PAT 498/598 (Fall 2024)

# Special Topics:
# Generative AI for Music and Audio Creation

## Lecture 5: Deep Learning Fundamentals II

Instructor: Hao-Wen Dong

# **Assignment 1**: AI Song Contest

- Please listen to the **ten finalists of AI Song Contest 2024** and **read the about pages** by clicking the cover arts

- **Vote for your favorites**

- **Answer the following questions** (in 10-20 sentences each)
  - Which is your favorite song? What did they do well? What can be improved?
  - What is one dimension that most finalists didn't look into or didn't do well on?
  - What tasks are easy for current AI? What are difficult?

aisongcontest.com/the-2024-finalists
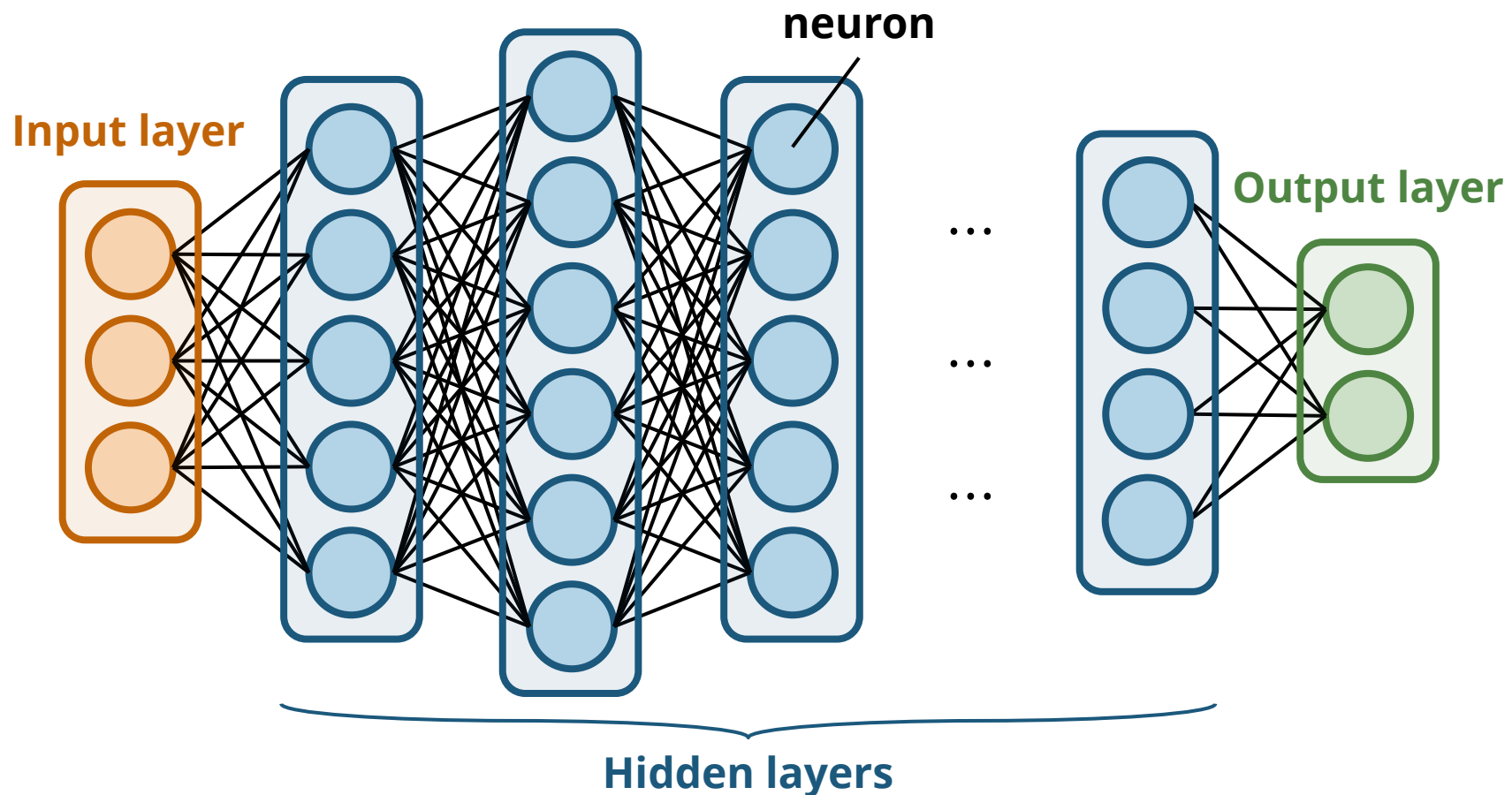
# Assignment 1: AI Song Contest

- Instructions will be released on Gradescope

- Due at **11:59pm ET** on **September 20**

- Late submissions:  **3 point deducted per day**

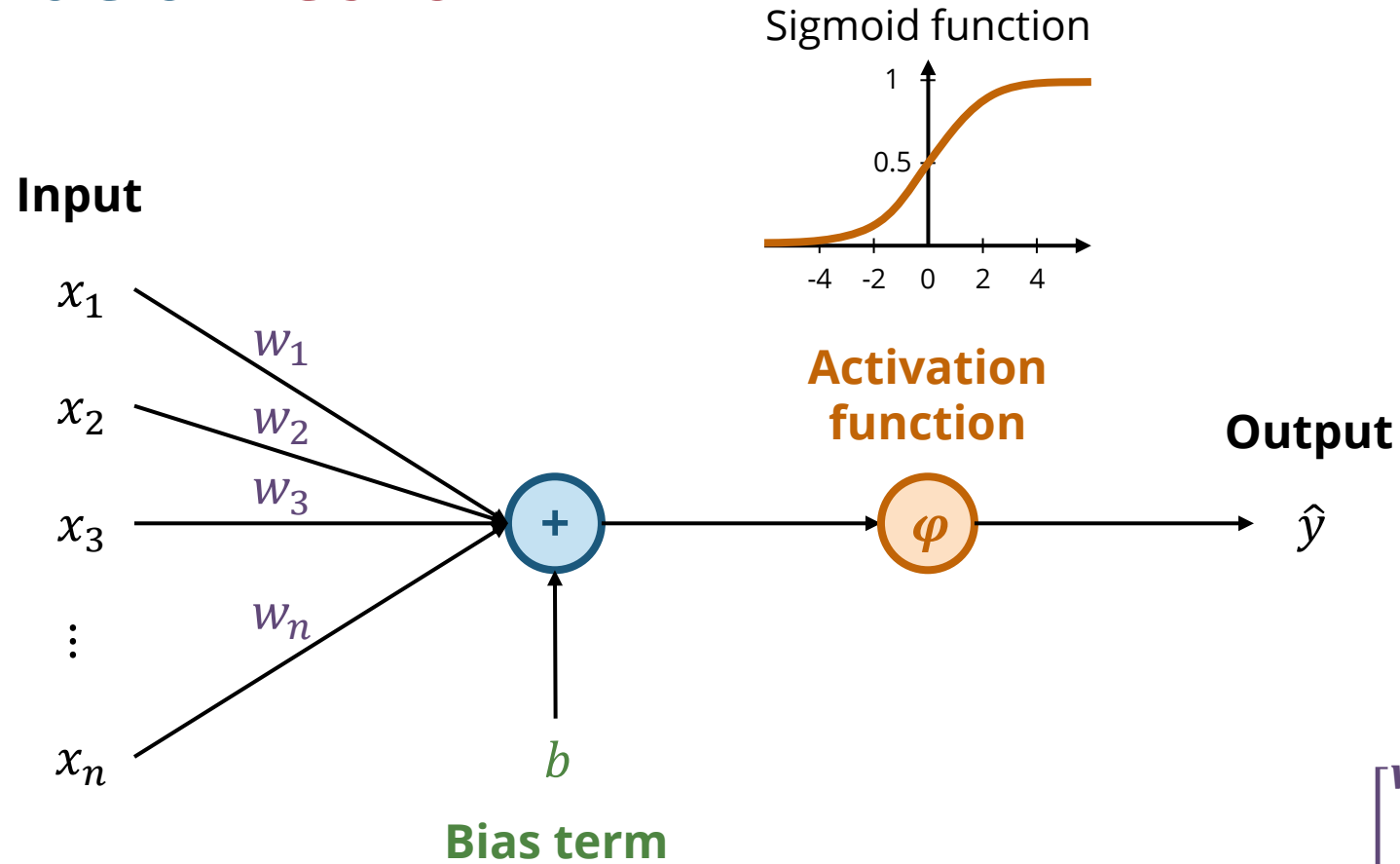[aisongcontest.com/the-2024-finalists](aisongcontest.com/the-2024-finalists)

# (Recap) What is Deep Learning?

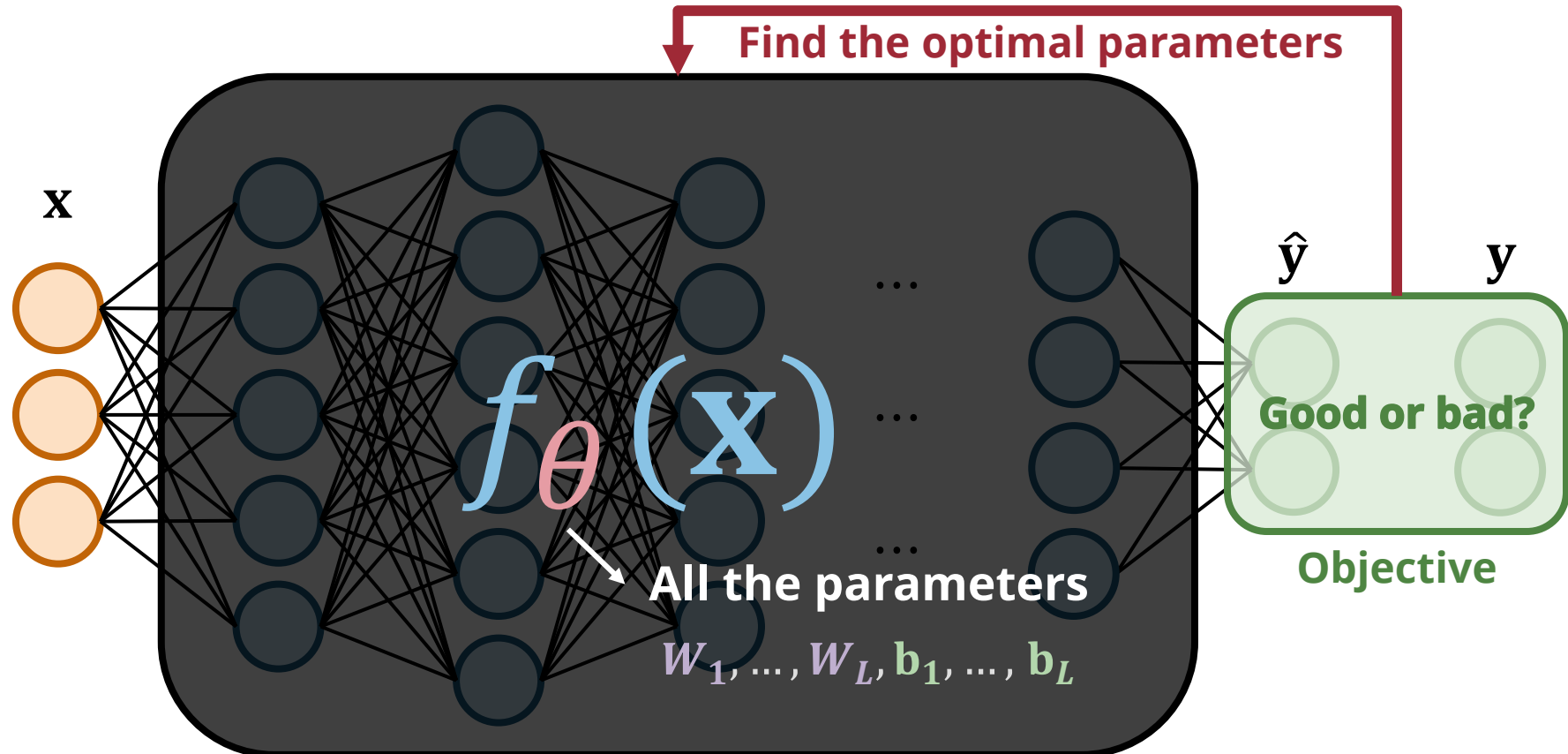- A type of machine learning that uses **deep neural networks**

# (Recap) Inside a Neuron

Sigmoid function

**Input**

$x_1$

$w_1$

$x_2$

$w_2$

$w_3$

$x_3$

$w_n$

$\vdots$

$x_n$

$+$

$b$

**Activation function**

$\varphi$

**Output**

$\hat{y}$

**Bias term**

$$\hat{y} = \varphi(w_1 x_1 + w_2 x_2 + \cdots + w_n x_n + b) = \varphi\left(\sum_{i=1}^{n} w_i x_i + b\right) = \varphi(\mathbf{w} \cdot \mathbf{x} + b)$$
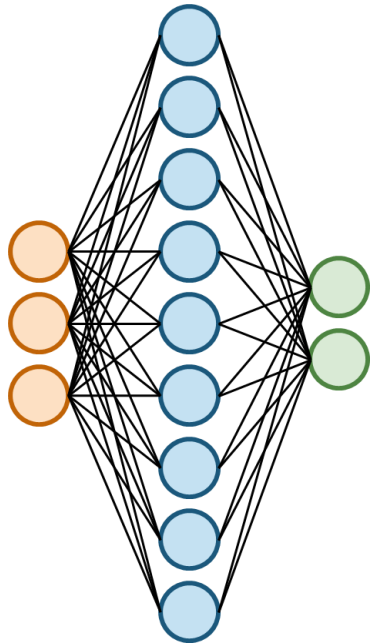
$$\begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

# (Recap) Neural Networks are Parameterized Functions

- A neural network represents **a set of functions**

**Find the optimal parameters**

$\hat{y}$      $y$

$f_\theta(\mathbf{X})$

**Good or bad?**

$\mathbf{x}$

**All the parameters**

$W_1, \dots, W_L, \mathbf{b}_1, \dots, \mathbf{b}_L$

**Objective**

# (Recap) Shallow vs Deep Neural Networks – In Practice

**Shallow neural nets**

**Deep neural nets**

**Less expressive**
(less parameter efficient)

**More expressive**
(more parameter efficient)

# Regression vs Classification
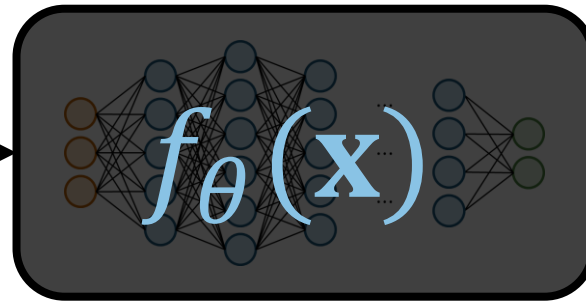
# Regression vs Classification

**Regression**



$f_\theta(\mathbf{x})$

**Age**

**5**

Output a number

**Classification**



$f_\theta(\mathbf{x})$

**Is human?**

**Yes / No**

Output a label

# Regression Example: Stock Price Prediction

$$y \in [0, \infty)$$



$$f\left( \quad \right) = 108.15$$



$$f\left( \quad \right) = 18.95$$

# Regression Example: Depth Estimation

$$\mathbf{y} \in [0, \infty)^{W \times H}$$

# Classification Example: Image Recognition

$y \in \{\text{cat}, \text{dog}, \text{bear}, \text{bird}\}$

$y \in \{0, 1, 2, \ldots, 9\}$

$f(\quad) = \text{cat}$

$f(\quad) = \text{dog}$

$f(\quad) = \text{bear}$

$f(\quad) = 8$

$f(\quad) = 6$

# Classification Example: Spam Filter

$$f \left( \begin{array}{c} \text{POWERBALL} \\ \text{CONGRATULATIONS!!} \\ \text{Your Email was selected in Powerball Lottery} \\ \text{Draw with the sum of 1.5million dollars.} \\ \text{Kindly send your Full Name, Address and} \\ \text{Phone Number for claims.} \\ \text{Yours Sincerely} \\ \text{Mr. James Hodges} \\ \text{Head Of Operations} \end{array} \right) = \text{spam}$$

$$y \in \{\text{spam}, \text{not spam}\}$$

$$f \left( \begin{array}{c} \text{Call for Panelists with Internship/work Experience for PAT Seminar @ Sep 13} \\ \text{Hi folks,} \\ \text{We are planning an internship panel for our PAT seminar this Friday...} \end{array} \right) = \text{not spam}$$

# How to Train a Neural Network?

# Training a Neural Network

**Build a neural network**
(which defines a set of functions)

**Define the objective**
(i.e., what is good for a function)

**Find the optimal parameters**
(which leads to the best function)

# Training a Neural Network



Build a neural network
(which defines a set of functions)

↓

Define the objective
(i.e., what is good for a function)

↓

Find the optimal parameters
(which leads to the best function)

# (Recap) Neural Networks are Parameterized Functions

- A neural network represents **a set of functions**



**Find the optimal parameters**

$$f_{\boldsymbol{\theta}}(\mathbf{X})$$

**All the parameters**

$$W_1, \ldots, W_L, \mathbf{b_1}, \ldots, \mathbf{b}_L$$

$\mathbf{x}$

$\hat{\mathbf{y}}$     $\mathbf{y}$

**Good or bad?**

**Objective**

$$L(\boldsymbol{\theta}) = L(\hat{\mathbf{y}}, \mathbf{y})$$

**Loss function**

# Loss Function

- Measure **how well the model perform** (in the opposite way)

- The choice of loss function depends on the task and the goals

$$L(\boldsymbol{\theta}) = L(\hat{\mathbf{y}}, \mathbf{y})$$

# Loss Function – The Many Names

- Sometimes called
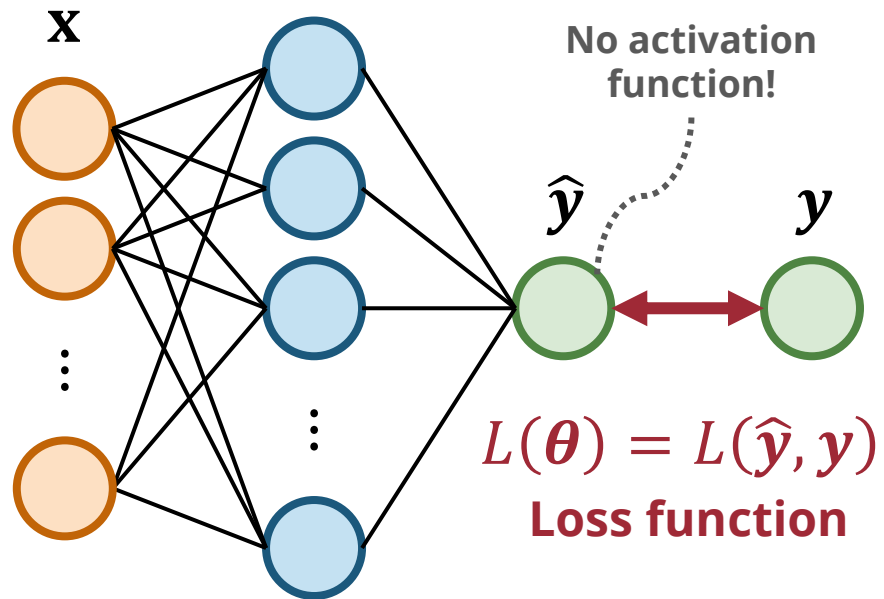    - **Cost** function
    - **Error** function

- The opposite is known as
    - **Objective** function
    - **Reward** function (reinforcement learning)
    - **Fitness** function (evolutionary algorithms & genetic algorithms)
    - **Utility** function (economics)
    - **Profit** function (economics)

# Example: Audio Codec

- What would be **a good objective to train a neural audio codec**?

- What do we care about for a codec?
  - Reconstruction quality     **Trainable**
  - Bit rate (compression rate)     **Likely not trainable but searchable**
  - Encoding/decoding speed     **Likely not trainable but searchable**

- How do we measure reconstruction quality?
  - Difference in raw waveforms?
  - Difference in spectrograms?
  - Perceptual quality (psychoacoustics)?

# Common Loss Functions for Regression

**x**

No activation function!

$\widehat{y}$     $y$

$L(\boldsymbol{\theta}) = L(\widehat{\boldsymbol{y}}, \boldsymbol{y})$

**Loss function**

**Why not** $L(\widehat{y}, y) = \widehat{y} - y$**?**

$$L(\widehat{y}, y) = |\widehat{y} - y|$$

**L1 loss**

$$L(\widehat{y}, y) = (\widehat{y} - y)^2$$

**L2 loss**

# L1 vs L2 Losses

**L1 loss**

$$L(\hat{y}, y) = |\hat{y} - y|$$

$$L(\hat{\mathbf{y}}, \mathbf{y}) = \mathbf{MAE}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^{n} |\hat{y}_i - y_i|$$

**Mean Absolute Error (MAE)**
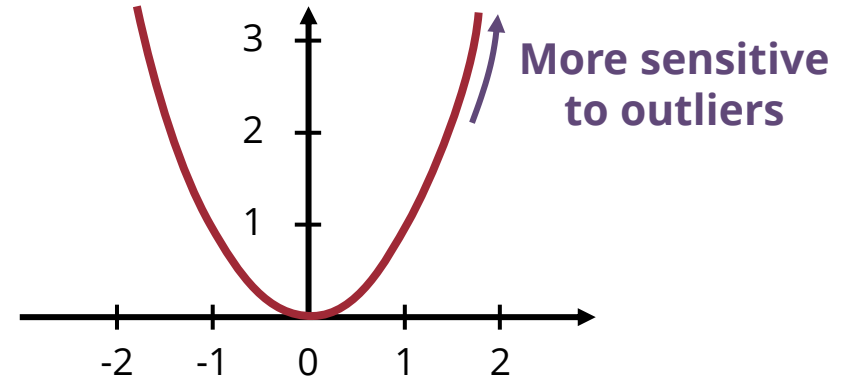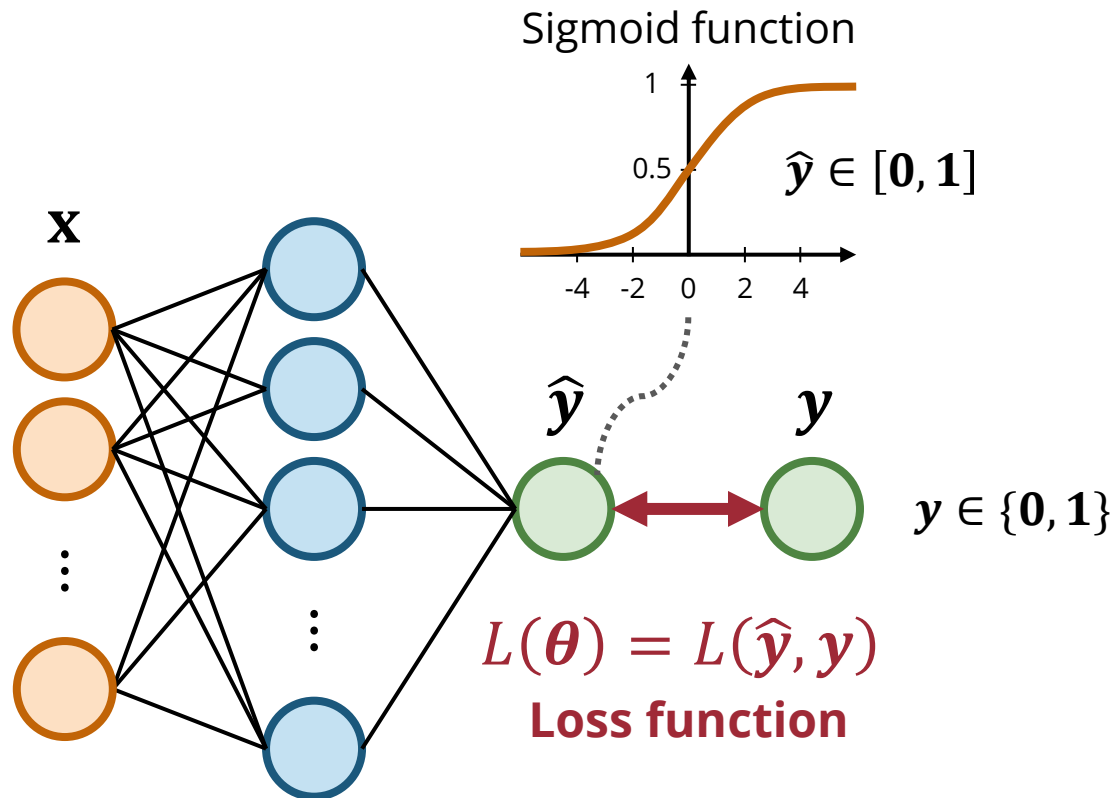
**L2 loss**

More sensitive to outliers

$$L(\hat{y}, y) = (\hat{y} - y)^2$$

$$L(\hat{\mathbf{y}}, \mathbf{y}) = \mathbf{MSE}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2$$

**Mean Squared Error (MSE)**

# Binary Cross Entropy for Binary Classification

- **Logistic regression** approaches classification like regression



Sigmoid function

$\hat{y} \in [0, 1]$

$\hat{y}$     $y$

$y \in \{0, 1\}$

$L(\boldsymbol{\theta}) = L(\hat{y}, y)$

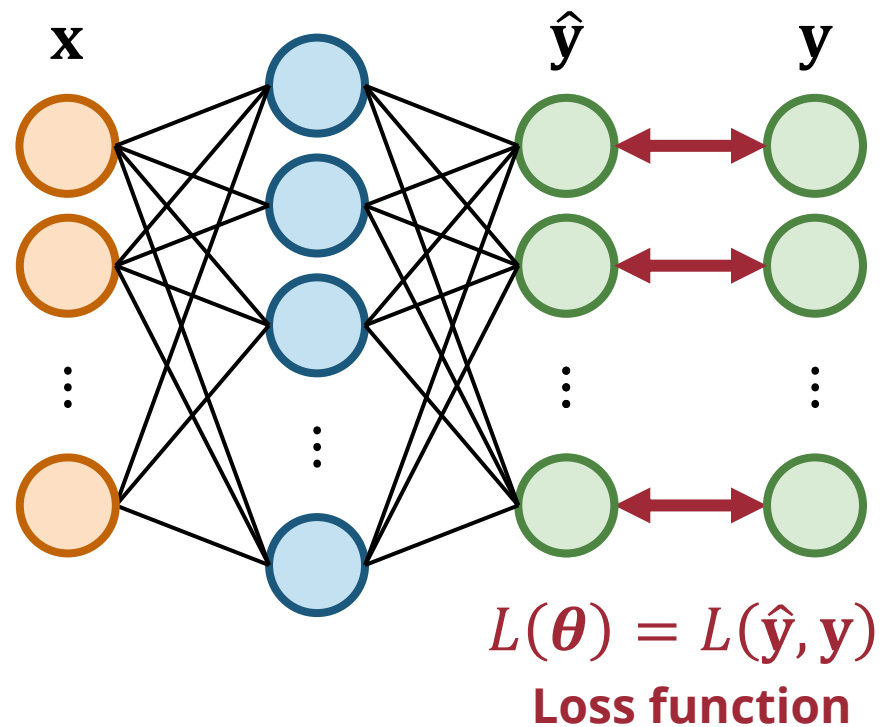**Loss function**

**Binary cross entropy**

**(Also called log loss)**

$$L(\hat{y}, y) = \begin{cases} -\log \hat{y}, & \text{if } y = 1 \\ -\log(1 - \hat{y}), & \text{if } y = 0 \end{cases}$$

$$= -y \log \hat{y} + (1 - y) \log(1 - \hat{y})$$

if $y = 1$     if $y = 0$

# Cross Entropy for Multiclass Classification



$$L(\boldsymbol{\theta}) = L(\hat{\mathbf{y}}, \mathbf{y})$$

**Loss function**

# Cross Entropy for Multiclass Classification



Real-valued numbers to probability-like numbers

$\widetilde{y}_i \in \mathbb{R}$　　$\widehat{y}_i \in [0, 1]$　$y_i \in \{0, 1\}$

$\mathbf{x}$　　$\widetilde{\mathbf{y}}$　　$\widehat{\mathbf{y}}$　$\mathbf{y}$

Softmax

$L(\boldsymbol{\theta}) = L(\widehat{\mathbf{y}}, \mathbf{y})$

Loss function

**Softmax**

$$\widehat{y}_i = \frac{e^{\widetilde{y}_i}}{\sum_{j=1}^{n} e^{\widetilde{y}_j}}$$

# Softmax

- **Intuition**: Map several numbers to $[0, 1]$ while **keeping their relative magnitude**
  - Softmax is like the **multivariate version of sigmoid**
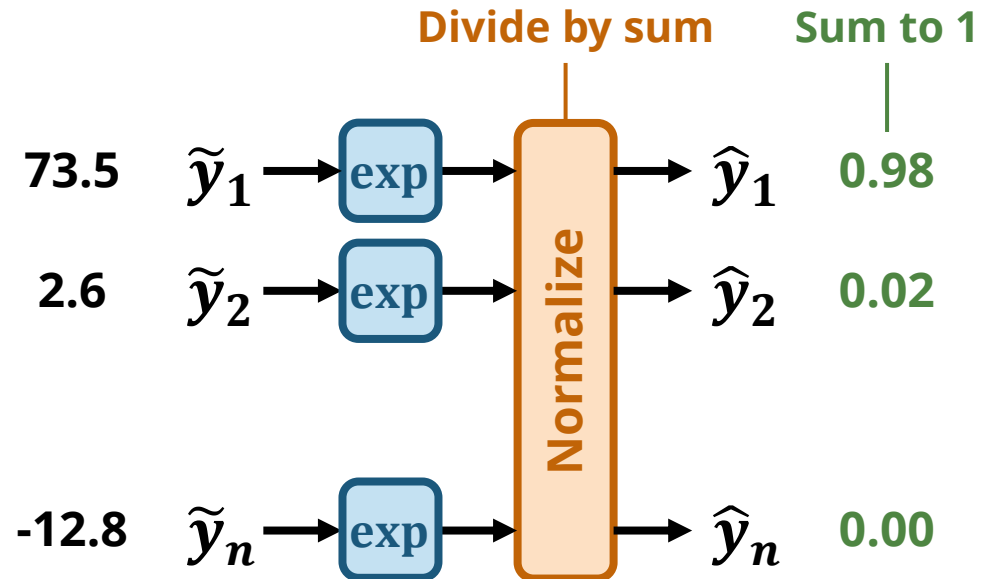
# Cross Entropy for Multiclass Classification

**Binary Cross Entropy**

**Only one of them will be one!**

$$L(\hat{y}, y) = -y \log \hat{y} + (1 - y) \log(1 - \hat{y})$$

**Cross Entropy**

**Only one of them will be one!**

$$L(\hat{\mathbf{y}}, \mathbf{y}) = -y_1 \log \hat{y}_1 - y_2 \log \hat{y}_2 - \cdots - y_i \log \hat{y}_n$$

$$= -\sum_i^n y_i \log \hat{y}_i$$

**Log likelihood**

# Why *not* MSE Loss for Classification?

- **Minimizing cross-entropy** is equivalent to **maximizing likelihood**!

- However, no one prohibits you from using an MSE loss on Softmax output
    - In fact, it will still train the model


- While loss functions can have the same global minima, they might have led to different **training dynamics** and **weights for different types of errors**
    - For example, MSE is more sensitive to MAE due to the quadratic term even though they have the same global minima

# Loss Functions vs Output Space

- Oftentimes, we change the **output space** instead of the loss function

- For example,

  $\mathcal{F}$: spectrogram → mel spectrogram

  - MSE on spectrograms → MSE on mel spectrograms
  - MSE of magnitude in raw values → MSE of magnitude in dB

  $\mathcal{F}$: raw value → db

- What's the difference?

|  | Model | Loss |
|---|---|---|
| Setup A | $f: x \to y$ | $L(\mathcal{F}(y), \mathcal{F}(\hat{y}))$ |
| Setup B | $f: x \to \mathcal{F}(y)$ | $L(y, \hat{y})$ |

# Optimization

# Training a Neural Network



**Build a neural network**
(which defines a set of functions)

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x})$$

**Define the objective**
(i.e., what is good for a function)

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

**Find the optimal parameters**
(which leads to the best function)

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

# Training a Neural Network

```
┌─────────────────────────────────────┐
│        Build a neural network        │          $\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x})$
│    (which defines a set of functions) │
└─────────────────────────────────────┘
                   │
                   ▼
┌─────────────────────────────────────┐
│        Define the objective          │          $L(\boldsymbol{\theta})$
│   (i.e., what is good for a function) │
└─────────────────────────────────────┘
                   │
                   ▼
┌─────────────────────────────────────┐
│      Find the optimal parameters      │      $\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$
│   (which leads to the best function)  │
└─────────────────────────────────────┘
```
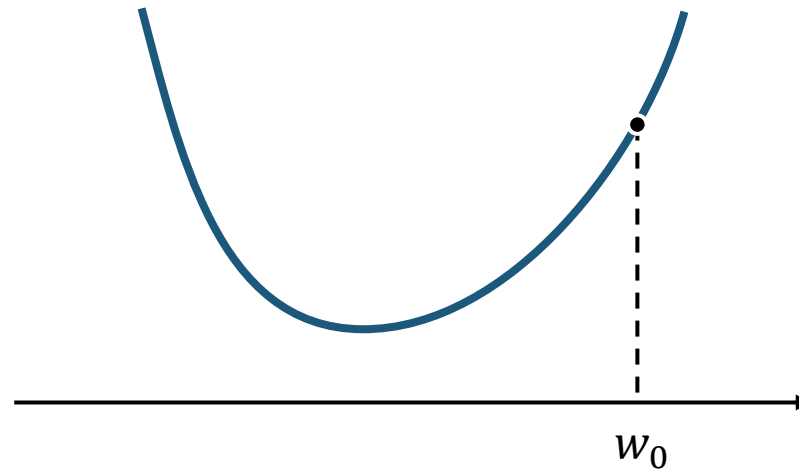
# Optimizing the Parameters of a Neural Network

- Many, many ways…

- Most commonly through **gradient descent** in deep learning

- Alternatively, we can use search or genetic algorithm

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$
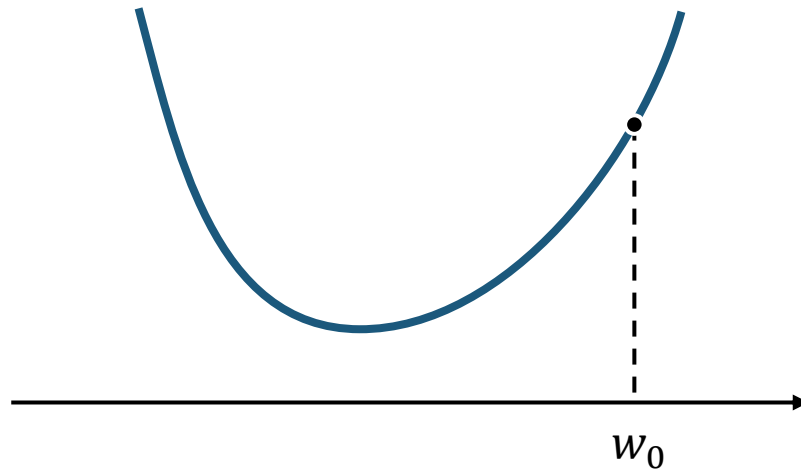
# Gradient Descent

- **Intuition**: Gradient can suggest a good direction to tune the parameters
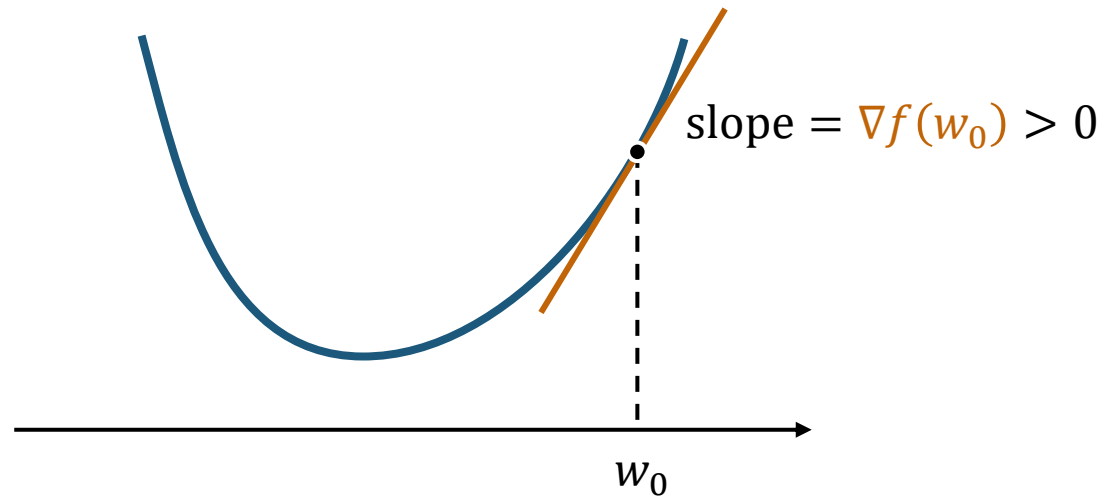
Derivative for a vector, matrix or tensor

# Gradient Descent – Pseudocode

- Pick an **initial weight vector** $w_0$ and **learning rate** $\eta$

- Repeat until convergence: $w_{t+1} = w_t - \eta \boxed{\nabla f(w_t)}$    **Gradient of function $f$ with respect to weight $w$**

# Gradient Descent – Pseudocode

- Pick an initial weight vector $w_0$ and learning rate $\eta$

- Repeat until convergence: $w_{t+1} = w_t - \eta \nabla f(w_t)$

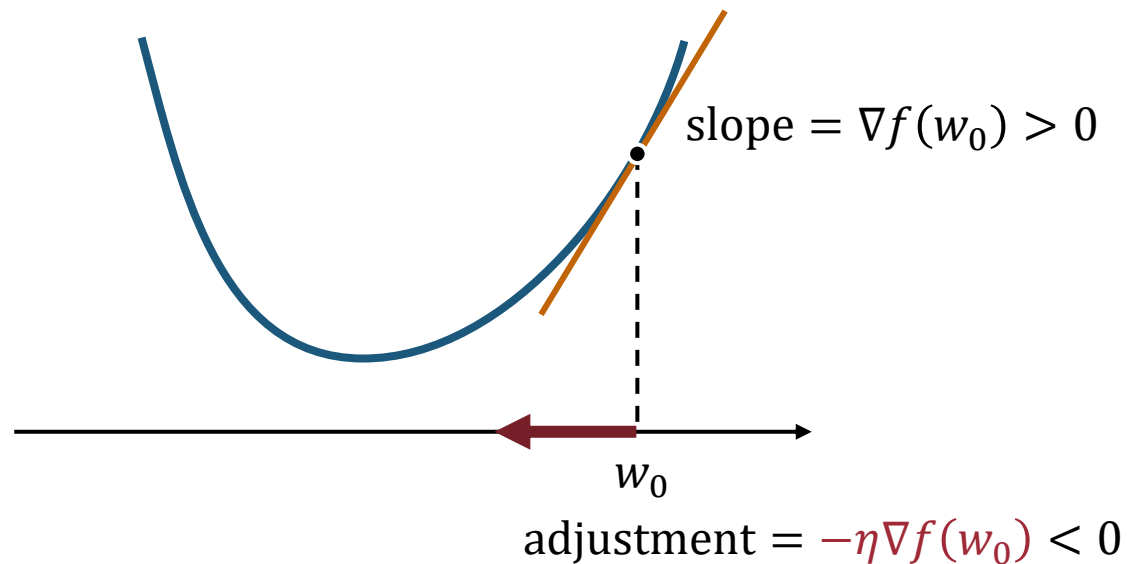slope $= \nabla f(w_0) > 0$

$w_0$

# Gradient Descent – Pseudocode

- Pick an initial weight vector $w_0$ and learning rate $\eta$

- Repeat until convergence:  $w_{t+1} = w_t - \eta \nabla f(w_t)$

slope $= \nabla f(w_0) > 0$

$w_0$

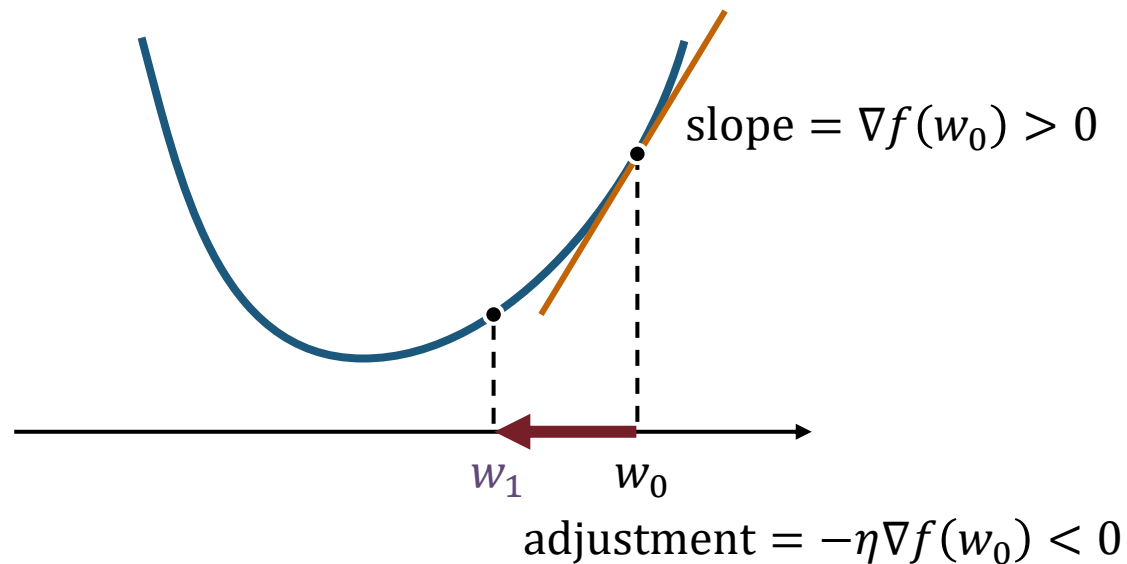adjustment $= -\eta \nabla f(w_0) < 0$

# Gradient Descent – Pseudocode

- Pick an initial weight vector $w_0$ and learning rate $\eta$

- Repeat until convergence: $w_{t+1} = w_t - \eta \nabla f(w_t)$

slope $= \nabla f(w_0) > 0$

$w_1$     $w_0$
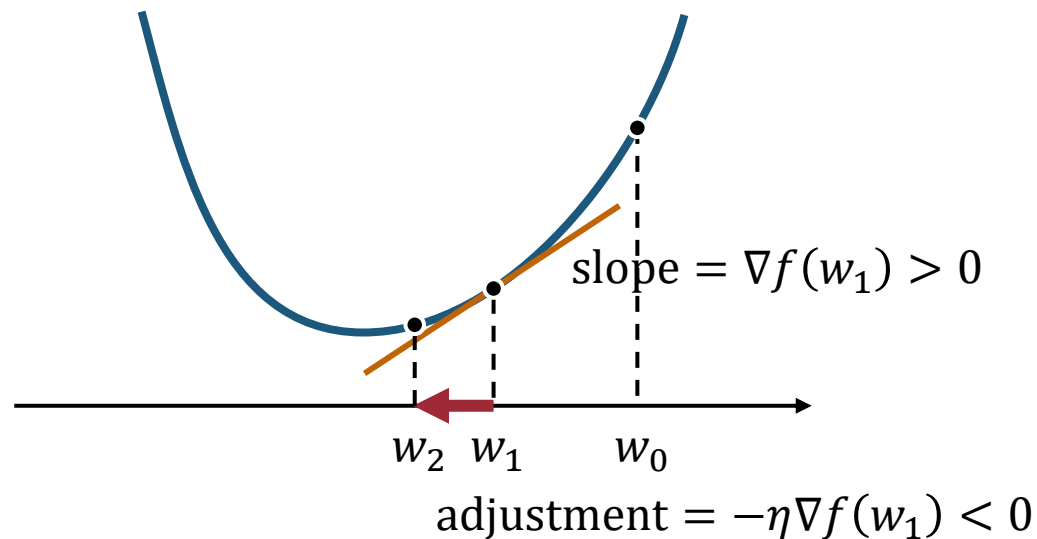
adjustment $= -\eta \nabla f(w_0) < 0$

# Gradient Descent – Pseudocode

- Pick an initial weight vector $w_0$ and learning rate $\eta$

- Repeat until convergence:  $w_{t+1} = w_t - \eta \nabla f(w_t)$



$$\text{slope} = \nabla f(w_1) > 0$$

$w_2 \quad w_1 \qquad w_0$
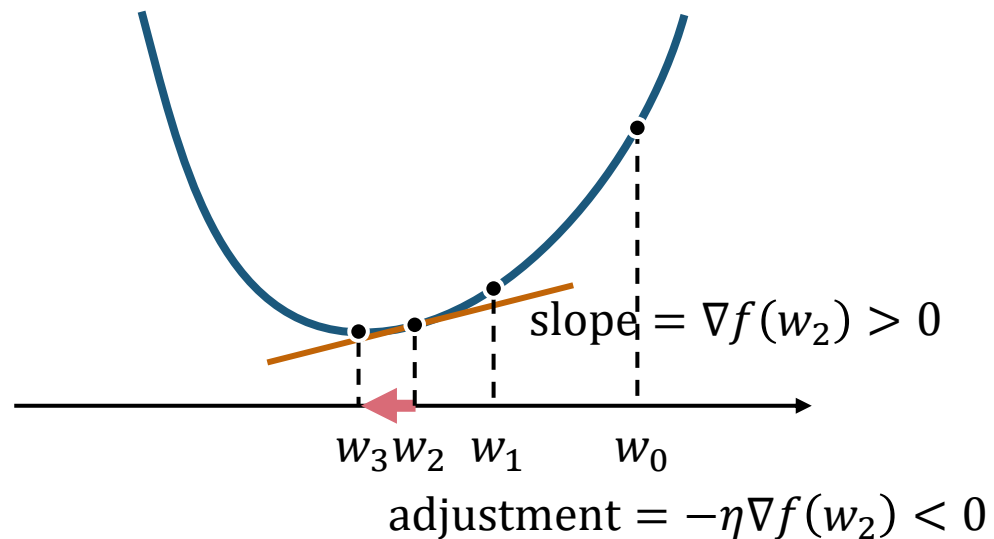
$$\text{adjustment} = -\eta \nabla f(w_1) < 0$$

# Gradient Descent – Pseudocode

- Pick an initial weight vector $w_0$ and learning rate $\eta$

- Repeat until convergence: $w_{t+1} = w_t - \eta \nabla f(w_t)$

slope $= \nabla f(w_2) > 0$

$w_3 w_2 \quad w_1 \qquad w_0$

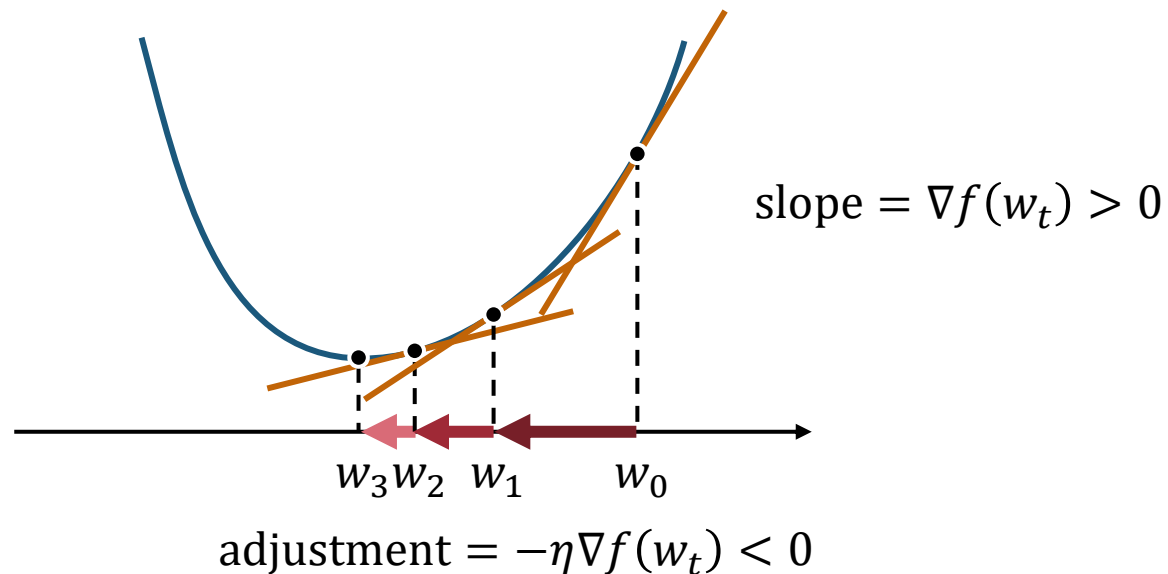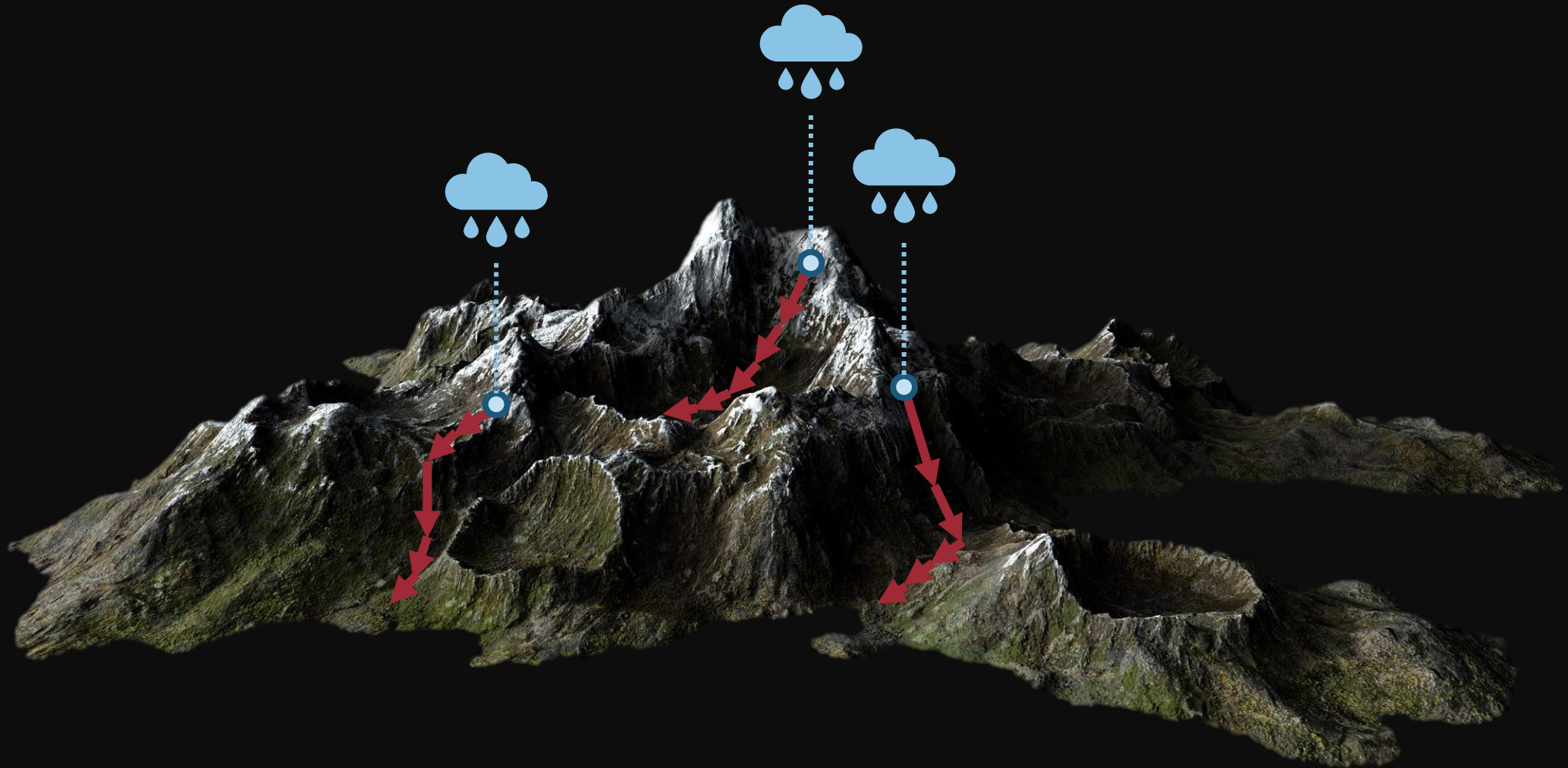adjustment $= -\eta \nabla f(w_2) < 0$

# Gradient Descent – Pseudocode

- Pick an initial weight vector $w_0$ and learning rate $\eta$

- Repeat until convergence: $w_{t+1} = w_t - \eta \nabla f(w_t)$

slope $= \nabla f(w_t) > 0$

$w_3 w_2 \quad w_1 \qquad w_0$

adjustment $= -\eta \nabla f(w_t) < 0$
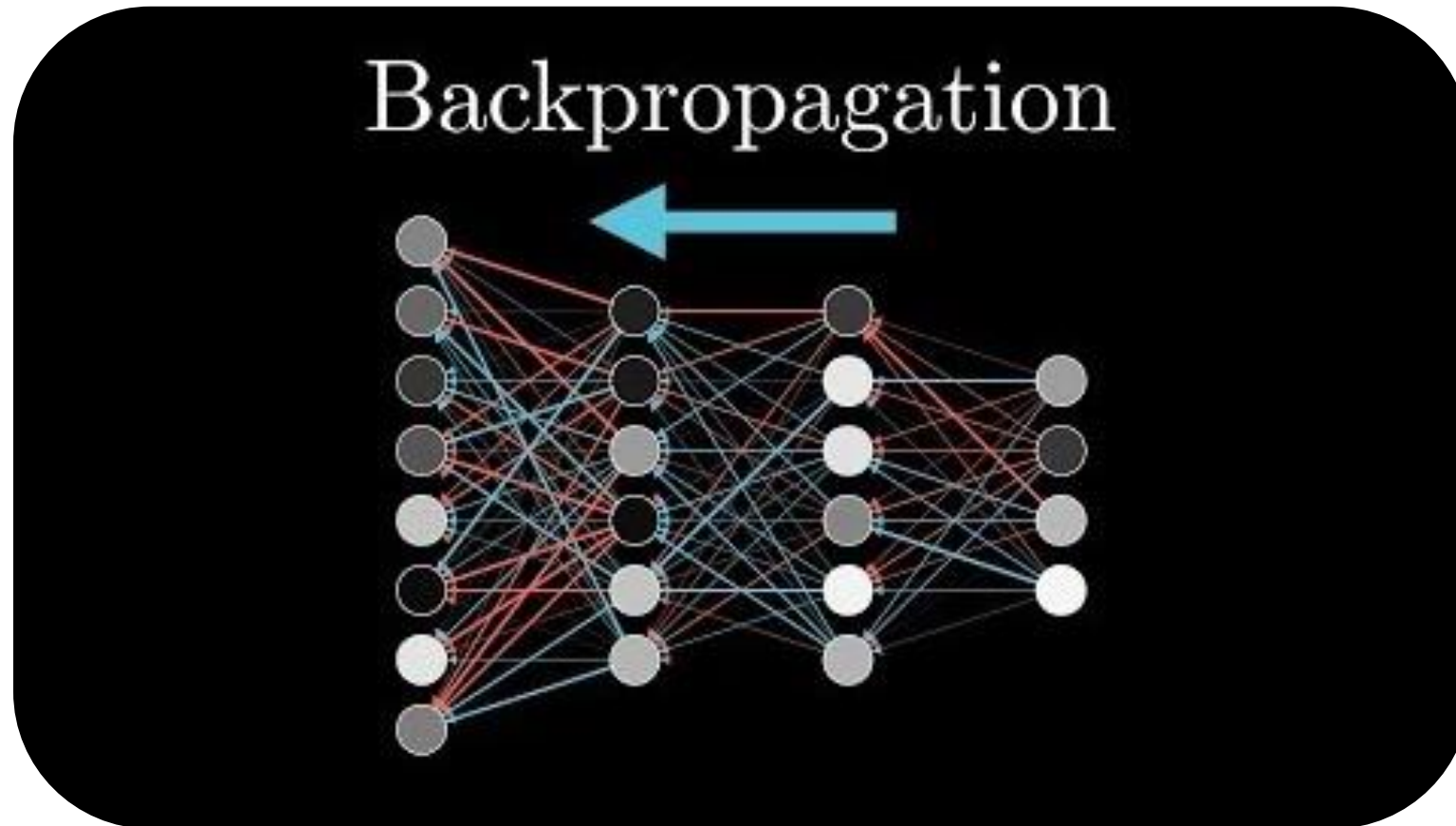
# Gradient Descent – 3D Case

# Backpropagation: Efficiently Computing the Gradients

- An efficient way of **computing gradients** using chain rule

- The reason why we want **everything to be differentiable** in deep learning
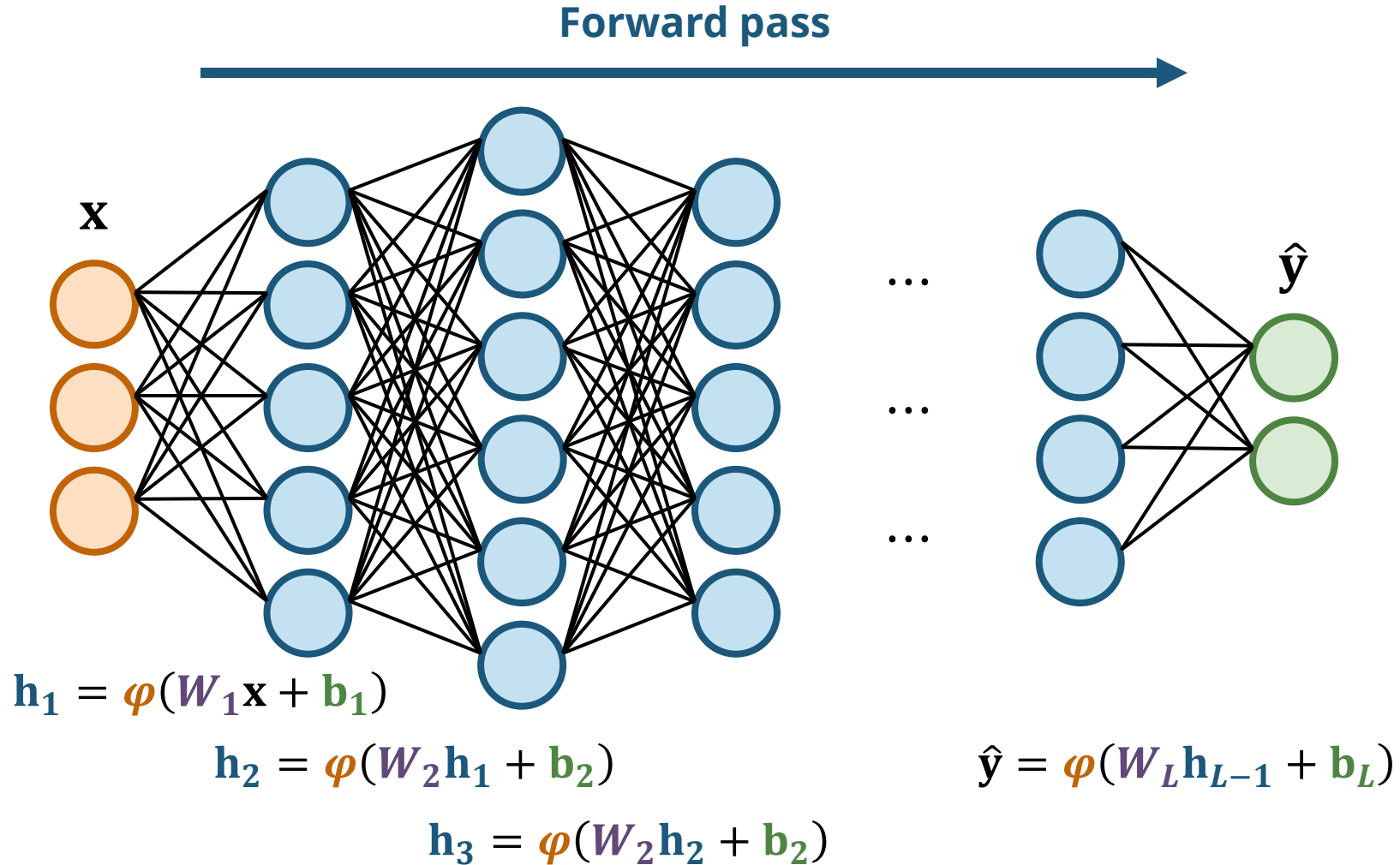
$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

# Backpropagation: Efficiently Computing the Gradients



youtu.be/Ilg3gGewQ5U?t=196

# Forward Pass & Backward Pass

**Forward pass**



$$\mathbf{h}_1 = \boldsymbol{\varphi}(W_1\mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{h}_2 = \boldsymbol{\varphi}(W_2\mathbf{h}_1 + \mathbf{b}_2)$$

$$\mathbf{h}_3 = \boldsymbol{\varphi}(W_2\mathbf{h}_2 + \mathbf{b}_2)$$

$$\hat{\mathbf{y}} = \boldsymbol{\varphi}(W_L\mathbf{h}_{L-1} + \mathbf{b}_L)$$

# Forward Pass & Backward Pass

$$\frac{\partial L}{\partial \mathbf{x}} \qquad \frac{\partial L}{\partial \mathbf{h}_1} \qquad \frac{\partial L}{\partial \mathbf{h}_2} \qquad \frac{\partial L}{\partial \mathbf{h}_3} \qquad \qquad \frac{\partial L}{\partial \mathbf{h}_{L-1}}$$



**Backward pass**

**loss.backward()**