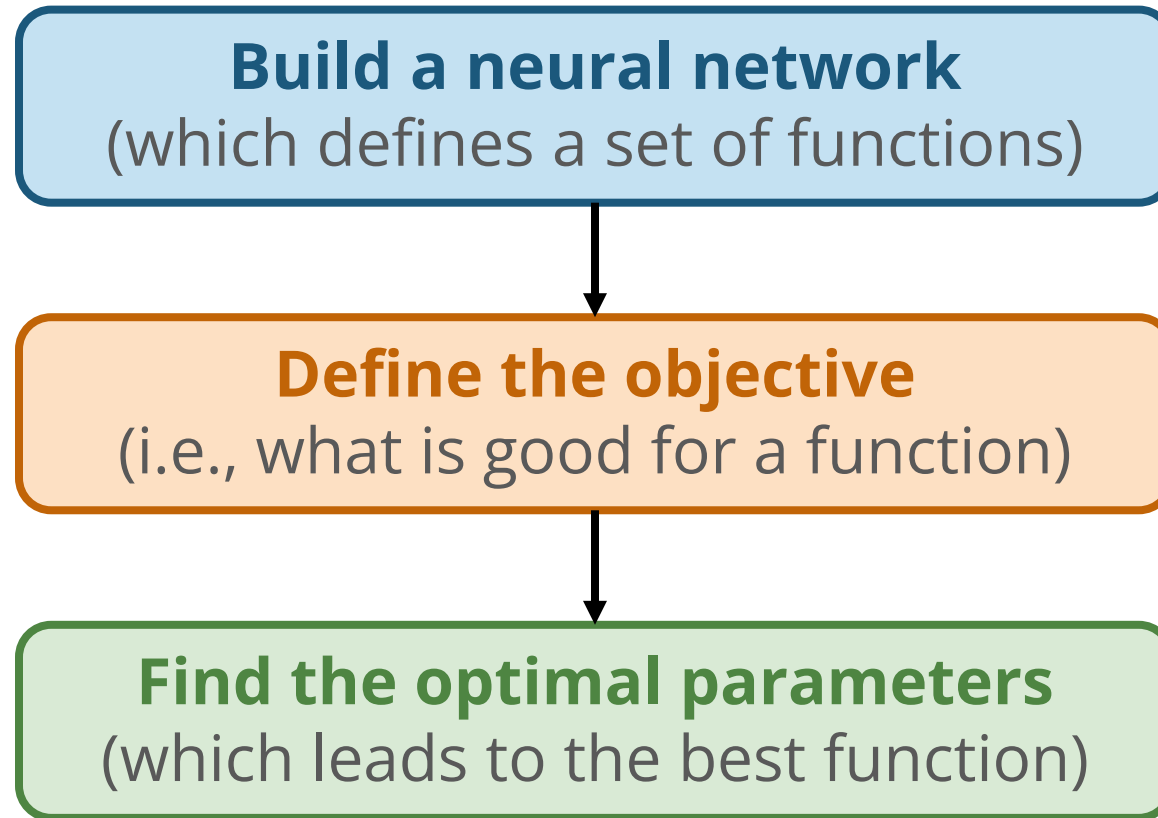PAT 463/563 (Fall 2025)

# Music & AI

**Lecture 7: Deep Learning Fundamentals II**
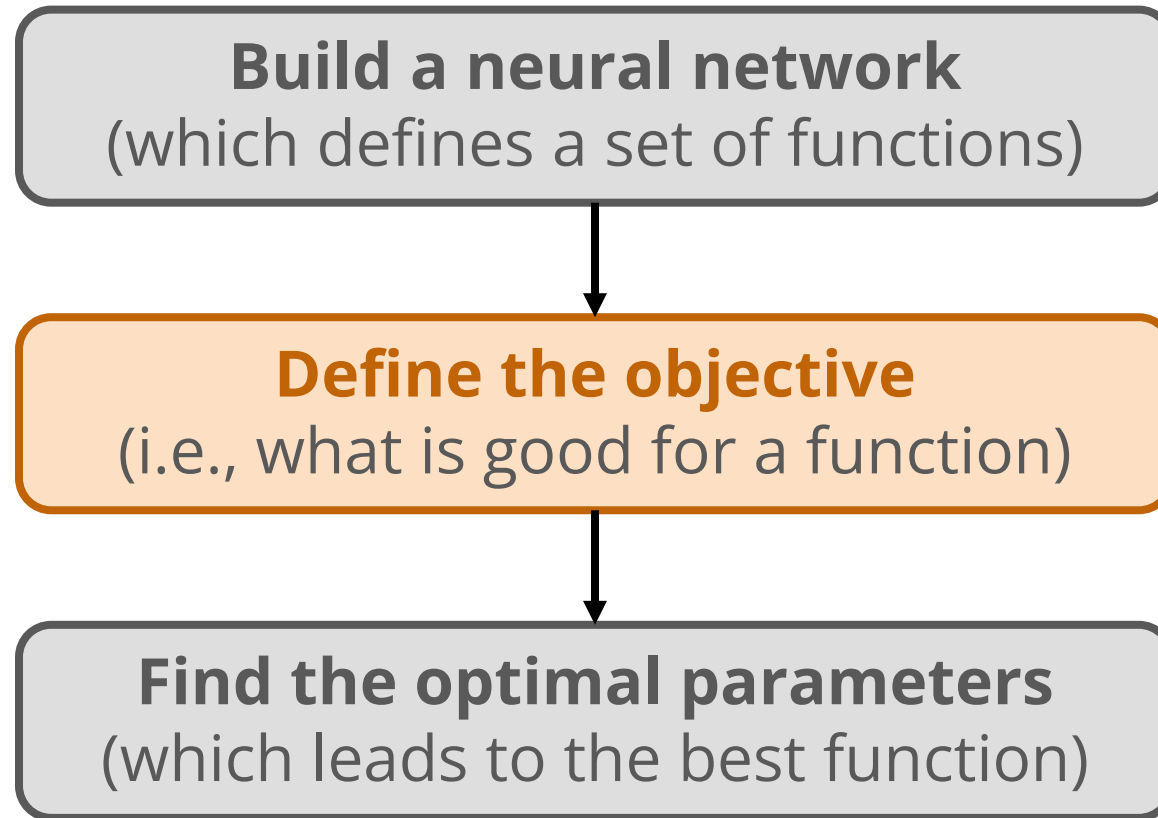
Instructor: Hao-Wen Dong

# Training a Neural Network

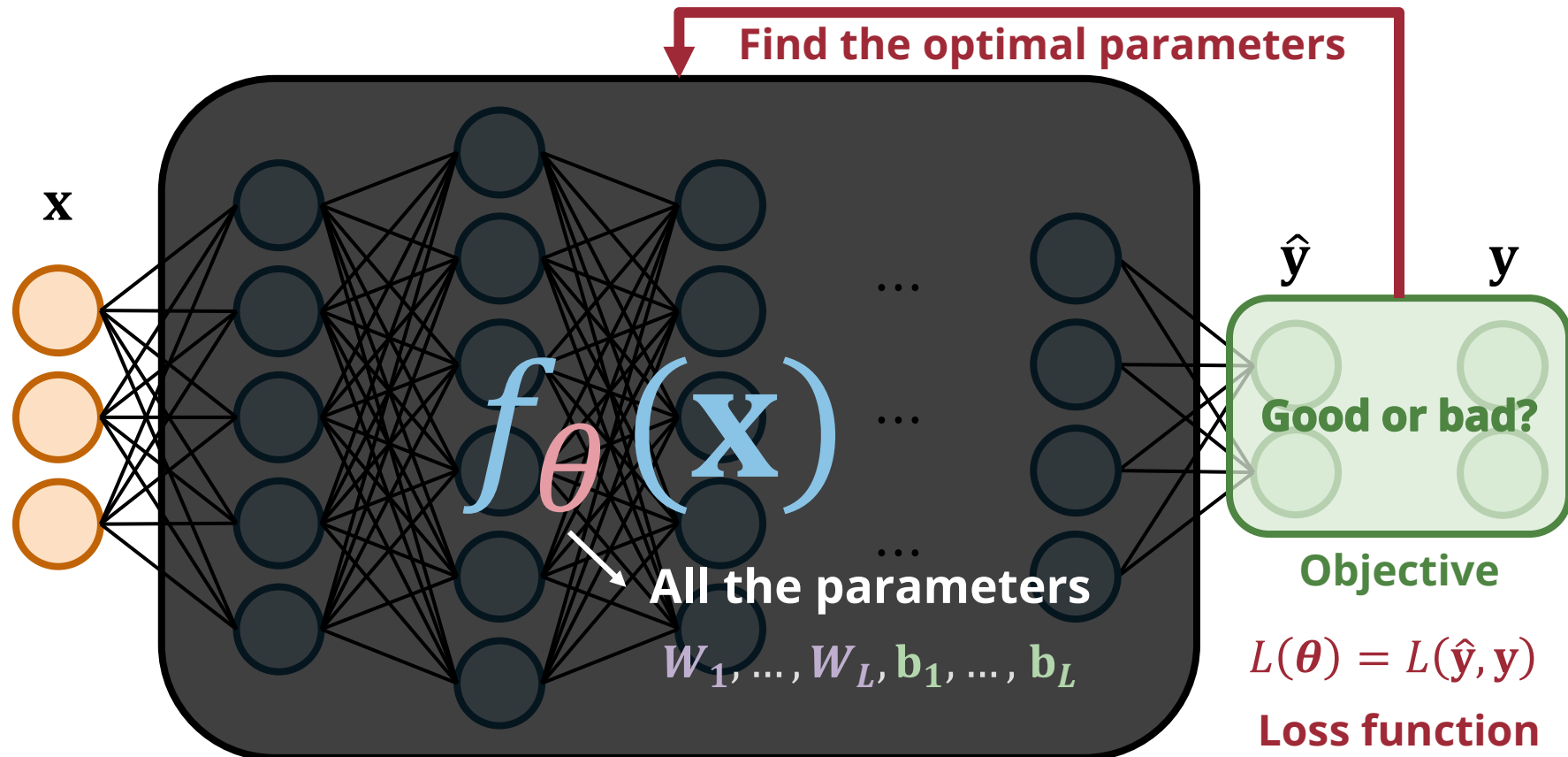# Training a Neural Network

**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)         │
└─────────────────────────────────────────────┘
```

# Neural Networks are Parameterized Functions

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



**Find the optimal parameters**

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

**All the parameters**

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

$\mathbf{x}$

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

**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

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

**L1 loss**

**x**

**No activation function!**

$\hat{y}$

$y$

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

**Loss function**

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

**L2 loss**

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

# L1 vs L2 Losses

**L1 loss**

**L2 loss**

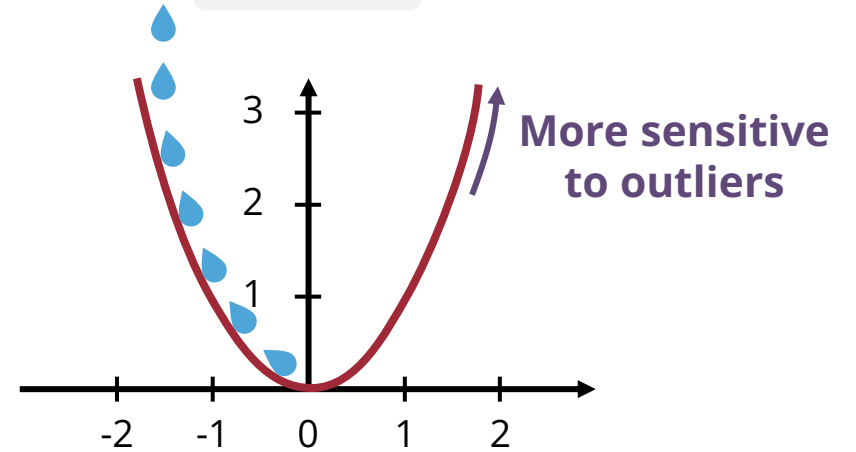**More sensitive to outliers**

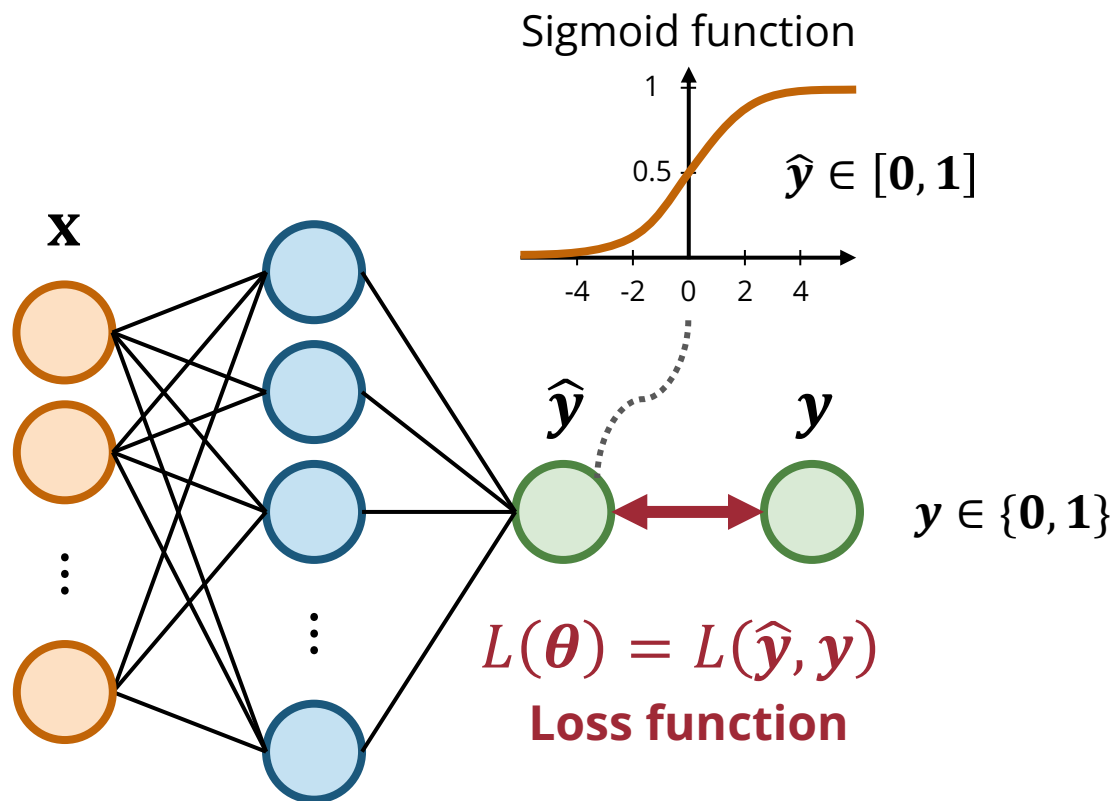$$L(\hat{y}, y) = |\hat{y} - y|$$

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

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

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

**Mean Absolute Error (MAE)**

**Mean Squared Error (MSE)**

# L1 vs L2 Losses

**L1 loss**

**L2 loss**

**More sensitive to outliers**

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

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

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

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

**Mean Absolute Error (MAE)**

**Mean Squared Error (MSE)**

# Binary Cross Entropy for Binary Classification

- **Logistic regression** approaches classification like regression



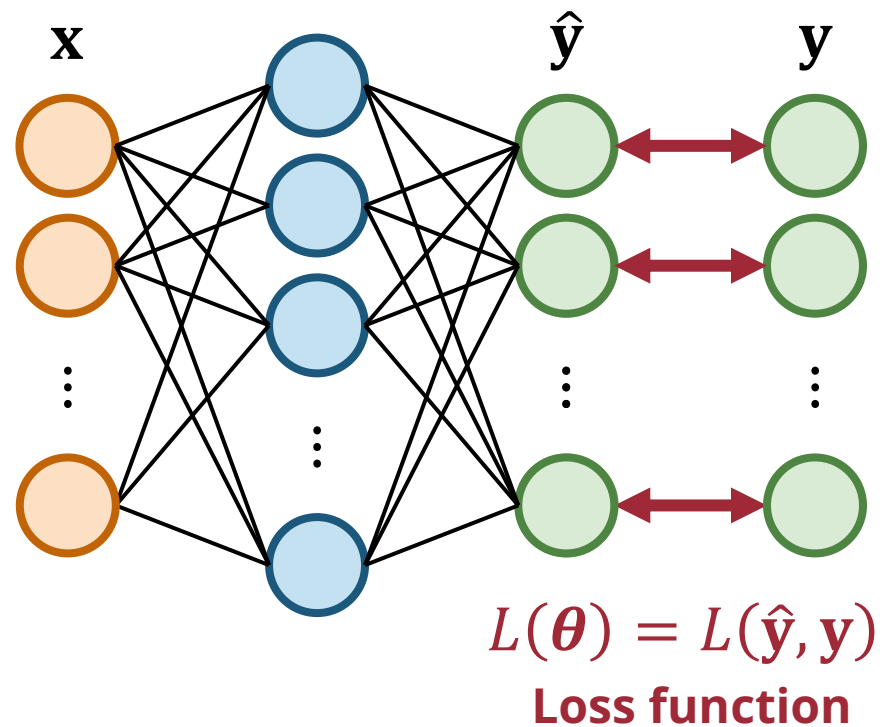Sigmoid function

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

$\mathbf{x}$

$\widehat{\mathbf{y}}$     $\mathbf{y}$

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

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

**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**

$\tilde{y}_i \in \mathbb{R}$  $\hat{y}_i \in [0, 1]$  $y_i \in \{0, 1\}$

$\mathbf{x}$  $\tilde{\mathbf{y}}$  Softmax  $\hat{\mathbf{y}}$  $\mathbf{y}$

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

**Loss function**

**Softmax**

$$\hat{y}_i = \frac{e^{\tilde{y}_i}}{\sum_{j=1}^{n} e^{\tilde{y}_j}}$$

# Softmax

- **Intuition**: Map several numbers to $[0, 1]$ while **keeping their relative magnitude**

  - Softmax is like the **multivariate version of sigmoid**

**Real-valued numbers to probability-like numbers**

**Divide by sum**     **Sum to 1**

$\tilde{y}_1 \rightarrow$ [Softmax] $\rightarrow \hat{y}_1$

$\tilde{y}_2 \rightarrow$ [Softmax] $\rightarrow \hat{y}_2$

$\tilde{y}_n \rightarrow$ [Softmax] $\rightarrow \hat{y}_n$

$73.5 \quad \tilde{y}_1 \rightarrow$ [exp] $\rightarrow$ [Normalize] $\rightarrow \hat{y}_1 \quad$ **0.98**

$2.6 \quad \tilde{y}_2 \rightarrow$ [exp] $\rightarrow$ [Normalize] $\rightarrow \hat{y}_2 \quad$ **0.02**

$-12.8 \quad \tilde{y}_n \rightarrow$ [exp] $\rightarrow$ [Normalize] $\rightarrow \hat{y}_n \quad$ **0.00**

# 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**

# 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**
(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)

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

# 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

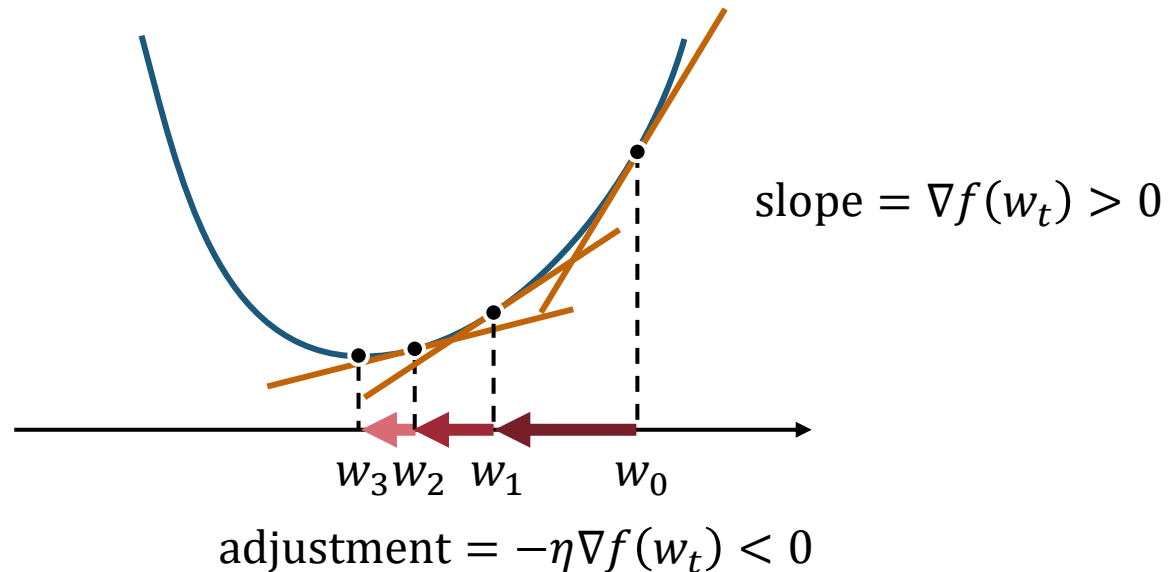Derivative for a vector,
matrix or tensor



$w_0$

# Gradient Descent: Pseudocode

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

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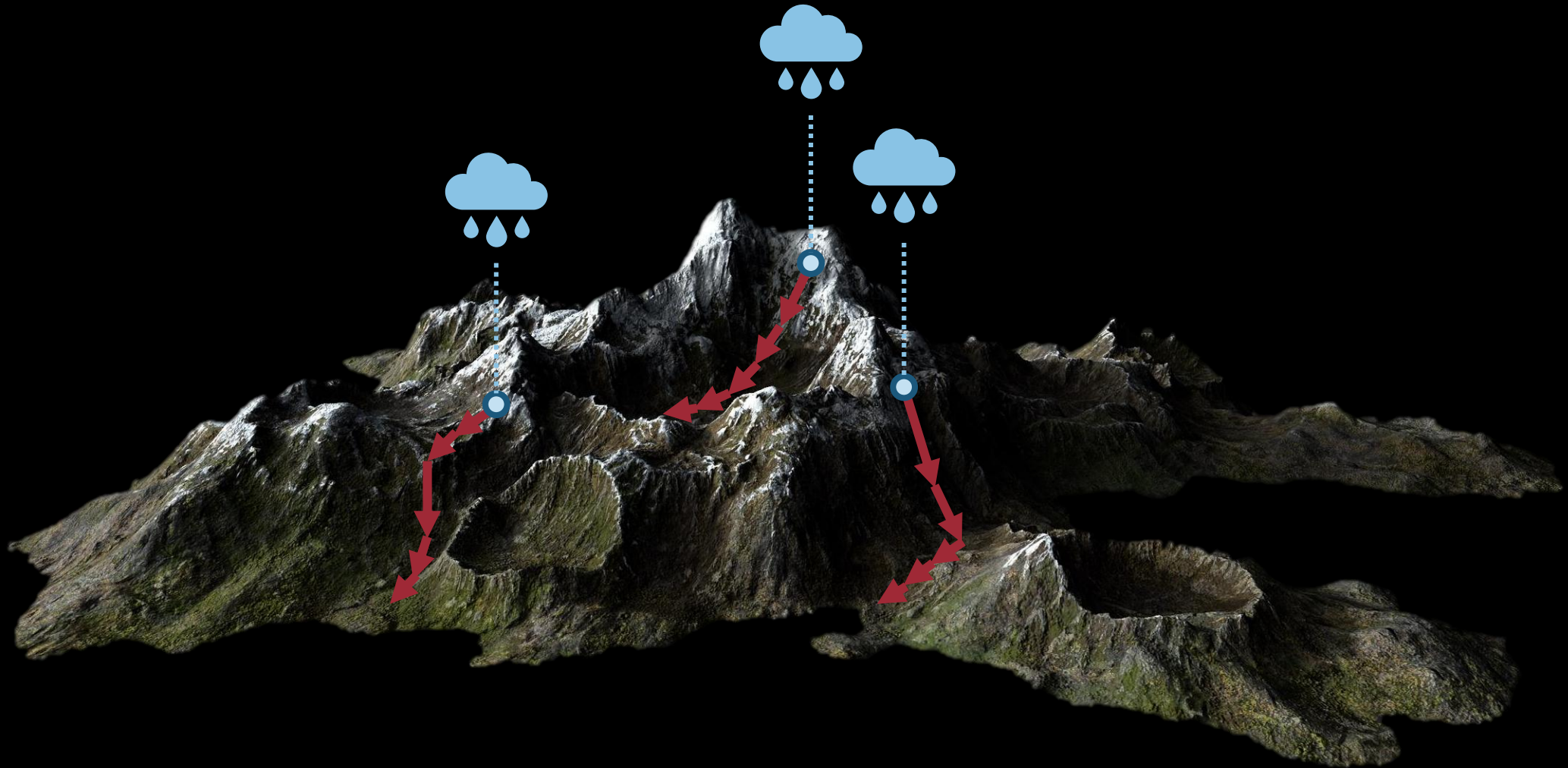
# Gradient Descent: Pseudocode

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

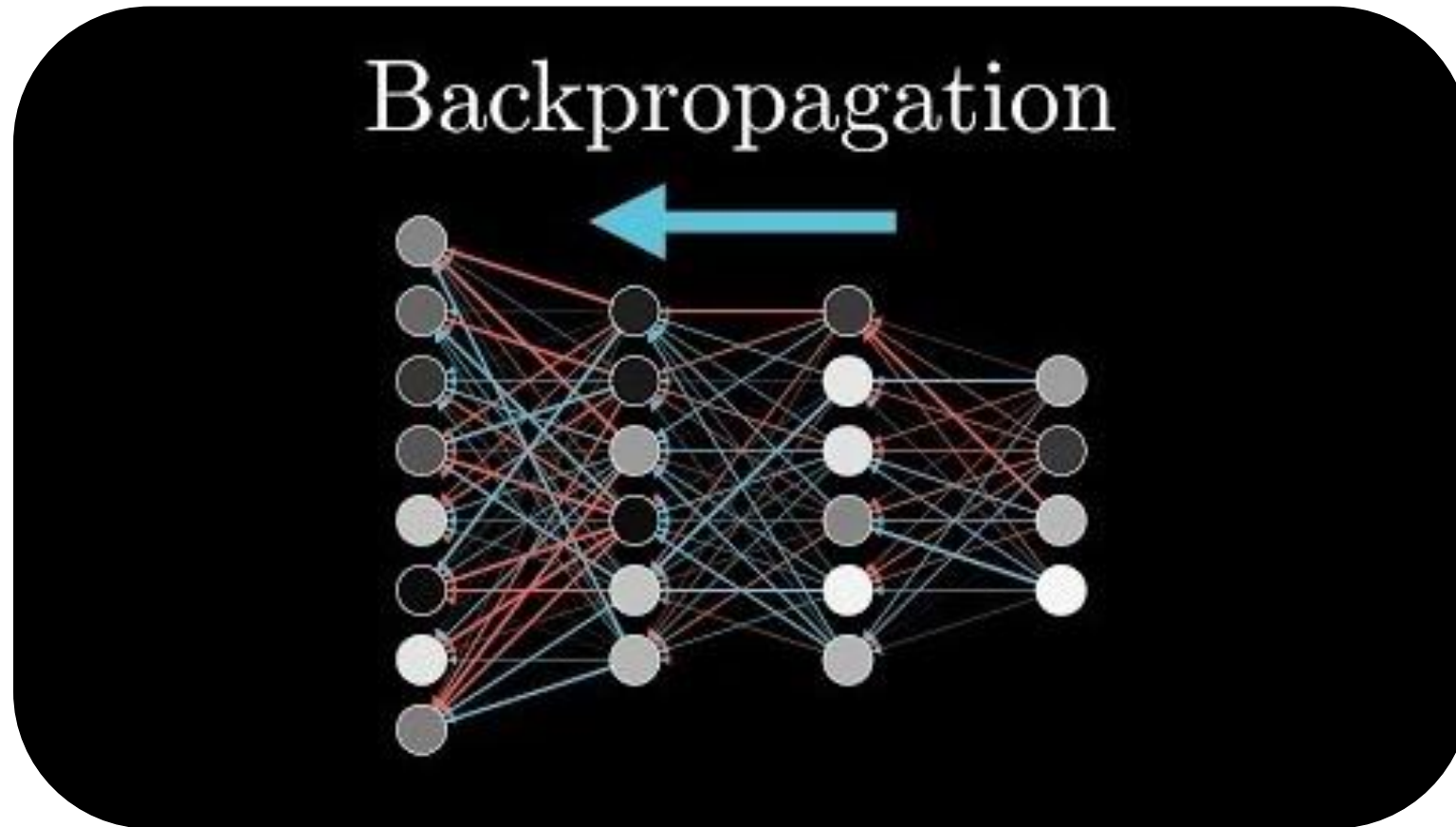- Repeat until convergence:  $w_{t+1} = w_t - \eta \nabla f(w_t)$

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 $= \nabla f(w_0) > 0$

$w_1 \quad w_0$

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 $= \nabla f(w_1) > 0$

$w_2$ $w_1$ $w_0$

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 $= \nabla f(w_t) > 0$

$w_3\, w_2\ \ w_1\qquad\ \ w_0$

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**

$\mathbf{x}$

$\hat{\mathbf{y}}$

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

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

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

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

# Forward Pass & Backward Pass



**Backward pass**

**loss.backward()**

# Training–Validation–Test

# In-distribution vs Out-of-distribution

# In-distribution vs Out-of-distribution

# In-distribution vs Out-of-distribution

**Training**

# In-distribution vs Out-of-distribution

- **Key**: Make the training distribution closer to the target distribution

- First, we need to **define our target distribution**

- Then, we can try to
  - Collect a **diverse** dataset covering that covers different parts of the target distribution
  - Apply **data augmentation** to fill the gaps in the distribution

# In-distribution vs Out-of-distribution

- What do we really want?

  ▪ Good performance on the **training samples**   **We already have their answers**

  ▪ Good performance on **unseen samples in the target distribution**   **Yep, we can do this!**

  ▪ Good performance on **out-of-distribution samples**   **Hopefully, but not guaranteed**

**How to achieve good performance on unseen samples in the target distribution**

# Overfitting & Underfitting

# Overfitting & Underfitting

**Underfitting**

**Good fit!**

**Overfitting**



**Model too inexpressive**

**Model too expressive**

# Train–Test Split

- **Goal**: Good performance on **unseen samples in the target distribution**

# Train–Test Split

- **Goal**: Good performance on **unseen samples in the target distribution**

**Training**

**Test**

# Test Set is an Estimation of the Test Distribution

- We create a test set because we want to **estimate the performance when the model is applied to an interested distribution**

# Train–Validation–Test Split

**Training**

**Test**

# Train–Validation–Test Split

**Training**

**Validation**

**Test**

# Training–Validation–Test Pipeline

**Training**

**Validation**

**Test**



**Optimize**

**Select**

# Training vs Validation Losses

# Training vs Validation Losses
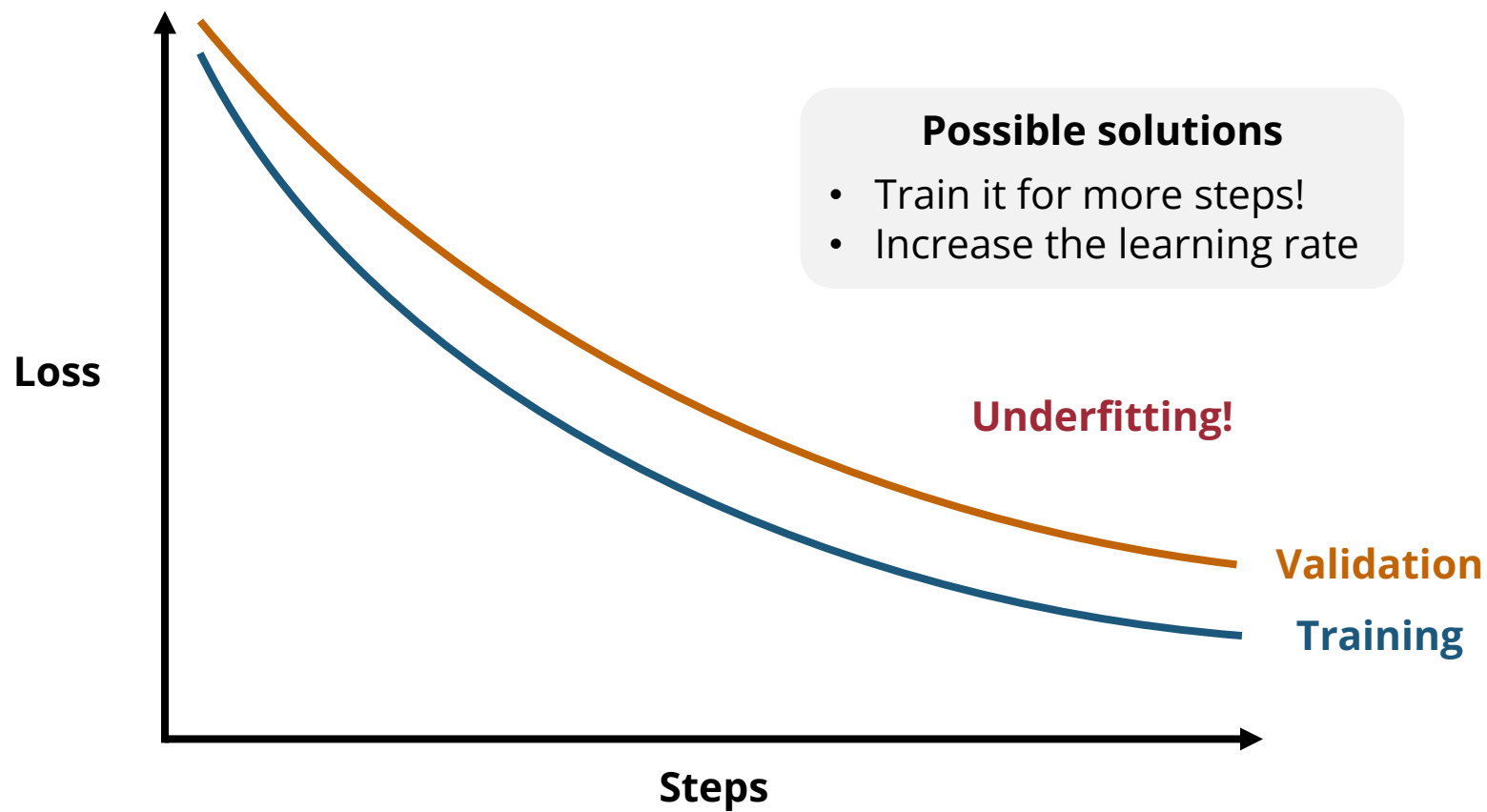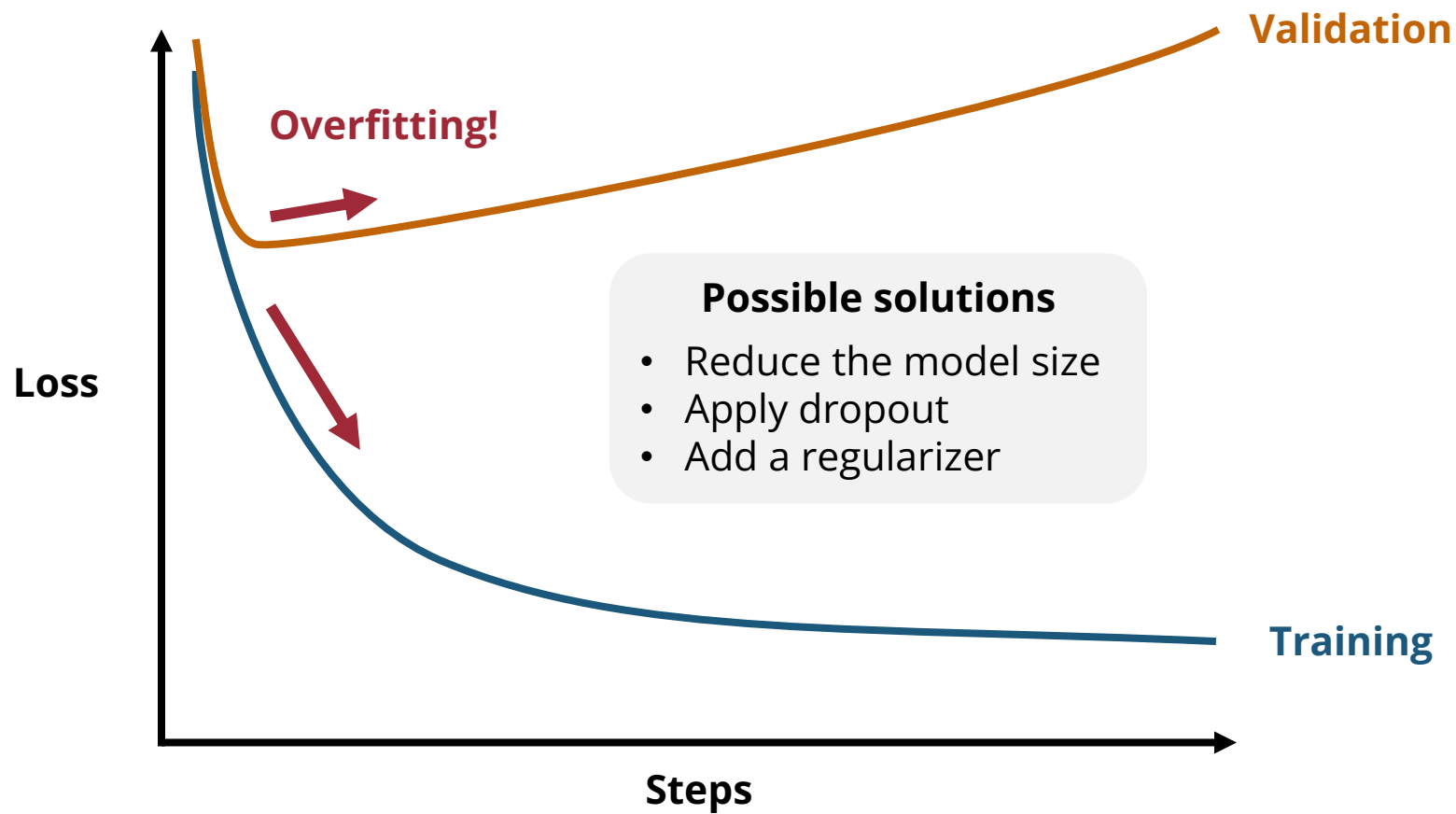
# Training vs Validation Losses



Loss

Pick the model with the
lowest validation loss

Validation

Training

Steps

# Training vs Validation Losses



**Possible solutions**
- Increase the size and diversity of the validation set
- Apply cross validation

**Validation**

**Unrepresentative validation samples**

**Training**

**Loss**

**Steps**

# Training vs Validation Losses



Loss

**Possible solutions**
- Train it for more steps!
- Increase the learning rate

**Underfitting!**

**Validation**

**Training**

Steps

# Training vs Validation Losses



**Validation**

**Overfitting!**

**Loss**

**Possible solutions**

- Reduce the model size
- Apply dropout
- Add a regularizer

**Training**

**Steps**

# Train–Validation–Test Split

- **Keys**
  - **Never train or select your model on test samples!**
  - Don't over-select your model on the validation set

- What's the **best ratio**?
  - Most common: **8:1:1** or 9:0.5:0.5
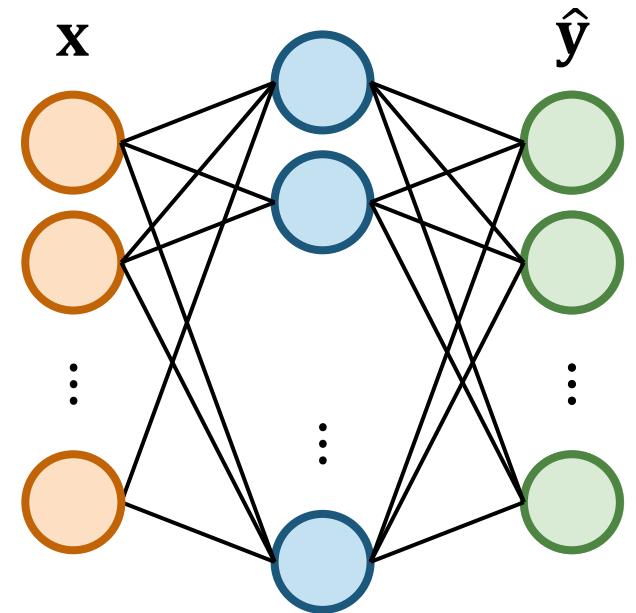  - For smaller dataset, you might even want 6:2:2

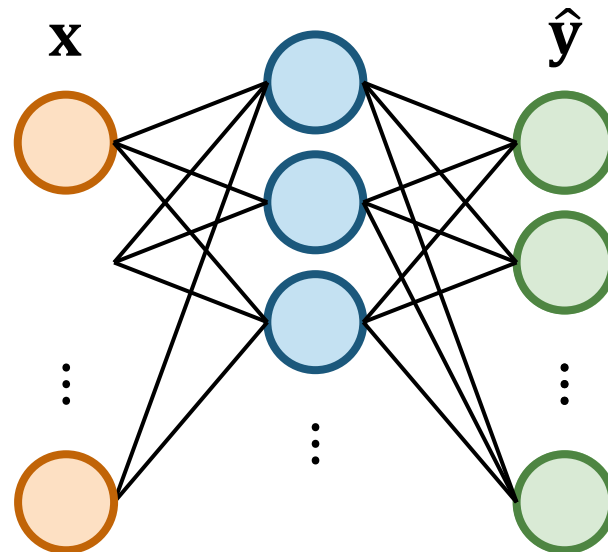# Overcoming Overfitting
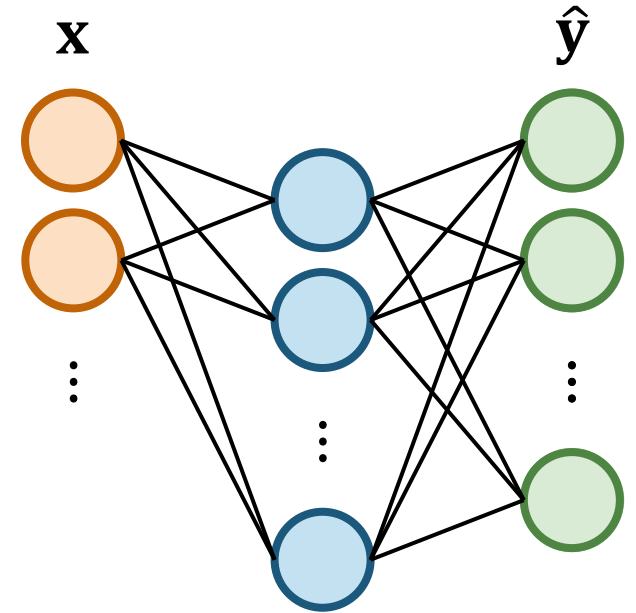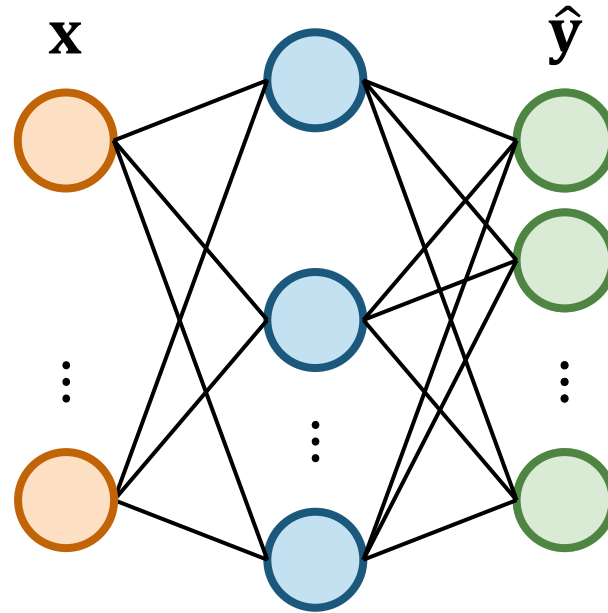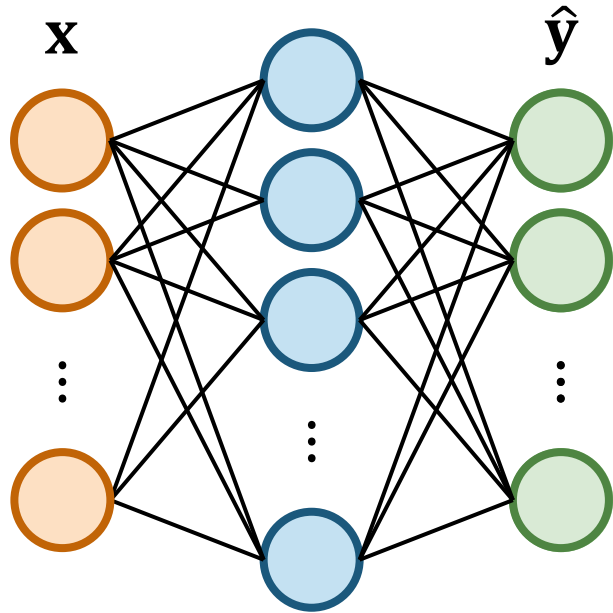
# Early Stopping

# Dropout



**Each neuron may be removed with probability $p$ during training**
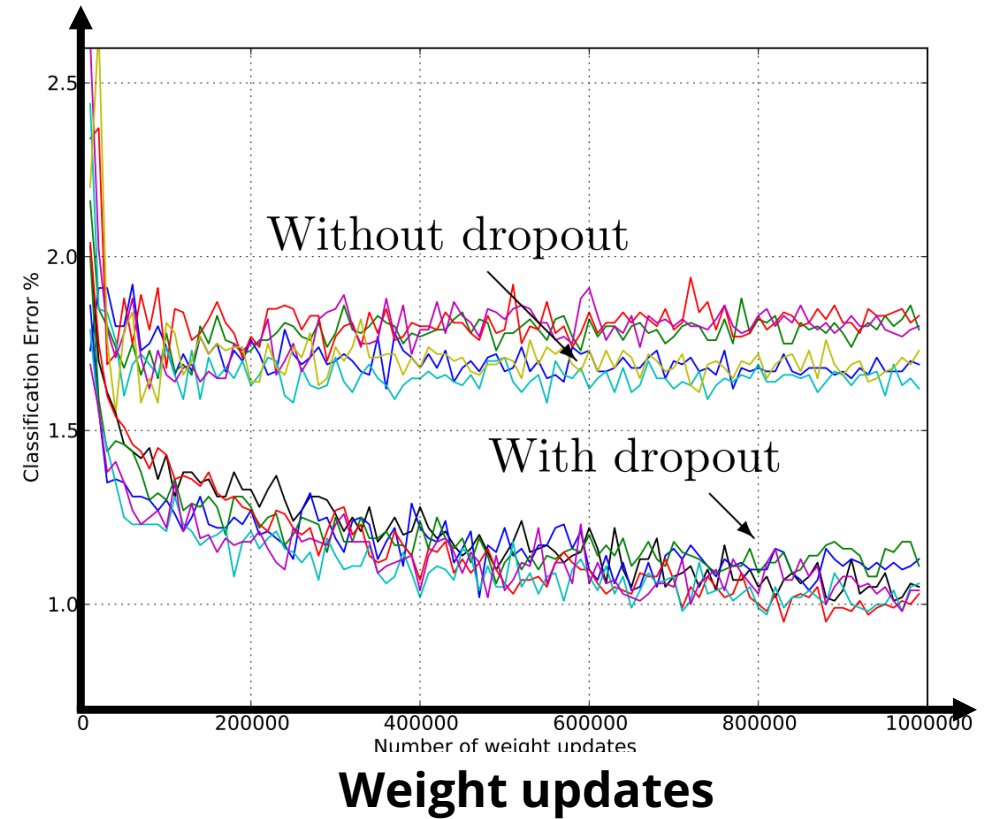
Dropout rate

# Dropout



**Each neuron may be removed**
**with probability $p$ during training**

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," *JMLR*, 2014.\

# Regularization Term

- A regularization term can help alleviate overfitting
  - **L1 regularization** (LASSO)
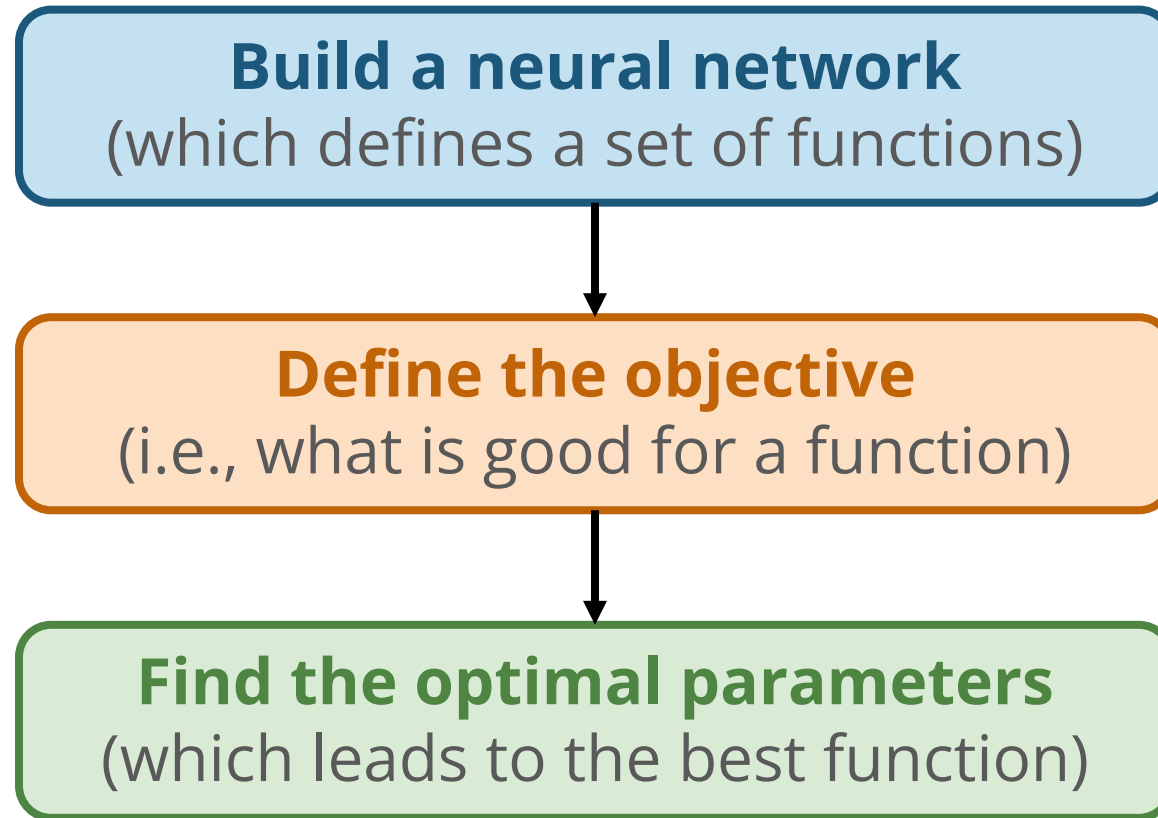
$$L' = L + \lambda(|w_1| + |w_2| + \cdots + |w_K|)$$

  - **L2 regularization** (ridge regression)

$$L' = L + \lambda\left(w_1^2 + w_2^2 + \cdots + w_K^2\right)$$

**Both L1 and L2 regularizations encourage smaller weights**

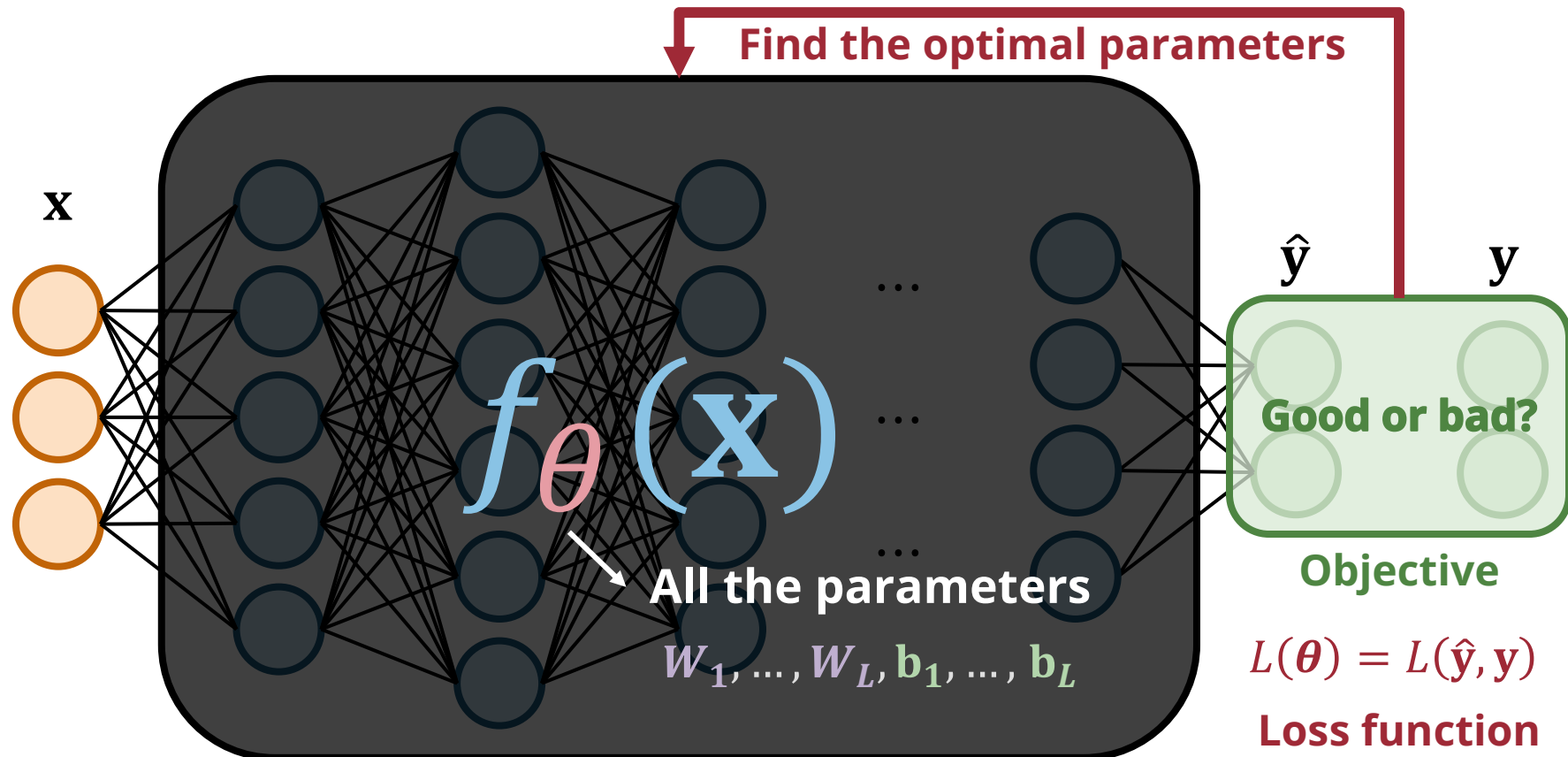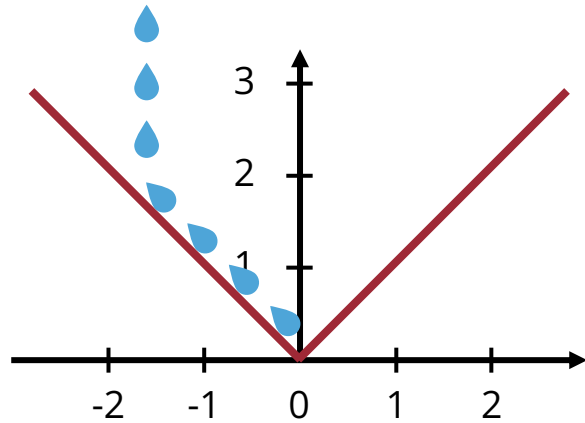# Recap

# Training a Neural Network

# Neural Networks are Parameterized Functions

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

**Find the optimal parameters**

$x$

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

**All the parameters**

$W_1, \ldots, W_L, \mathbf{b_1}, \ldots, \mathbf{b_L}$

$\hat{\mathbf{y}}$　$\mathbf{y}$

**Good or bad?**

**Objective**

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

**Loss function**
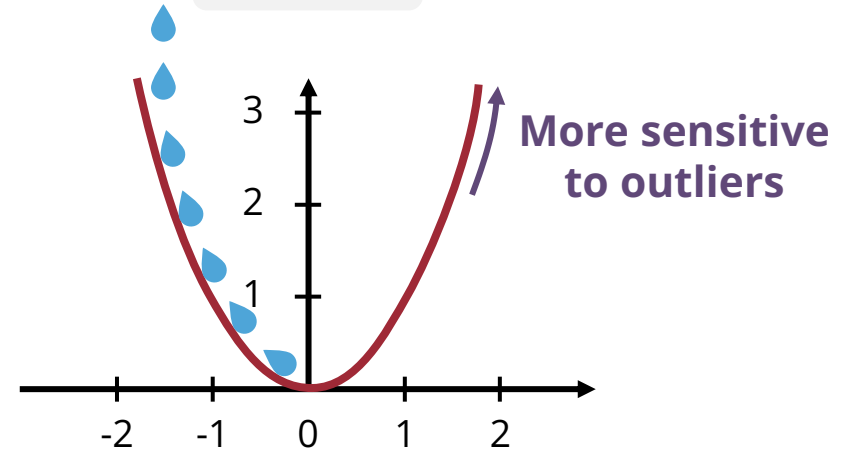
# L1 vs L2 Losses

**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)**

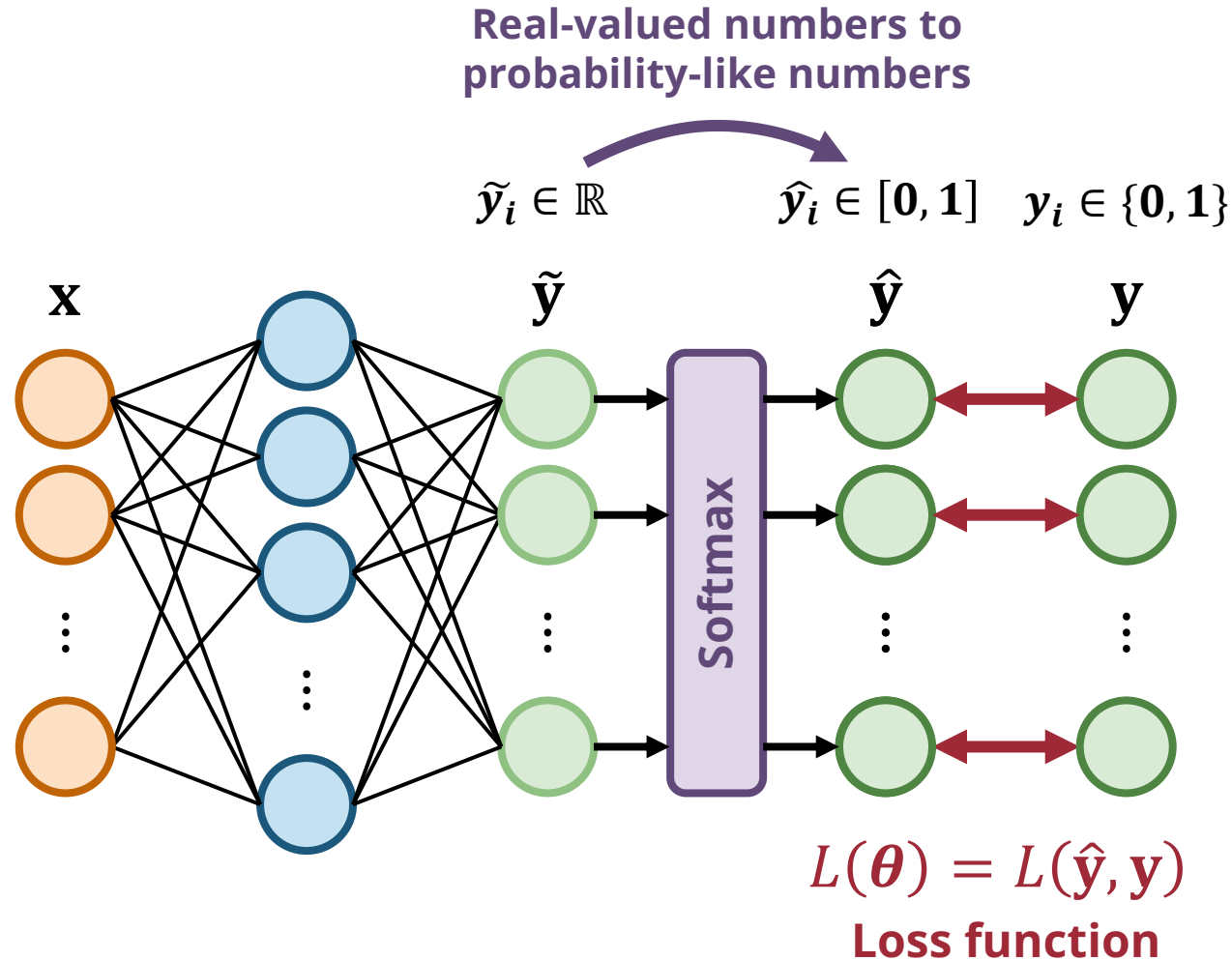**L2 loss**



**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)**
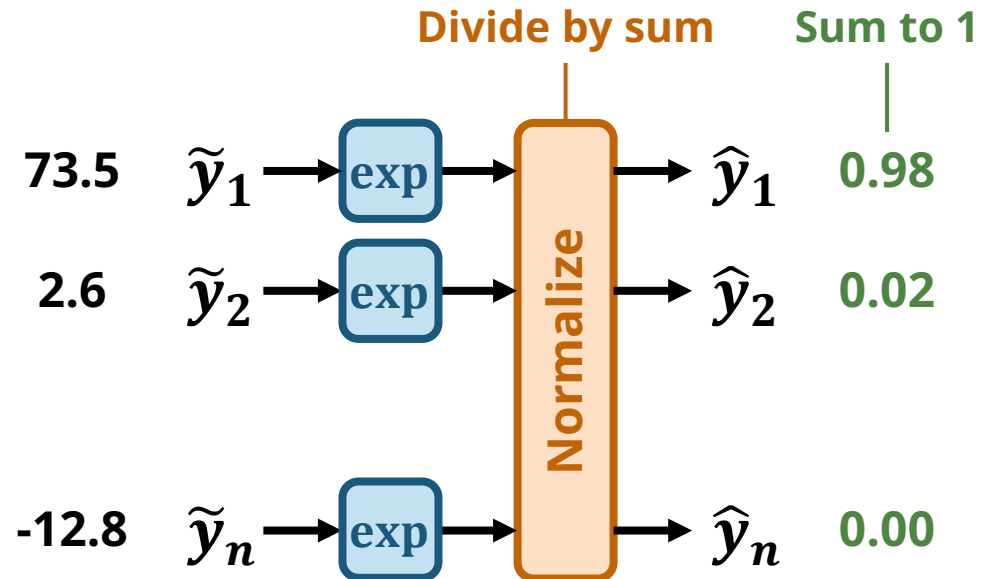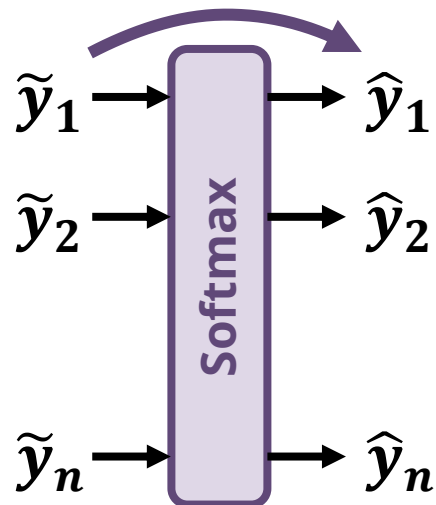
# Cross Entropy for Multiclass Classification

**Real-valued numbers to probability-like numbers**

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

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

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

**Loss function**

**Softmax**

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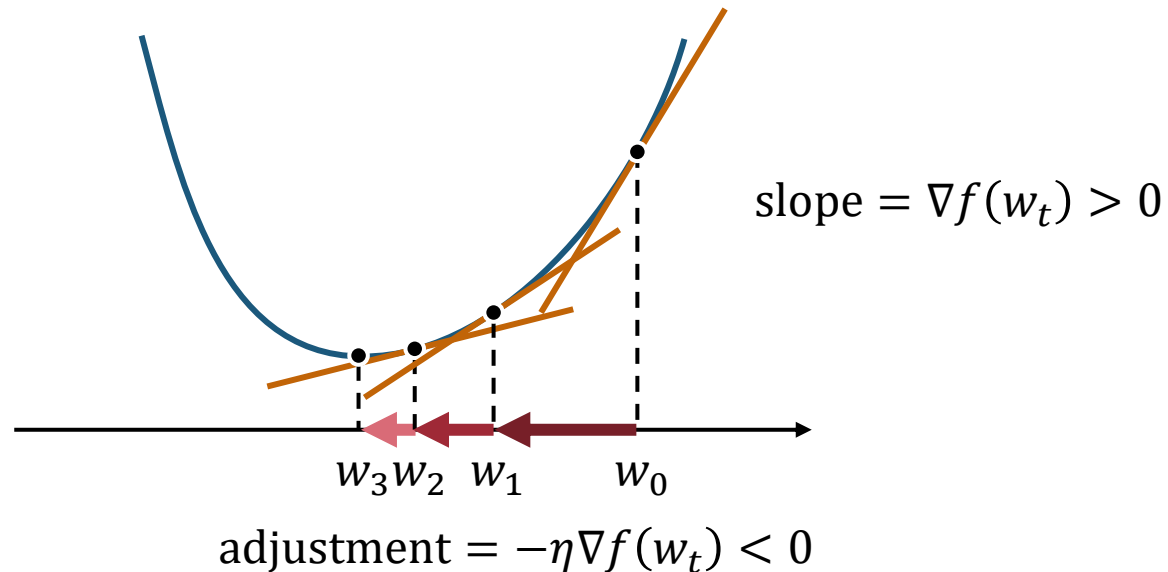
# Softmax

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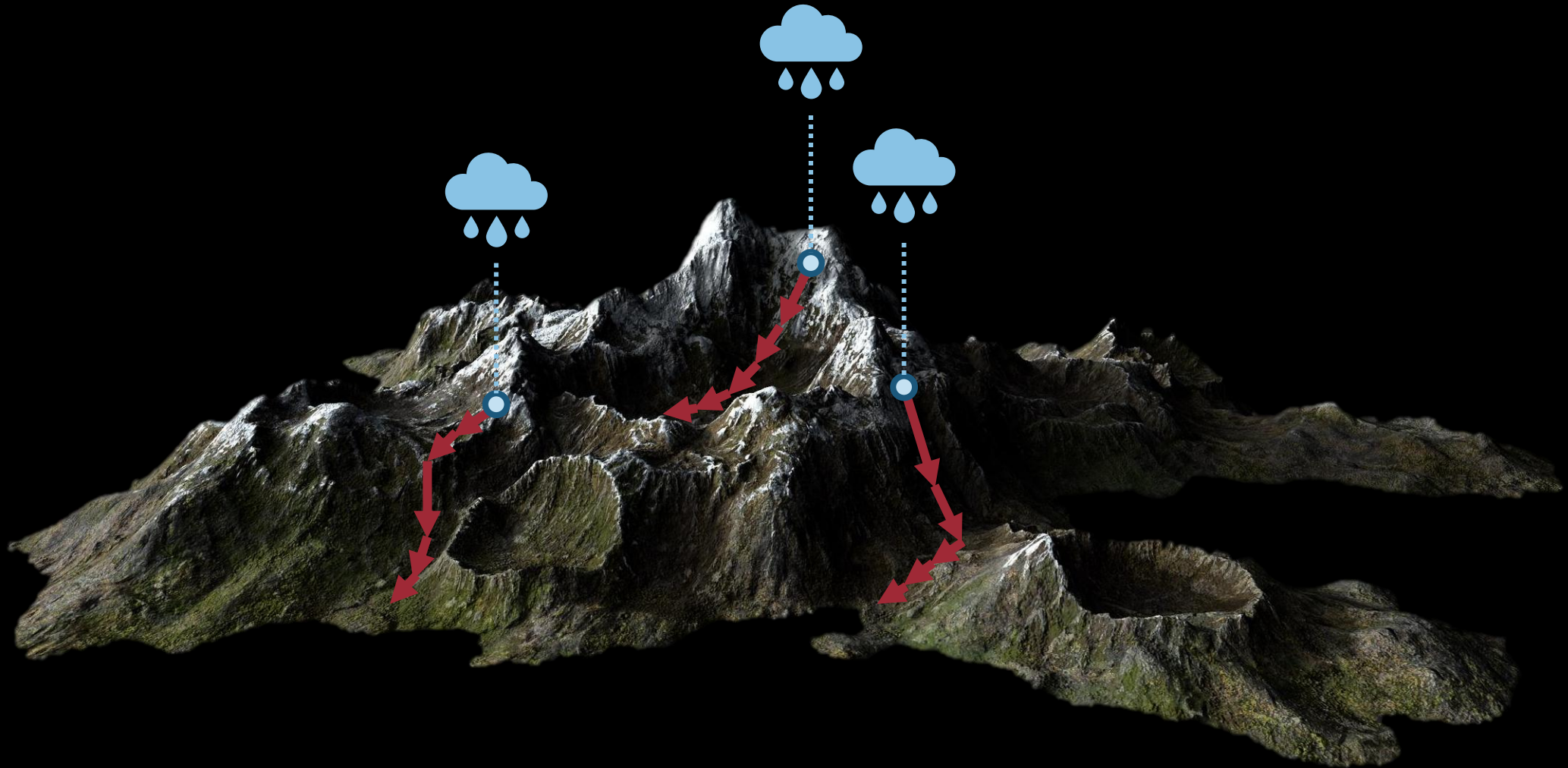  - Softmax is like the **multivariate version of sigmoid**

# Gradient Descent: Pseudocode

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

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

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

$$w_3\ w_2\quad w_1\qquad w_0$$

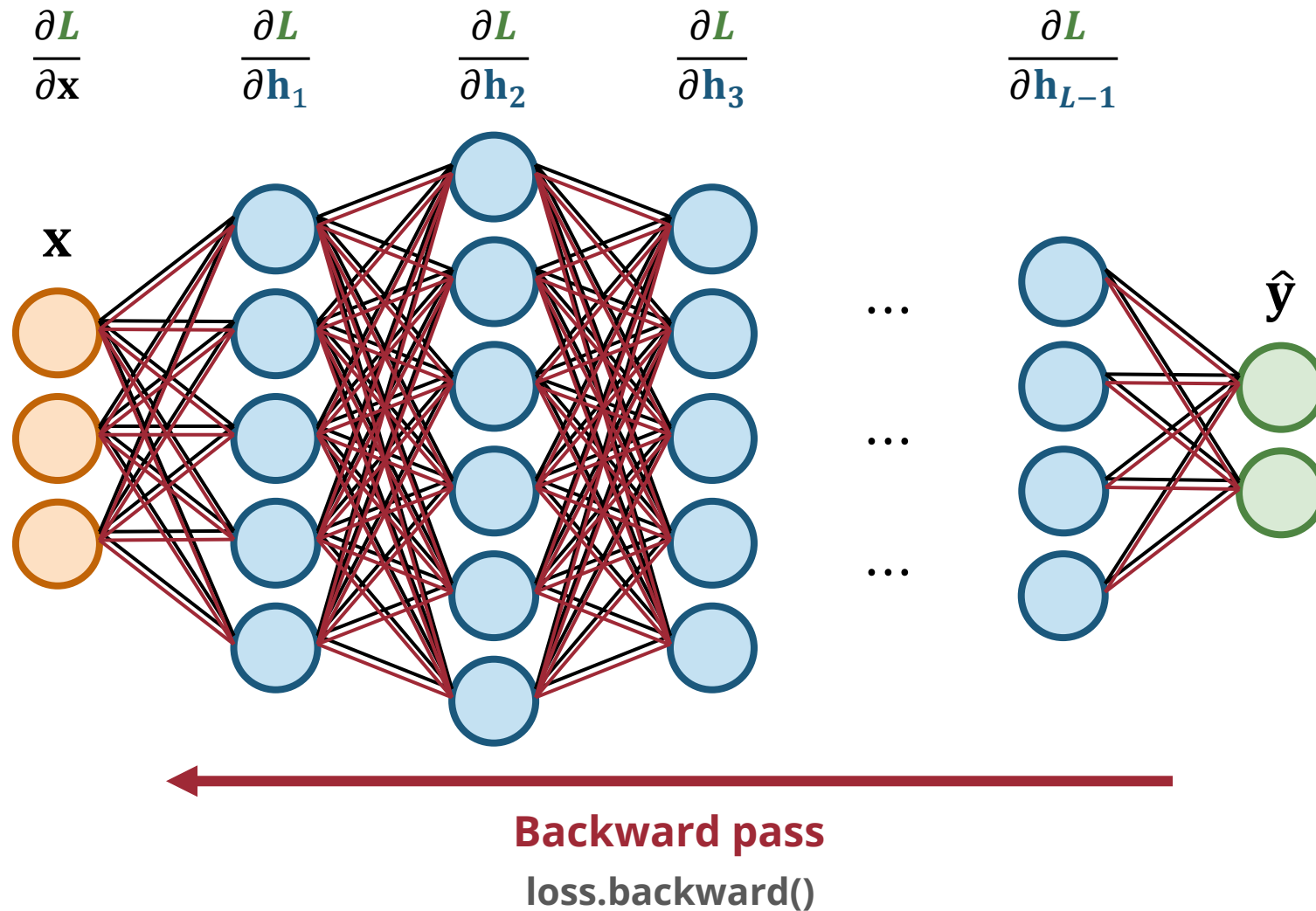$$\text{adjustment} = -\eta \nabla f(w_t) < 0$$

# Gradient Descent: 3D Case

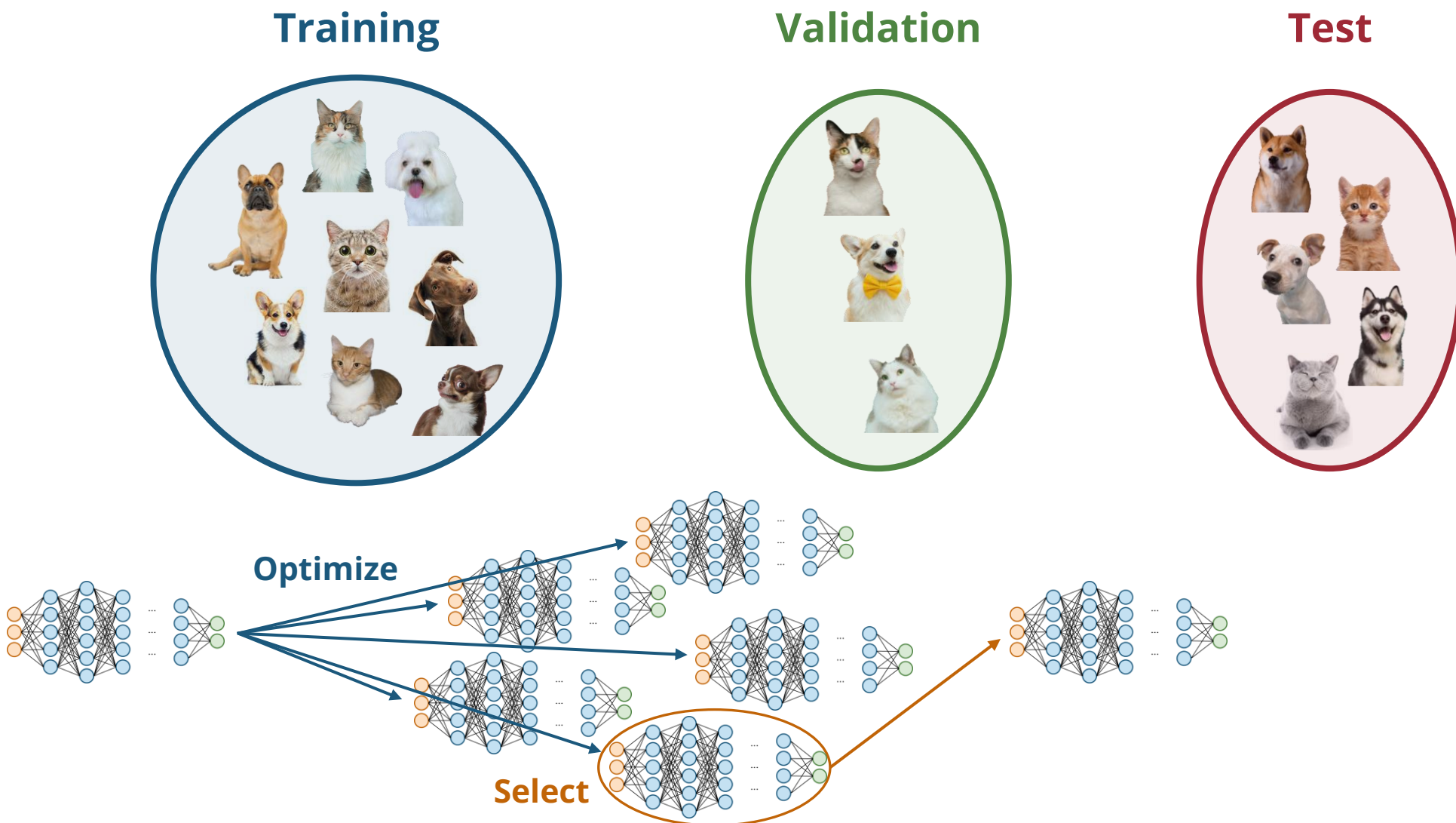# Backpropagation: Efficiently Computing the Gradients

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

- The reason why we want **everything to be differentiable** in deep learning

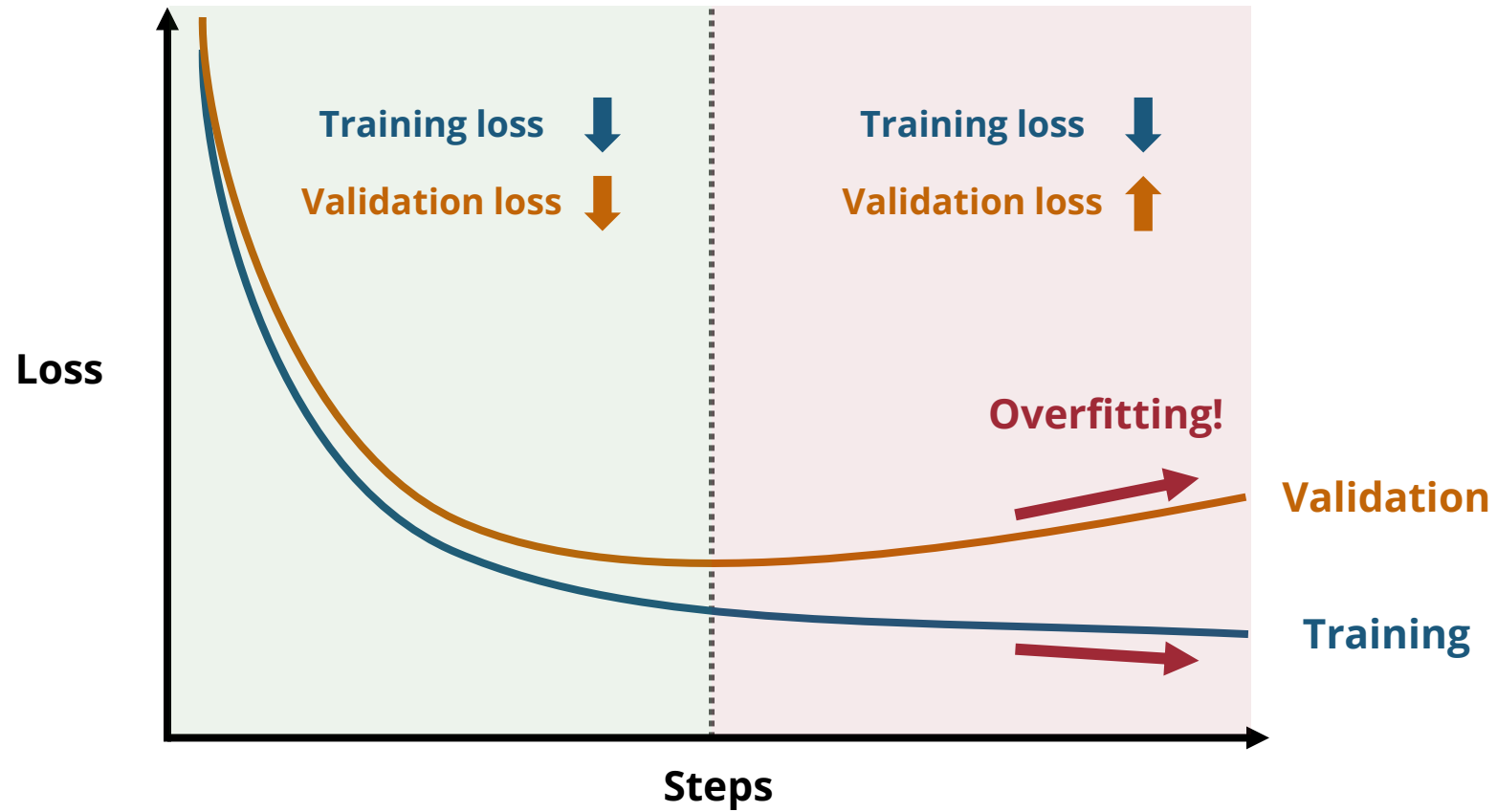$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

# Forward Pass & Backward Pass



$$\frac{\partial L}{\partial \mathbf{x}} \qquad \frac{\partial L}{\partial \mathbf{h_1}} \qquad \frac{\partial L}{\partial \mathbf{h_2}} \qquad \frac{\partial L}{\partial \mathbf{h_3}} \qquad \qquad \frac{\partial L}{\partial \mathbf{h}_{L-1}}$$

$\mathbf{x}$

$\hat{\mathbf{y}}$

**Backward pass**

**loss.backward()**

# Training–Validation–Test Pipeline

**Training**

**Validation**

**Test**

**Optimize**

**Select**

# Training vs Validation Losses

# Training vs Validation Losses



**Loss**

**Validation**

**Training**

**Steps**

# Training vs Validation Losses



Loss

Pick the model with the lowest validation loss

Validation

Training

Steps

# Next Lecture

# Source Separation



(Source: "Like Before" by Bessonn&sa)

(Other)

UNIVERSITY OF MICHIGAN