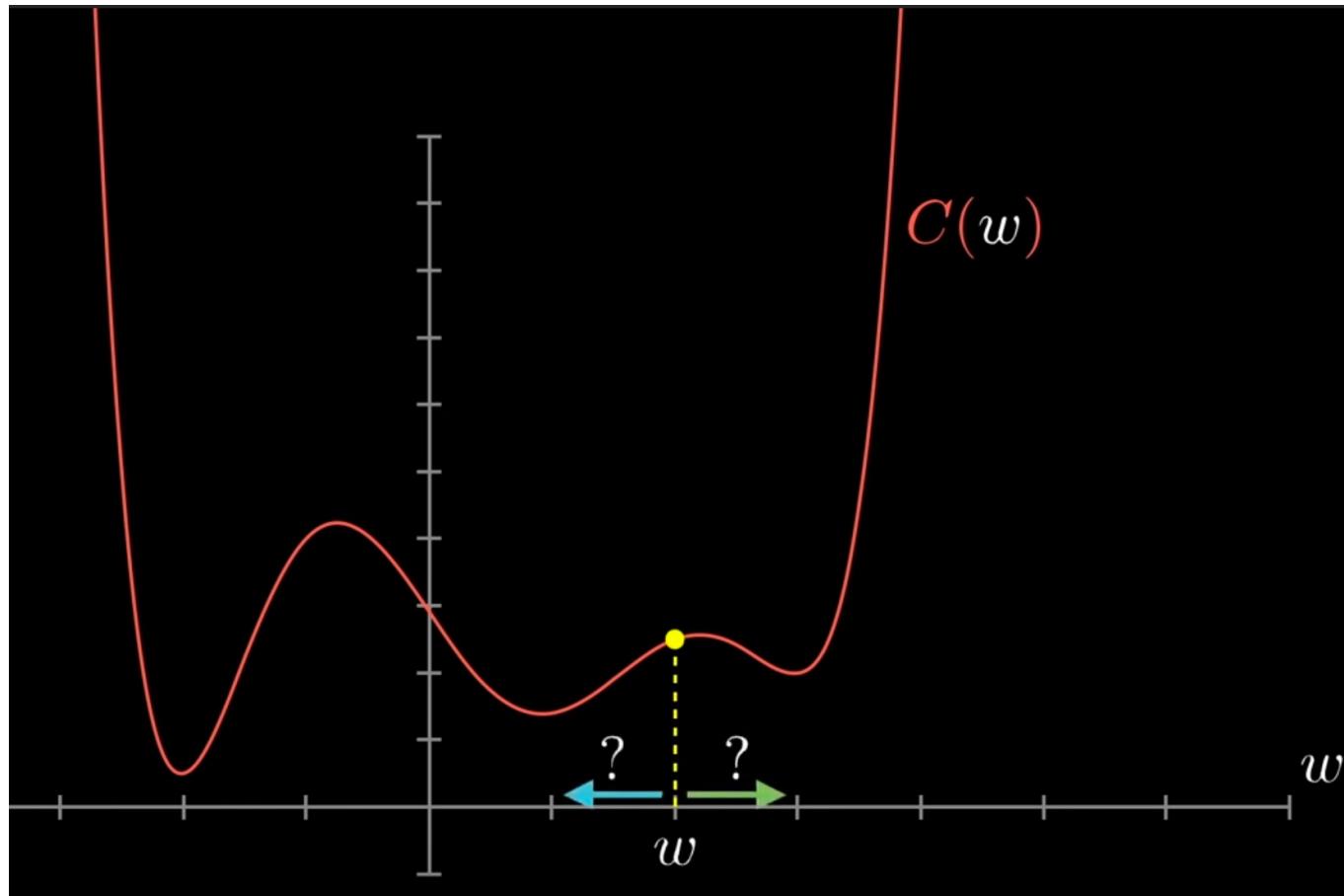
Calculating the Minimum



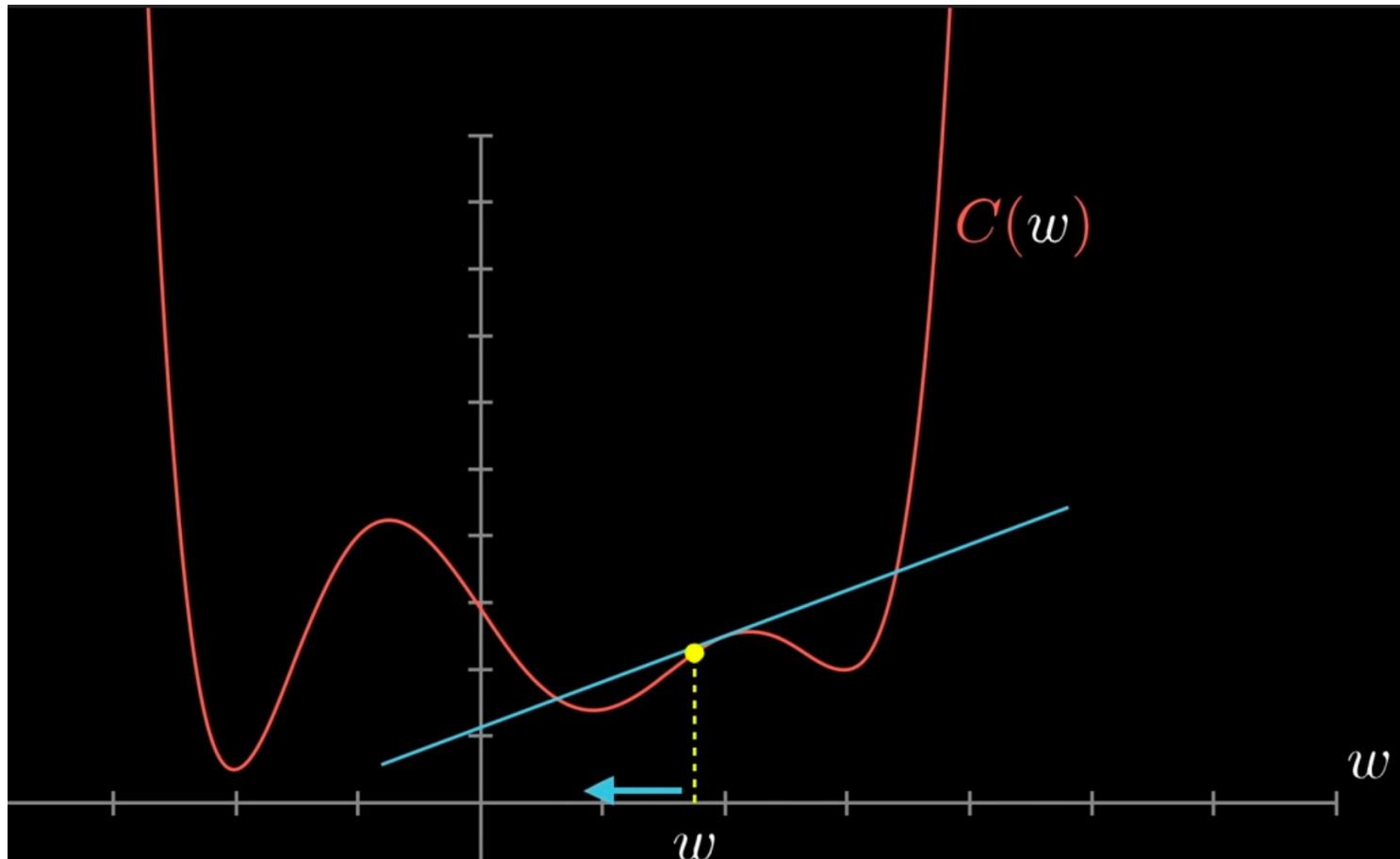
Calculating the minimum is a bit more complex now for this function



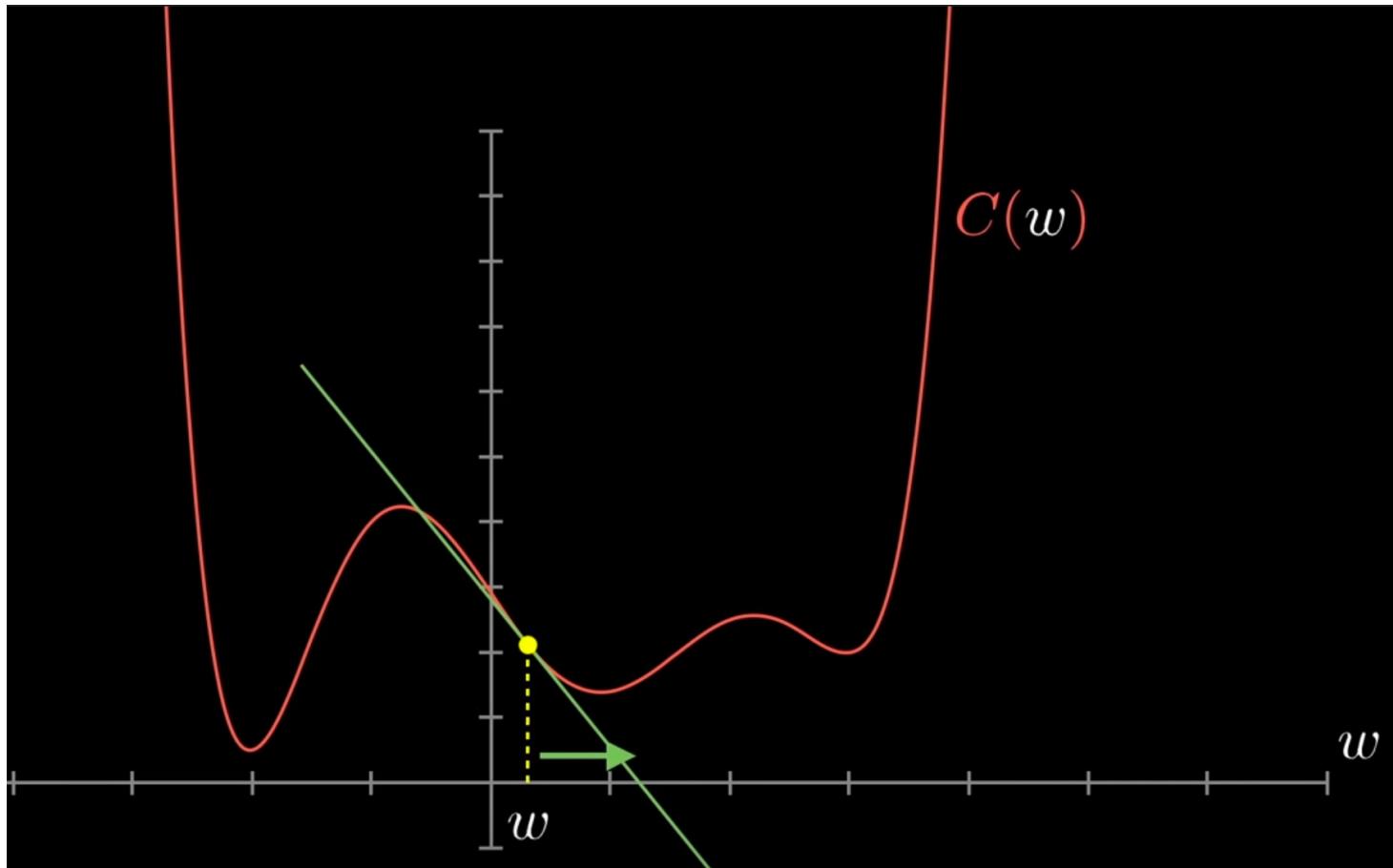
**Alternative: Start at a specific point
Check in which direction to go to make output lower**



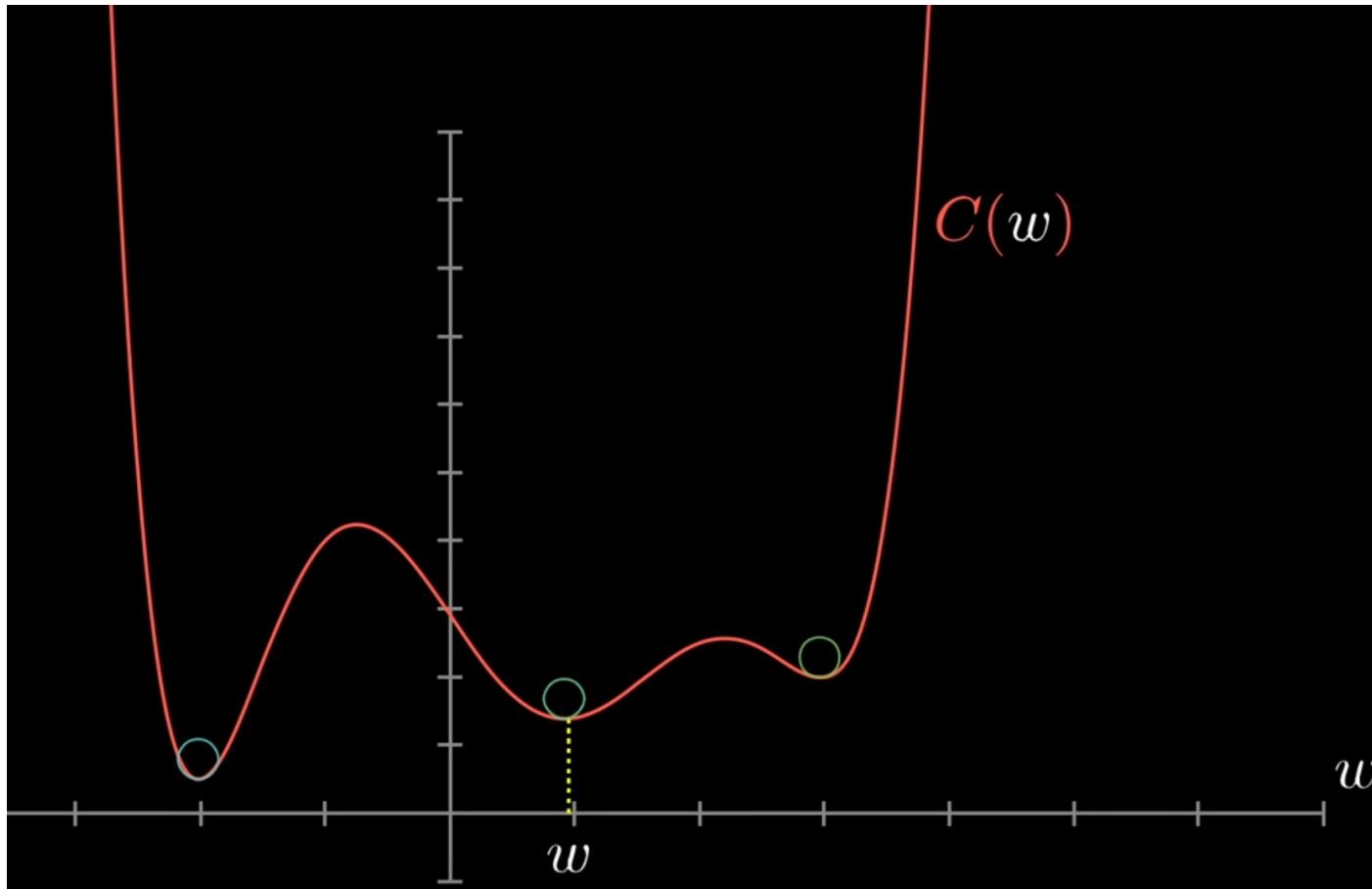
**Slope tells you in which direction to go:
If positive, go left**



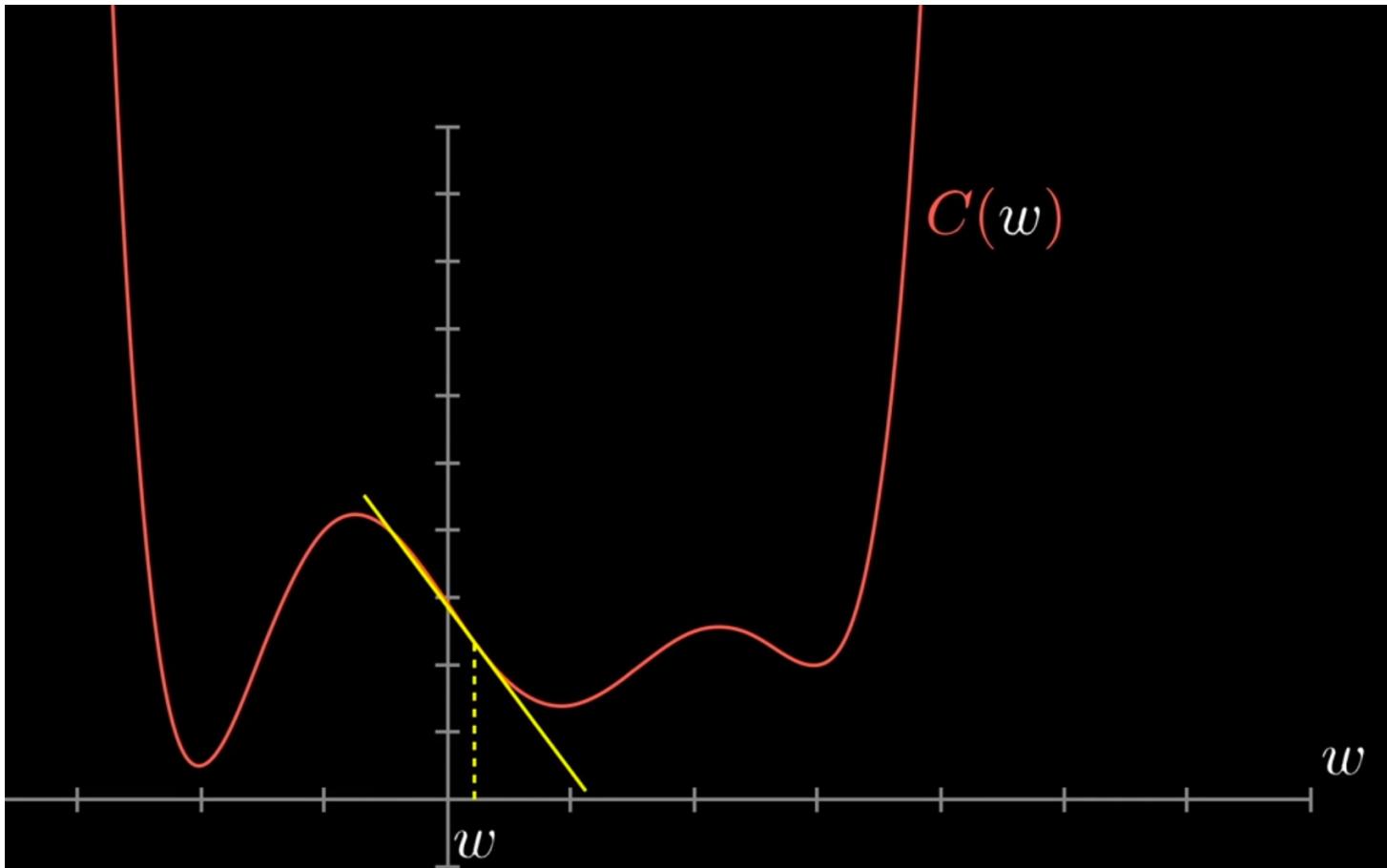
**Slope tells you in which direction to go:
If negative, go right**



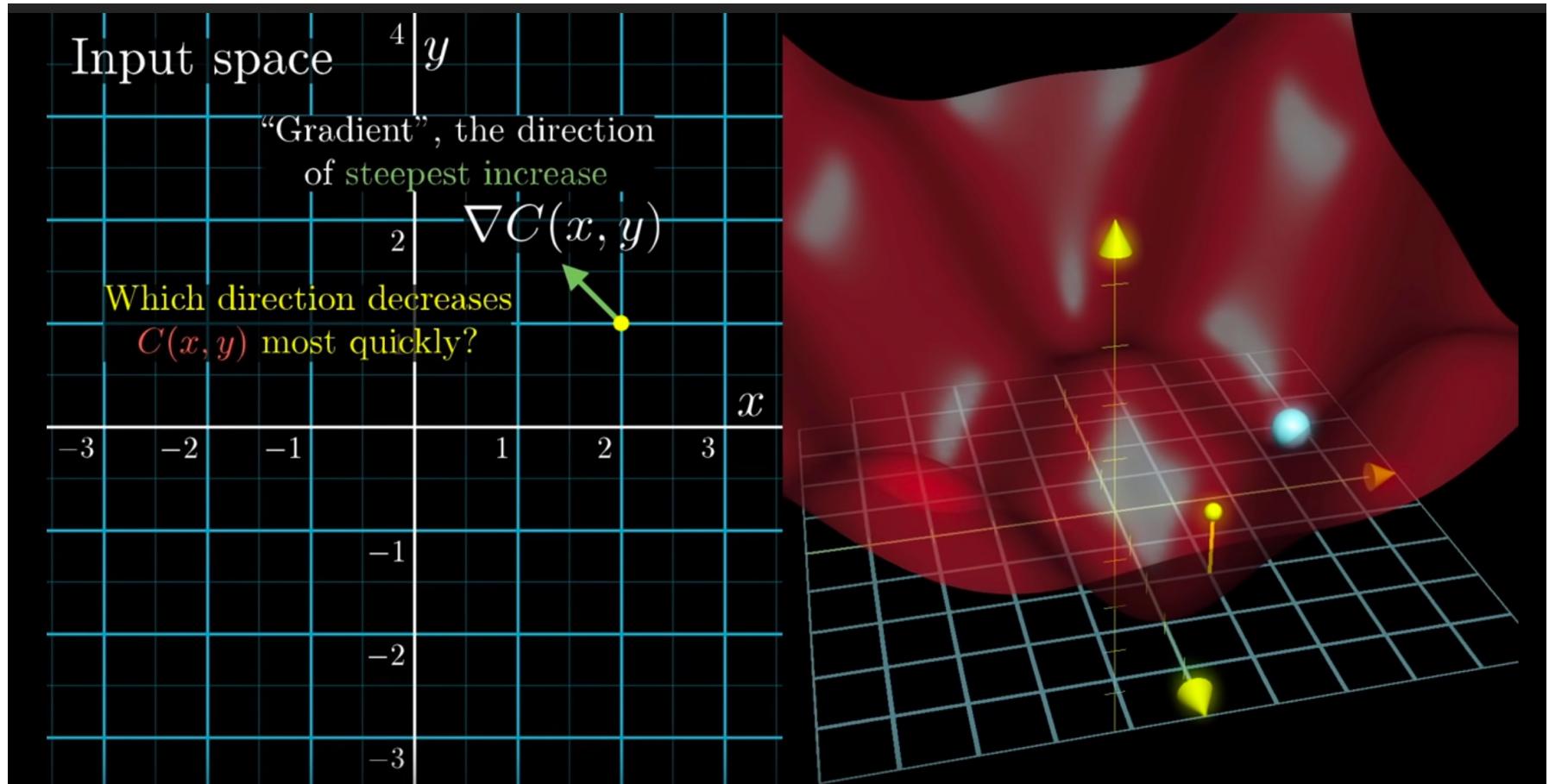
Various (local) minima possible:
No guarantee of finding global minimum



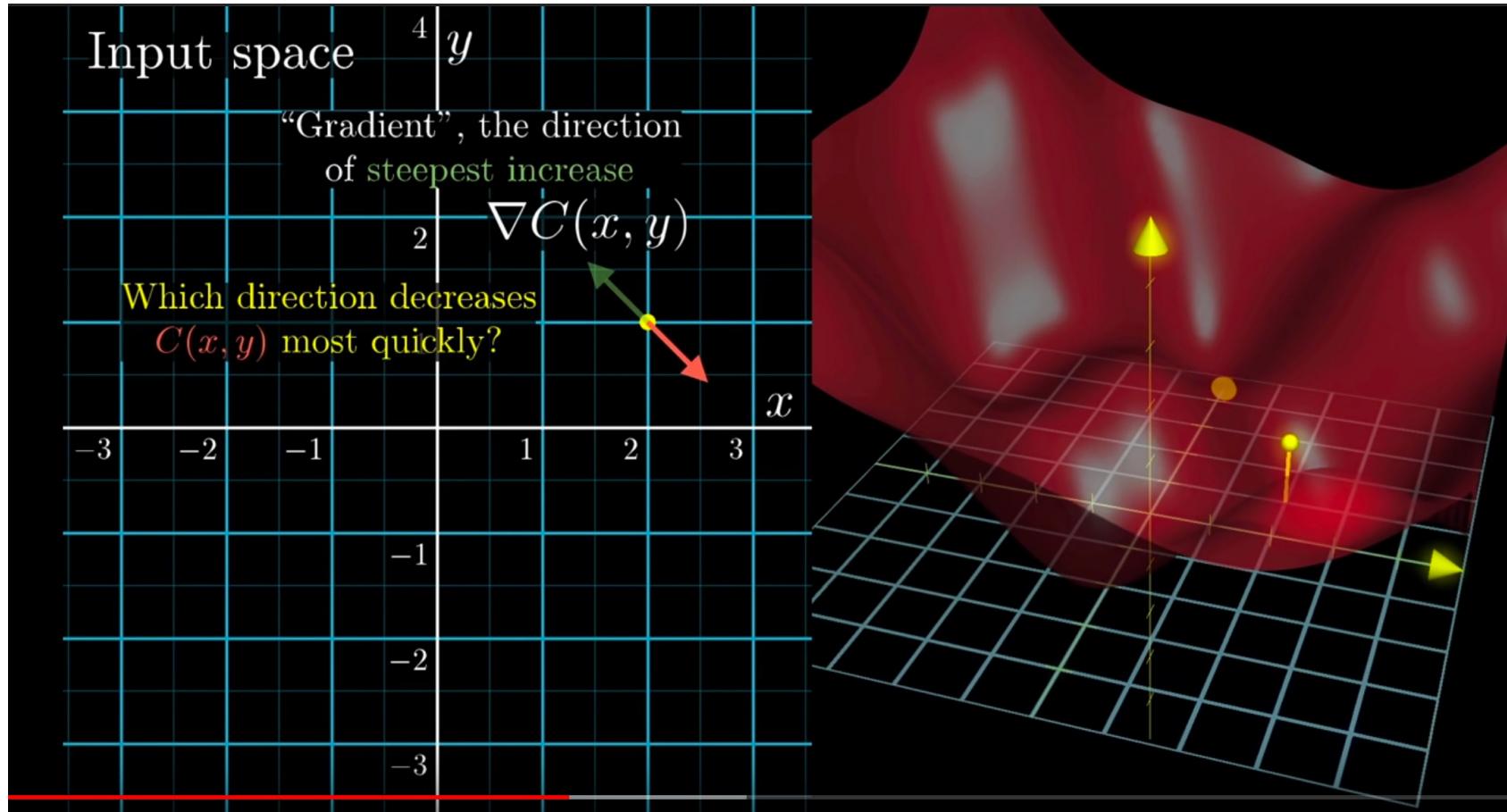
Make step size proportional to slope →
Avoids “overshooting”



Gradient of function gives direction of steepest ascent



Length of gradient vector = steepness of slope

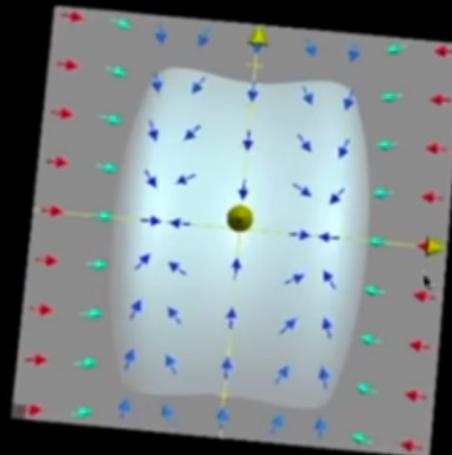


Direction of steepest descent can be calculated
But: don't worry about details

$$f(x, y) = x^2 + y^2$$

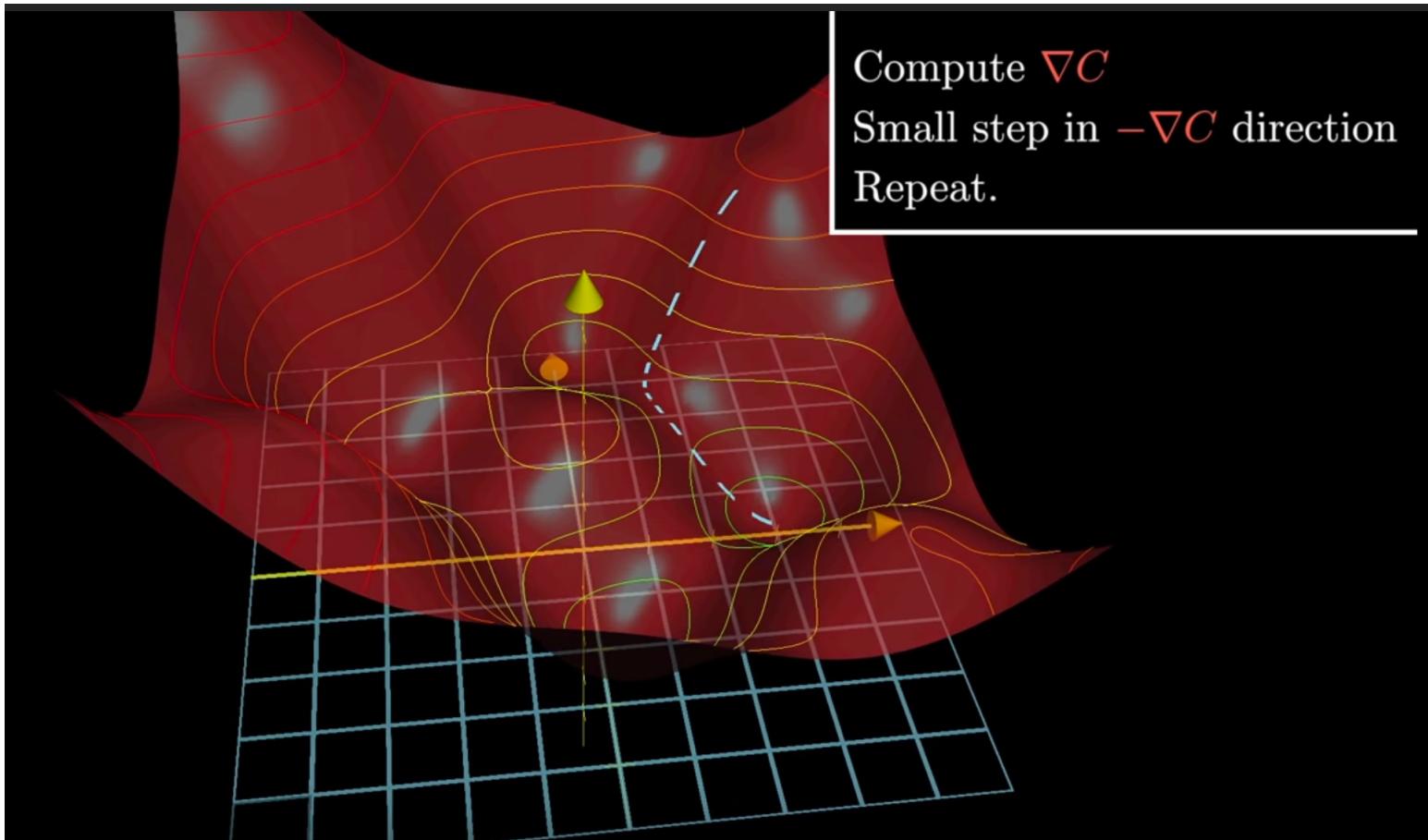
$$\nabla f(x, y) = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$$

$$= \begin{bmatrix} z_x \\ z_y \end{bmatrix}$$



direction of
steepest
ascent

Gradient descent algorithm: (1) compute cost (2) take small step (3) repeat



Negative gradient of cost function = vector of weights

13,002 weights and biases

$$\vec{W} = \begin{bmatrix} 2.43 \\ -1.12 \\ 1.47 \\ \vdots \\ -0.76 \\ 3.50 \\ 2.03 \end{bmatrix}$$

How to nudge all
weights and biases

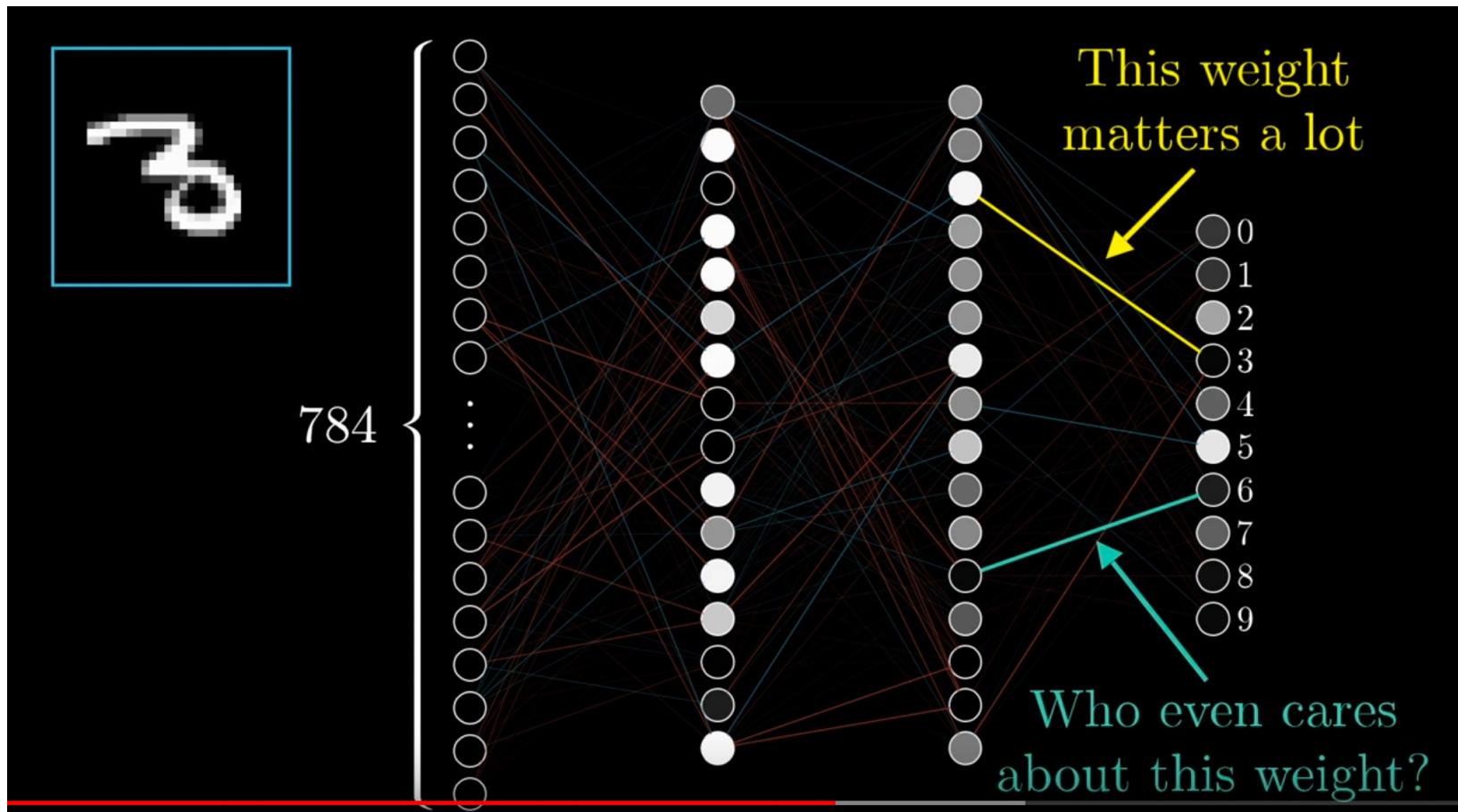
$$-\nabla C(\vec{W}) = \begin{bmatrix} 0.18 \\ 0.45 \\ -0.51 \\ \vdots \\ 0.40 \\ -0.32 \\ 0.82 \end{bmatrix}$$

Sign of gradient says if weight should be pushed up or down

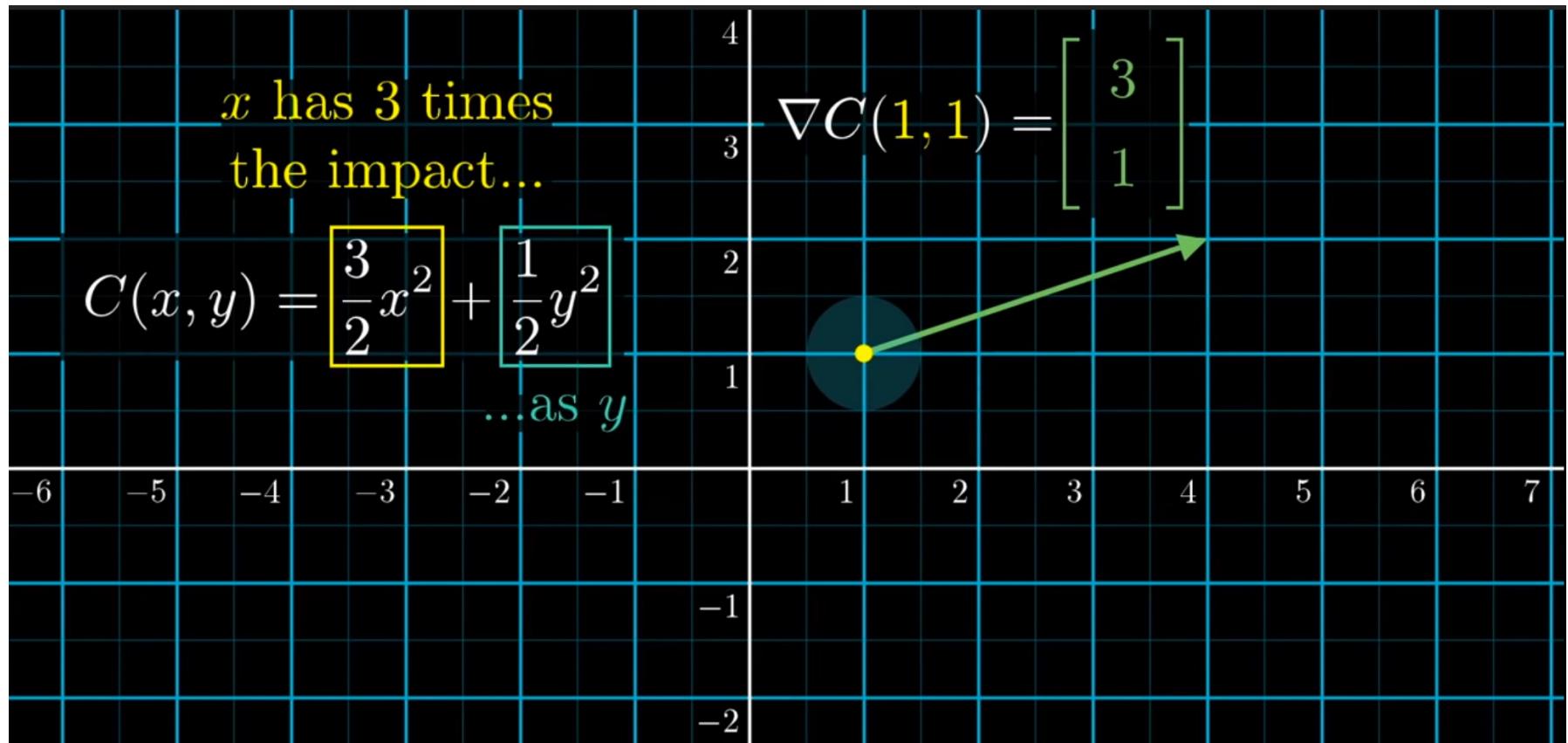
$$\vec{\mathbf{W}} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_{13,000} \\ w_{13,001} \\ w_{13,002} \end{bmatrix}$$

$$-\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -1.25 \\ \vdots \\ 0.78 \\ -0.37 \\ 0.16 \end{bmatrix} \quad \begin{array}{l} w_0 \text{ should increase} \\ w_1 \text{ should increase} \\ w_2 \text{ should decrease} \\ \\ w_{13,000} \text{ should increase} \\ w_{13,001} \text{ should decrease} \\ w_{13,002} \text{ should increase} \end{array}$$

Higher weights make some neurons “more important”



x value has 3 times more importance
than y value



Summary

- Neural network consists of input, (hidden) and output layers
- The goal of neural networks is to minimize a cost function:
 - Error between expected and observed value
- Neural networks have a large parameter space of weights and biases
 - Need to find a good parameter setting
- Gradient descent helps you find the (local) minimum
 - The slope of the gradient tells you in which direction to go
 - Step size is proportional to the size of the gradient