

RSM338: Machine Learning in Finance

Lecture 9: Neural Networks & Deep Learning | March 18–19, 2026

Kevin Mott

Rotman School of Management

Today's Goal

Today we introduce the last and most flexible model in the course: neural networks. To get there, we need to reconnect to a framework from Lecture 5 that has been quietly running behind every model we've studied since.

Today's roadmap:

1. **Deep learning = regression:** Connecting back to the Lecture 5 framework
2. **Opening the hood:** What's inside a neural network?
3. **The universal approximation theorem:** Why neural networks can represent any function
4. **Training neural networks:** Loss, backpropagation, optimizers, and the training loop
5. **Regularization:** Preventing overfitting in overparameterized models
6. **Demo:** Building a neural network in PyTorch

Part I: Deep Learning = Regression

The Regression Toolkit (Lecture 5 Recap)

In Lecture 5 we generalized OLS into a “choose your ingredients” framework. Here was the slide:

$$\hat{\theta} = \arg \min_{\theta} \left\{ \underbrace{\sum_{i=1}^n \square(y_i, f_{\theta}(\mathbf{x}_i))}_{\text{loss function}} + \underbrace{\lambda \cdot \text{Penalty}(\theta)}_{\text{regularization}} \right\}$$

where f_{θ} is the prediction function (parameterized by θ), \square is the loss measuring how wrong our predictions are, and $\lambda \cdot \text{Penalty}(\theta)$ is an optional regularization term that penalizes complex models.

Component	OLS Choice	Alternatives
Function f_{θ}	Linear: $\mathbf{x}^T \boldsymbol{\beta}$	Polynomial, tree, neural network
Loss \square	Squared error	Absolute error, Huber, cross-entropy
Penalty	None ($\lambda = 0$)	Ridge (L2), Lasso (L1), Elastic Net

$$\hat{\theta} = \arg \min_{\theta} \left\{ \underbrace{\sum_{i=1}^n \square(y_i, f_{\theta}(\mathbf{x}_i))}_{\text{loss function}} + \underbrace{\lambda \cdot \text{Penalty}(\theta)}_{\text{regularization}} \right\}$$

where f_{θ} is the prediction function (parameterized by θ), \square is the loss measuring how wrong our predictions are, and $\lambda \cdot \text{Penalty}(\theta)$ is an optional regularization term that penalizes complex models.

Component	OLS Choice	Alternatives
Function f_{θ}	Linear: $\mathbf{x}^T \boldsymbol{\beta}$	Polynomial, tree, neural network
Loss \square	Squared error	Absolute error, Huber, cross-entropy
Penalty	None ($\lambda = 0$)	Ridge (L2), Lasso (L1), Elastic Net

At the time, we said “non-linear f comes in later lectures.” Every lecture since has been filling in a specific choice for f_{θ} from the “Alternatives” column — we just didn’t always call it that.

How Every Model Fits the Framework

Here's where each model sits in that table:

Lecture	Model	What $f_{\theta}(\mathbf{x})$ is	Parameters θ
5	Linear regression	$\mathbf{x}^T \boldsymbol{\beta}$	Slopes and intercept
5	Ridge / Lasso	$\mathbf{x}^T \boldsymbol{\beta}$ (with penalty)	Slopes and intercept
7	Logistic regression	$\sigma(\mathbf{x}^T \mathbf{w})$	Weights and bias
8	Decision trees / ensembles	Piecewise-constant regions	Split rules and leaf values

In every case, the parameters θ had clear interpretations. You could look at a regression coefficient and say “ β_3 means a one-unit increase in X_3 is associated with a 0.5-unit increase in Y .” Even tree splits are interpretable: “if

Here's where each model sits in that table:

Lecture	Model	What $f_{\theta}(\mathbf{x})$ is	Parameters θ
5	Linear regression	$\mathbf{x}^T \boldsymbol{\beta}$	Slopes and intercept
5	Ridge / Lasso	$\mathbf{x}^T \boldsymbol{\beta}$ (with penalty)	Slopes and intercept
7	Logistic regression	$\sigma(\mathbf{x}^T \mathbf{w})$	Weights and bias
8	Decision trees / ensembles	Piecewise-constant regions	Split rules and leaf values

In every case, the parameters θ had clear interpretations. You could look at a regression coefficient and say “ $\beta_3 = 0.5$ means a one-unit increase in X_3 is associated with a 0.5-unit increase in Y .” Even tree splits are interpretable: “if momentum > 0.02 , go left.”

The Neural Network as f_{θ}

Today we add one more row to the table:

Lecture	Model	What $f_{\theta}(\mathbf{x})$ is	Parameters θ
9	Neural network	A composition of layers	Thousands (millions? billions?) of weights

A neural network is just another choice of f_{θ} — same loss functions, same regularization ideas. What's different:

- ▶ f_{θ} is built by **composing many simple operations** — layer after layer of weighted sums and nonlinear transformations
- ▶ θ contains thousands or millions of weights and biases
- ▶ The parameters have **no individual interpretation** — you cannot look at weight #4,817 and say what it “means”
- ▶ The model is a **black box**: features go in, predictions come out

Why Give Up Interpretability?

If we can't interpret the parameters, why would we use a neural network?

Because the other models on our list have **structural limitations**:

- ▶ Linear regression can only learn linear relationships
- ▶ Logistic regression can only learn linear decision boundaries
- ▶ Decision trees partition the feature space into rectangles
- ▶ Even Random Forests and XGBoost are combinations of piecewise-constant functions

Neural networks trade interpretability for **flexibility**:

- ▶ By composing many layers, they can approximate essentially any continuous function (we'll make this precise later)
- ▶ The cost: you lose the ability to explain *why* the model makes a particular prediction
- ▶ Worth it when prediction accuracy matters more than explanation (fraud detection, algorithmic trading)

Part II: Opening the Hood

A Single Neuron

The building block of every neural network is a **neuron** (also called a **unit** or **node**). A single neuron does two things:

1. Compute a weighted sum of its inputs, plus a bias term
2. Apply a nonlinear function (called an **activation function**) to the result

$$a = g\left(\sum_{j=1}^p w_j x_j + b\right) = g(\mathbf{w}^\top \mathbf{x} + b)$$

where:

- ▶ x_1, x_2, \dots, x_p are the inputs (our features)
- ▶ w_1, w_2, \dots, w_p are the **weights** — how much each input contributes
- ▶ b is the **bias** — an offset (like an intercept)
- ▶ $g(\cdot)$ is the **activation function** — introduces nonlinearity

The weighted sum $\mathbf{w}^\top \mathbf{x} + b$ should look familiar: it's the same linear combination from logistic regression. The

A Neuron Is Logistic Regression

If we choose the sigmoid as our activation function, a single neuron is exactly logistic regression:

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

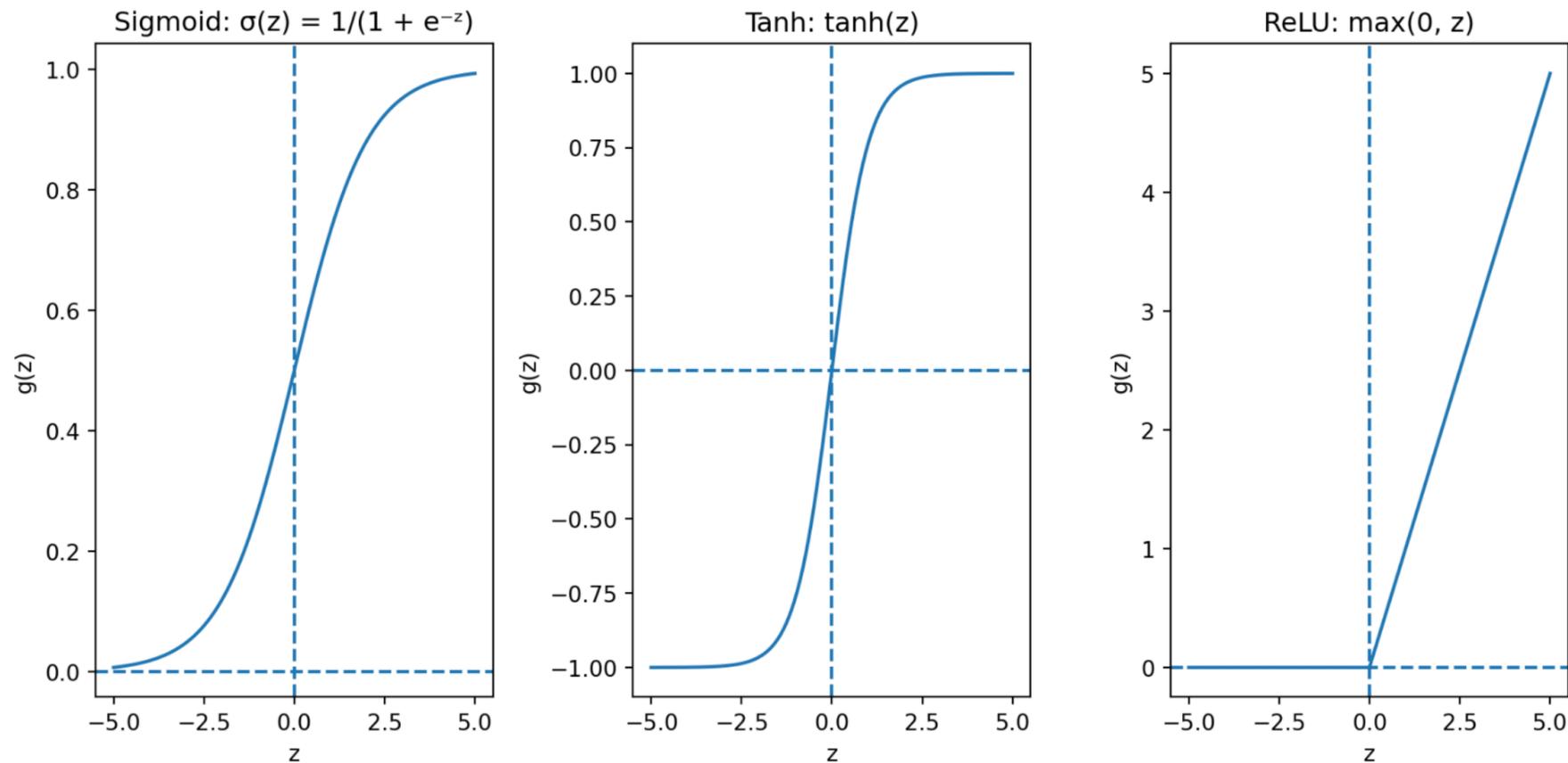
This is the same formula from the classification lecture. One neuron with a sigmoid activation = logistic regression. A neural network is what happens when we stack many of these neurons together.

If we choose no activation function at all (the “identity” activation, $g(z) = z$), a single neuron is just linear regression:
 $\hat{y} = \mathbf{w}^\top \mathbf{x} + b$.

So neural networks are not a departure from what we’ve learned — they’re a **generalization**. The models we already know are special cases.

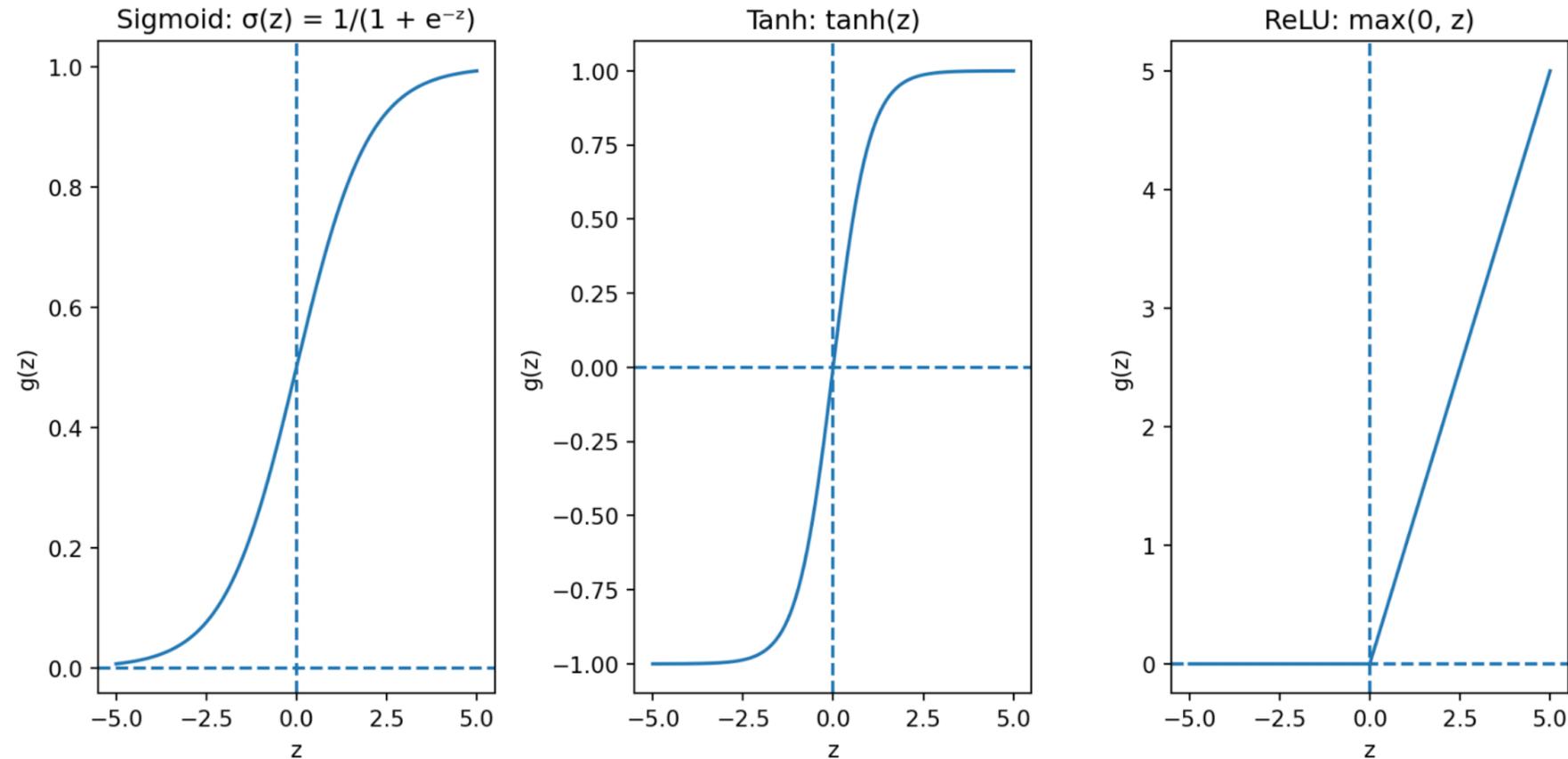
Activation Functions

The activation function g controls what kind of nonlinearity each neuron introduces. Here are the most common choices:



- ▶ **Sigmoid (σ):** Squashes to (0, 1). Used for binary classification output layers. Suffers from vanishing gradients in deep networks.
- ▶ **Tanh:** Squashes to (-1, 1). Centred at zero (better for hidden layers). Still has vanishing gradients.

The activation function g controls what kind of nonlinearity each neuron introduces. Here are the most common choices:



- ▶ **Sigmoid (σ):** Squashes to (0, 1). Used for binary classification output layers. Suffers from vanishing gradients in deep networks.
- ▶ **Tanh:** Squashes to (-1, 1). Centred at zero (better for hidden layers). Still has vanishing gradients.
- ▶ **ReLU (Rectified Linear Unit):** $g(z) = \max(0, z)$. The modern default. Simple, fast, avoids vanishing gradients for positive inputs. Downside: outputs exactly zero for negative inputs (“dead neurons”).

Why Activation Functions Matter

Without activation functions, stacking layers doesn't help. Suppose we have two layers of linear transformations:

$$\text{Layer 1: } \mathbf{h} = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$$

$$\text{Layer 2: } \hat{y} = \mathbf{W}_2 \mathbf{h} + \mathbf{b}_2$$

Substituting:

$$\hat{y} = \mathbf{W}_2 (\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2 = \underbrace{(\mathbf{W}_2 \mathbf{W}_1)}_{\mathbf{W}'} \mathbf{x} + \underbrace{(\mathbf{W}_2 \mathbf{b}_1 + \mathbf{b}_2)}_{\mathbf{b}'}$$

This is still a linear function of \mathbf{x} . No matter how many linear layers we stack, the result is always equivalent to a single linear transformation. The depth buys us nothing.

Activation functions break this collapse. When we insert g between layers:

$$\hat{y} = \mathbf{W}_2 g(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$$

the g prevents us from simplifying. Each layer now genuinely transforms the representation. This is why depth matters.

From One Neuron to a Layer

A **layer** is a collection of neurons that all receive the same inputs but have different weights. With d neurons and p input features, the whole layer is one matrix operation:

$$\mathbf{h} = g(\mathbf{W}\mathbf{x} + \mathbf{b})$$

where:

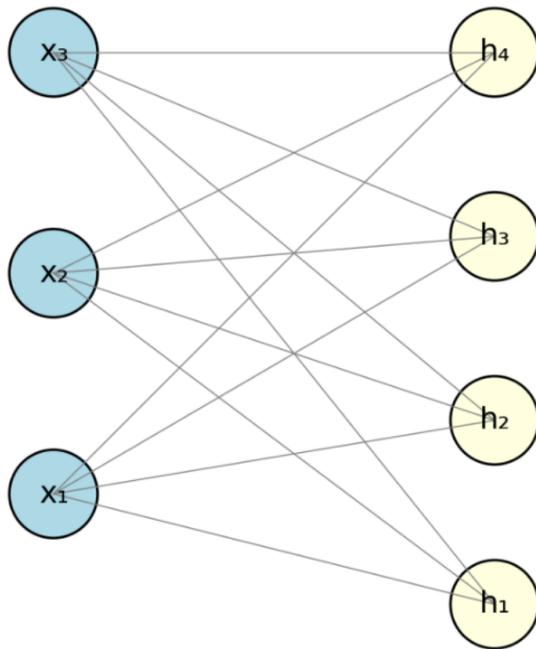
- ▶ \mathbf{W} is a $d \times p$ **weight matrix** — row k contains the weights for neuron k
- ▶ \mathbf{b} is a $d \times 1$ **bias vector** — one bias per neuron
- ▶ g is applied element-wise to each neuron's output
- ▶ \mathbf{h} is the $d \times 1$ **output vector** — one value per neuron

One layer: 3 inputs \rightarrow 4 hidden neurons (12 weights + 4 biases = 16 parameters)



- ▶ \mathbf{W} is a $d \times p$ weight matrix — row k contains the weights for neuron k
- ▶ \mathbf{b} is a $d \times 1$ bias vector — one bias per neuron
- ▶ g is applied element-wise to each neuron's output
- ▶ \mathbf{h} is the $d \times 1$ output vector — one value per neuron

One layer: 3 inputs \rightarrow 4 hidden neurons (12 weights + 4 biases = 16 parameters)



Each line in the diagram represents one weight. With 3 inputs and 4 neurons, that's $3 \times 4 = 12$ weights plus 4 biases = 16 parameters. Already more than the 4 parameters in a 3-variable linear regression.

**Rotman
Commerce**

Stacking Layers: The Feed-Forward Network

A feed-forward neural network (or multilayer perceptron / MLP) stacks layers — each layer's output feeds into the next:

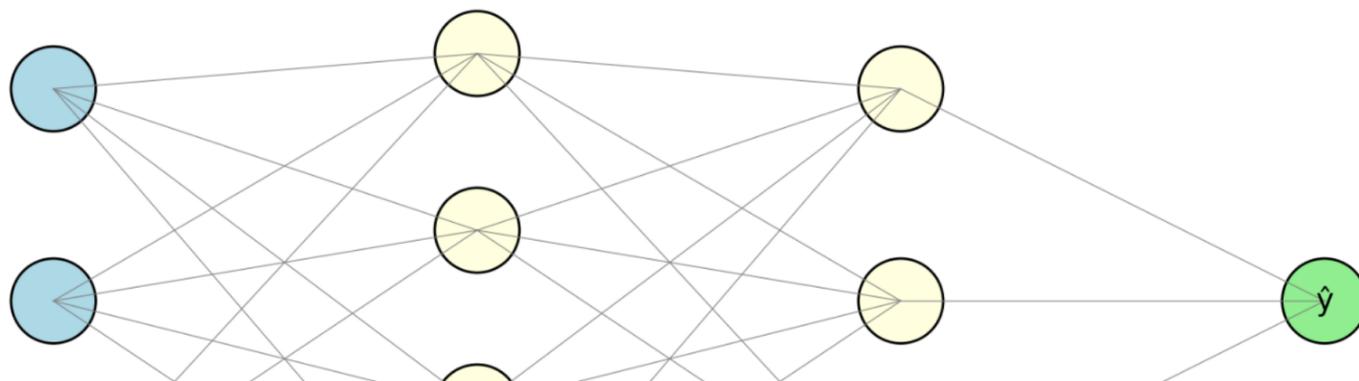
$$\mathbf{h}^{(1)} = g(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

$$\mathbf{h}^{(2)} = g(\mathbf{W}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}^{(2)})$$

$$\hat{y} = \mathbf{W}^{(3)}\mathbf{h}^{(2)} + \mathbf{b}^{(3)}$$

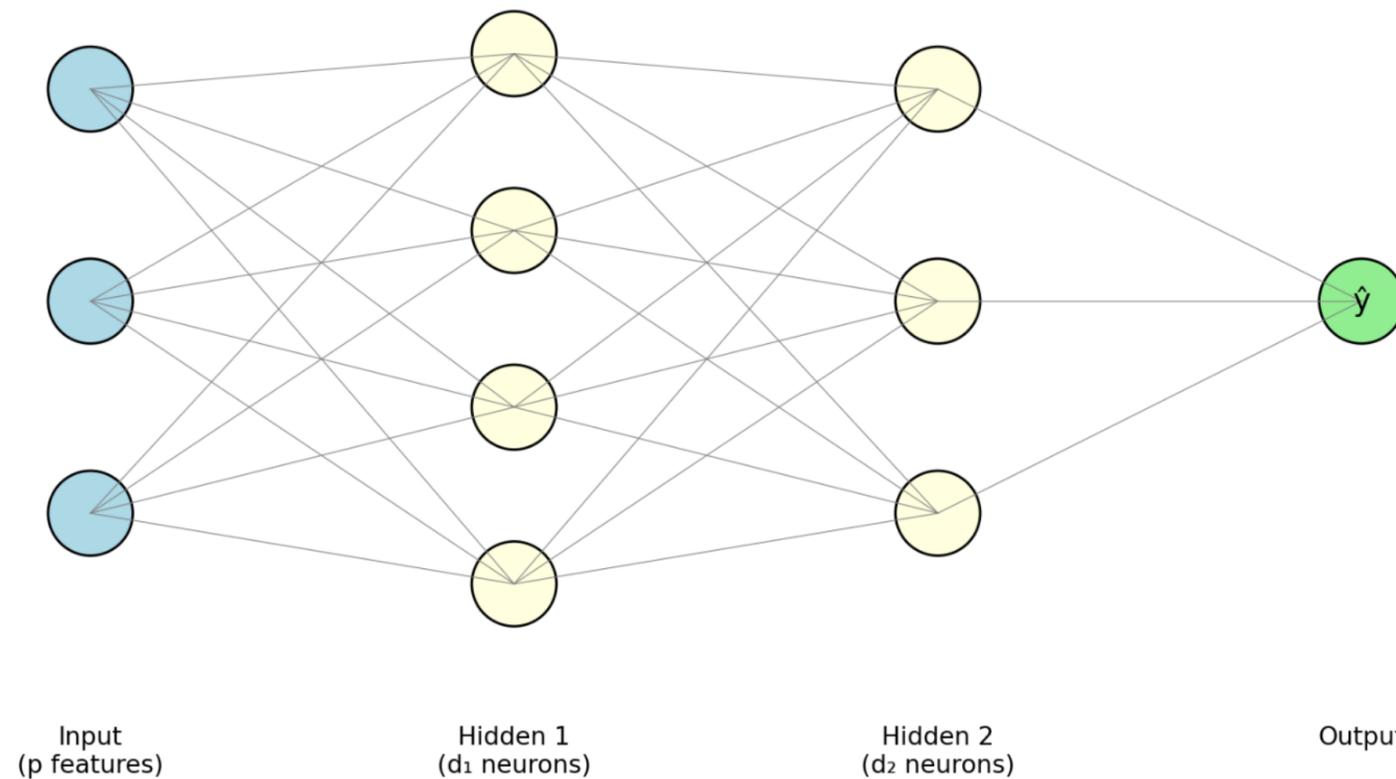
- ▶ $\mathbf{h}^{(\ell)}$ = **hidden layer** — “hidden” because we never observe these values directly
- ▶ The final layer is the **output layer**, producing the prediction

Feed-Forward Neural Network: Input → Hidden 1 → Hidden 2 → Output



- ▶ $\mathbf{h}^{(\ell)}$ = **hidden layer** — “hidden” because we never observe these values directly
- ▶ The final layer is the **output layer**, producing the prediction

Feed-Forward Neural Network: Input → Hidden 1 → Hidden 2 → Output



- ▶ Information flows one direction: input → hidden layers → output (no loops or feedback)
- ▶ **Depth** = number of hidden layers. Modern “deep learning” models may have dozens or hundreds.

The Parameter Explosion

Let's count parameters for a concrete example. Suppose we have:

- ▶ $p = 10$ input features
- ▶ Hidden layer 1: 64 neurons
- ▶ Hidden layer 2: 32 neurons
- ▶ Output: 1 neuron (regression) or K neurons (classification)

Connection	Weights	Biases	Total
Input → Hidden 1	$10 \times 64 = 640$	64	704
Hidden 1 → Hidden 2	$64 \times 32 = 2,048$	32	2,080
Hidden 2 → Output	$32 \times 1 = 32$	1	33
Total			2,817

- ▶ Hidden layer 1: 64 neurons
- ▶ Hidden layer 2: 32 neurons
- ▶ Output: 1 neuron (regression) or K neurons (classification)

Connection	Weights	Biases	Total
Input → Hidden 1	$10 \times 64 = 640$	64	704
Hidden 1 → Hidden 2	$64 \times 32 = 2,048$	32	2,080
Hidden 2 → Output	$32 \times 1 = 32$	1	33
Total			2,817

- ▶ 10 features, two modest hidden layers → nearly **3,000 parameters** (vs. 11 for linear regression)
- ▶ ResNet-50 (image classifier): **25 million** parameters
- ▶ Large language models: **billions**
- ▶ Parameters typically far exceed training observations — a regime where OLS wouldn't even have a unique solution

The Output Layer

The final layer of the network depends on the task:

Regression (predicting a continuous value like returns):

- ▶ One output neuron with **no activation** (identity): $\hat{y} = \mathbf{w}^\top \mathbf{h} + b$
- ▶ The output can be any real number

Binary classification (predicting default / no default):

- ▶ One output neuron with **sigmoid activation**: $\hat{y} = \sigma(\mathbf{w}^\top \mathbf{h} + b)$
- ▶ The output is a probability between 0 and 1
- ▶ Same as logistic regression's output, but \mathbf{h} is a learned representation instead of raw features

Multi-class classification (predicting which of K sectors):

- ▶ K output neurons with **softmax activation**:

$$\hat{y}_k = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

Regression (predicting a continuous value like returns):

- ▶ One output neuron with **no activation** (identity): $\hat{y} = \mathbf{w}^\top \mathbf{h} + b$
- ▶ The output can be any real number

Binary classification (predicting default / no default):

- ▶ One output neuron with **sigmoid activation**: $\hat{y} = \sigma(\mathbf{w}^\top \mathbf{h} + b)$
- ▶ The output is a probability between 0 and 1
- ▶ Same as logistic regression's output, but \mathbf{h} is a learned representation instead of raw features

Multi-class classification (predicting which of K sectors):

- ▶ K output neurons with **softmax activation**:

$$\hat{y}_k = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$

- ▶ Outputs sum to 1 — a probability distribution over classes

Hidden layers learn a useful representation; the output layer translates it into the prediction format we need.

Putting It All Together: The Full Picture

A neural network is just a function $f_{\theta}(\mathbf{x})$ that we plug into the same regression framework from Lecture 5:

$$\hat{\theta} = \arg \min_{\theta} \left\{ \sum_{i=1}^n \square(y_i, f_{\theta}(\mathbf{x}_i)) + \lambda \cdot \text{Penalty}(\theta) \right\}$$

What makes it different:

- ▶ f_{θ} is a composition of layers: weighted sums \rightarrow activations \rightarrow weighted sums \rightarrow activations $\rightarrow \dots \rightarrow$ output
- ▶ θ = all the weights and biases across all layers (thousands or millions of numbers)
- ▶ There is no closed-form solution. We must use **iterative optimization** (gradient descent) to find the best θ
- ▶ The parameters θ have no individual interpretation — the model is a black box

The loss \square , the regularization penalty, the train/test split, cross-validation — all the machinery from earlier lectures carries over. The only new ingredient is the architecture of f_{θ} .

Part III: The Universal Approximation Theorem

Can Neural Networks Actually Learn Anything?

How flexible is a neural network, exactly? Could it represent *any* relationship?

Yes — with a caveat.

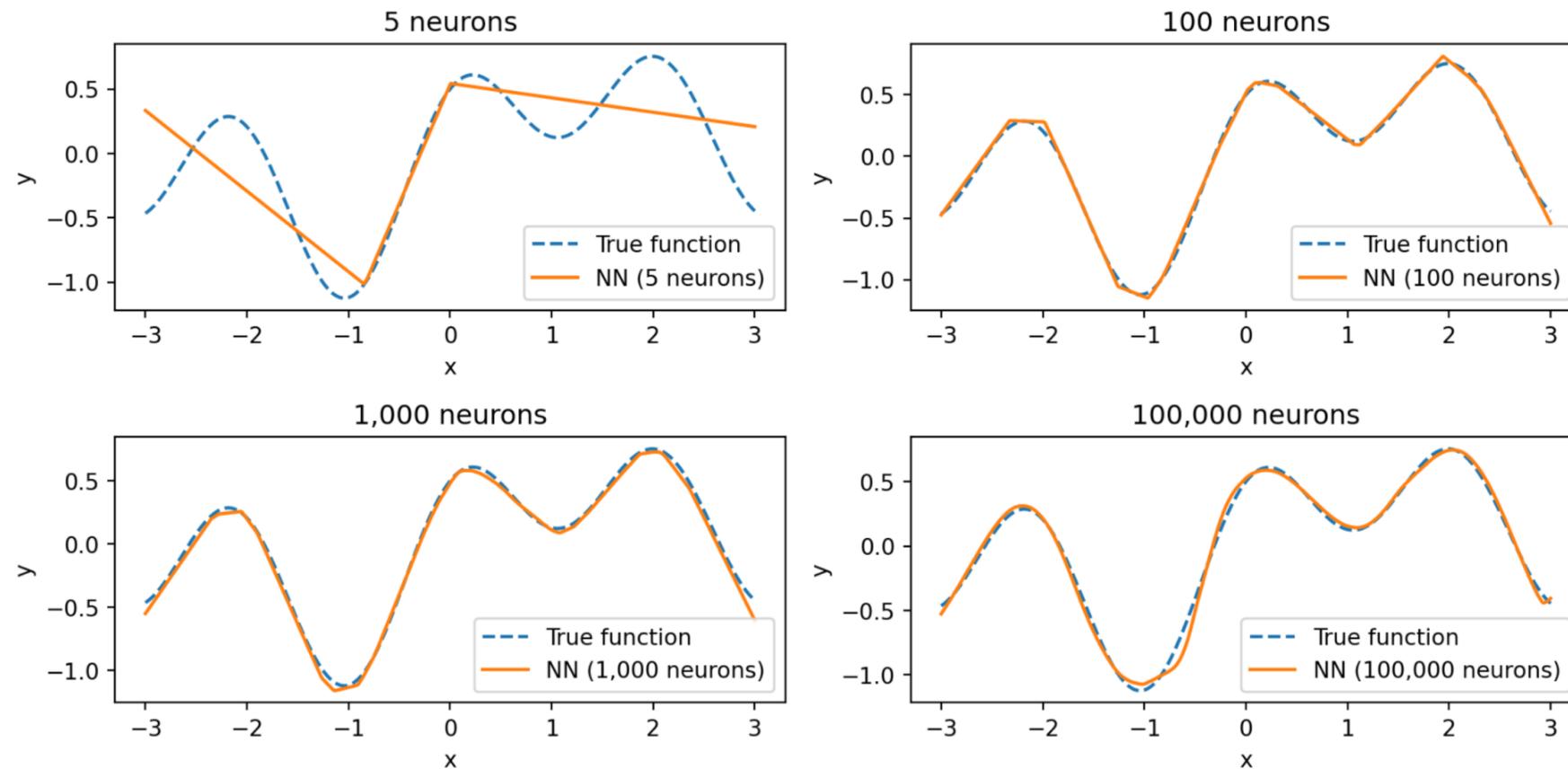
The Universal Approximation Theorem (Cybenko, 1989; Hornik, 1991):

A feed-forward neural network with a single hidden layer containing a sufficient number of neurons can approximate any continuous function on a compact domain to arbitrary accuracy.

In plain language: given enough neurons, a one-hidden-layer network can get as close as you want to any smooth function.

What the Theorem Promises

The theorem says a wide enough network can represent extremely complex functions:



With 5 neurons, the network captures only the broad shape. With 100, it gets most of the detail. With 1,000 and 100,000 neurons, the fit becomes nearly perfect — the network has enough capacity to match every wiggle of the true function. The universal approximation theorem guarantees that we can get as close as we want — if we use enough neurons.

What the Theorem Does NOT Promise

The theorem is an **existence** result, not a practical guarantee. It says a solution *exists*, not that we can *find* it.

It doesn't promise:

- 1. That gradient descent will find the right weights.** The loss landscape of a neural network is non-convex — it has many local minima and saddle points. The optimization could get stuck.
- 2. How many neurons you need.** The theorem says “a sufficient number,” but that number could be impractically large.
- 3. That the network will generalize.** A network with enough parameters can memorize the training data perfectly (just like a high-degree polynomial). The theorem says nothing about performance on new data.
- 4. That one hidden layer is the best architecture.** The theorem proves one wide layer is enough, but in practice, **deeper networks with fewer neurons per layer** are more efficient. They can represent certain functions with far fewer total parameters than a single wide layer would need.

This last point is the practical motivation for **deep** learning. Depth is not about theoretical power (one layer is enough) — it's about **efficiency** and **learning useful representations** at each level.

Part IV: Training Neural Networks

The Training Problem

We need to find the parameters θ (all weights and biases) that minimize the loss:

$$\theta^* = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \square(y_i, f_{\theta}(\mathbf{x}_i))$$

- ▶ Linear regression: closed-form solution
- ▶ Neural network: **no closed-form** — the loss is **non-convex** (many bumps and valleys) because of nonlinear activations

We solve it with **gradient descent**: start with random weights, compute the gradient, step in the opposite direction:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \square$$

where η (eta) is the learning rate — how large a step we take.

Loss Functions

Same losses from earlier lectures — they carry over directly.

Regression:

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

Mean squared error — same as OLS.

Binary classification:

$$\square_{\text{BCE}} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

Binary cross-entropy — same as logistic regression. Heavily penalizes confident wrong predictions.

Multi-class classification (K classes):

$$\square_{\text{CE}} = -\frac{1}{n} \sum \sum_{k=1}^K y_{ik} \log(\hat{y}_{ik})$$

Regression:

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

Mean squared error — same as OLS.

Binary classification:

$$\square_{\text{BCE}} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

Binary cross-entropy — same as logistic regression. Heavily penalizes confident wrong predictions.

Multi-class classification (K classes):

$$\square_{\text{CE}} = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K y_{ik} \log(\hat{y}_{ik})$$

where $y_{ik} = 1$ if observation i belongs to class k and 0 otherwise.

Backpropagation: How Does the Network Learn?

With thousands of parameters, how do we compute the gradient $\nabla_{\theta} \square$?

Backpropagation (Rumelhart, Hinton & Williams, 1986) computes the gradient of the loss with respect to every weight. It works because a neural network is a **nested composition** of functions. For a 3-layer network:

$$\begin{aligned}\mathbf{h}^{(1)} &= g(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) \\ \mathbf{h}^{(2)} &= g(\mathbf{W}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}^{(2)}) \\ \hat{y} &= \mathbf{W}^{(3)}\mathbf{h}^{(2)} + \mathbf{b}^{(3)}\end{aligned}$$

Substituting, the prediction is a nested composition:

$$\hat{y} = f^{(3)}\left(f^{(2)}\left(f^{(1)}(\mathbf{x}) \right) \right)$$

where $f^{(\ell)}(\cdot) = g(\mathbf{W}^{(\ell)}(\cdot) + \mathbf{b}^{(\ell)})$. The **chain rule** differentiates through each layer:

$$\frac{\partial \square}{\partial \mathbf{W}^{(1)}} = \underbrace{\frac{\partial \square}{\partial \hat{y}}}_{\text{}} \cdot \underbrace{\frac{\partial \hat{y}}{\partial \mathbf{h}^{(2)}}}_{\text{}} \cdot \underbrace{\frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}}}_{\text{}} \cdot \underbrace{\frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{W}^{(1)}}}_{\text{}}$$

$$\mathbf{h}^{(2)} = \mathbf{g}(\mathbf{W}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}^{(2)})$$

$$\hat{y} = \mathbf{W}^{(3)}\mathbf{h}^{(2)} + \mathbf{b}^{(3)}$$

Substituting, the prediction is a nested composition:

$$\hat{y} = f^{(3)}\left(f^{(2)}\left(f^{(1)}(\mathbf{x})\right)\right)$$

where $f^{(\ell)}(\cdot) = \mathbf{g}(\mathbf{W}^{(\ell)}(\cdot) + \mathbf{b}^{(\ell)})$. The **chain rule** differentiates through each layer:

$$\frac{\partial \square}{\partial \mathbf{W}^{(1)}} = \underbrace{\frac{\partial \square}{\partial \hat{y}}}_{\text{output}} \cdot \underbrace{\frac{\partial \hat{y}}{\partial \mathbf{h}^{(2)}}}_{\text{layer 3}} \cdot \underbrace{\frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{h}^{(1)}}}_{\text{layer 2}} \cdot \underbrace{\frac{\partial \mathbf{h}^{(1)}}{\partial \mathbf{W}^{(1)}}}_{\text{layer 1}}$$

- ▶ Each factor involves only one layer's local operation — easy to compute
- ▶ **Forward pass:** compute the loss
- ▶ **Backward pass:** multiply these local derivatives from right to left
- ▶ Weights that contributed more to the error get larger updates
- ▶ Modern frameworks (PyTorch, JAX) do this automatically

The Vanishing Gradient Problem

The chain rule multiplies many factors together. If each is small, the gradient **shrinks exponentially**:

- ▶ If each layer's gradient is 0.25, then after 4 layers: $0.25^4 \approx 0.004$ — nearly zero
- ▶ Early layers barely learn because the error signal is too weak by the time it reaches them

Why this happens with sigmoid/tanh: Their gradients are always < 1 (sigmoid's max is 0.25), so the chain rule multiplies many small numbers together.

Why ReLU helps: Its gradient is either 0 or 1. For positive inputs, the gradient passes through unchanged — no shrinkage.

Stochastic Gradient Descent and Mini-Batches

Standard gradient descent uses the **entire** training set at each step — slow with millions of observations.

Stochastic Gradient Descent (SGD) uses a small random subset (a **mini-batch**) instead:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \square_{\text{batch}}$$

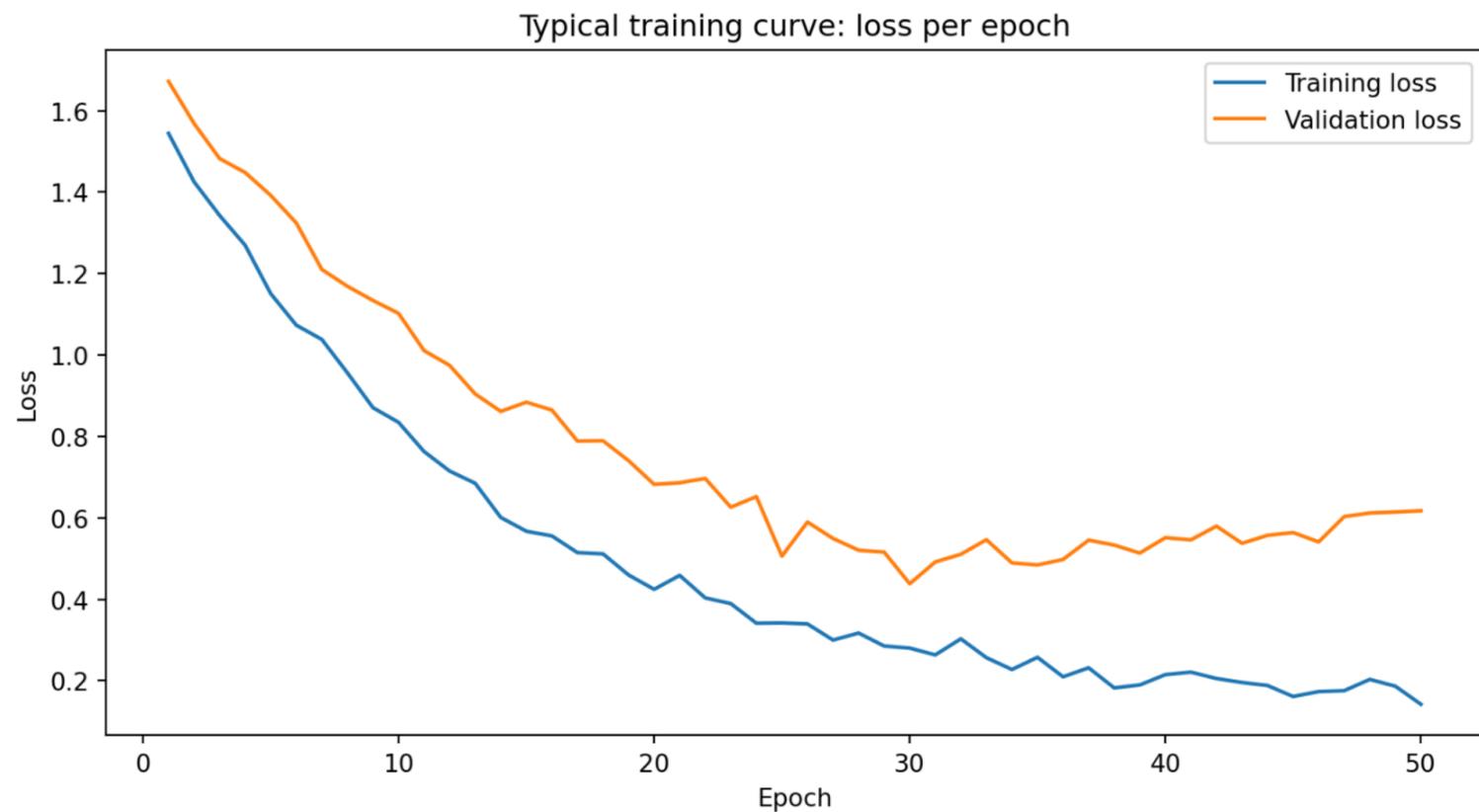
- ▶ The mini-batch gradient is a noisy estimate of the true gradient — on average, it points the right way
- ▶ The noise actually helps: it can bounce the optimizer out of shallow local minima
- ▶ Common batch sizes: 32, 64, 128, 256
- ▶ Smaller batches → more noise, slower convergence
- ▶ Larger batches → smoother gradients, more memory

Epochs

An **epoch** is one complete pass through the entire training set.

If you have 10,000 training observations and a batch size of 100, each epoch consists of 100 mini-batch updates. After one epoch, every observation has been used exactly once.

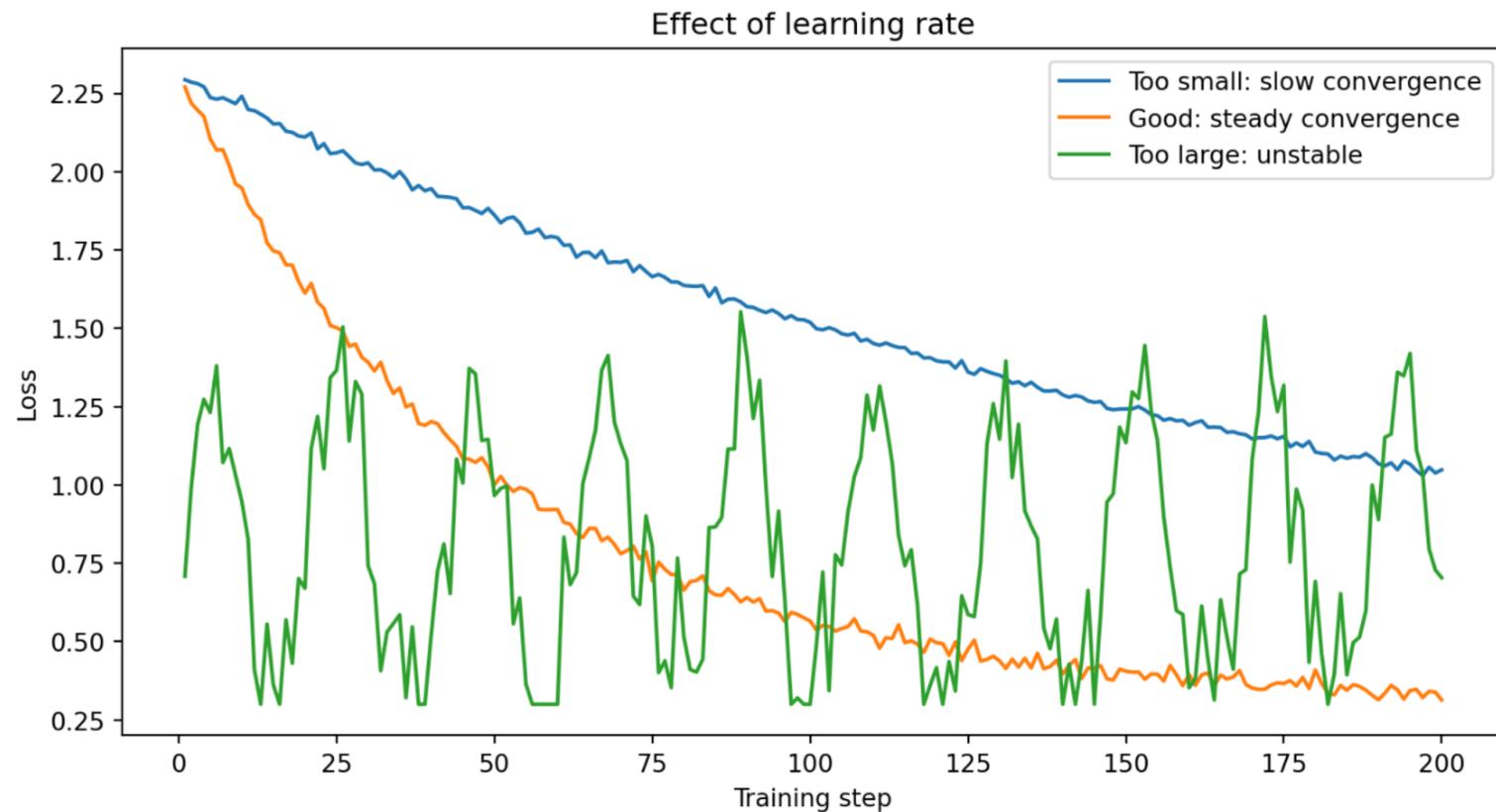
Training typically runs for many epochs — the optimizer sees the same data repeatedly, refining the weights each time.



Training loss decreases steadily. Validation loss decreases initially, then starts to **increase** — the model begins overfitting. The gap between the curves signals that the model is memorizing training data rather than learning generalizable patterns.

The Learning Rate

The learning rate η is usually the first hyperparameter to tune when training a neural network.

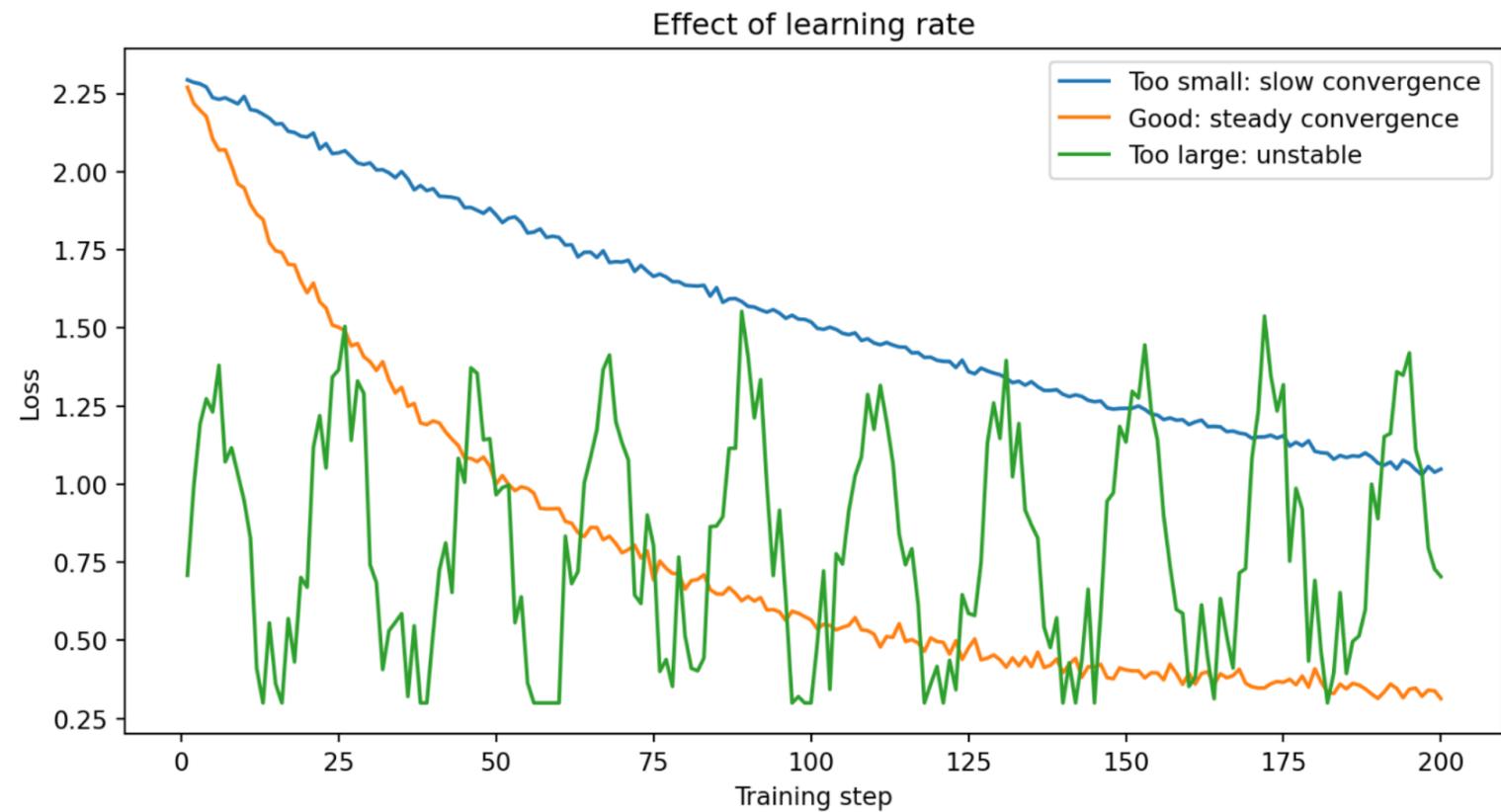


Too small: The model learns very slowly. You might need thousands of epochs to converge — wasteful and sometimes impractical.

Too large: The model overshoots and bounces around, never settling into a good minimum. The loss oscillates or even diverges.

Just right: The loss decreases steadily and converges to a low value.

The learning rate η is usually the first hyperparameter to tune when training a neural network.



Too small: The model learns very slowly. You might need thousands of epochs to converge — wasteful and sometimes impractical.

Too large: The model overshoots and bounces around, never settling into a good minimum. The loss oscillates or even diverges.

Just right: The loss decreases steadily and converges to a low value.

Typical starting values: $\eta = 0.001$ for Adam optimizer, $\eta = 0.01$ for SGD. Often reduced during training (learning rate scheduling).

Optimizers: Beyond Basic Gradient Descent

Plain SGD uses the same learning rate for every parameter and doesn't account for the history of past gradients. Modern optimizers improve on this.

SGD with Momentum adds a “velocity” term \mathbf{v}_t that accumulates past gradients:

$$\begin{aligned}\mathbf{v}_t &= \gamma \mathbf{v}_{t-1} + \eta \nabla_{\theta} \square \\ \theta_{t+1} &= \theta_t - \mathbf{v}_t\end{aligned}$$

- ▶ γ (typically 0.9) controls how much history to keep
- ▶ Consistent gradient direction \rightarrow velocity builds up, optimizer moves faster
- ▶ Oscillating gradients \rightarrow velocity averages out the noise
- ▶ Analogy: a ball rolling downhill — picks up speed on consistent slopes, dampens jitter

Adam (Adaptive Moment Estimation; Kingma & Ba, 2015) tracks the mean \mathbf{m}_t and variance \mathbf{v}_t of past gradients, adapting the learning rate per parameter:

$$\begin{aligned}\mathbf{m}_t &= \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla_{\theta} \square && \text{(mean of gradients)} \\ \mathbf{v}_t &= \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla_{\theta} \square)^2 && \text{(variance of gradients)}\end{aligned}$$

- ▶ Consistent gradient direction → velocity builds up, optimizer moves faster
- ▶ Oscillating gradients → velocity averages out the noise
- ▶ Analogy: a ball rolling downhill — picks up speed on consistent slopes, dampens jitter

Adam (Adaptive Moment Estimation; Kingma & Ba, 2015) tracks the mean \mathbf{m}_t and variance \mathbf{v}_t of past gradients, adapting the learning rate per parameter:

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla_{\theta} \square \quad (\text{mean of gradients})$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla_{\theta} \square)^2 \quad (\text{variance of gradients})$$

$$\theta_{t+1} = \theta_t - \eta \frac{\hat{\mathbf{m}}_t}{\sqrt{\hat{\mathbf{v}}_t + \epsilon}}$$

- ▶ $\hat{\mathbf{m}}_t, \hat{\mathbf{v}}_t$ are bias-corrected versions; ϵ avoids division by zero
- ▶ Defaults: $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$
- ▶ Division by $\sqrt{\hat{\mathbf{v}}_t}$ makes it adaptive: large-gradient parameters get smaller learning rates, noisy parameters get larger ones
- ▶ Adam is the default for most applications — less tuning than SGD, faster convergence

The Training Workflow

Putting it all together, training a neural network involves these steps:

- 1. Choose the architecture:** How many hidden layers? How many neurons per layer? What activation functions?
- 2. Choose the loss function:** MSE for regression, cross-entropy for classification.
- 3. Choose the optimizer:** Adam is usually the default starting point.
- 4. Set hyperparameters:** Learning rate, batch size, number of epochs.
- 5. Train:** Feed mini-batches through the network, compute the loss, backpropagate, update weights. Repeat for many epochs.
- 6. Monitor:** Track training and validation loss each epoch. Stop when validation loss stops improving.
- 7. Evaluate:** Test on held-out data to estimate real-world performance.

This is the same workflow as fitting any ML model — define the model, choose the loss, optimize, evaluate. The difference is that neural networks have more architectural choices and take longer to train.

Part V: Regularization and Practical Considerations

The Overfitting Problem

Neural networks typically have far more parameters than training observations (e.g., 50,000 parameters, 5,000 observations). Without regularization:

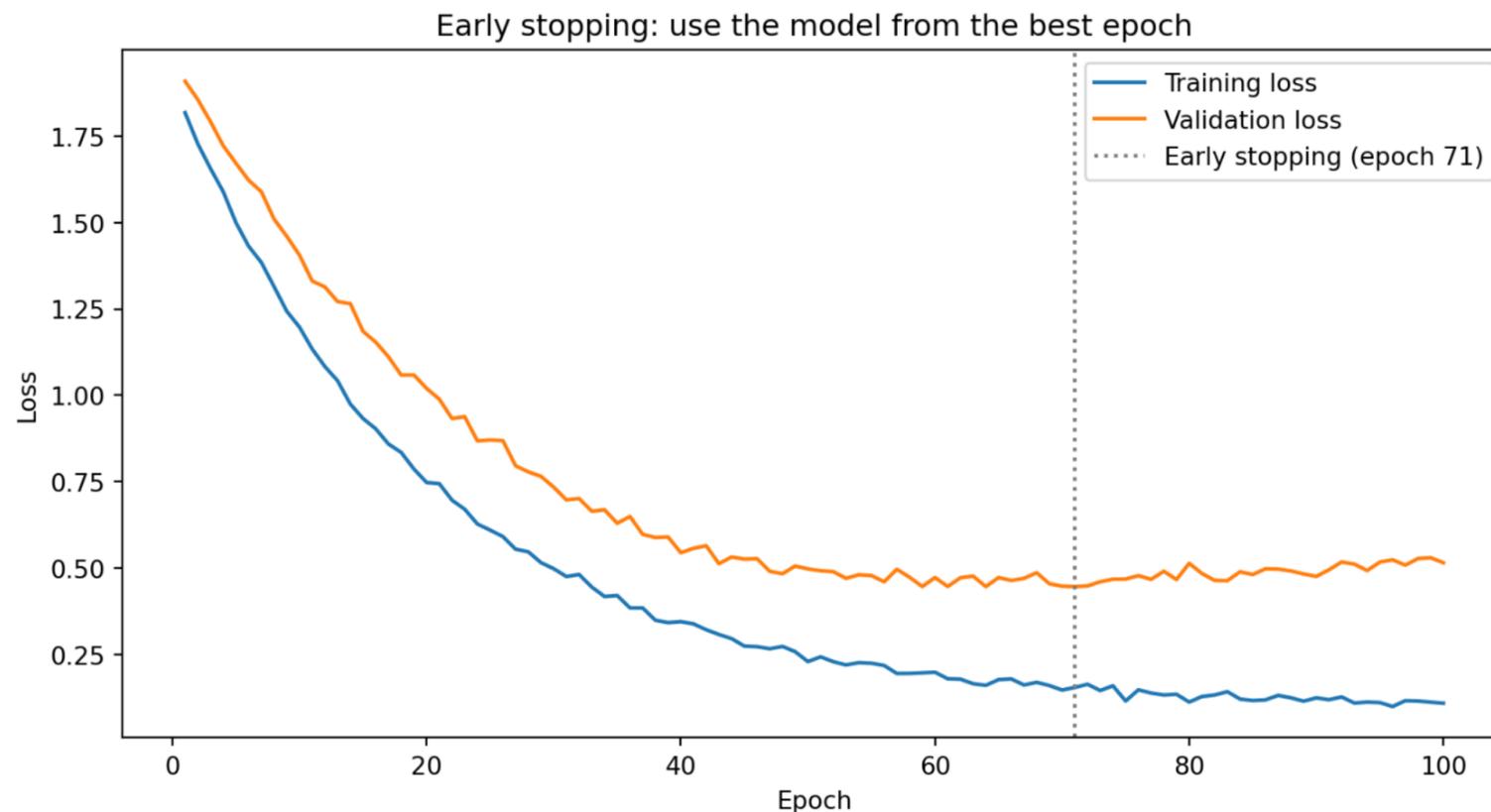
- ▶ The model can *potentially* memorize the training data, including the noise
- ▶ Training loss \rightarrow zero, test loss stays high

Unlike Ridge/Lasso (one knob: λ), neural networks use a **toolkit of complementary strategies**.

Early Stopping

The simplest and most effective regularization technique: **stop training before the model overfits.**

Monitor validation loss during training. When it stops improving (or starts increasing), stop training and use the weights from the best epoch.

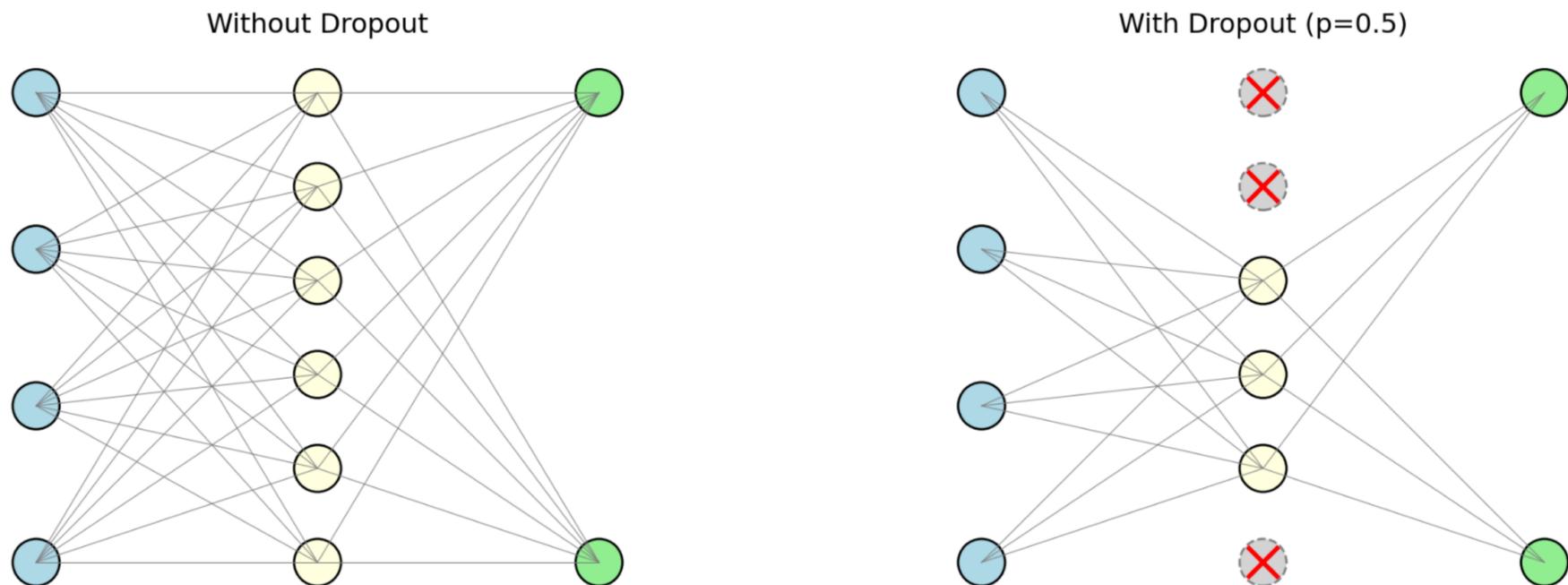


Early stopping acts as implicit regularization: by limiting the number of gradient descent steps, we prevent the model from fully fitting the training noise. It's analogous to how reducing the number of boosting rounds controls overfitting in XGBoost.

Dropout

Dropout (Srivastava et al., 2014): at each mini-batch, randomly set a fraction of neurons to zero.

- ▶ Each neuron is turned off independently with probability p (typically 0.2 to 0.5)
- ▶ Surviving neurons must learn to predict without relying on any single neuron



- ▶ **Ensemble interpretation:** each mini-batch trains a slightly different sub-network → the final model is an average over all of them (like Random Forests averaging trees)
- ▶ **At test time:** dropout is turned off, all neurons are used, weights are scaled accordingly

Weight Decay (L2 Regularization)

Same L2 penalty as Ridge regression, applied to neural network weights:

$$\square_{\text{total}} = \square_{\text{data}} + \lambda \sum_{\ell} \|\mathbf{w}^{(\ell)}\|_2^2$$

The gradient update now includes a term that pulls weights toward zero:

$$w_{t+1} = w_t - \eta \frac{\partial \square_{\text{data}}}{\partial w_t} - \underbrace{\eta \lambda w_t}_{\text{decay}}$$

- ▶ The $-\eta \lambda w_t$ term shrinks each weight at every step → called **weight decay**
- ▶ Encourages small, diffuse weights rather than a few strong connections
- ▶ Same logic as Ridge: constrain parameters to prevent fitting noise

Deep vs. Wide Networks

Two ways to add capacity to a network:

Wider: More neurons per layer (e.g., one hidden layer with 1,000 neurons)

Deeper: More layers with fewer neurons each (e.g., five hidden layers with 64 neurons each)

The universal approximation theorem says one wide layer is sufficient. In practice, deeper is often better:

- ▶ **Hierarchical representations:** early layers detect simple patterns, later layers combine them (edges → shapes → objects)
- ▶ **Parameter-efficient:** some functions need exponentially many neurons in one layer but only a few layers deep
- ▶ **Harder to train:** vanishing gradients, more hyperparameters, longer training

For tabular financial data: 2–4 hidden layers usually suffice. 50+ layers are for images, text, and structured data.

Batch Normalization and Layer Normalization

As preceding layers change during training, the distribution of each layer's inputs shifts (**internal covariate shift**). Normalization techniques fix this.

Batch Normalization (Ioffe & Szegedy, 2015) — normalize across the mini-batch:

$$\hat{z} = \frac{z - \mu_{\text{batch}}}{\sigma_{\text{batch}}}$$

- ▶ The network learns two extra parameters per neuron (γ, β) to scale and shift if needed

Layer Normalization (Ba, Lei & Hinton, 2016) — normalize across features within each observation:

- ▶ Works better for small batch sizes
- ▶ Standard for transformer architectures

Both stabilize training, allow higher learning rates, and act as mild regularizers.

Summary of Regularization Techniques

Technique	What it does	Analogy
Early stopping	Stop before the model overfits	Limiting boosting rounds in XGBoost
Dropout	Randomly disable neurons during training	Ensemble averaging (like Random Forests)
Weight decay	Penalize large weights (L2 penalty)	Ridge regression's $\lambda \ \beta\ _2^2$
Batch / Layer norm	Normalize intermediate values	Standardizing features before regression

In practice, these techniques are often combined. A typical setup: ReLU activation + Adam optimizer + dropout (0.2–0.5) + early stopping + some weight decay.

Part VI: Demo — Neural Networks in PyTorch

PyTorch: The Standard Framework

PyTorch is the most widely used deep learning framework in research and increasingly in industry. It handles all the backpropagation, gradient computation, and optimization automatically — you define the architecture and the training loop.

We'll build a simple classifier for the same synthetic data we used with Random Forests and XGBoost.

```
1 import numpy as np
2 import torch
3 import torch.nn as nn
4 from sklearn.datasets import make_classification
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7
8 # Generate data
9 X, y = make_classification(
10     n_samples=1000, n_features=10, n_informative=5,
11     random_state=42
12 )
13
14 # Split into train, validation, and test
15 X_temp, X_test, y_temp, y_test = train_test_split(
16     X, y, test_size=0.2, random_state=42
17 )
18 X_train, X_val, y_train, y_val = train_test_split(
```

```
Training set: 640 observations
Validation set: 160 observations
Test set: 200 observations
```

Building the Model

```

1 # Build a feed-forward neural network
2 model = nn.Sequential(
3     nn.Linear(10, 64),    # hidden layer 1: 64 neurons
4     nn.ReLU(),
5     nn.Linear(64, 32),   # hidden layer 2: 32 neurons
6     nn.ReLU(),
7     nn.Linear(32, 1),    # output layer: 1 neuron
8     nn.Sigmoid()        # sigmoid for binary classification
9 )
10
11 # Choose optimizer and loss function
12 optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
13 loss_fn = nn.BCELoss() # binary cross-entropy
14
15 print(model)
16 print(f"\nTotal parameters: {sum(p.numel() for p in model.parameters()):,}")

```

```

Sequential(
  (0): Linear(in_features=10, out_features=64, bias=True)
  (1): ReLU()
  (2): Linear(in_features=64, out_features=32, bias=True)
  (3): ReLU()
  (4): Linear(in_features=32, out_features=1, bias=True)
  (5): Sigmoid()
)

```

Total parameters: 2,817

```

1 # Build a feed-forward neural network
2 model = nn.Sequential(
3     nn.Linear(10, 64), # hidden layer 1: 64 neurons
4     nn.ReLU(),
5     nn.Linear(64, 32), # hidden layer 2: 32 neurons
6     nn.ReLU(),
7     nn.Linear(32, 1), # output layer: 1 neuron
8     nn.Sigmoid()      # sigmoid for binary classification
9 )
10
11 # Choose optimizer and loss function
12 optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
13 loss_fn = nn.BCELoss() # binary cross-entropy
14
15 print(model)
16 print(f"\nTotal parameters: {sum(p.numel() for p in model.parameters()):,}")

```

```

Sequential(
  (0): Linear(in_features=10, out_features=64, bias=True)
  (1): ReLU()
  (2): Linear(in_features=64, out_features=32, bias=True)
  (3): ReLU()
  (4): Linear(in_features=32, out_features=1, bias=True)
  (5): Sigmoid()
)

```

Total parameters: 2,817

Each **Linear** layer is a fully-connected layer (every neuron connects to every neuron in the previous layer). **ReLU** activations introduce nonlinearity between layers.

Training the Model

```

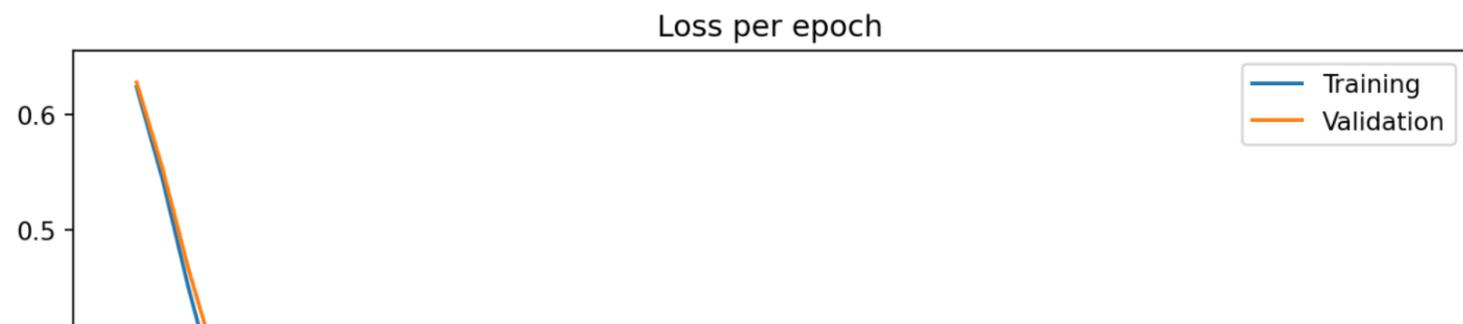
1 from torch.utils.data import TensorDataset, DataLoader
2
3 # Create data loader for mini-batching
4 train_loader = DataLoader(TensorDataset(X_train_t, y_train_t), batch_size=32, shuffle=True)
5
6 # Training loop
7 train_losses = []
8 val_losses = []
9
10 torch.manual_seed(42)
11
12 for epoch in range(50):
13     # Training phase
14     for X_batch, y_batch in train_loader:
15         optimizer.zero_grad()           # reset gradients
16         y_pred = model(X_batch)         # forward pass
17         loss = loss_fn(y_pred, y_batch) # compute loss
18         loss.backward()                 # backpropagation

```

```

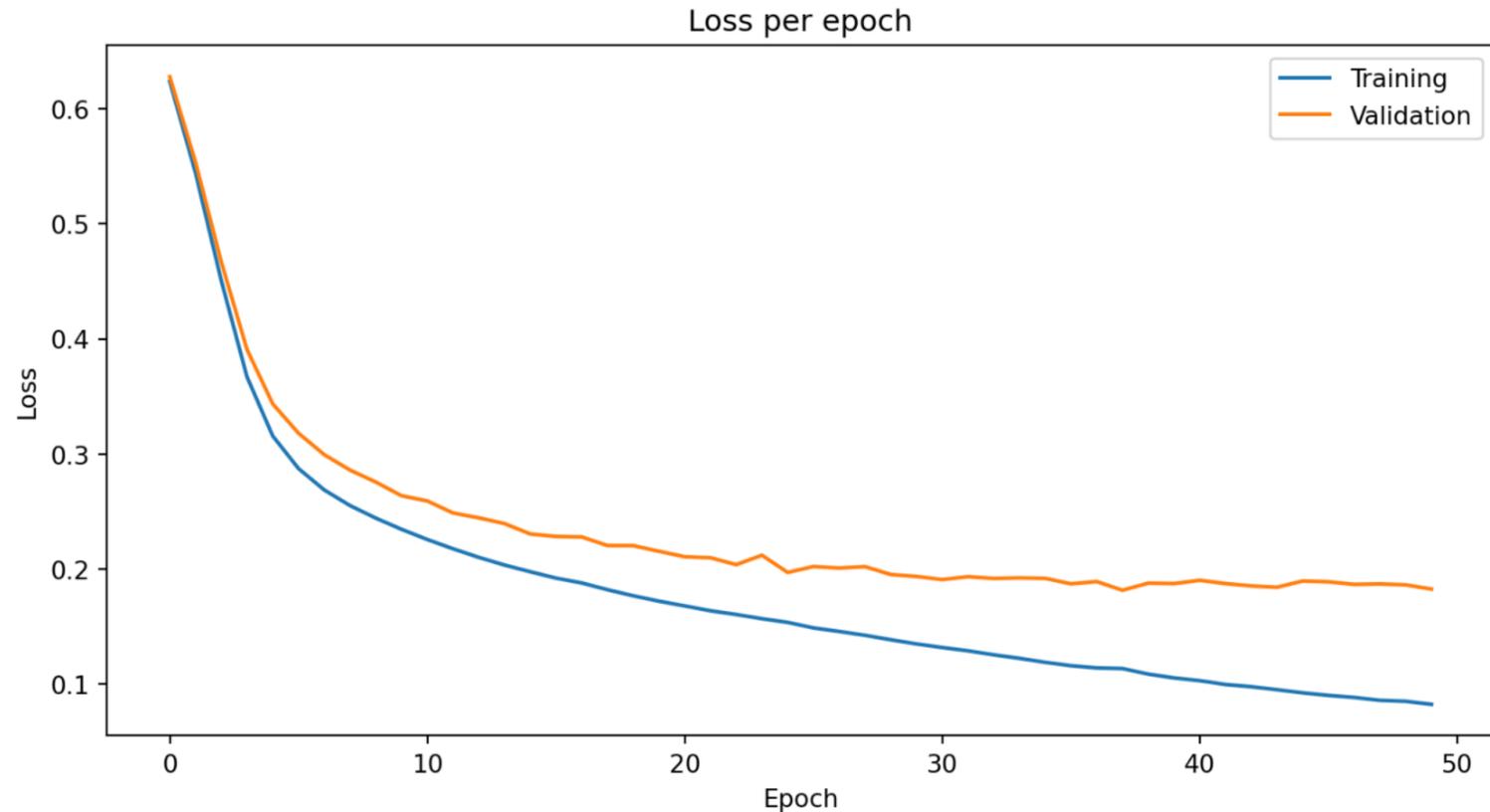
Final training loss: 0.0821
Final validation loss: 0.1823

```



```
11
12 for epoch in range(50):
13     # Training phase
14     for X_batch, y_batch in train_loader:
15         optimizer.zero_grad()           # reset gradients
16         y_pred = model(X_batch)         # forward pass
17         loss = loss_fn(y_pred, y_batch) # compute loss
18         loss.backward()                  # backpropagation
```

```
Final training loss: 0.0821
Final validation loss: 0.1823
```



Training loss is consistently below validation loss — the model fits the training data better than unseen data, as expected. The gap is small here, meaning the model generalizes reasonably well.

Evaluating on the Test Set

```

1 # Evaluate on the held-out test set
2 with torch.no_grad():
3     test_pred = model(X_test_t)
4     test_accuracy = ((test_pred > 0.5) == y_test_t).float().mean().item()
5
6 print(f"Test accuracy: {test_accuracy:.3f}")

```

Test accuracy: 0.950

```

1 # Compare with Random Forest and XGBoost on the same data
2 from sklearn.ensemble import RandomForestClassifier
3 from xgboost import XGBClassifier
4
5 rf = RandomForestClassifier(n_estimators=200, random_state=42)
6 rf.fit(X_train, y_train)
7
8 xgb = XGBClassifier(n_estimators=200, learning_rate=0.1, max_depth=3, random_state=42)
9 xgb.fit(X_train, y_train)
10
11 print(f"\nModel comparison on test set:")
12 print(f"Neural Network: {test_accuracy:.3f}")
13 print(f"Random Forest: {rf.score(X_test, y_test):.3f}")
14 print(f"XGBoost: {xgb.score(X_test, y_test):.3f}")

```

Model comparison on test set:
 Neural Network: 0.950
 Random Forest: 0.950

```

4     test_accuracy = ((test_pred > 0.5) == y_test_t).float().mean().item()
5
6     print(f"Test accuracy: {test_accuracy:.3f}")

```

```
Test accuracy: 0.950
```

```

1 # Compare with Random Forest and XGBoost on the same data
2 from sklearn.ensemble import RandomForestClassifier
3 from xgboost import XGBClassifier
4
5 rf = RandomForestClassifier(n_estimators=200, random_state=42)
6 rf.fit(X_train, y_train)
7
8 xgb = XGBClassifier(n_estimators=200, learning_rate=0.1, max_depth=3, random_state=42)
9 xgb.fit(X_train, y_train)
10
11 print(f"\nModel comparison on test set:")
12 print(f"Neural Network: {test_accuracy:.3f}")
13 print(f"Random Forest: {rf.score(X_test, y_test):.3f}")
14 print(f"XGBoost: {xgb.score(X_test, y_test):.3f}")

```

```

Model comparison on test set:
Neural Network: 0.950
Random Forest: 0.950
XGBoost: 0.945

```

On small tabular datasets like this, neural networks typically perform **comparably** to tree-based methods — sometimes slightly better, sometimes slightly worse. The real advantage of neural networks emerges with large datasets and non-tabular data (images, text, sequences).

Neural Networks in Finance: Where They Shine

Application	Why neural networks?
High-frequency trading	Large datasets, complex nonlinear patterns in order flow
Factor models	Gu, Kelly & Xiu (2020): NNs outperform linear/tree models for stock return prediction — but need large panels
Alternative data	Images, news, earnings calls, social media require specialized architectures (CNNs, transformers)
Risk management	Estimating loss distributions, stress testing, credit risk

Where tree-based models still win: small-to-medium tabular datasets (most day-to-day finance). XGBoost/RF are competitive, faster to train, easier to tune, and interpretable.

Neural Network Hyperparameters: A Summary

Hyperparameter	What it controls	Typical choices
Number of layers	Network depth	2–4 for tabular data; dozens for images/text
Neurons per layer	Network width	32, 64, 128, 256
Activation function	Nonlinearity type	ReLU (hidden), sigmoid/softmax (output)
Learning rate	Step size in gradient descent	0.001 (Adam), 0.01 (SGD)
Batch size	Observations per gradient update	32–256
Epochs	Passes through training data	Until validation loss stops improving
Dropout rate	Fraction of neurons dropped	0.2–0.5
Weight decay	L2 penalty strength	10^{-4} to 10^{-2}

images/text

Neurons per layer	Network width	32, 64, 128, 256
Activation function	Nonlinearity type	ReLU (hidden), sigmoid/softmax (output)
Learning rate	Step size in gradient descent	0.001 (Adam), 0.01 (SGD)
Batch size	Observations per gradient update	32–256
Epochs	Passes through training data	Until validation loss stops improving
Dropout rate	Fraction of neurons dropped	0.2–0.5
Weight decay	L2 penalty strength	10^{-4} to 10^{-2}
Optimizer	Gradient descent variant	Adam (default), SGD with momentum

Start with a simple architecture (2 hidden layers, 64 neurons each, ReLU, Adam, dropout 0.3) and adjust based on validation performance.

Tuning Hyperparameters: A Practical Guide

- 1. Start simple.** 2 hidden layers, 64 neurons each. Can the model overfit the training data? If not, fix the architecture or learning rate first.
- 2. Get the learning rate right.** Loss barely moves \rightarrow increase. Loss explodes \rightarrow decrease. Try 10^{-2} , 10^{-3} , 10^{-4} .
- 3. Regularize.** Once you can overfit, add dropout (0.2–0.3), early stopping, and/or weight decay to close the train–val gap.
- 4. Scale up if needed.** Val loss high and close to train loss \rightarrow underfitting (go wider/deeper). Big gap \rightarrow overfitting (more regularization or smaller network).

Common Pitfalls

- ▶ **Tuning too many things at once.** Change one hyperparameter at a time — otherwise you won't know what helped.
- ▶ **Not standardizing inputs.** Feature scales matter. A feature in billions and one in decimals → gradients dominated by the large one. Always standardize first.
- ▶ **No early stopping.** The model *will* memorize. Monitor validation loss.
- ▶ **Too deep too soon.** For tabular data, 2–4 layers usually suffice. Start shallow.
- ▶ **Ignoring the baseline.** Always compare against logistic regression or Random Forest. If the NN doesn't beat it, the complexity isn't worth it.
- ▶ **Forgetting `model.eval()`.** Dropout and batch norm behave differently during training vs. evaluation. Always switch modes.

Part VII: Advanced Deep Learning

Beyond Feed-Forward: A Preview

- ▶ The feed-forward network processes each input independently — no notion of **order** or **sequence**
- ▶ But many problems involve sequential data: sentences, time series, order books
- ▶ Specialized architectures handle this — we won't go deep, but you should know what each one solves

Recurrent Neural Networks (RNNs) and LSTMs

A recurrent neural network adds a connection from a layer back to itself across time steps:

- ▶ At each step t : receives current input x_t and previous hidden state h_{t-1}

$$h_t = g(\mathbf{W}_h h_{t-1} + \mathbf{W}_x x_t + \mathbf{b})$$

- ▶ h_t summarizes everything the network has seen up to time t
- ▶ Output is computed from h_t (at each step or the final step)
- ▶ **Problem:** for long sequences, the gradient travels through many multiplications of $\mathbf{W}_h \rightarrow$ vanishing gradients again

LSTM (Hochreiter & Schmidhuber, 1997) fixes this by adding a **cell state** c_t — a separate memory channel with three learned **gates**:

- ▶ **Forget gate:** What fraction of the previous cell state to keep ($f_t \in [0, 1]$)
- ▶ **Input gate:** What fraction of new information to write into the cell ($i_t \in [0, 1]$)
- ▶ **Output gate:** What fraction of the cell state to expose as the hidden state ($o_t \in [0, 1]$)

LSTM (Hochreiter & Schmidhuber, 1997) fixes this by adding a **cell state** C_t — a separate memory channel with three learned **gates**:

- ▶ **Forget gate:** What fraction of the previous cell state to keep ($f_t \in [0, 1]$)
- ▶ **Input gate:** What fraction of new information to write into the cell ($i_t \in [0, 1]$)
- ▶ **Output gate:** What fraction of the cell state to expose as the hidden state ($o_t \in [0, 1]$)

$$f_t = \sigma(\mathbf{W}_f[h_{t-1}, \mathbf{x}_t] + \mathbf{b}_f)$$

$$i_t = \sigma(\mathbf{W}_i[h_{t-1}, \mathbf{x}_t] + \mathbf{b}_i)$$

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$$

$$o_t = \sigma(\mathbf{W}_o[h_{t-1}, \mathbf{x}_t] + \mathbf{b}_o)$$

$$h_t = o_t \odot \tanh(C_t)$$

where \odot is element-wise multiplication, $\tilde{C}_t = \tanh(\mathbf{W}_c[h_{t-1}, \mathbf{x}_t] + \mathbf{b}_c)$ is candidate new memory.

- ▶ The cell state update $C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$ is the key
- ▶ When the forget gate ≈ 1 , the gradient flows through almost unchanged \rightarrow solves vanishing gradients over long sequences

Transformers and Attention

The problem with RNNs:

- ▶ Read one word at a time, left to right — step 50 waits for steps 1–49
- ▶ Word 1’s meaning must survive through every intermediate hidden state to interact with word 50 (vanishing gradients)
- ▶ But language puts related words far apart: “The executive **who** the board hired last quarter after a lengthy search **resigned.**” — subject and verb are nine words apart

The fix: process the **entire sequence at once**, letting any word attend to any other word directly.

The Attention Mechanism

The **transformer** (Vaswani et al., 2017) projects each word into three learned roles:

- ▶ **Query (q)**: “What am I looking for?”
- ▶ **Key (k)**: “What do I have to offer?”
- ▶ **Value (v)**: “Here is my actual content.”

Stacking into matrices **Q**, **K**, **V**:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \mathbf{V}$$

- ▶ $\mathbf{Q}\mathbf{K}^\top$ = score matrix — one score for every pair of words (dot product of query with key)
- ▶ Softmax turns each row into weights that sum to 1 (how much to attend to each word)
- ▶ Multiply by \mathbf{V} → each word gets a new representation = weighted average of the whole sentence
- ▶ $\sqrt{d_k}$ = scaling factor to keep dot products from getting too large

The **transformer** (Vaswani et al., 2017) projects each word into three learned roles:

- ▶ **Query (q)**: “What am I looking for?”
- ▶ **Key (k)**: “What do I have to offer?”
- ▶ **Value (v)**: “Here is my actual content.”

Stacking into matrices **Q**, **K**, **V**:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d_k}}\right) \mathbf{V}$$

- ▶ $\mathbf{Q}\mathbf{K}^{\top}$ = score matrix — one score for every pair of words (dot product of query with key)
- ▶ Softmax turns each row into weights that sum to 1 (how much to attend to each word)
- ▶ Multiply by \mathbf{V} → each word gets a new representation = weighted average of the whole sentence
- ▶ $\sqrt{d_k}$ = scaling factor to keep dot products from getting too large
- ▶ All pairwise scores computed at once (matrix multiplication) → runs in parallel on GPUs
- ▶ This is why transformers scale to billions of parameters — and the architecture behind GPT, Claude

Generative Adversarial Networks (GANs)

- ▶ Every model so far: **discriminative** ($\mathbf{x} \rightarrow y$)
- ▶ A **generative** model instead learns to produce new data that resembles the training data

A **GAN** (Goodfellow et al., 2014) trains two networks against each other:

- ▶ **Generator G** : takes random noise \mathbf{z} as input and produces a synthetic data point $G(\mathbf{z})$
- ▶ **Discriminator D** : takes a data point (real or synthetic) and outputs the probability that it is real

The two networks play a minimax game:

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z} [\log(1 - D(G(\mathbf{z})))]$$

- ▶ **Discriminator** maximizes: assign high $D(\mathbf{x})$ to real data, low $D(G(\mathbf{z}))$ to fakes
- ▶ **Generator** minimizes: produce $G(\mathbf{z})$ that fools D

- ▶ **Generator G** : takes random noise \mathbf{z} as input and produces a synthetic data point $G(\mathbf{z})$
- ▶ **Discriminator D** : takes a data point (real or synthetic) and outputs the probability that it is real

The two networks play a minimax game:

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z} [\log(1 - D(G(\mathbf{z})))]$$

- ▶ **Discriminator maximizes**: assign high $D(\mathbf{x})$ to real data, low $D(G(\mathbf{z}))$ to fakes
- ▶ **Generator minimizes**: produce $G(\mathbf{z})$ that fools D
- ▶ Training alternates between updating D and G
- ▶ **At convergence** (in theory): generator produces data indistinguishable from real, D outputs 0.5 for everything
- ▶ **In practice**: GAN training is notoriously unstable — the two networks can oscillate rather than converge

My Research: Finance-Informed Neural Networks

Every model we've seen today uses **data** to train: we have (x_i, y_i) pairs and minimize a loss like MSE. But what if we don't have data — we have an **equation** we know the solution must satisfy?

- ▶ Many problems in finance boil down to: find a function f such that some condition like $f(x) = 0$ holds
- ▶ We know from theory what the equation *is*, but we can't solve it by hand — especially when the problem is high-dimensional
- ▶ The idea: let f_θ be a neural network, and use the equation itself as the loss function. If the equation says $f(x) = 0$, then $\square = \sum_i f_\theta(x_i)^2$ — this is zero when f_θ satisfies the equation, positive when it doesn't. Gradient descent does the rest.

This idea originated in physics (Raissi, Perdikaris & Karniadakis, 2019), where it's called **Physics-Informed Neural Networks** (PINNs). My dissertation applied it to finance — **Finance-Informed Neural Networks** (FINNs) — in two settings.

FINNs Application 1: Macroeconomic Policy (Chapters 1 & 2)

Model 60 overlapping generations of households (ages 20–79)–3 types: 0-, low-, high-LTI–each choosing how much to consume vs. save. A household of age i maximizes lifetime utility — the discounted sum of utility from consumption over its remaining life:

$$\max \quad \mathbb{E}_t \sum_{s=0}^{59} \beta^s u(c_{i,t+i})$$

subject to a budget constraint each period:

$$c_{i,t} + k_{i+1,t+1} = (1 - \tau_t) w_t \ell_i + (1 + r_t) k_{i,t} - a_{i,t}$$

where $c_{i,t}$ is consumption, $k_{i+1,t+1}$ is savings carried into next period, $w_t \ell_i$ is after-tax labour income, r_t is the return on savings, and $a_{i,t}$ is the student loan payment. This budget constraint is what links all the households together — everyone's savings become the economy's capital, which determines next period's wages and returns for everyone else.

When you solve this optimization (take the first-order condition), you get a condition that must hold at every age: the marginal benefit of consuming one more dollar today must equal the expected marginal benefit of saving that dollar and consuming it tomorrow:

$$c_{i,t} + k_{i+1,t+1} - (1+r_t)k_{i,t} = w_t \ell_i - a_{i,t}$$

where $c_{i,t}$ is consumption, $k_{i+1,t+1}$ is savings carried into next period, $w_t \ell_i$ is after-tax labour income, r_t is the return on savings, and $a_{i,t}$ is the student loan payment. This budget constraint is what links all the households together — everyone's savings become the economy's capital, which determines next period's wages and returns for everyone else.

When you solve this optimization (take the first-order condition), you get a condition that must hold at every age: the marginal benefit of consuming one more dollar today must equal the expected marginal benefit of saving that dollar and consuming it tomorrow:

$$u'(c_{i,t}) = \beta \mathbb{E}_t [u'(c_{i+1,t+1}) \cdot (1 + r_{t+1})]$$

where $u'(\cdot)$ is marginal utility, β is the discount factor (patience), and r_{t+1} is the return on savings. In words: if this condition is violated, the household can do better by shifting consumption between today and tomorrow.

That's 59×3 of these conditions (one per age per type) that must all hold simultaneously. The state space — tracking everyone's wealth — is 178-dimensional. Traditional methods that use grids are hopeless here: a grid with just 10 points per dimension would have 10^{178} grid points.

Instead, a neural network f_θ learns the optimal saving decision for all 60 age groups at once, trained by penalizing violations of these optimality conditions. To avoid grids entirely, I use the neural network's own policy to simulate the economy forward in time — generating a time series of states the economy actually visits — and evaluate the optimality conditions only at those points. No grid, no curse of dimensionality.

FINNs Application 2: Pricing Interest Rate Derivatives (Chapter 3)

A caplet is an option on a future interest rate — it pays off if the rate exceeds a strike L_E . From RSM332, the price of any derivative is the expected discounted payoff under no-arbitrage:

$$V(t) = \mathbb{E} \left[\exp \left(- \int_t^T r(s) ds \right) \cdot \text{payoff at } T \right]$$

where $r(s)$ is the short-term interest rate along the path. The standard approach is Monte Carlo simulation: simulate thousands of random interest rate paths, compute the payoff on each, discount back, and average. This is slow — seconds per contract — and you need to re-simulate for every change in contract terms.

Instead, there is an equivalent way to express exactly the same price, not as an expectation over random paths, but as a deterministic equation that the price function V must satisfy:

$$-\frac{\partial V}{\partial \tau} + \mu' \nabla_f V + \frac{1}{2} \sum_{n=1}^N \sigma_n' \nabla_f^2 V \sigma_n - r \cdot V = 0$$

where τ is time to maturity, f is the current forward rate curve, and μ, σ come from the term structure model. You don't need to understand every symbol here — the point is that this is an equation the price *must* satisfy, and it involves

where $r(s)$ is the short-term interest rate along the path. The standard approach is Monte Carlo simulation: simulate thousands of random interest rate paths, compute the payoff on each, discount back, and average. This is slow — seconds per contract — and you need to re-simulate for every change in contract terms.

Instead, there is an equivalent way to express exactly the same price, not as an expectation over random paths, but as a deterministic equation that the price function V must satisfy:

$$-\frac{\partial V}{\partial \tau} + \mu' \nabla_f V + \frac{1}{2} \sum_{n=1}^N \sigma_n' \nabla_f^2 V \sigma_n - r \cdot V = 0$$

where τ is time to maturity, f is the current forward rate curve, and μ, σ come from the term structure model. You don't need to understand every symbol here — the point is that this is an equation the price *must* satisfy, and it involves derivatives of V with respect to its inputs.

A neural network V_θ approximates the price. The loss function squares the left-hand side of that equation: if V_θ satisfies the equation, the loss is zero. Backpropagation — the same machinery used to train any neural network — computes all the partial derivatives of V_θ that appear in the equation. No random simulation at all — the network learns the price by satisfying the equation directly.

The result: once trained, the network prices any contract instantly. Prices are 300,000–4,500,000× faster than Monte Carlo simulation, with accuracy to 0.04 cents per dollar.

Summary and Preview

What We Learned Today

Deep learning = regression. A neural network is a choice of f_θ in the same framework from Lecture 5. The loss functions, regularization, and evaluation methods all carry over.

The architecture. Neurons compute weighted sums plus nonlinear activations. Layers stack neurons together. Feed-forward networks chain layers sequentially. Activation functions (especially ReLU) are what make depth meaningful — without them, multiple layers collapse to a single linear transformation.

Universal approximation. A wide enough single-layer network can approximate any continuous function. But this is an existence theorem, not a practical recipe. Deep networks are more parameter-efficient and learn hierarchical representations.

Training. No closed-form solution — we use gradient descent. Backpropagation efficiently computes gradients by applying the chain rule backwards through the network. SGD with mini-batches makes training practical for large datasets. Adam is the default optimizer.

Regularization. Neural networks are massively overparameterized and prone to overfitting. Early stopping, dropout, weight decay, and batch normalization are the standard toolkit.

Modern architectures. RNNs and LSTMs handle sequential data by passing hidden states through time. Transformers replace recurrence with self-attention, processing entire sequences in parallel. GANs train two networks adversarially to generate synthetic data.

In practice. On tabular financial data, neural networks compete with but rarely dominate tree-based methods. Their real

Next Lecture: Text and Natural Language Processing

Neural networks unlock a new kind of data for finance: **text**.

Earnings calls, news articles, analyst reports, SEC filings, social media — vast amounts of financially relevant information is locked in natural language.

Lecture 10: Text & Natural Language Processing

- ▶ How to represent text as numbers
- ▶ Bag-of-words, TF-IDF, and word embeddings
- ▶ Sentiment analysis and topic modelling
- ▶ Transformers and large language models

References

- ▶ Ba, J. L., Kiros, J. R., & Hinton, G. E. (2016). Layer normalization. *arXiv preprint arXiv:1607.06450*.
- ▶ Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4), 303–314.
- ▶ Gu, S., Kelly, B., & Xiu, D. (2020). Empirical asset pricing via machine learning. *Review of Financial Studies*, 33(5), 2223–2273.
- ▶ Hornik, K. (1991). Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2), 251–257.
- ▶ Ioffe, S. & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *Proceedings of ICML*, 448–456.
- ▶ Kingma, D. P. & Ba, J. (2015). Adam: A method for stochastic optimization. *Proceedings of ICLR*.
- ▶ Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533–536.
- ▶ Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15, 1929–1958.
- ▶ Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27.
- ▶ Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning* (2nd ed.). Springer. Chapter 11.

- ▶ Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4), 303–314.
- ▶ Gu, S., Kelly, B., & Xiu, D. (2020). Empirical asset pricing via machine learning. *Review of Financial Studies*, 33(5), 2223–2273.
- ▶ Hornik, K. (1991). Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2), 251–257.
- ▶ Ioffe, S. & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *Proceedings of ICML*, 448–456.
- ▶ Kingma, D. P. & Ba, J. (2015). Adam: A method for stochastic optimization. *Proceedings of ICLR*.
- ▶ Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533–536.
- ▶ Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15, 1929–1958.
- ▶ Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27.
- ▶ Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning* (2nd ed.). Springer. Chapter 11.
- ▶ Hochreiter, S. & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.
- ▶ Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.