PAT 498/598 (Winter 2025)

# Music & AI

## Lecture 9: Deep Learning Fundamentals III
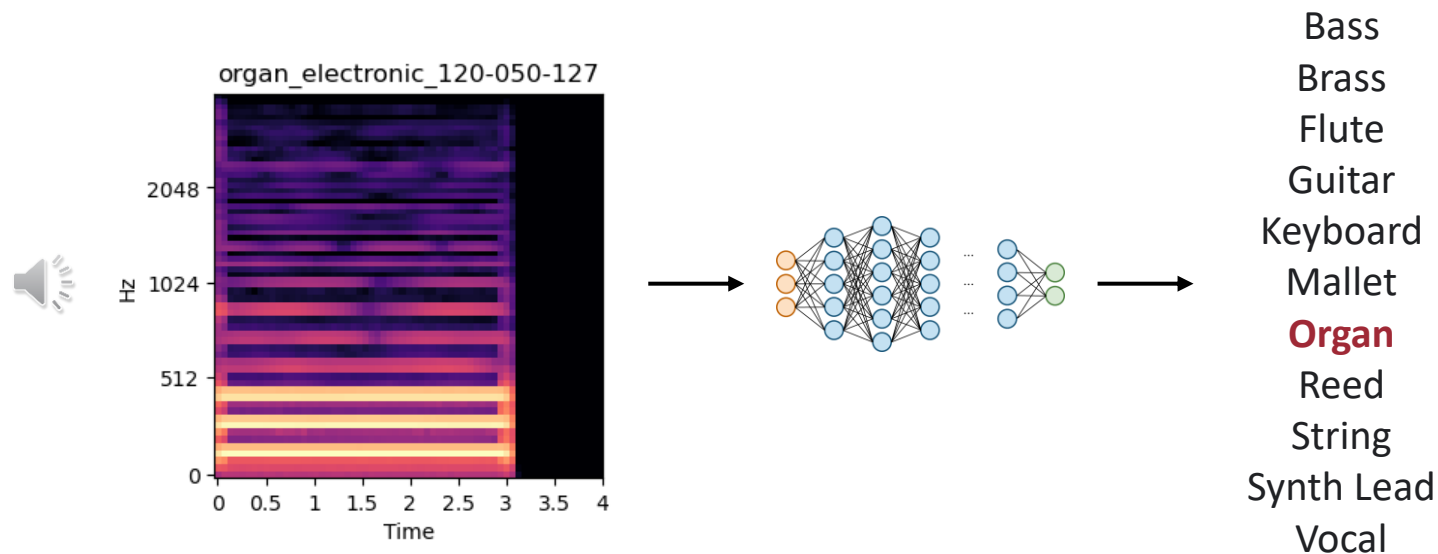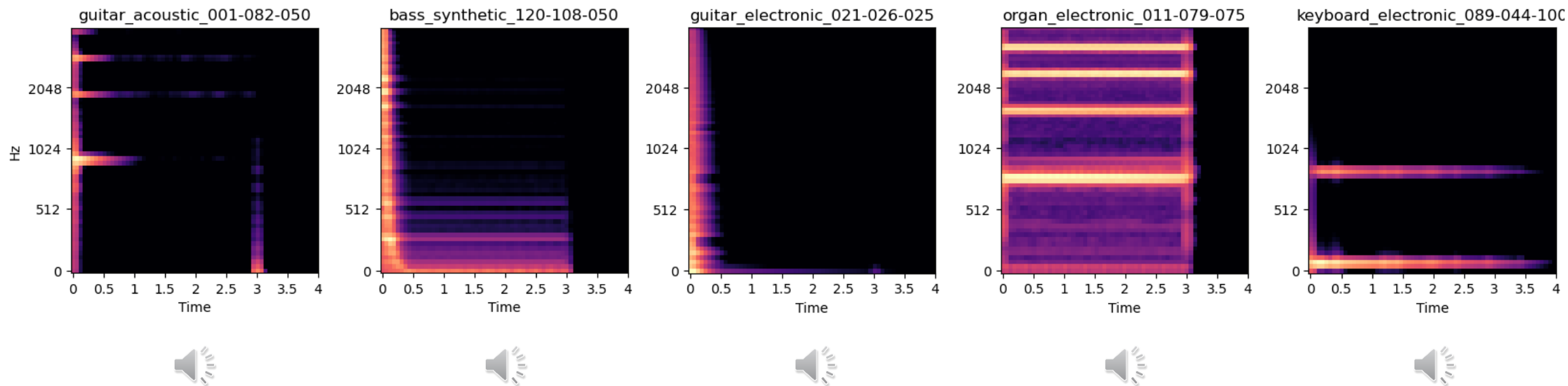
Instructor: Hao-Wen Dong

# Homework 3: Musical Note Classification using CNNs

- Train a CNN that can classify audio files into their **instrument families**
  - **Input**: 64x64 mel spectrogram
  - **Output**: 11 instrument classes
  - Using the **NSynth** dataset (Engel et al., 2017)



Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Douglas Eck, Karen Simonyan, and Mohammad Norouzi, "Neural Audio Synthesis of Musical Notes with WaveNet Autoencoders," *ICML*, 2017.
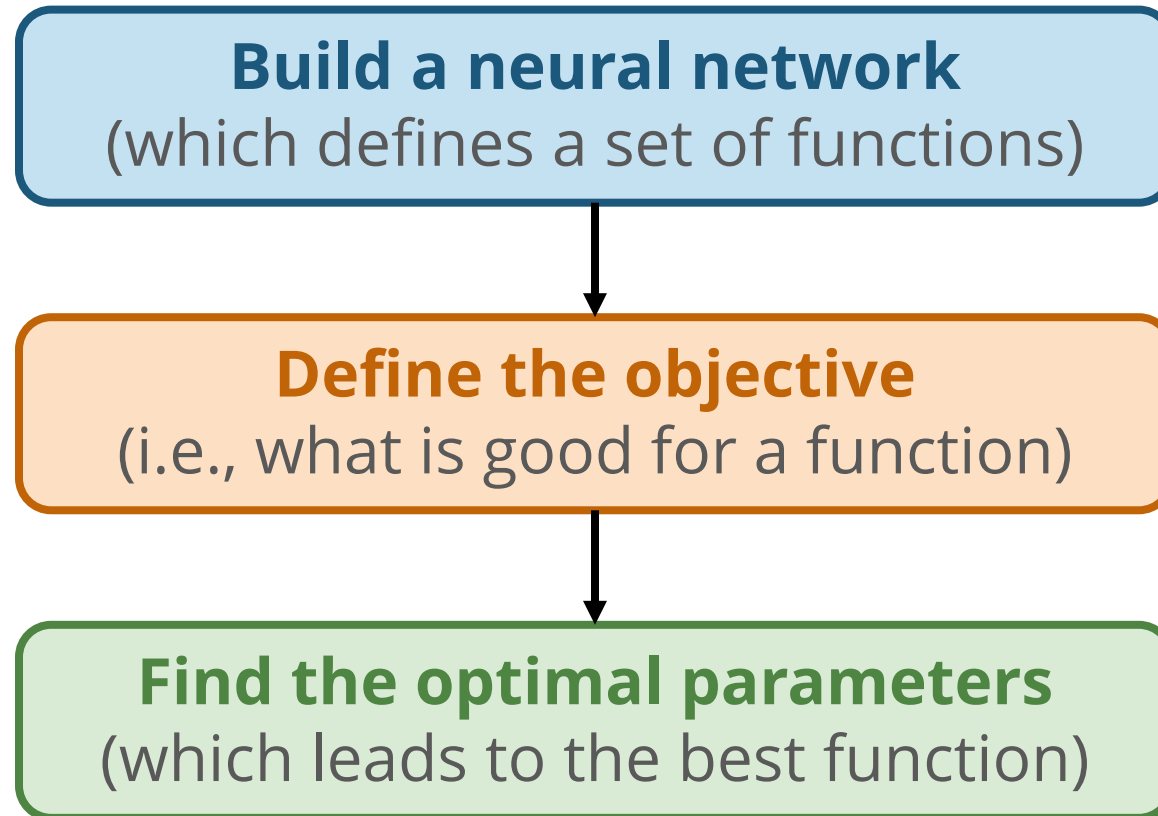
# NSynth Dataset

- A collection of 305,979 **single-shot musical notes** (Engel et al., 2017)
  - Produced from 1,006 **commercial sample libraries**
  - With different **MIDI pitches** (21–108) and **velocities** (25, 50, 75, 100, 127)

Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Douglas Eck, Karen Simonyan, and Mohammad Norouzi, "Neural Audio Synthesis of Musical Notes with WaveNet Autoencoders," *ICML*, 2017.

# Homework 3: Musical Note Classification using CNNs

- Instructions will be released on Gradescope

- Due at **11:59pm ET** on **February 17**

- Late submissions:  **1 point deducted per day**
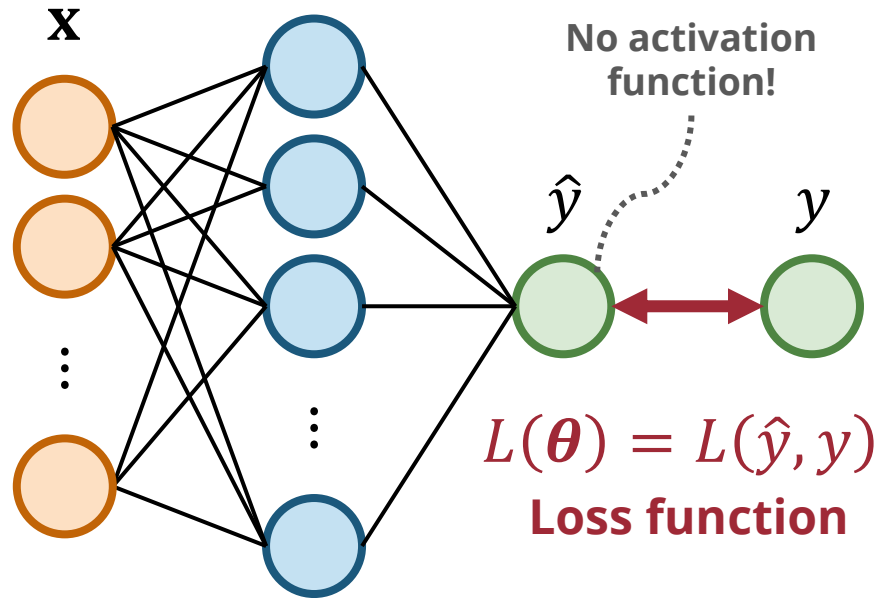
# (Recap) Training a Neural Network

**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)

$$Loss(\boldsymbol{\theta}) = \sum_{k}^{N} L(\hat{\mathbf{y}}_k, \mathbf{y}_k)$$

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

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

# (Recap) Common Loss Functions for Regression
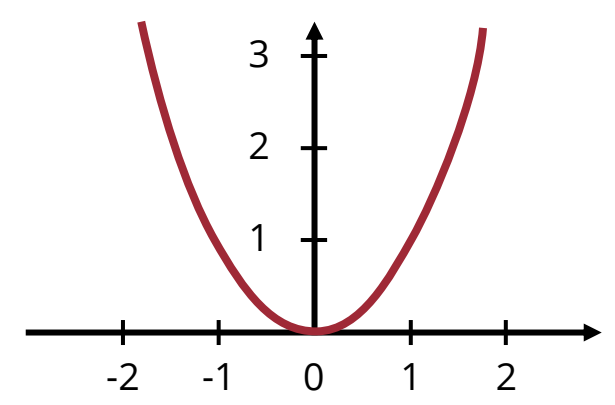
$$L(\hat{y}, y) = |\hat{y} - y|$$

**L1 loss**

**x**

**No activation function!**

$\hat{y}$          $y$

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

**Loss function**

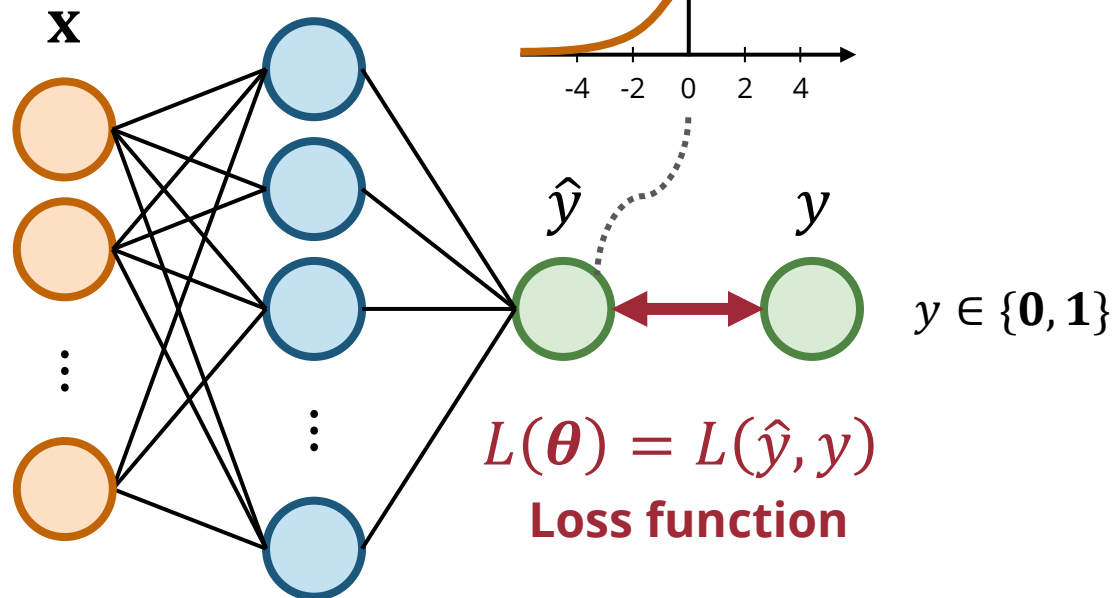$$Loss(\boldsymbol{\theta}) = \sum_{k}^{N} L(\hat{y}_k, y_k)$$
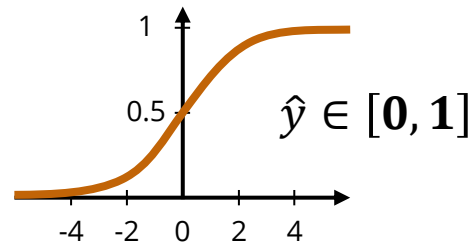
$$L(\hat{y}, y) = (\hat{y} - y)^2$$

**L2 loss**

# (Recap) Binary Cross Entropy for Binary Classification

- **Logistic regression** approaches classification like regression

$$Loss(\boldsymbol{\theta}) = \sum_{k}^{N} L(\hat{y}_k, y_k)$$

Sigmoid function



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

$\hat{y}$   $y$

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

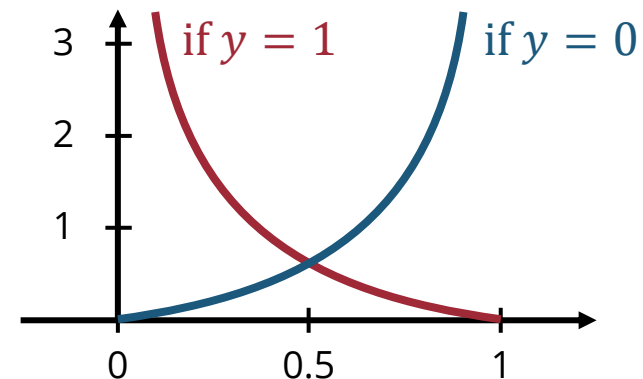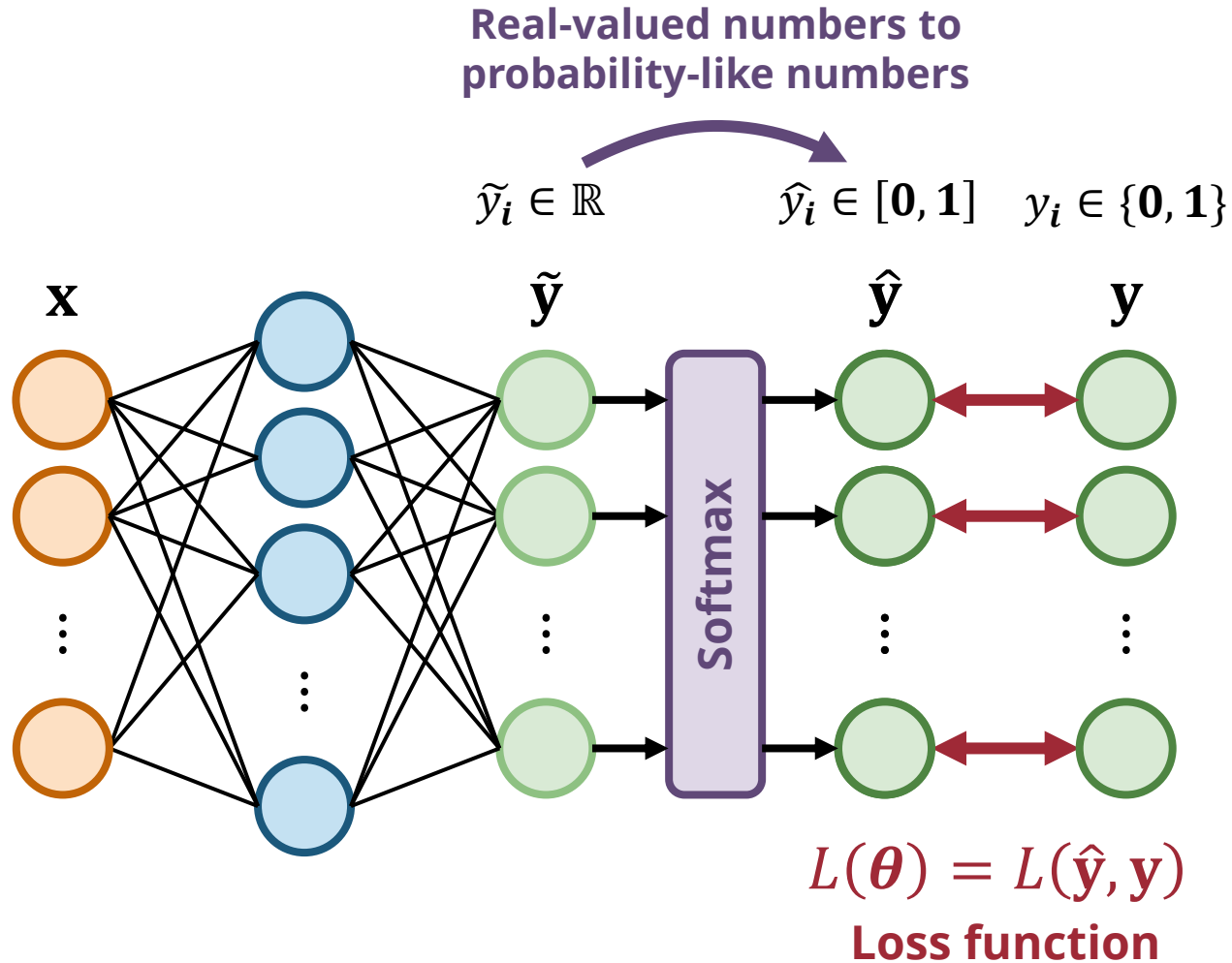$$L(\boldsymbol{\theta}) = L(\hat{y}, y)$$
**Loss function**

**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 $y = 0$

# (Recap) Cross Entropy for Multiclass Classification

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

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

$$\mathbf{x}$$

$$\tilde{\mathbf{y}} \qquad \hat{\mathbf{y}} \qquad \mathbf{y}$$

Softmax

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

**Loss function**

**Softmax**

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

**Cross entropy**

$$L(\hat{\mathbf{y}}, \mathbf{y}) = -\sum_{i}^{n} y_i \log \hat{y}_i$$

$$Loss(\boldsymbol{\theta}) = \sum_{k}^{N} L(\hat{\mathbf{y}}_k, \mathbf{y}_k)$$

# (Recap) Cross Entropy for Multiclass Classification

**Binary Cross Entropy**

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

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

**Cross Entropy**

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

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

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

**Log likelihood**

# (Recap) Training a Neural Network

<div style="border:2px solid #2e5e8c; background:#cfe2f3; border-radius:10px; padding:10px;">

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

</div>

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

<div style="border:2px solid #c87d1e; background:#fbe4cc; border-radius:10px; padding:10px;">

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

</div>

$$Loss(\boldsymbol{\theta}) = \sum_{k}^{N} L(\hat{\mathbf{y}}_k, \mathbf{y}_k)$$

<div style="border:2px solid #3e7d3e; background:#d9ead3; border-radius:10px; padding:10px;">

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

</div>

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

# (Recap) 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 $= -\eta \nabla f(w_t) < 0$

# (Recap) Forward Pass & Backward Pass



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

$\mathbf{x}$

$\hat{\mathbf{y}}$

**Backward pass**
**loss.backward()**

# Local Minima in Complex Loss Landscape



Local minima

Global minimum

**Solution 1**
**Use an optimizer with adaptive learning rate**

**Solution 2**
**Use a stochastic optimizer**

**Solution 3**
**Make the loss landscape smoother**

Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein, "Visualizing the Loss Landscape of Neural Nets," *NeurIPS*, 2018.

# (Recap) Comparison of Optimizers

- **Momentum**
  - Gets you out of spurious local minima
  - Allows the model to explore around

- **Gradient-based adaption**
  - Maintains steady improvement
  - Allows faster convergence



Optimizer Comparison

SDG
SGD with Momentum
AdaGrad
RMSprop
Adam

# (Recap) Mini-batch Gradient Descent

- **Intuition:** **Estimate** the gradient using **several random training samples**



**Loss**

| GD | SGD | Mini-batch GD |
| --- | --- | --- |
| batch size = $N$ | batch size = $1$ | $1 <$ batch size $< N$ |

# (Recap) Skip Connections

**Without skip connections**



**With skip connections**



Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein, "Visualizing the Loss Landscape of Neural Nets," *NeurIPS*, 2018.

# 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
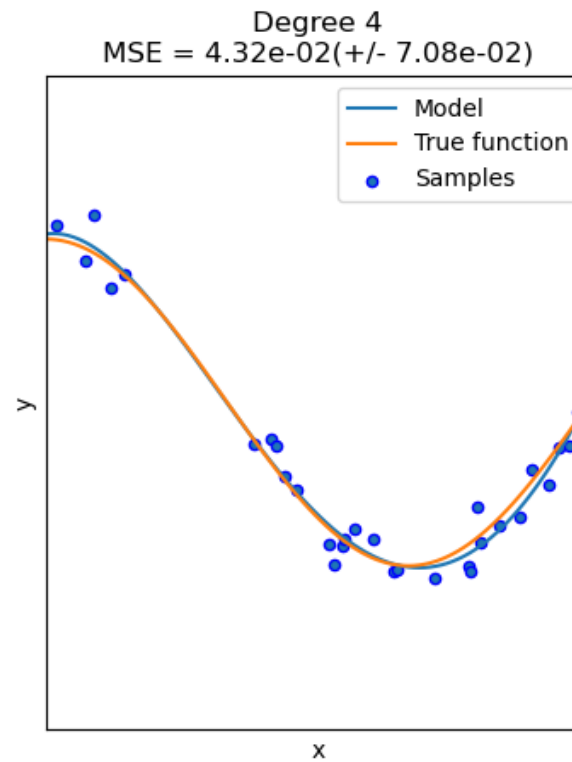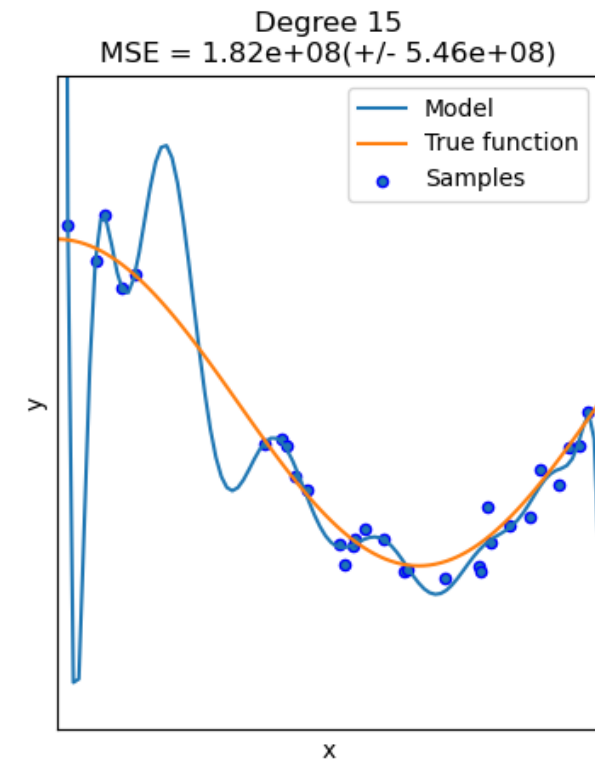  unseen samples in the target distribution**

# Overfitting & Underfitting



scikit-learn.org/stable/auto_examples/model_selection/plot_underfitting_overfitting.html

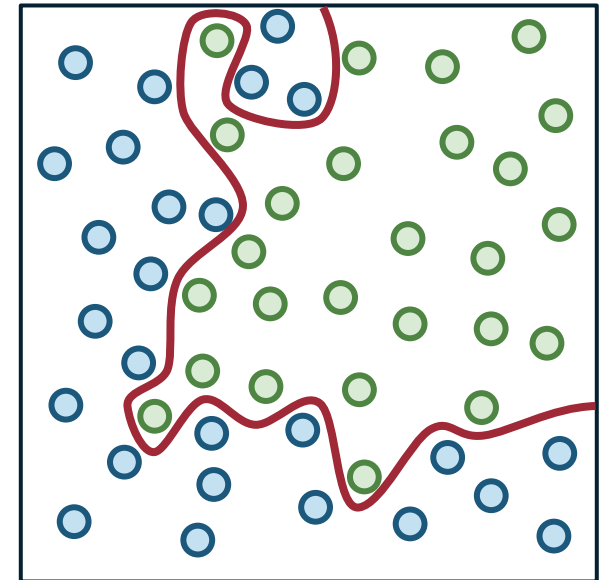# Overfitting & Underfitting

**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**

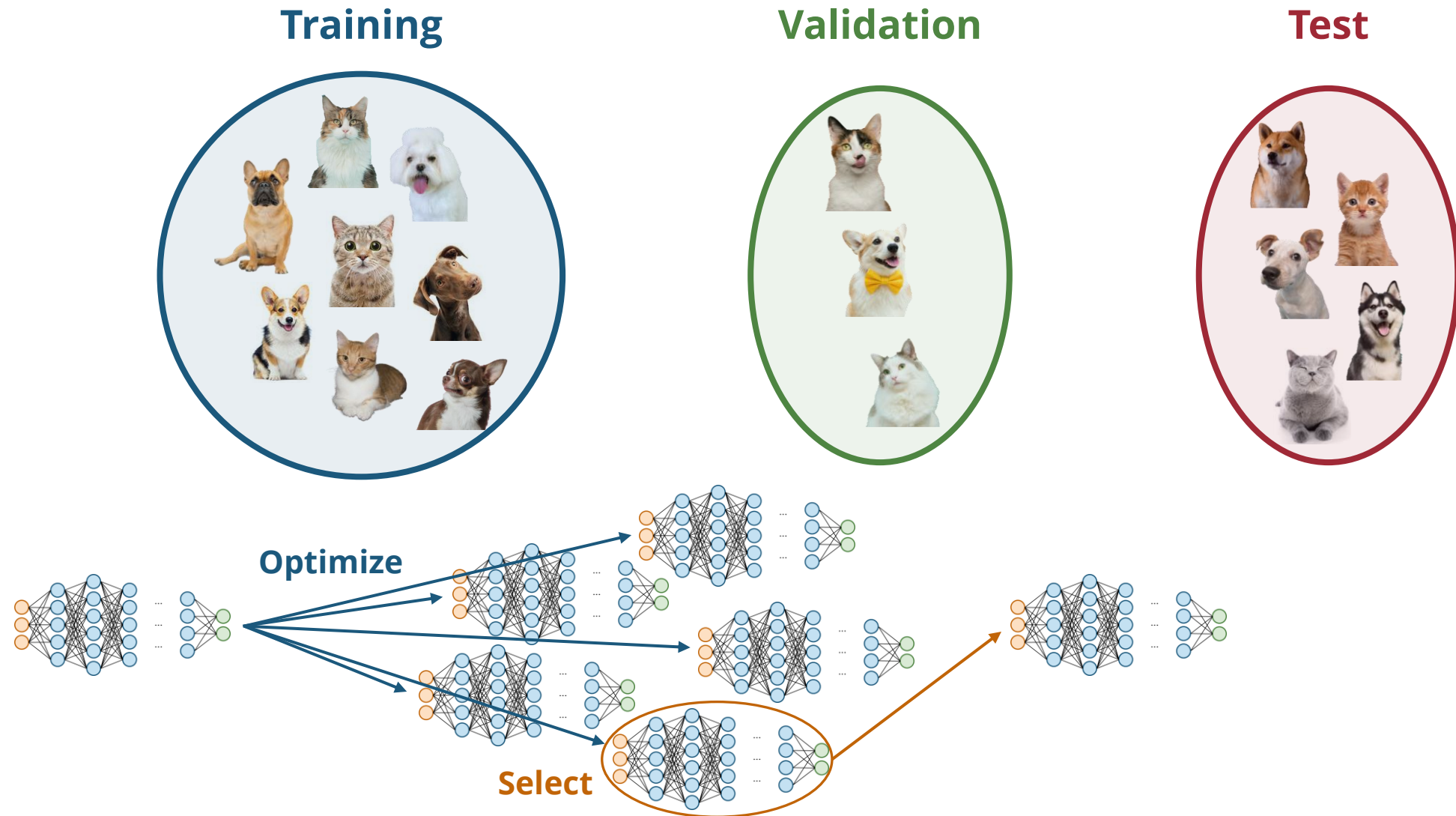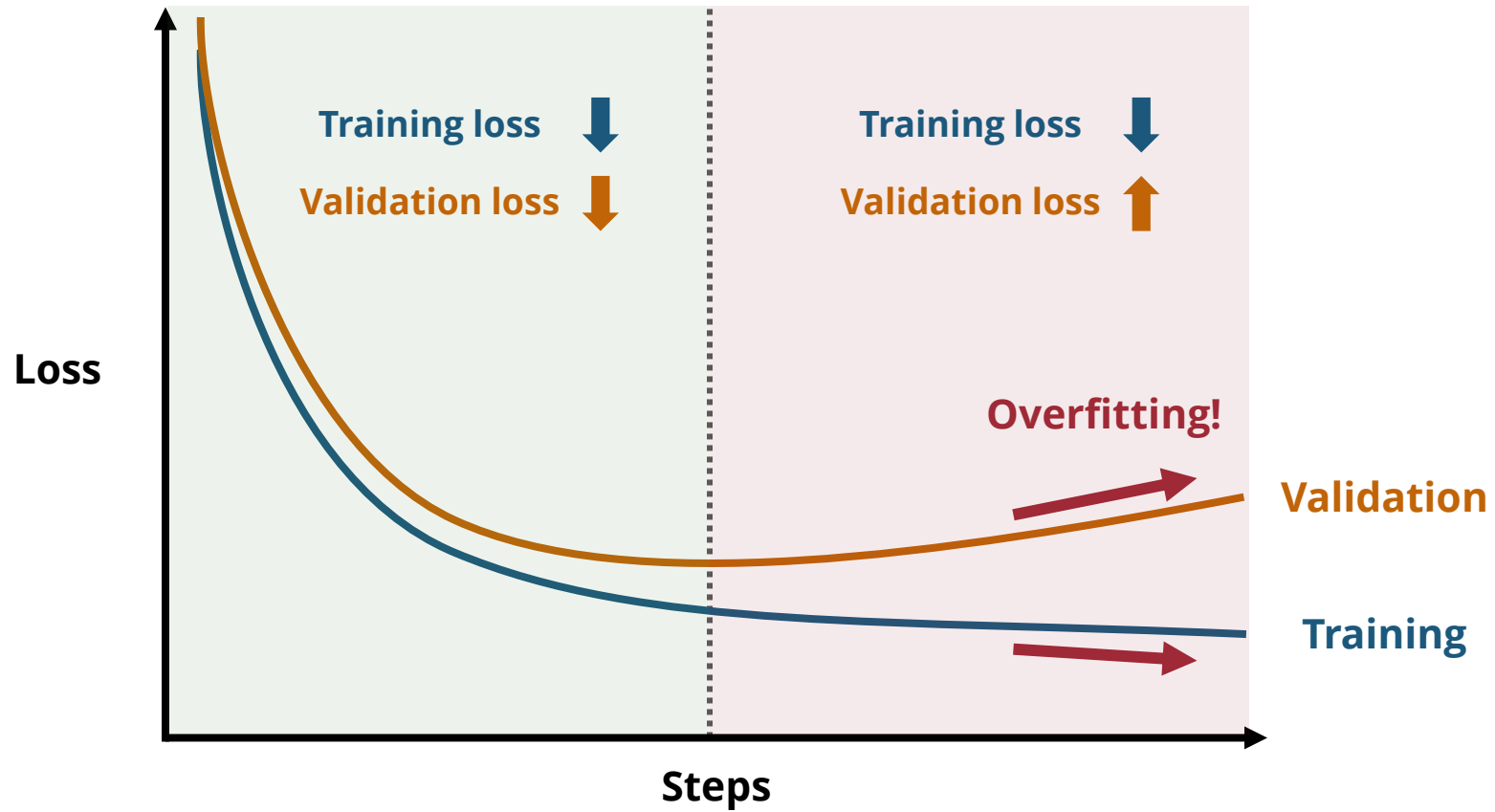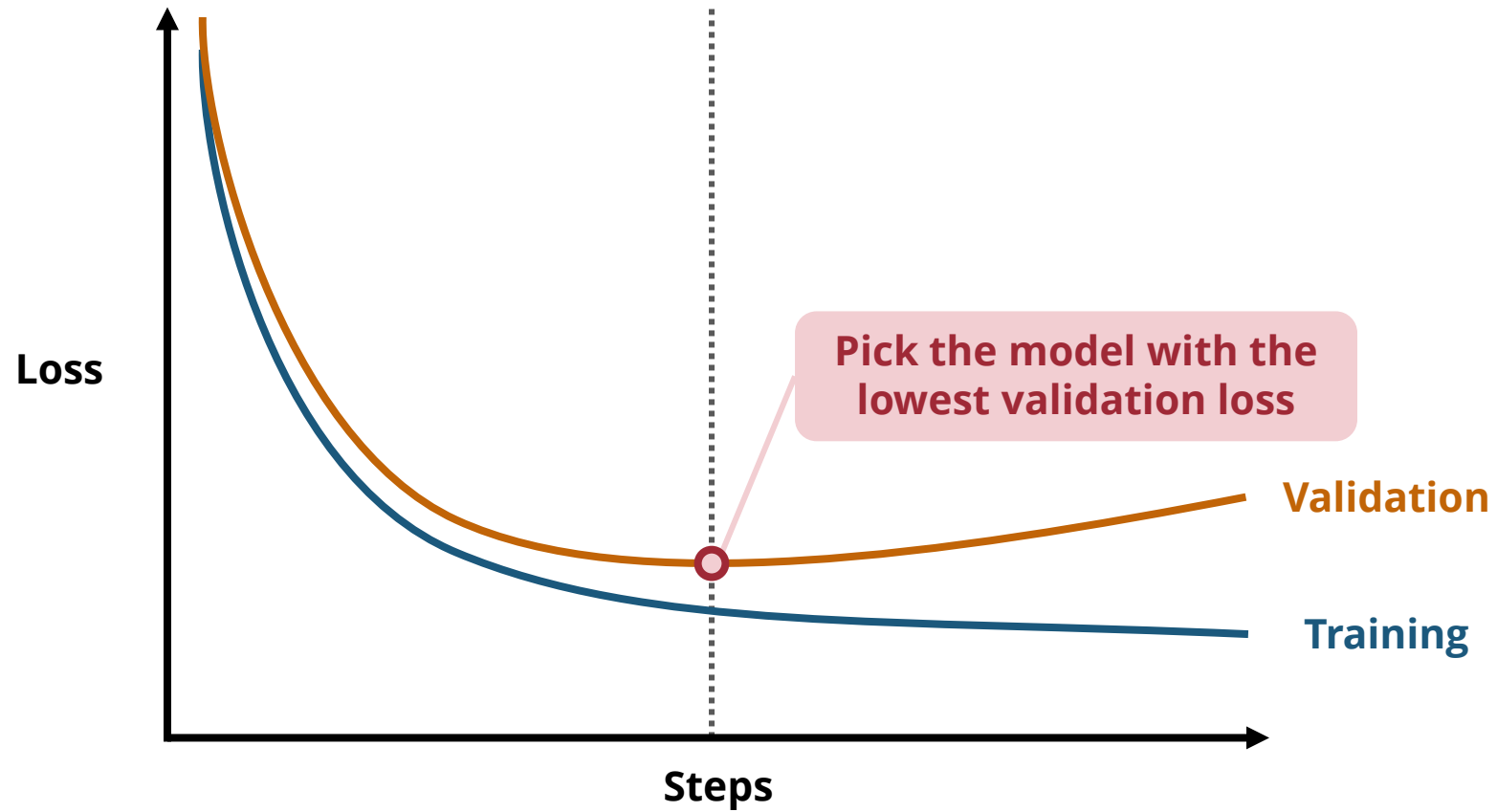**Test**

# Train–Validation–Test Split



**Training**

**Validation**

**Test**

# Training–Validation–Test Pipeline

**Training**

**Validation**

**Test**

Optimize

Select

# Training vs Validation Losses

# Training vs Validation Losses



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**

34

# Training vs Validation Losses



**Validation**

**Overfitting!**

**Loss**

**Possible solutions**
- Reduce the model size
- Apply dropout
- Add a regularizer
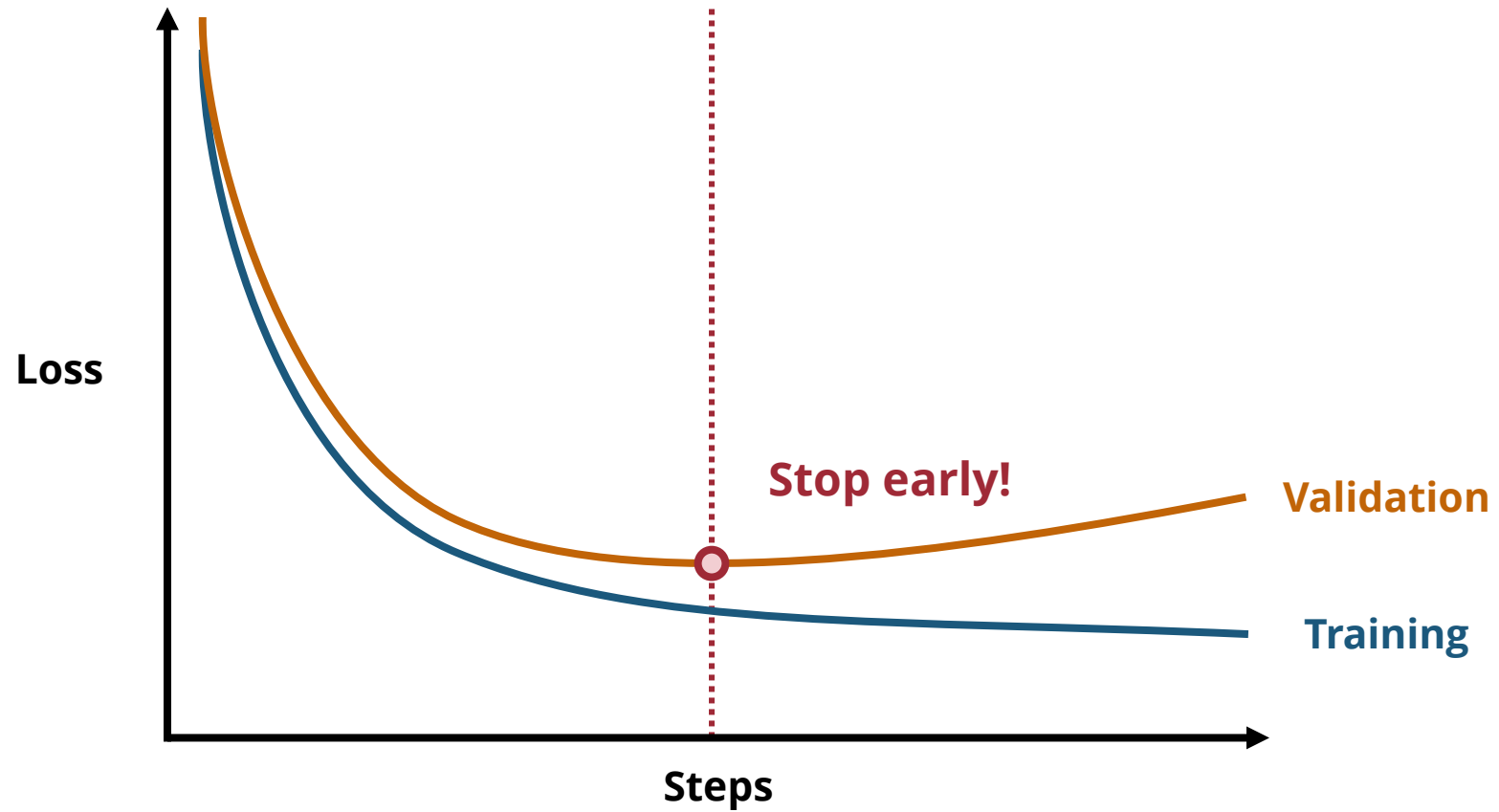
**Training**
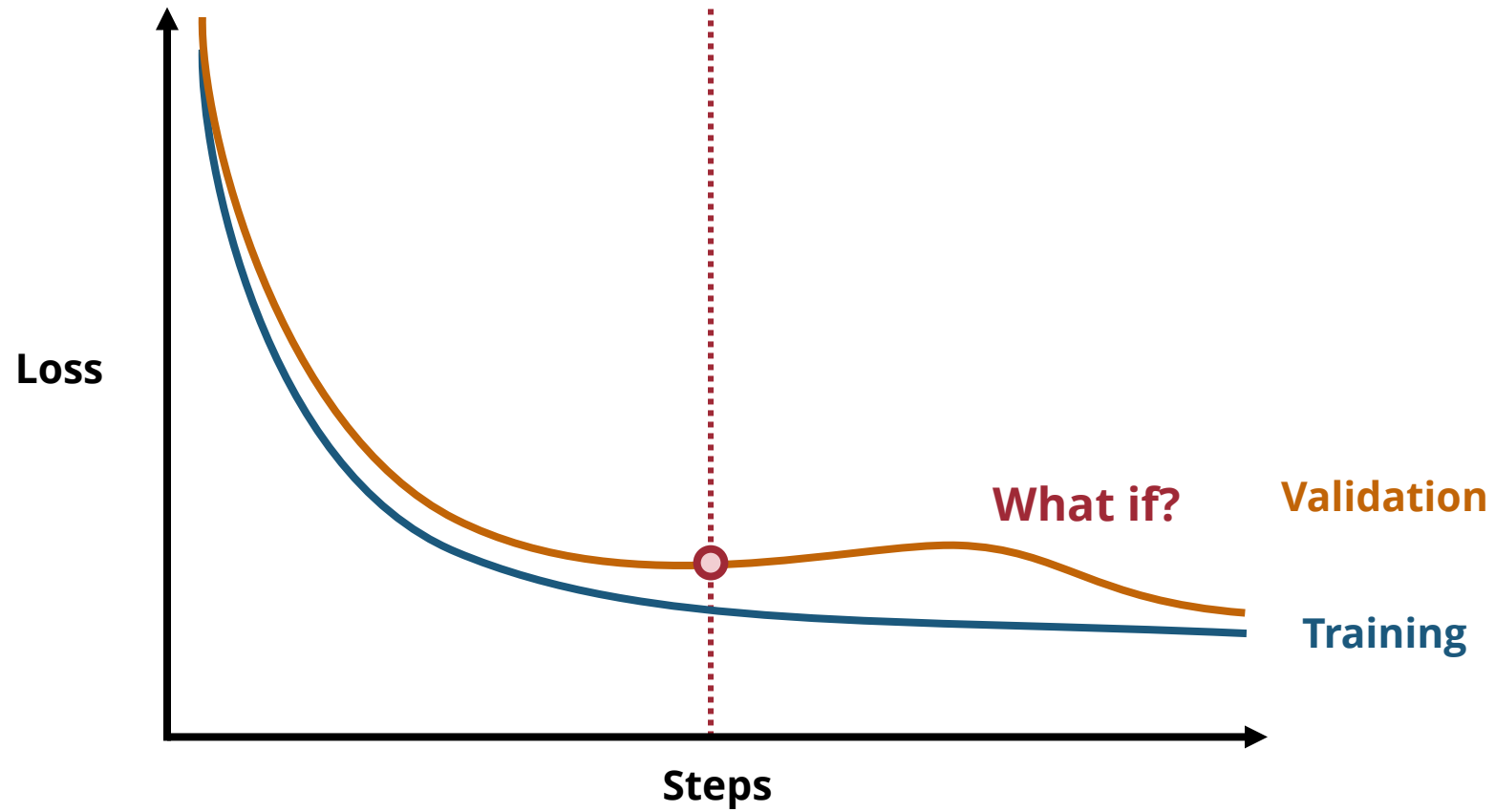
**Steps**

# 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
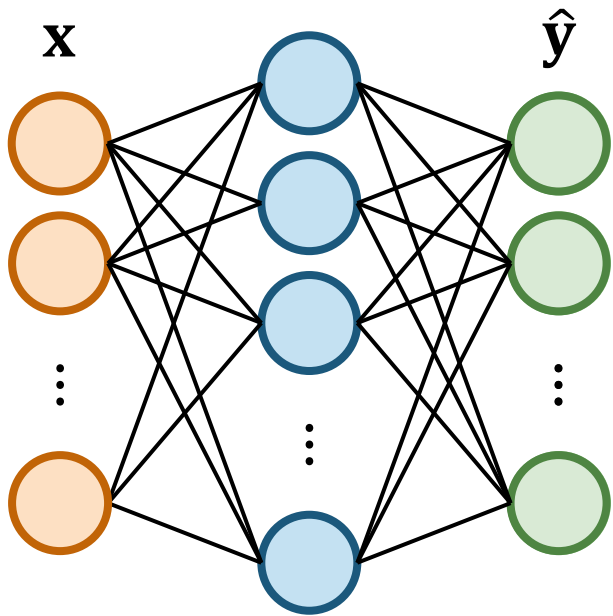
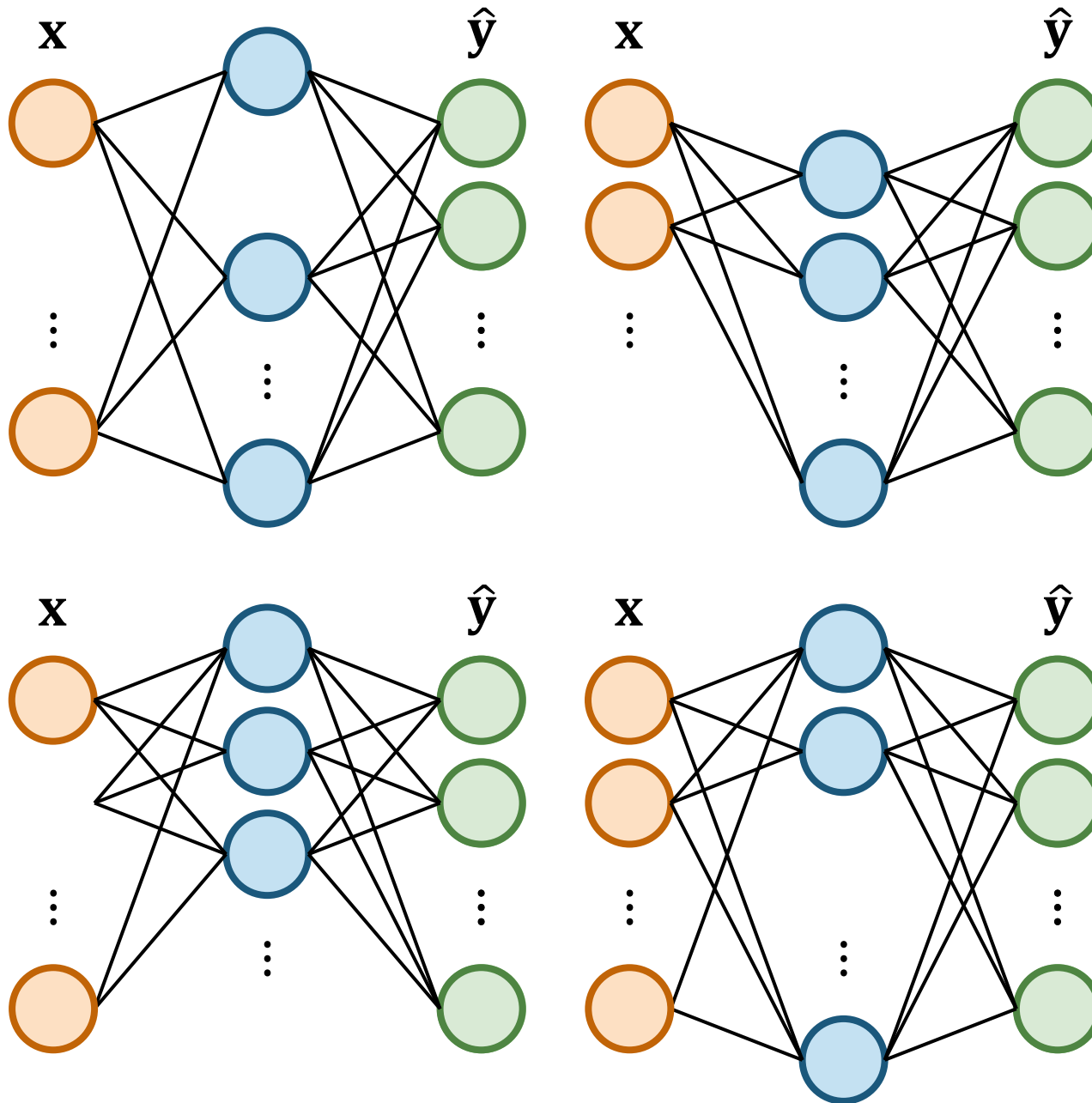# Overcoming Overfitting

# Early Stopping

# Early Stopping
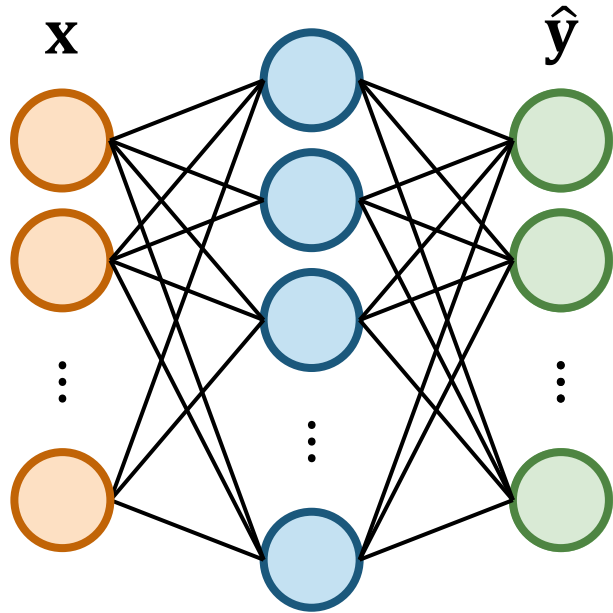
# Dropout



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

Dropout rate

# Dropout



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

Test
error
rate

Weight updates

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 regularization encourage smaller weights**