

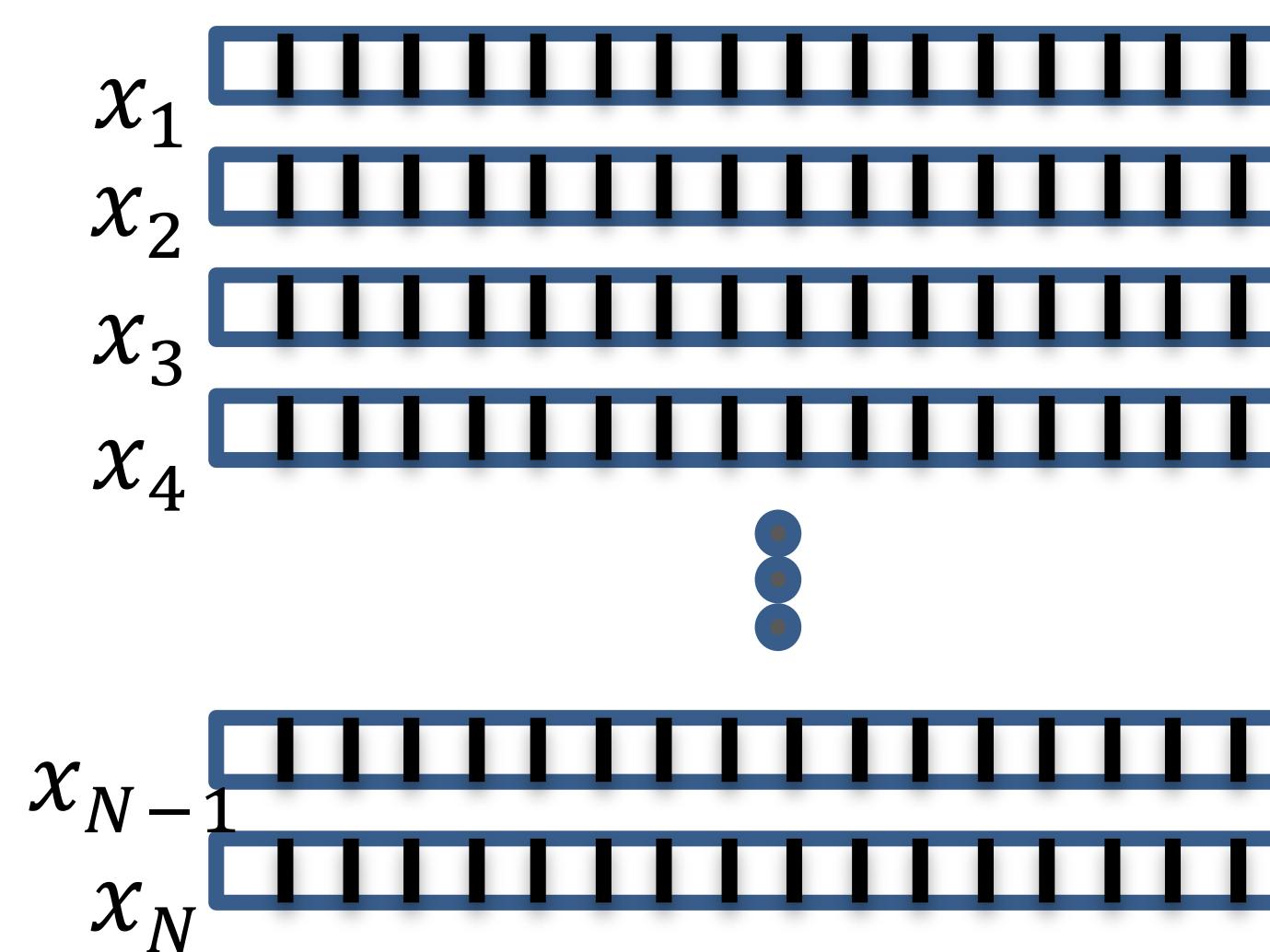
# Machine Learning Summer School: Day 1

David Carlson

# How do we learn a Deep Model?

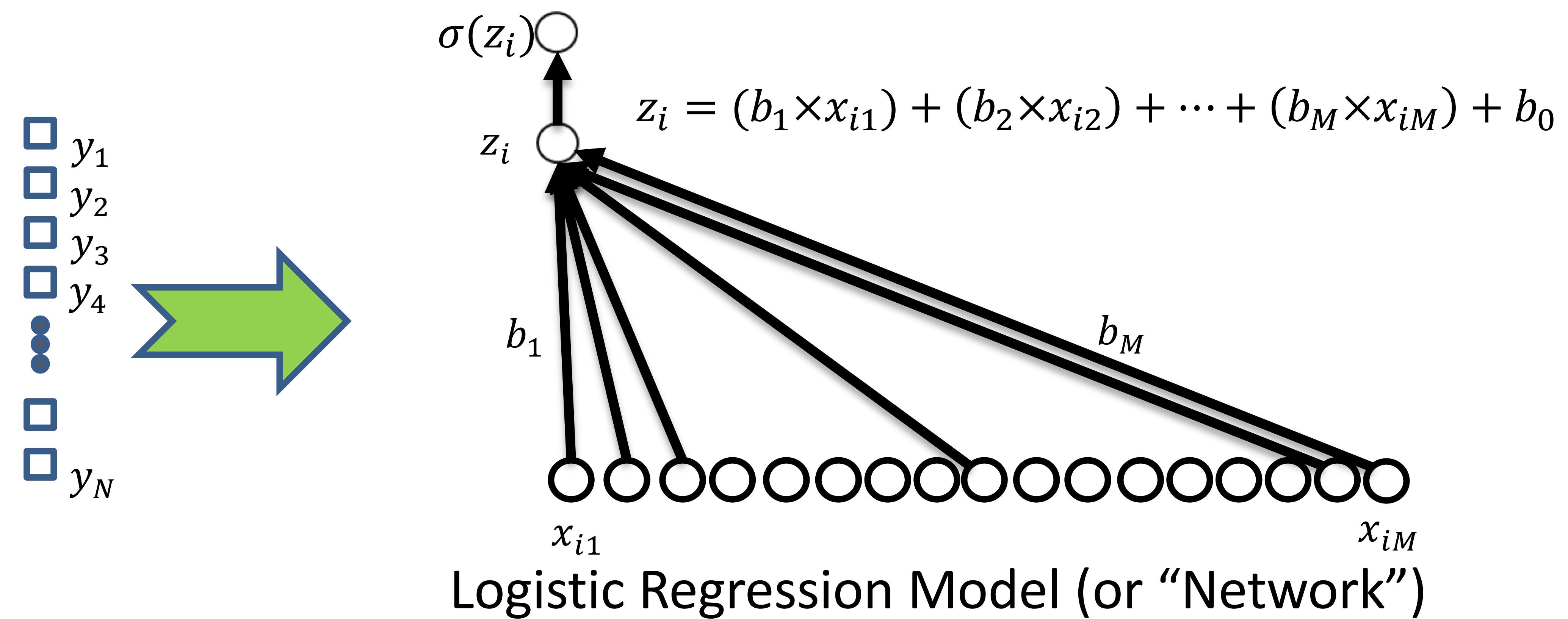
Part II

# (Recall) Learning Model Parameters



Training Set

How do we  
actually do this?



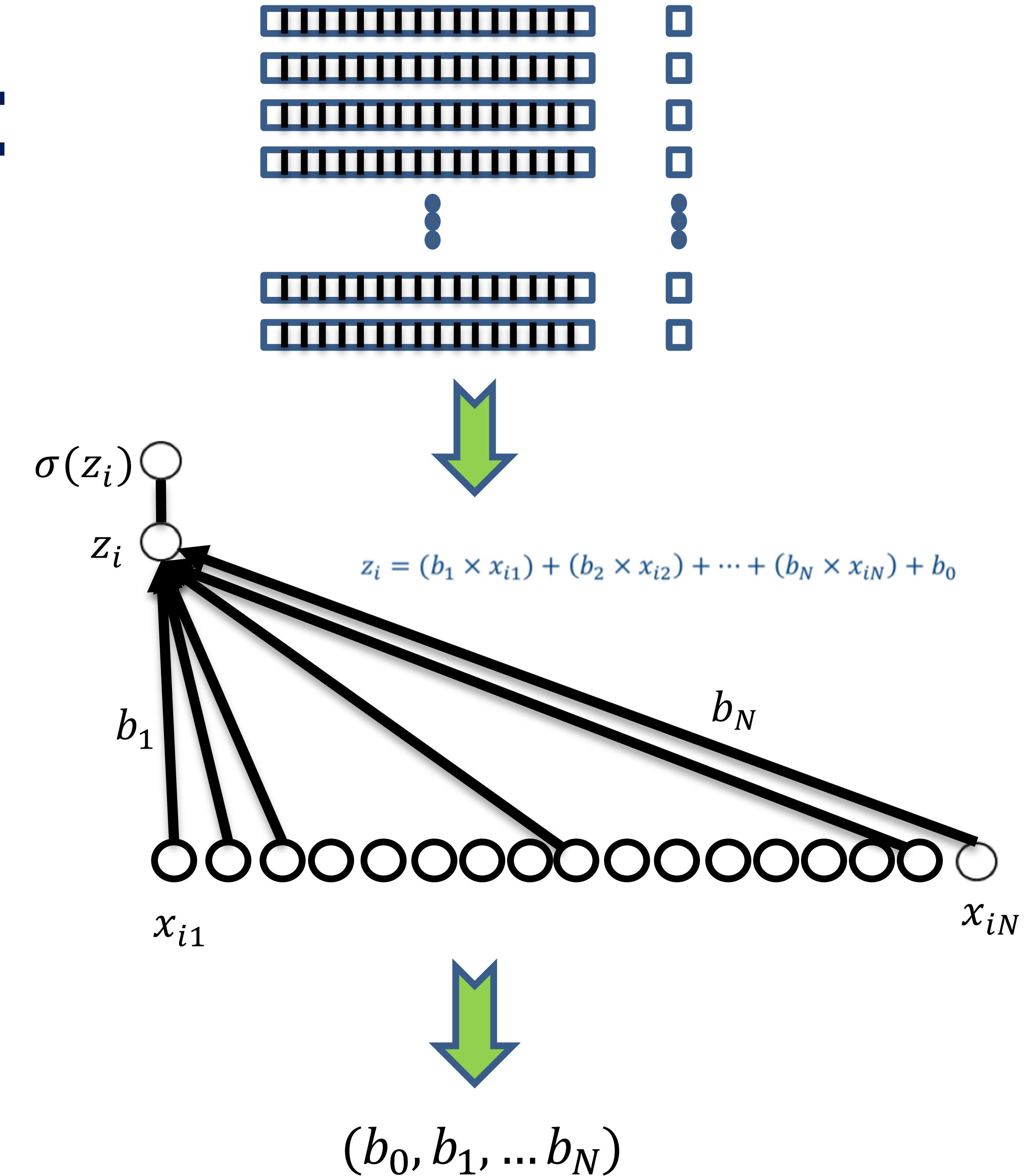
Learned  
Parameters  $(b_0, b_1, \dots, b_N)$

# Looking Forward

- This afternoon you will be implementing and learning a multilayer perceptron *yourself* (with help from the TAs!)
- In this second lecture, we want to set up the mathematical structure to learn the network

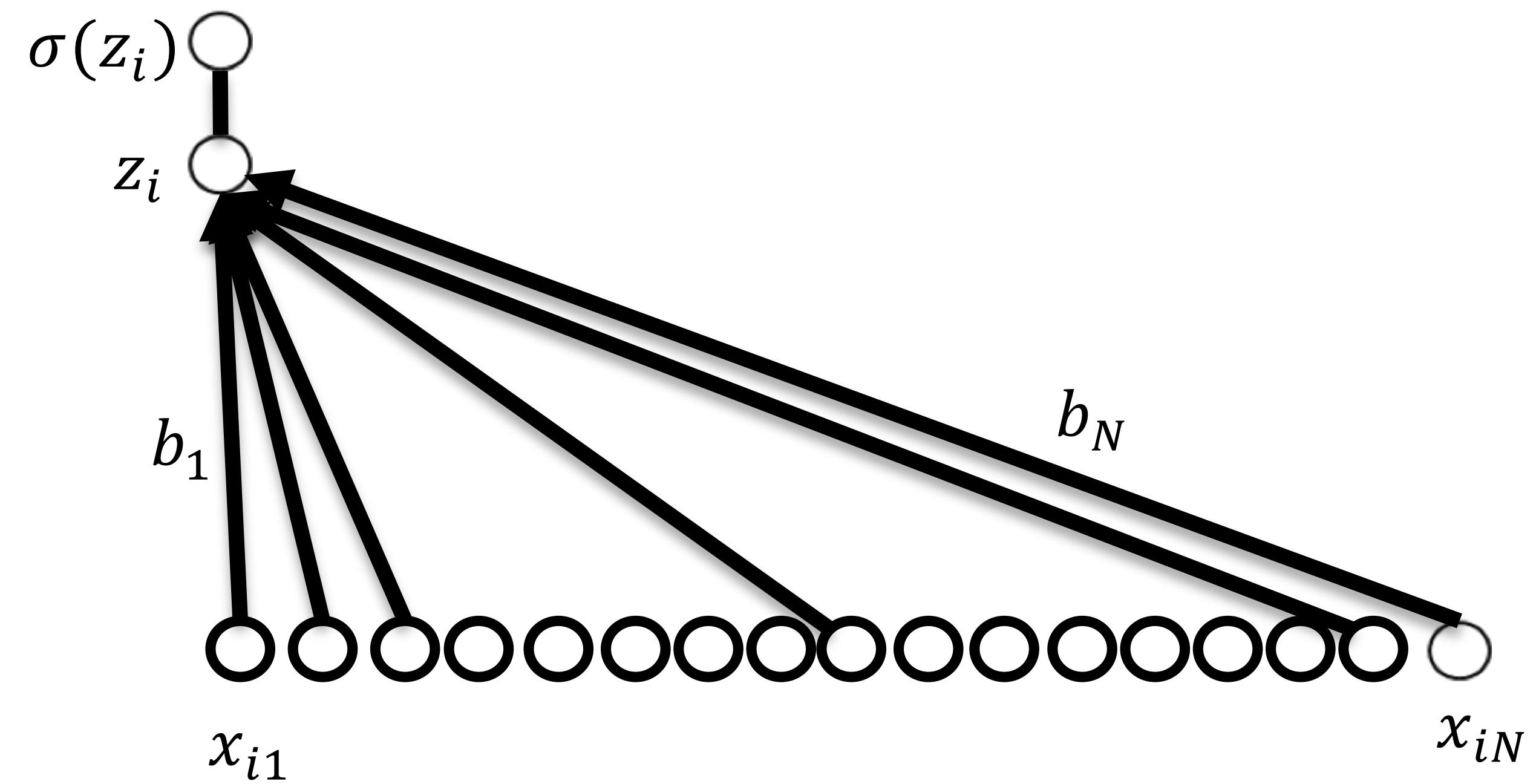
# Fundamental Question:

- Given a large amount of data and a model/network to fit, how to *efficiently* and *effectively* learn the model parameters?



# What are we trying to do?

- Want to learn parameters that give us the best performance
- Essentially: given data, find the “best”  $b$  for that data
- How do we define performance?

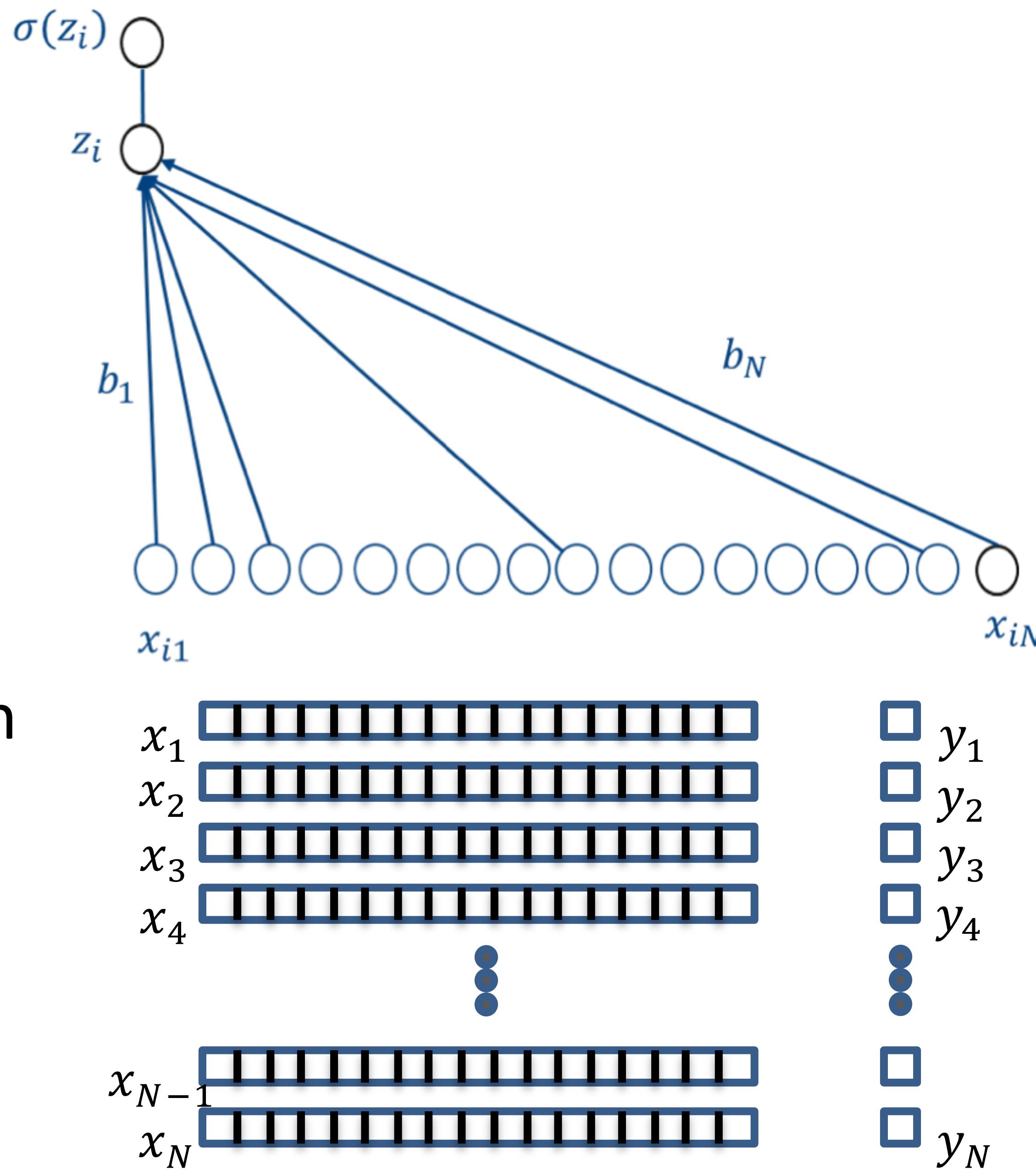


# “Empirical Risk Minimization”

- A *loss function*  $\ell(\text{true}, \text{prediction})$  defines a penalty for poor predictions
- Want to minimize average loss
- Mathematically, this can be stated as

$$\mathbf{b}^* = \arg \min_{\mathbf{b}} \frac{1}{N} \sum_i^N \ell(y_i, \sigma(z_i))$$

↑                      ↑  
True Label      Loss function  
                    Guess

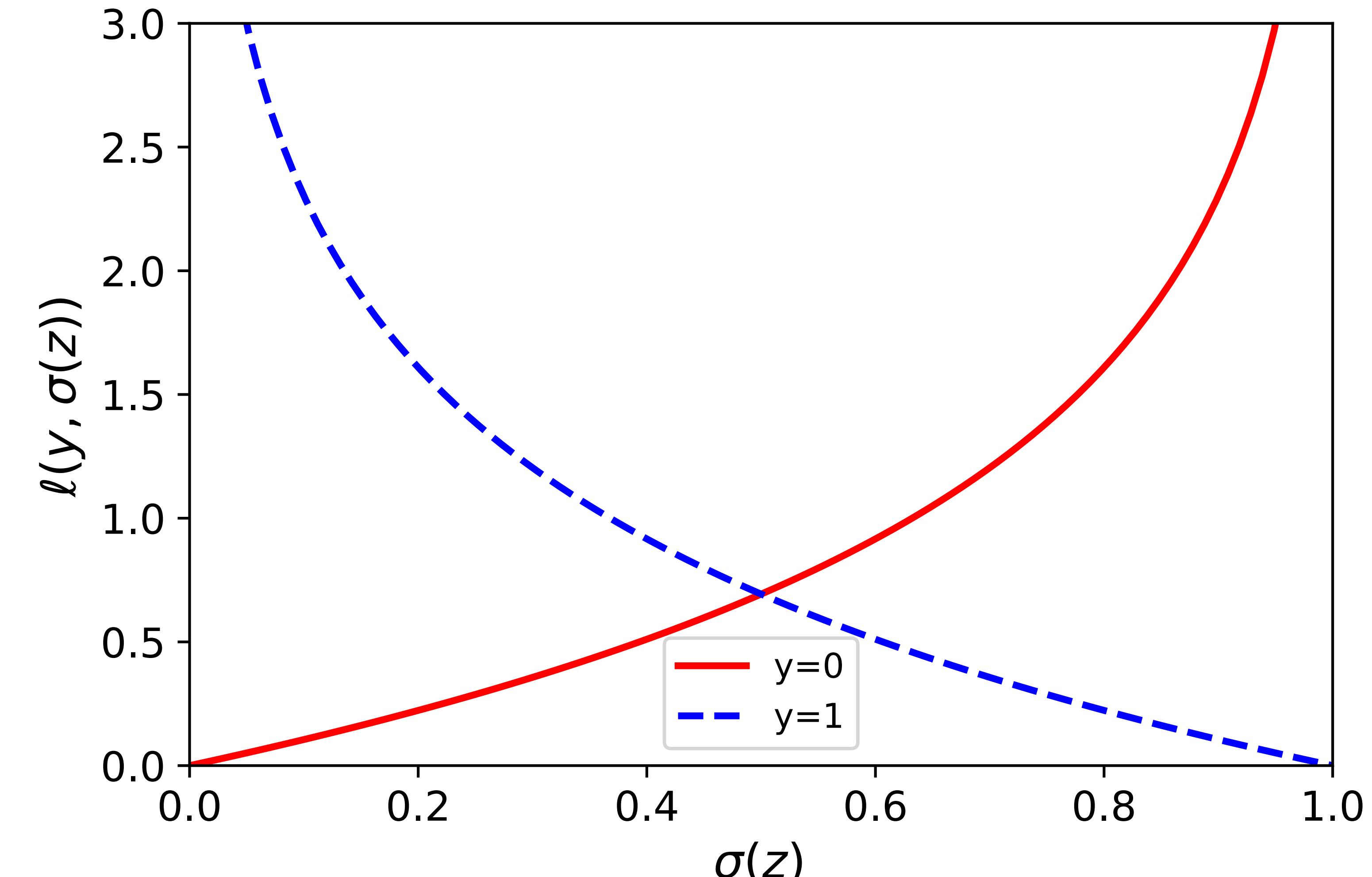


# What is the Loss Function?

- Define  $\sigma(z_i)$  as the predicted probability
- $y_i$  is our true label
- Can be viewed as the negative log-likelihood  
$$\ell(y_i, \sigma(z_i)) = -\log p(y_i | \sigma(z_i))$$
- Specific mathematical form is:  
$$\ell(y, \sigma(z)) = -y \log \sigma(z) - (1 - y) \log(1 - \sigma(z))$$

# Visualization of Logistics Loss Function

- The logistic/cross-entropy loss is:  
 $\ell(y, \sigma(z)) = -y \log \sigma(z) - (1 - y) \log(1 - \sigma(z))$
- Dashed blue line shows loss when the true prediction is positive
- When we give 100% of a 1, we pay no penalty
- If we are less confident or wrong, we pay an increasing penalty



Predicted Probability of a “one”

# Visualization of Logistics Loss Function

- The logistic/cross-entropy loss is:

$$\ell(y, \sigma(z)) = -y \log \sigma(z) - (1 - y) \log(1 - \sigma(z))$$

- Some example values:

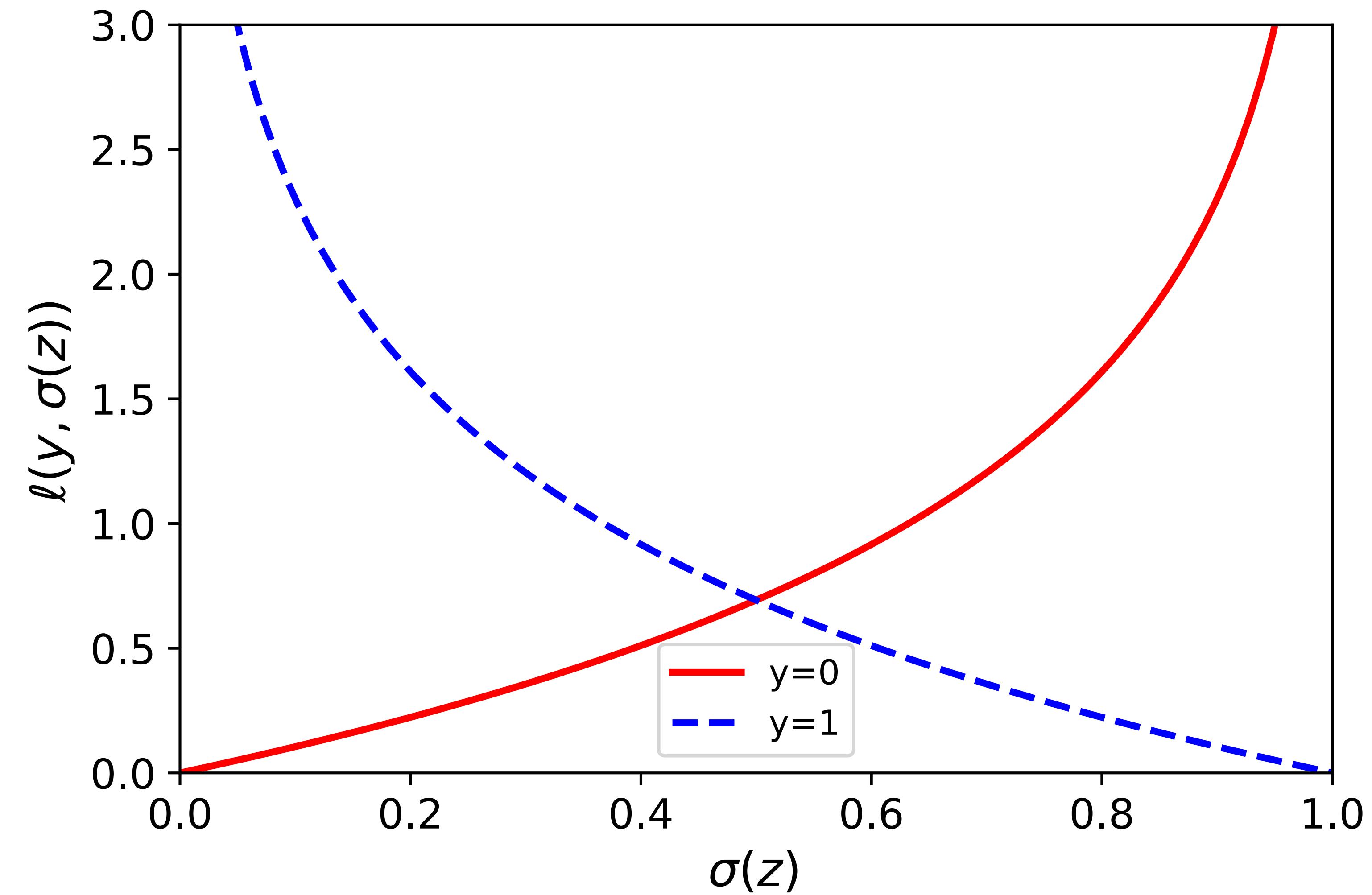
- $\ell\left(0, \frac{1}{2}\right) = \ell\left(1, \frac{1}{2}\right) = \log 2$

- $\ell\left(1, \frac{9}{10}\right) \approx .105$

- $\ell\left(0, \frac{9}{10}\right) \approx 2.303$

- $\ell\left(1, \frac{99}{100}\right) \approx 0.001$

- $\ell\left(0, \frac{99}{100}\right) \approx 4.605$



Predicted Probability of a “one”

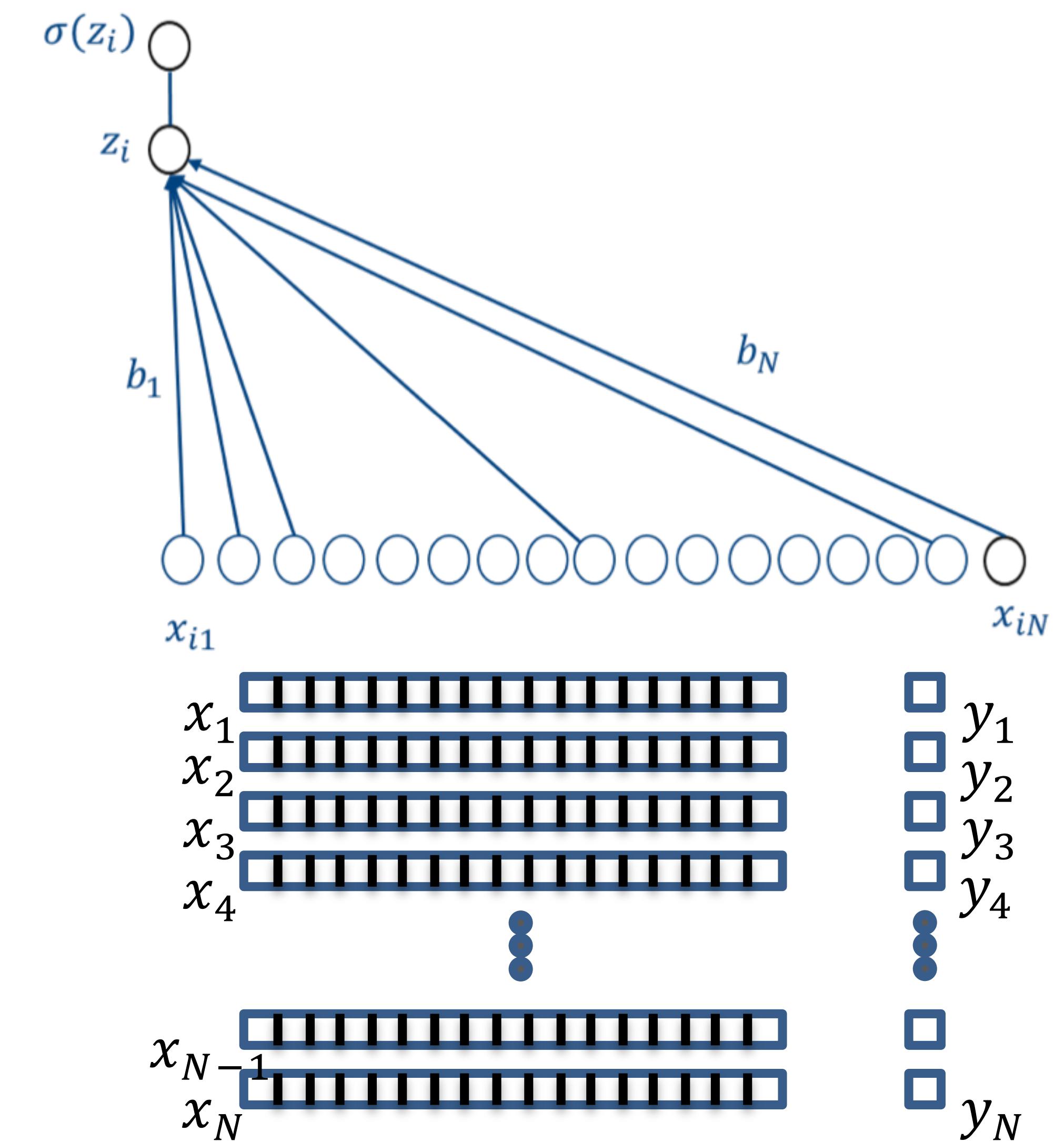
# Optimization Goal for Binary Classification

- The optimization goal is to minimize the average loss

$$\mathbf{b}^* = \arg \min_{\mathbf{b}} \frac{1}{N} \sum_i^N \ell(y_i, \sigma(z_i))$$

- For binary (0/1) problems, the logistic or cross-entropy loss is

$$\begin{aligned}\ell(y, \sigma(z)) &= -y \log \sigma(z) - (1 - y) \log(1 - \sigma(z))\end{aligned}$$

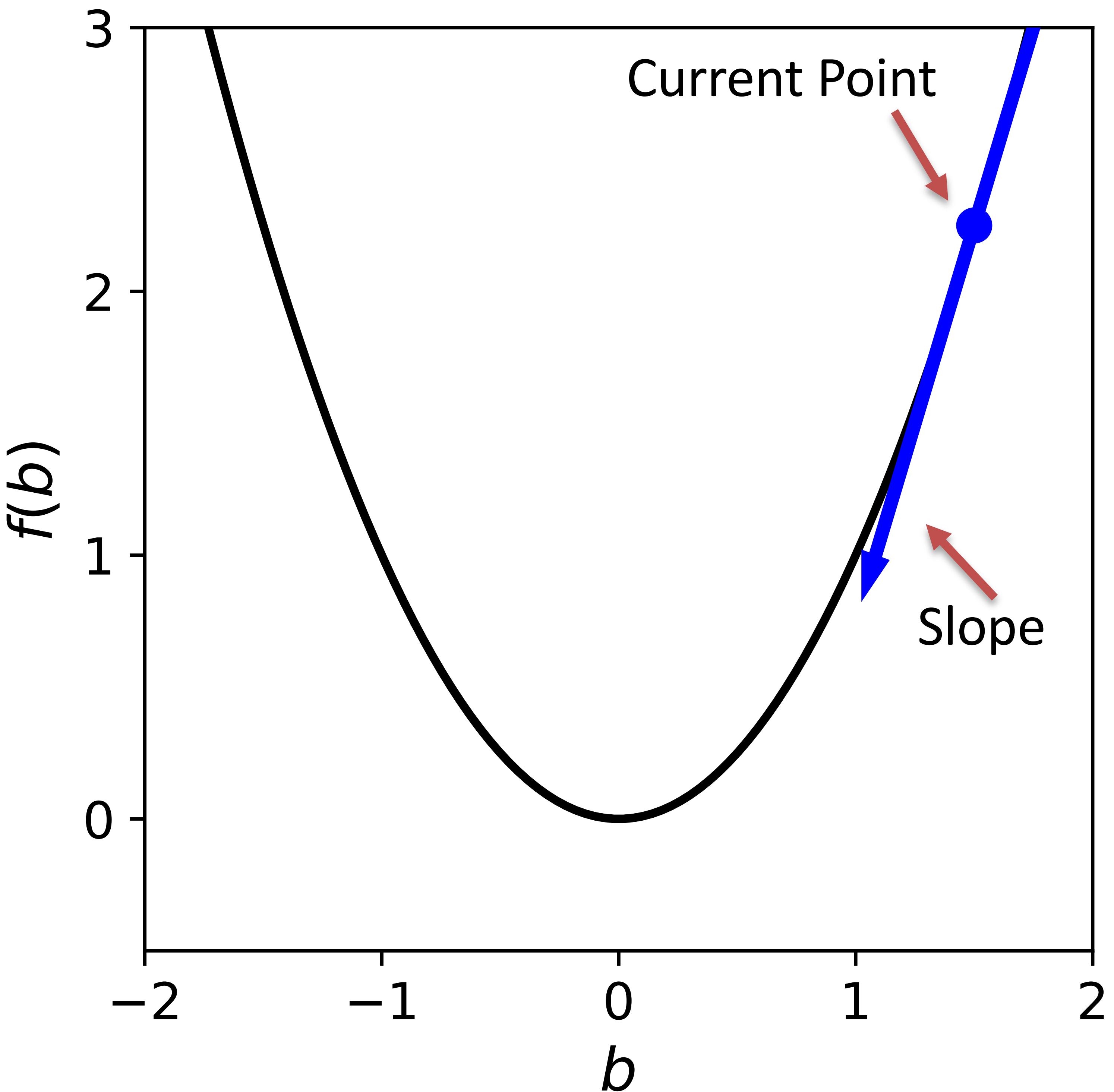


Brief Intro to Gradient Descent

# LEARNING THE PARAMETERS

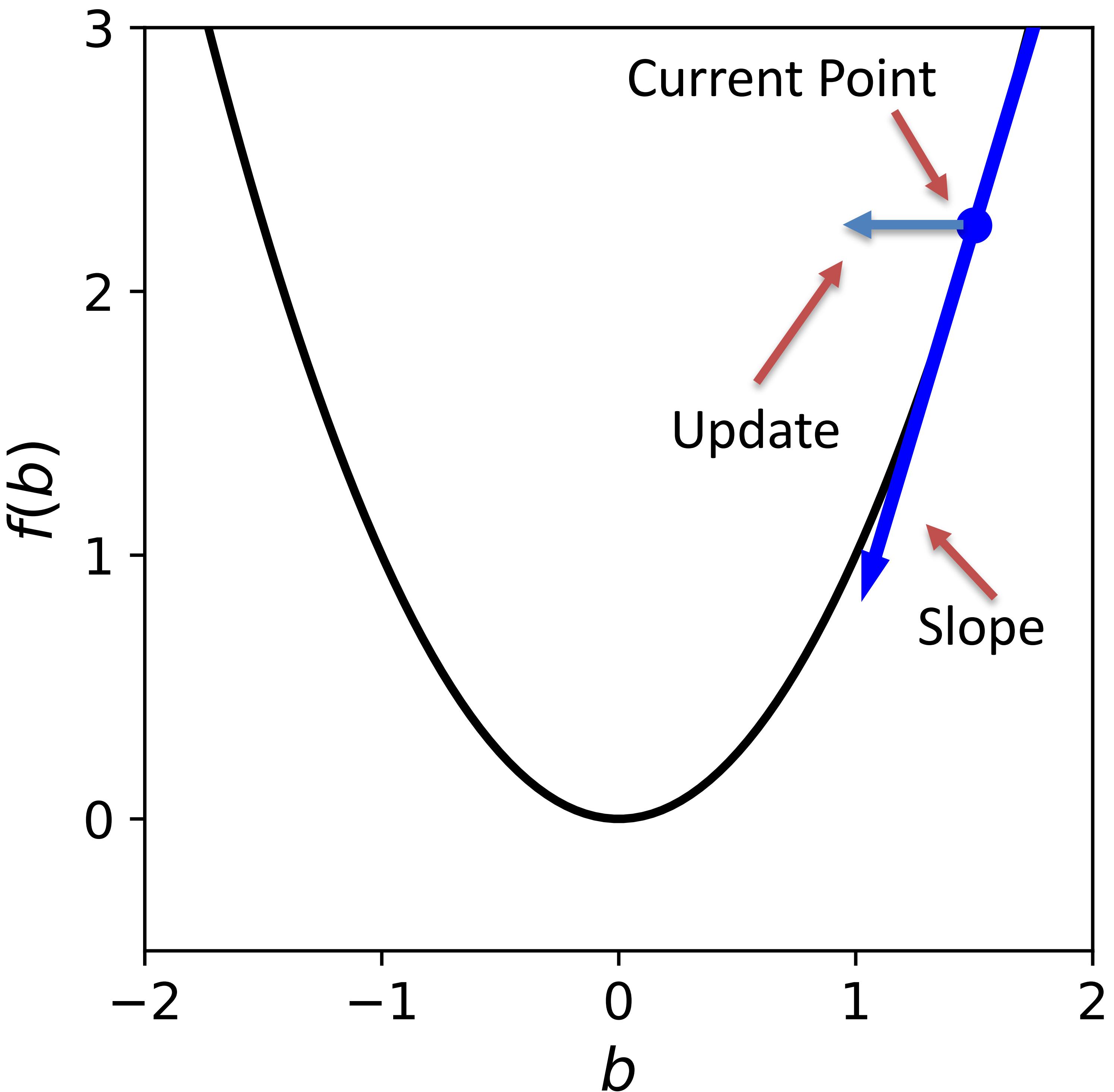
# Visualization of Optimization Method

- We want to minimize a mathematical function (i.e. our average loss function)
- One approach is to:
  1. Find the direction pointing “down the hill” (towards a smaller value)
  2. Move a bit in that direction
  3. Repeat 1-2 until satisfied



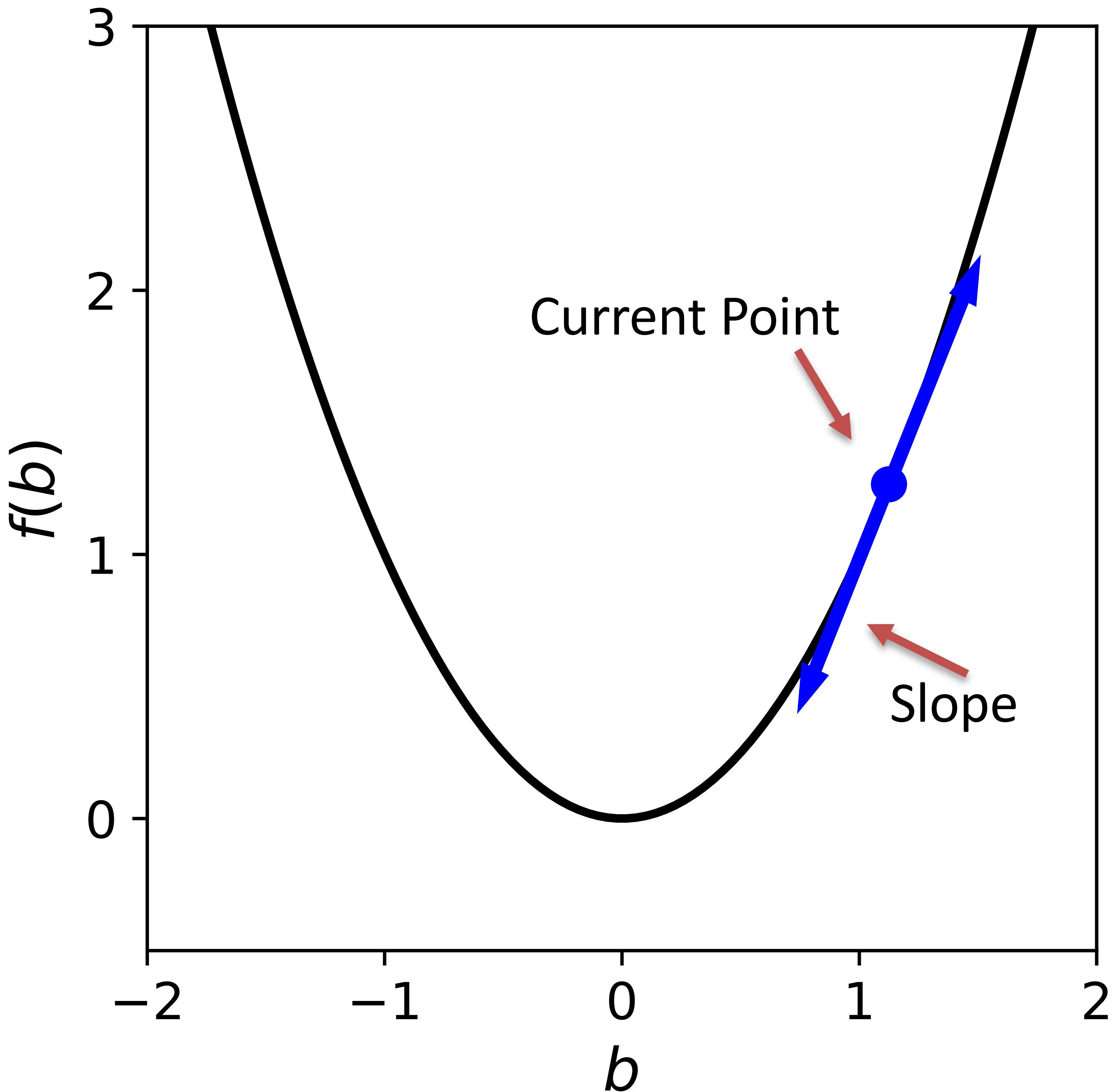
# Visualization of Optimization Method

- We want to minimize a mathematical function (i.e. our average loss function)
- One approach is to:
  1. Find the direction pointing “down the hill” (towards a smaller value)
  2. Move a bit in that direction
  3. Repeat 1-2 until satisfied



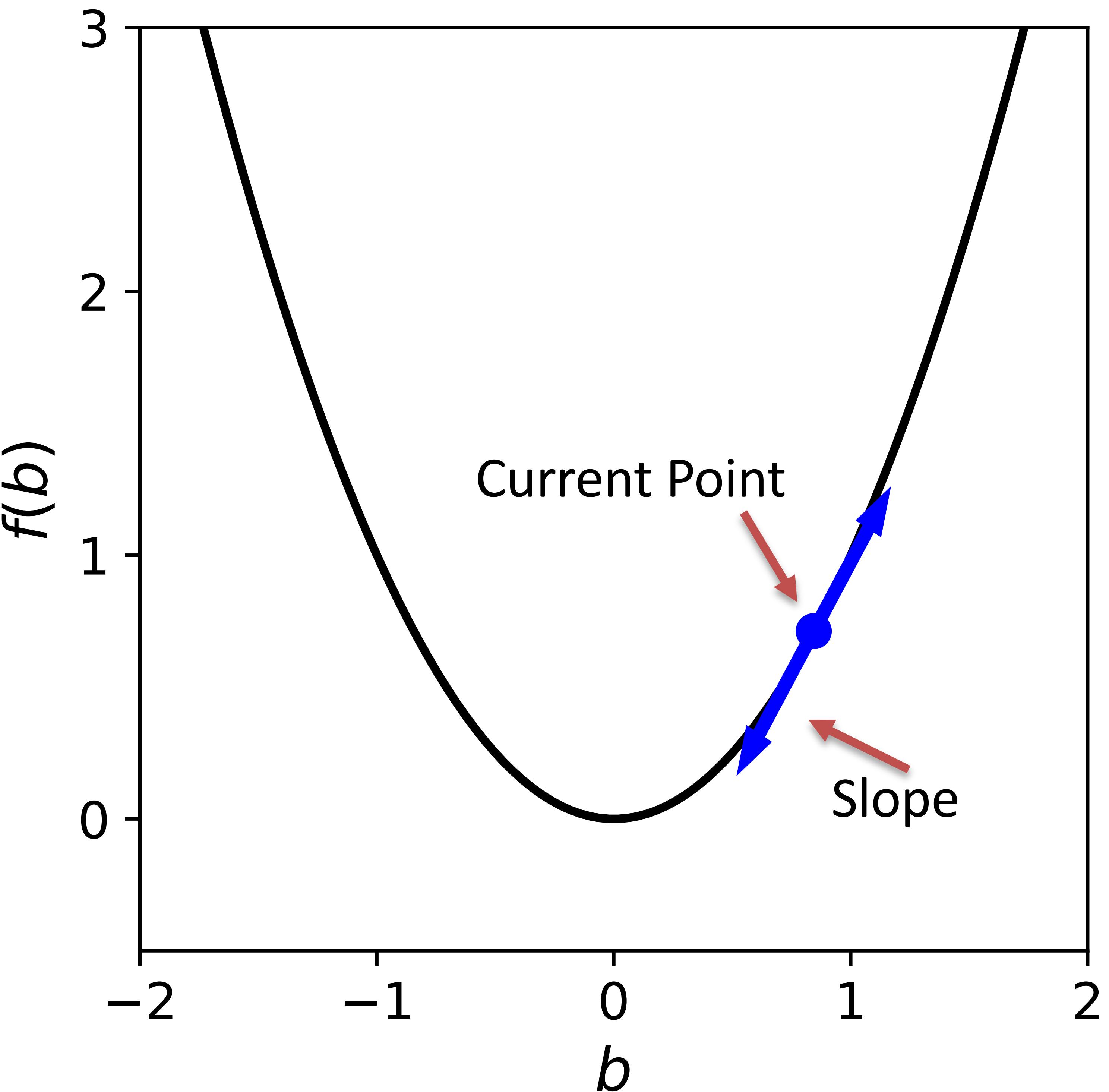
# Visualization of Optimization Method

- We want to minimize a mathematical function (i.e. our average loss function)
- One approach is to:
  1. Find the direction pointing “down the hill” (towards a smaller value)
  2. Move a bit in that direction
  3. Repeat 1-2 until satisfied
- This shows the **first** update



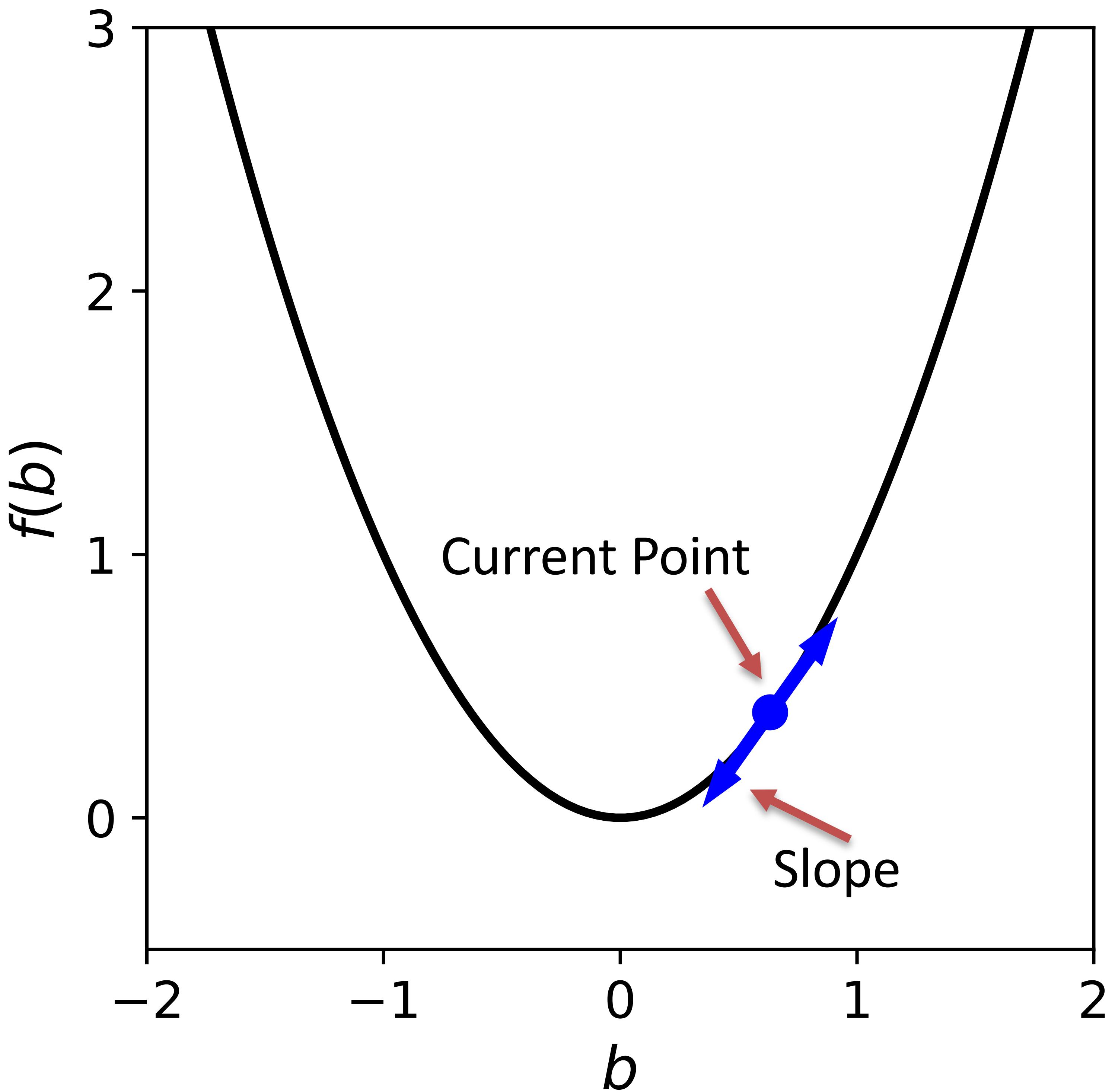
# Visualization of Optimization Method

- We want to minimize a mathematical function (i.e. our average loss function)
- One approach is to:
  1. Find the direction pointing “down the hill” (towards a smaller value)
  2. Move a bit in that direction
  3. Repeat 1-2 until satisfied
- This shows the **second** update



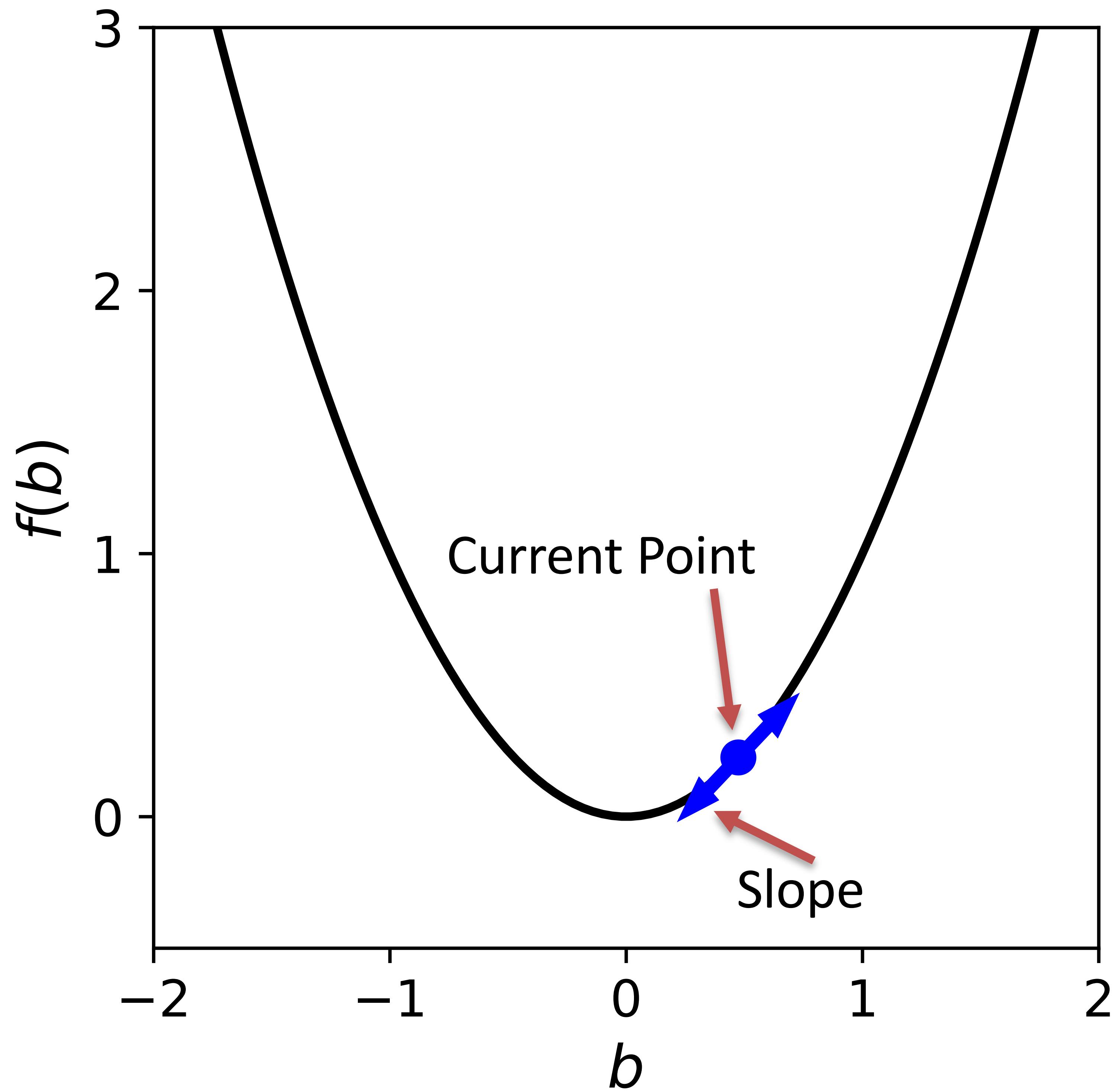
# Visualization of Optimization Method

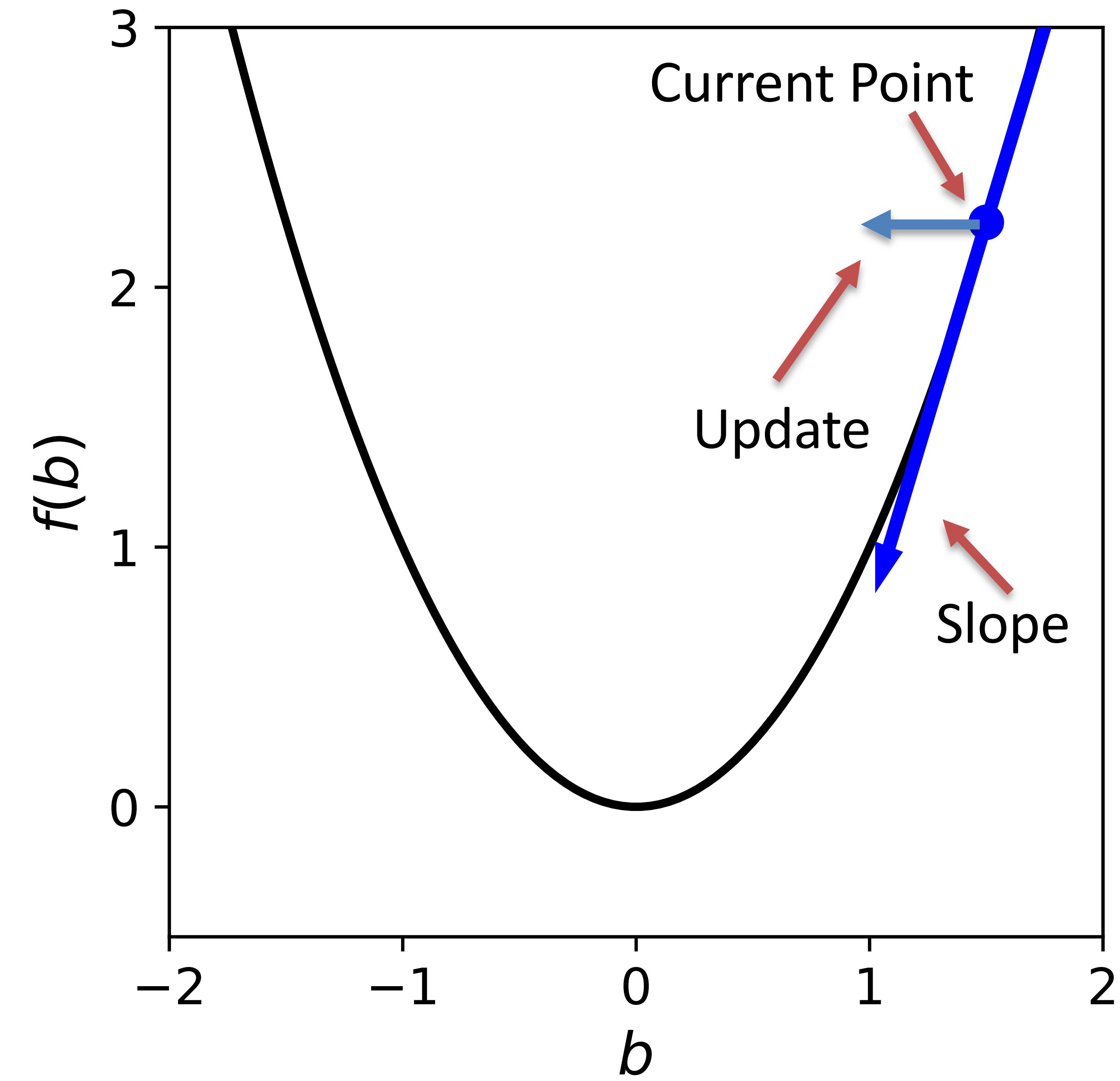
- We want to minimize a mathematical function (i.e. our average loss function)
- One approach is to:
  1. Find the direction pointing “down the hill” (towards a smaller value)
  2. Move a bit in that direction
  3. Repeat 1-2 until satisfied
- This shows the **third** update



# Visualization of Optimization Method

- We want to minimize a mathematical function (i.e. our average loss function)
- One approach is to:
  1. Find the direction pointing “down the hill” (towards a smaller value)
  2. Move a bit in that direction
  3. Repeat 1-2 until satisfied
- This shows the **fourth** update





## Mathematical Description of Gradient Descent

- We want to minimize a function  

$$b^* = \arg \min_b f(b)$$
- Start at an initial value  $b^0$
- We will run a series of updates to move from  $b^k$  to  $b^{k+1}$  (i.e. from  $b^0$  to  $b^1$ )
- Iteratively run the procedure:
  1. Calculate the slope at the current point (For one parameter, this is the derivative. For multiple parameters, this is the *gradient*.):  
 $\nabla f(b^k)$ ,  
 $\nabla$  means gradient or multidimensional slope
  2. Move in the direction of the negative gradient with step size  $\alpha^k$ :  

$$b^{k+1} = b^k - \alpha^k \nabla f(b^k)$$
  3. Repeat 1-2 until converged

Brief Introduction to *Stochastic* Gradient Descent (SGD)

# SCALING TO BIG DATA

# Return to our Optimization Goal

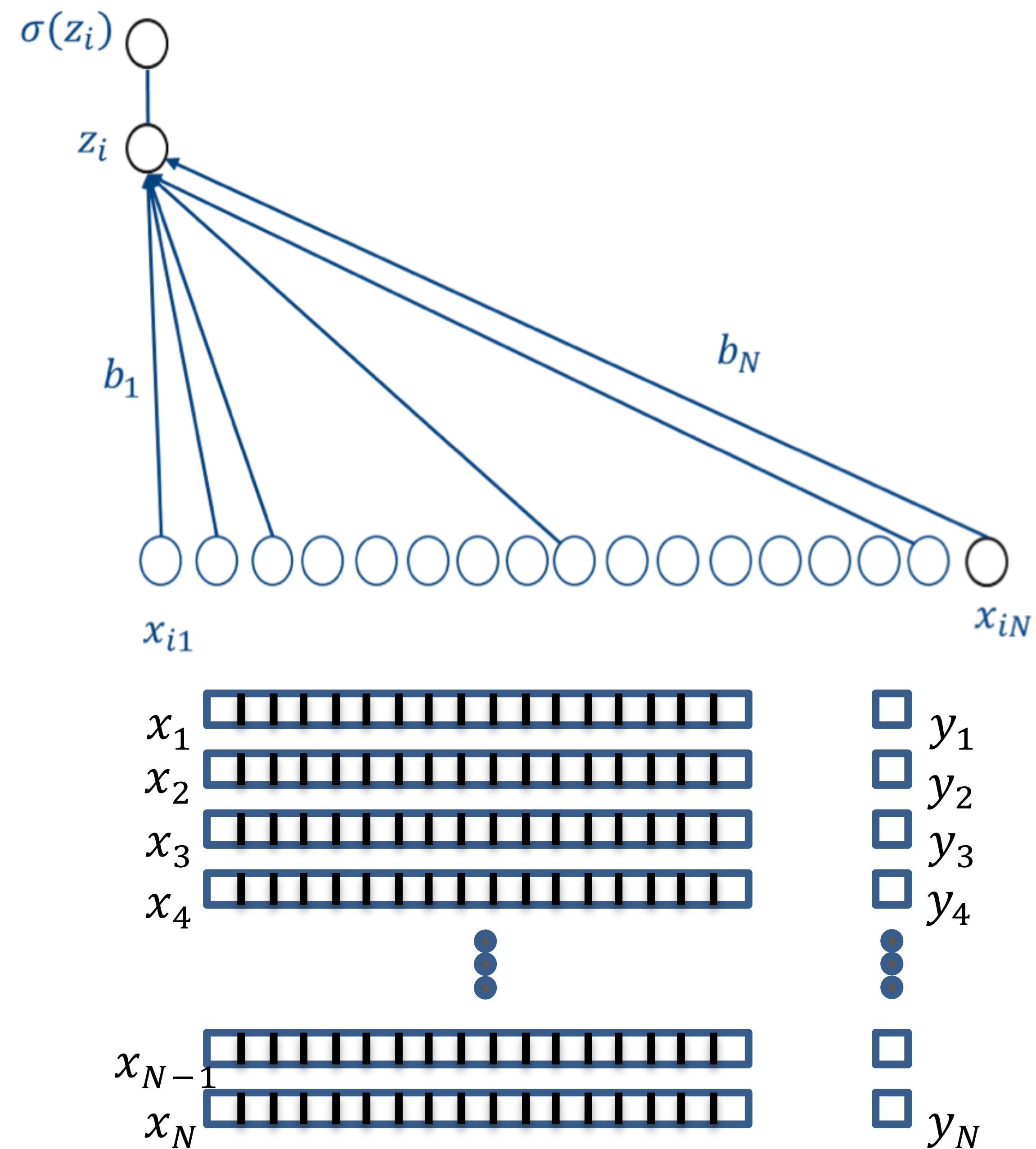
- We want to find the best parameters for all data points

$$\mathbf{b}^* = \arg \min_{\mathbf{b}} \frac{1}{N} \sum_i^N \ell(y_i, \sigma(z_i))$$

- However, calculating the gradient requires looking at every data point

$$\nabla \frac{1}{N} \sum_i^N \ell(y_i, \sigma(z_i)) = \frac{1}{N} \sum_i^N \nabla \ell(y_i, \sigma(z_i))$$

- This can be problematic in big data:
  - MNIST has ~60,000 images
  - ImageNet has ~1,000,000 images



# Let's Just Approximate the Gradient

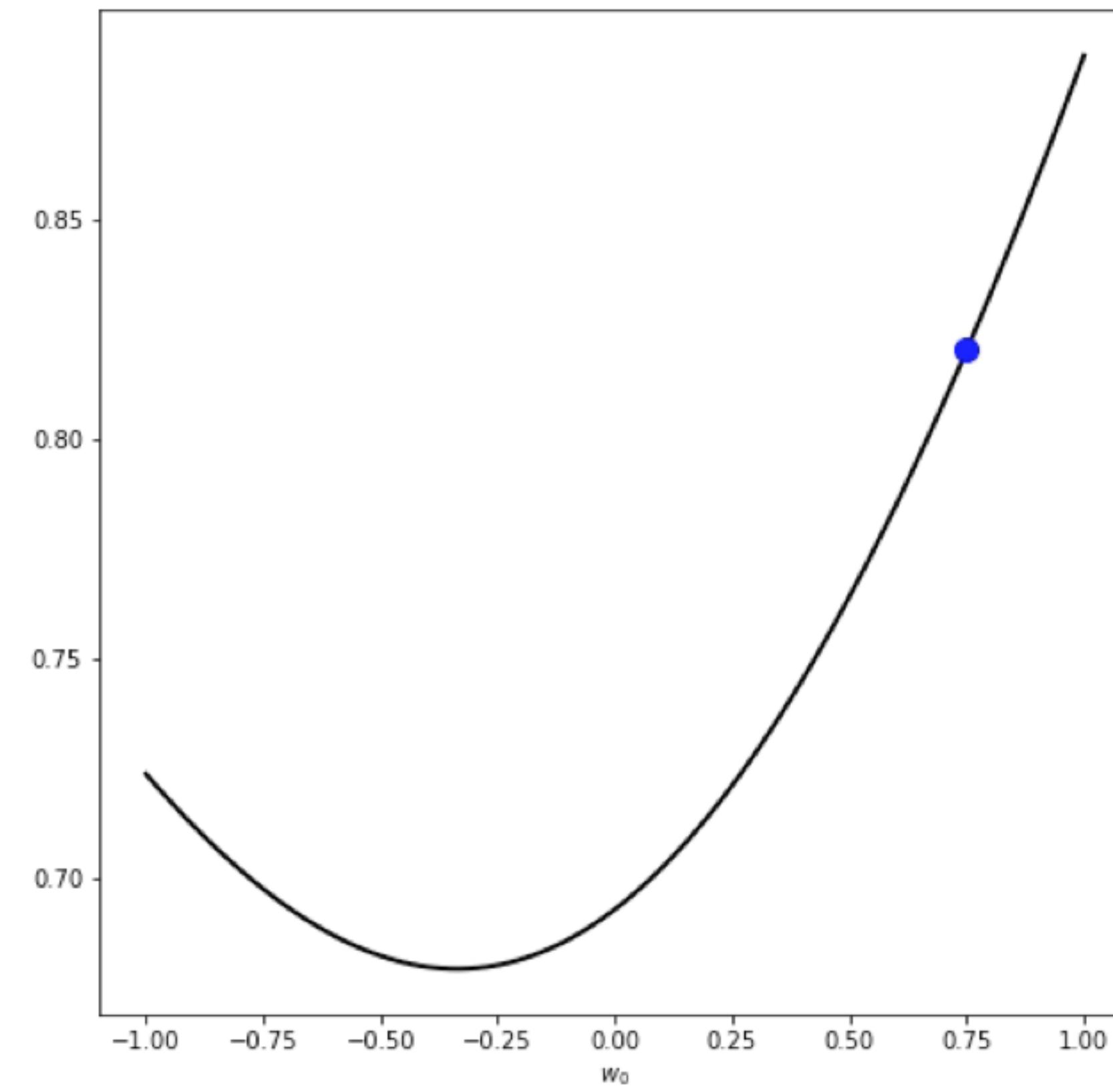
- We don't want to look at *every* data point to update our parameters
- Let's just take a single example  $j$  and use it to approximate the gradient

$$\nabla \ell(y_j, \sigma(z_j)) \simeq \frac{1}{N} \sum_{i=1}^N \nabla \ell(y_i, \sigma(z_i))$$

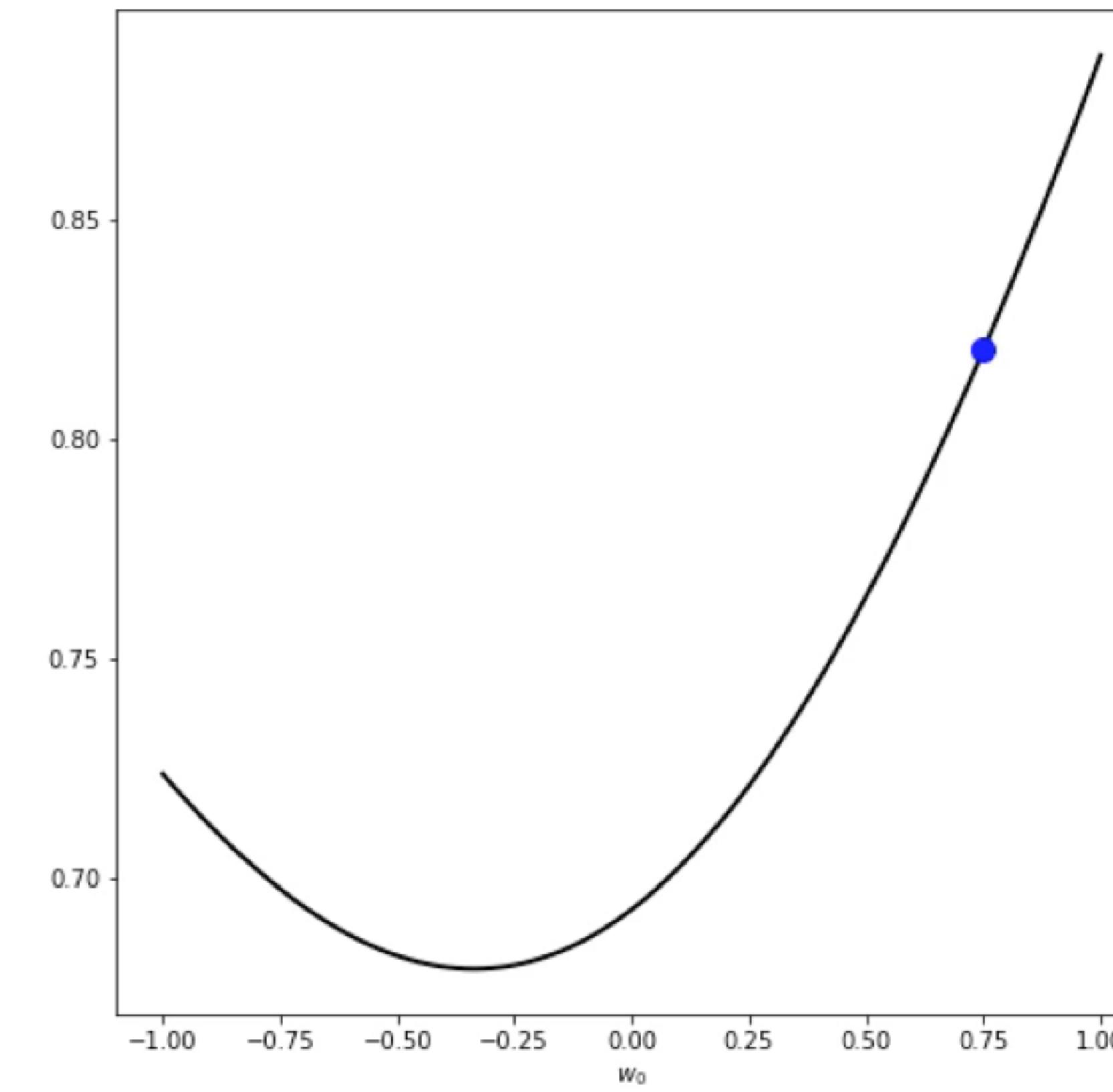
- In ImageNet, this is *1,000,000 times faster*
- Does this work? What would this look like?

# Visualizations

## Gradient Descent



## Stochastic Gradient Descent



# Data is Often Redundant



- 60,000 images, but really only 10 different types...

# Mathematical Descriptions

## Gradient Descent

1. Start with an initial  $\mathbf{b}^0$
2. Calculate gradient  $\nabla f(\mathbf{b}^k)$  over *all* data
3. Iteratively update:  
$$\mathbf{b}^{k+1} \leftarrow \mathbf{b}^k - \alpha_k \nabla f(\mathbf{b}^k)$$
4. Repeat 2-3 until solution is good enough

## Stochastic Gradient Descent

1. Start with an initial  $\mathbf{b}^0$
2. Choose a random data entry  $j$
3. Estimate gradient  $\widehat{\nabla f}(\mathbf{b}^k)$  by data point  $j$
4. Iteratively update:  
$$\mathbf{b}^{k+1} \leftarrow \mathbf{b}^k - \alpha_k \widehat{\nabla f}(\mathbf{b}^k)$$
5. Repeat 2-4 until solution is good enough

# Question:

- If the updates for gradient descent and stochastic gradient descent start at the same place, how do their resulting updates vary? How about the expectations of their updates?
- Reminder that:

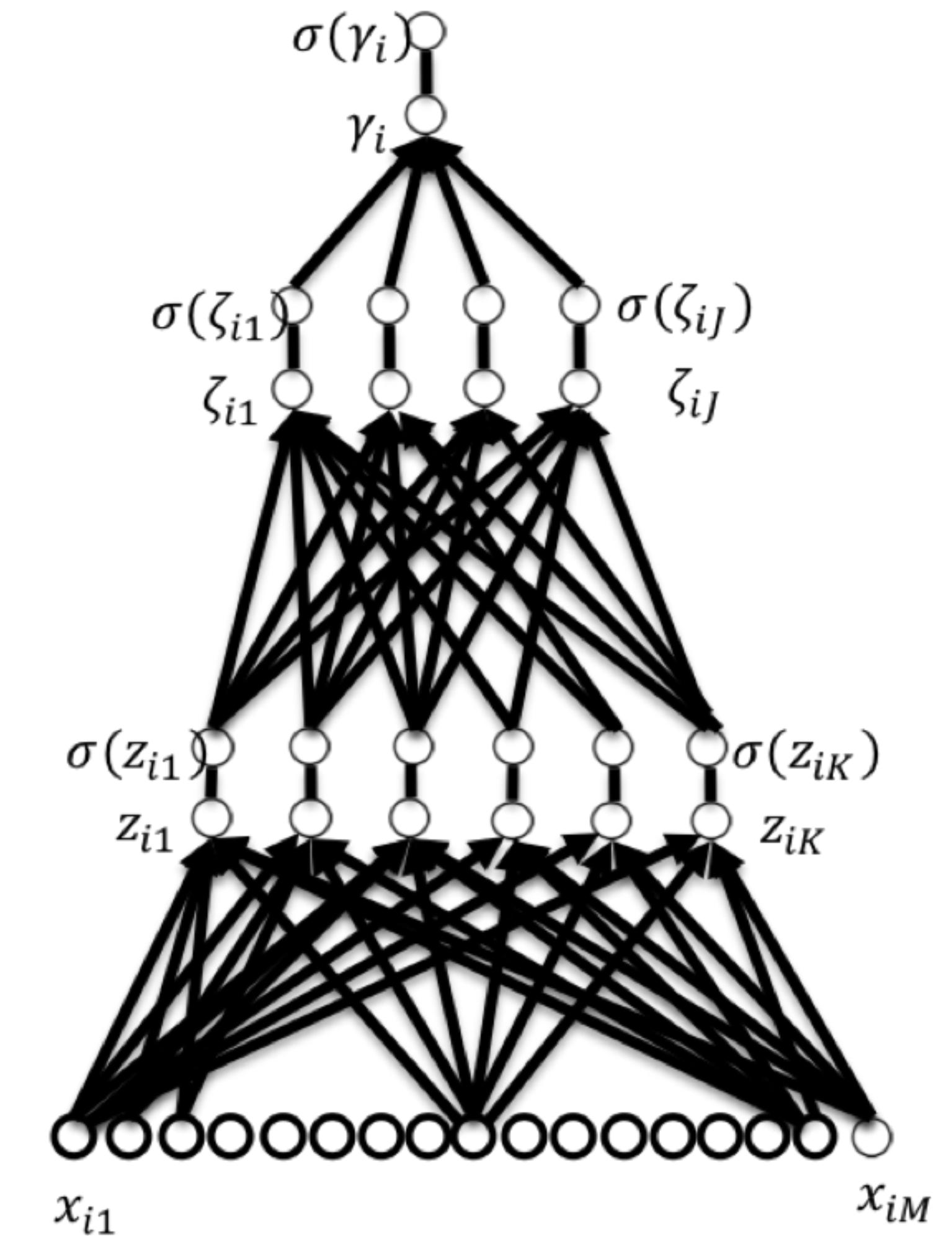
$$\mathbb{E}_{j \sim Unif(1, \dots, N)} [\nabla f_j(\mathbf{b})] = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{b})$$

# Initial Conclusions

- Stochastic Gradient Descent can update *many more times* than Gradient Descent
- Gets *near* the solution very quickly
- Allows scaling to *big data* (update time doesn't increase with the data size)
- In practice, we often use a minibatch, which uses a few data examples to estimate the gradient

# Calculating the Gradient

- Calculating the gradient (multi-dimensional slope) on a multi-layer perceptron may seem daunting
- Algorithms can automatically calculate the gradient (“backpropagation”)
- In code, we simply have to use a “gradient” function call once the network is defined



Building Intuition

# **SMALL-SCALE EXAMPLE OF STOCHASTIC UPDATES**

# Finding the mean

- Consider the collection of points  $\{1,2,3,4,5\}$  and we want to find the mean:

$$F(w) = \frac{1}{5} \sum_{i=1}^5 f_i(w), \quad f_i(w) = (i - w)^2$$

- There is no need or benefit to stochastic optimization here – minimizer is 3
- Useful to think concretely about why stochastic optimization can get us near the solution

# Finding the mean

Simple example:

$$f_i(w) = (i - w)^2$$

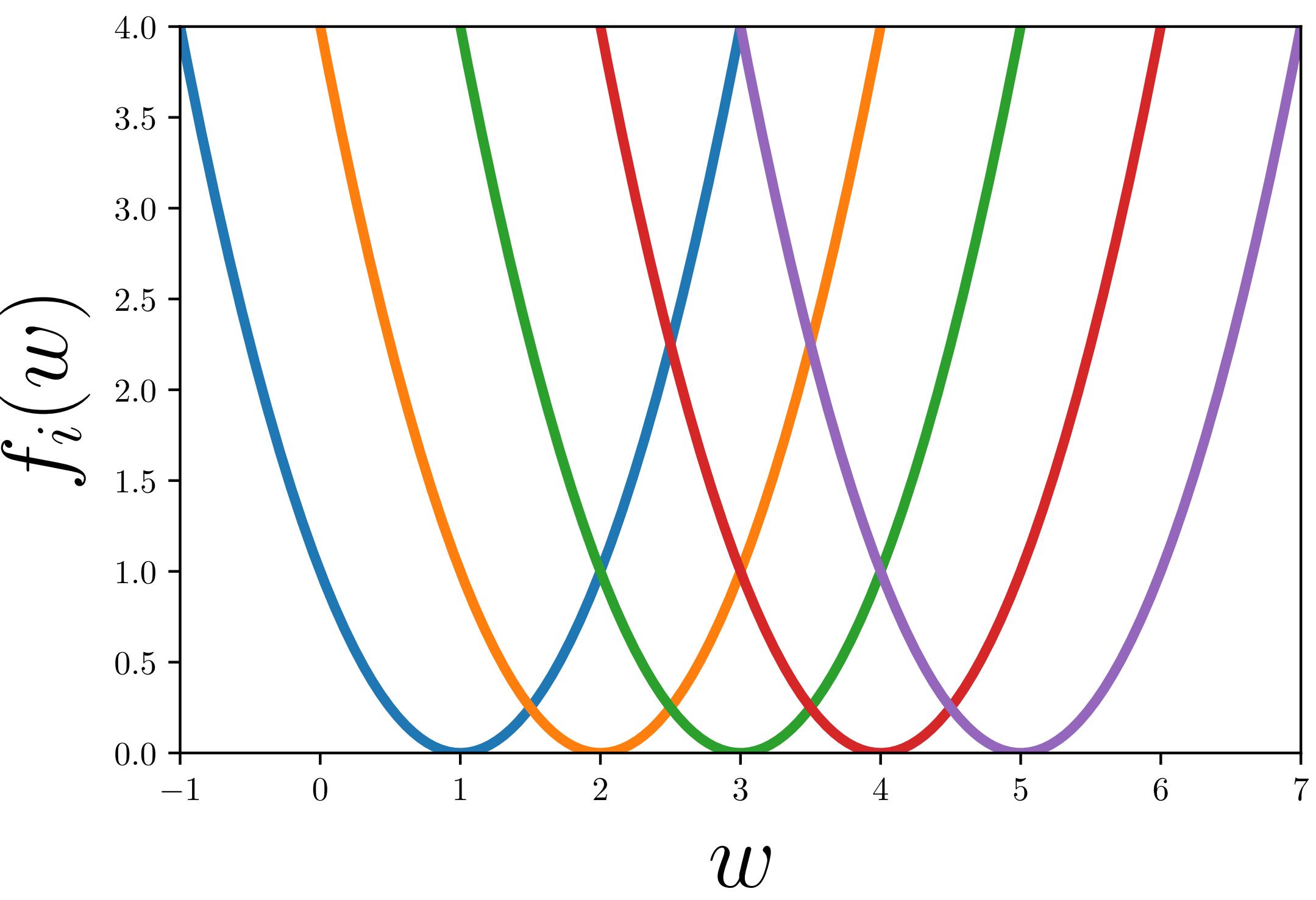
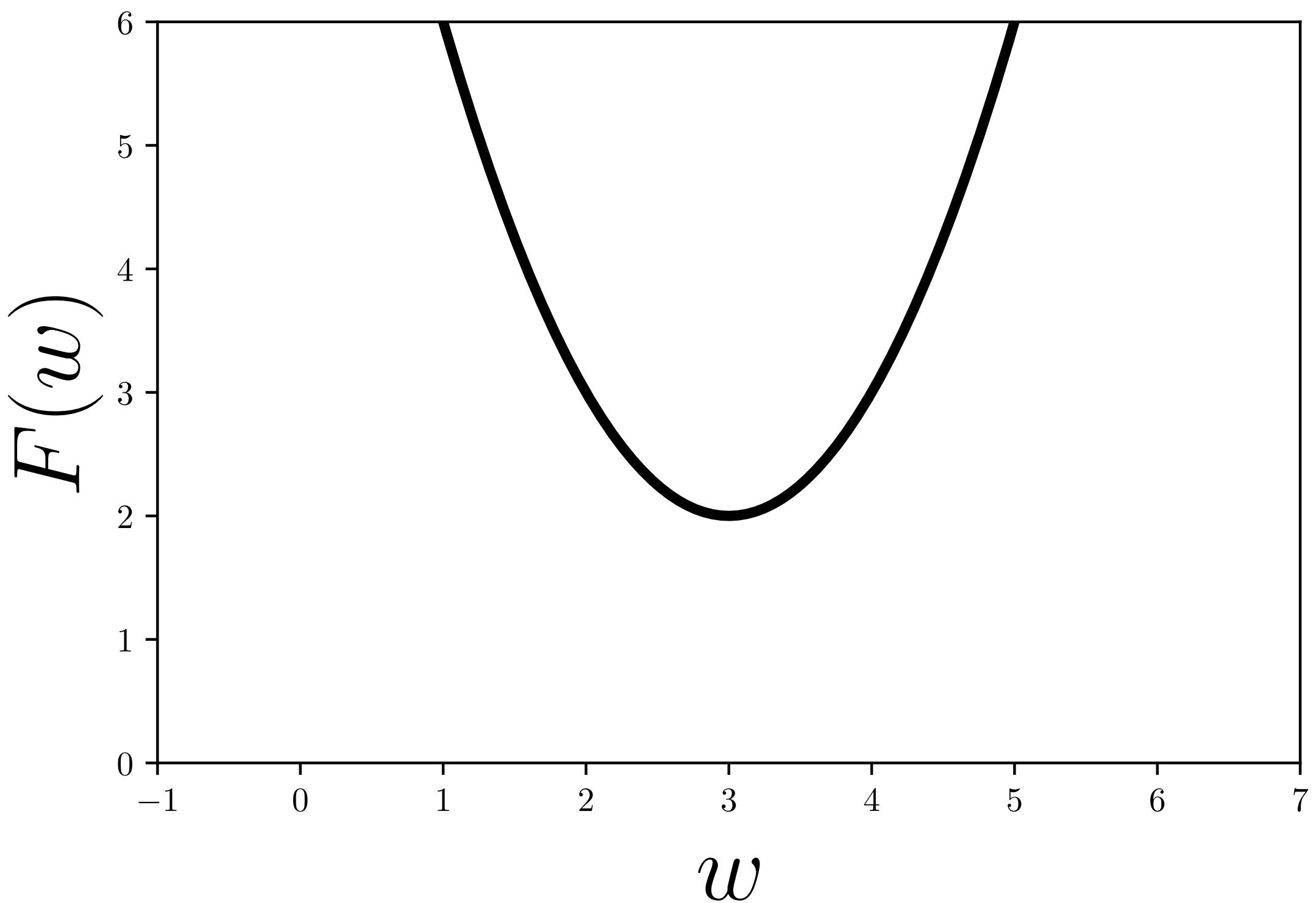
Can we find the mean with stochastic gradient descent?

What happens when  $w_0 = -100$ ?

What happens when  $w_0 = 0$ ?

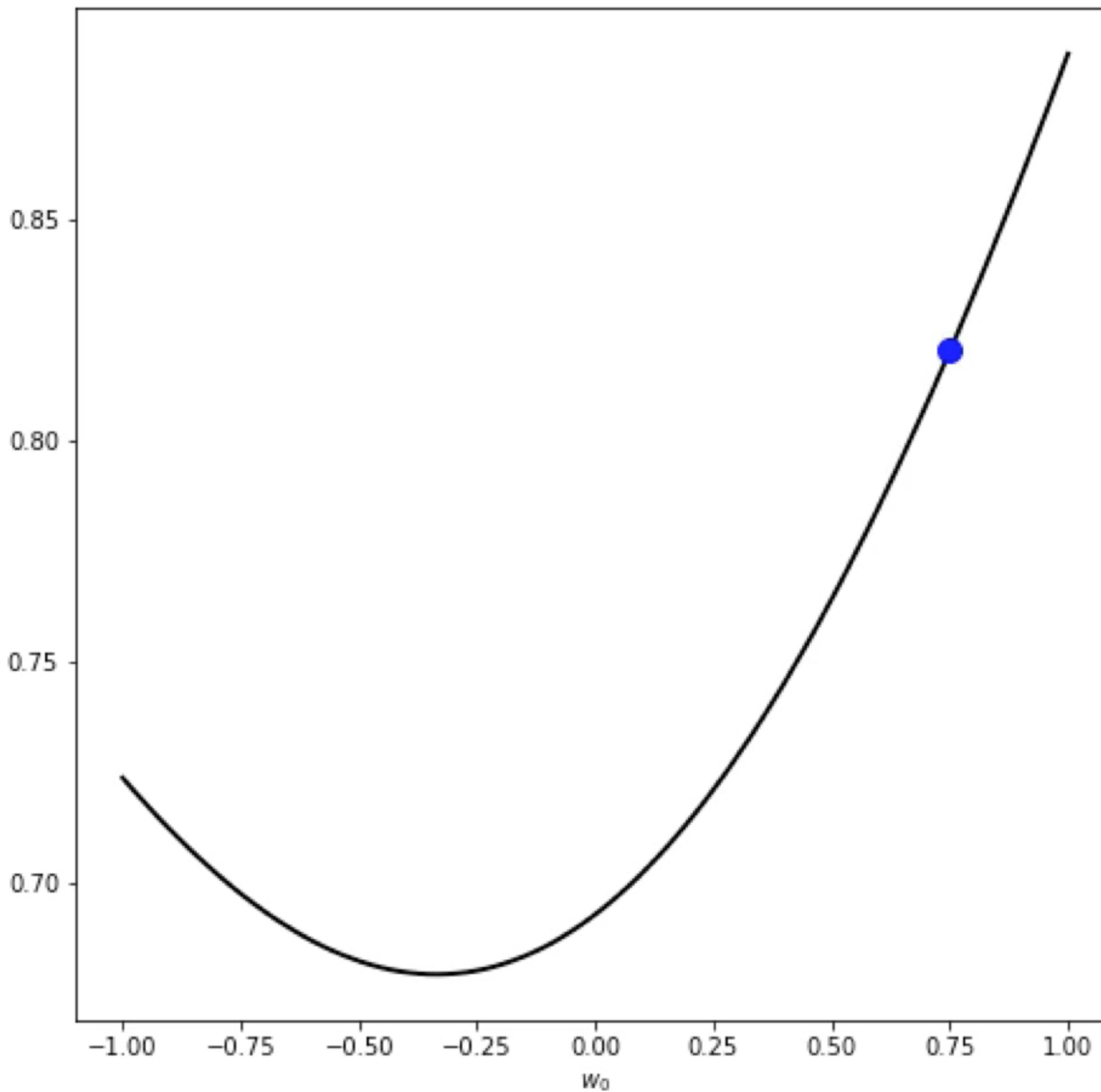
What happens when  $w_0 = 2$ ?

What happens when  $w_0 = 3$ ?

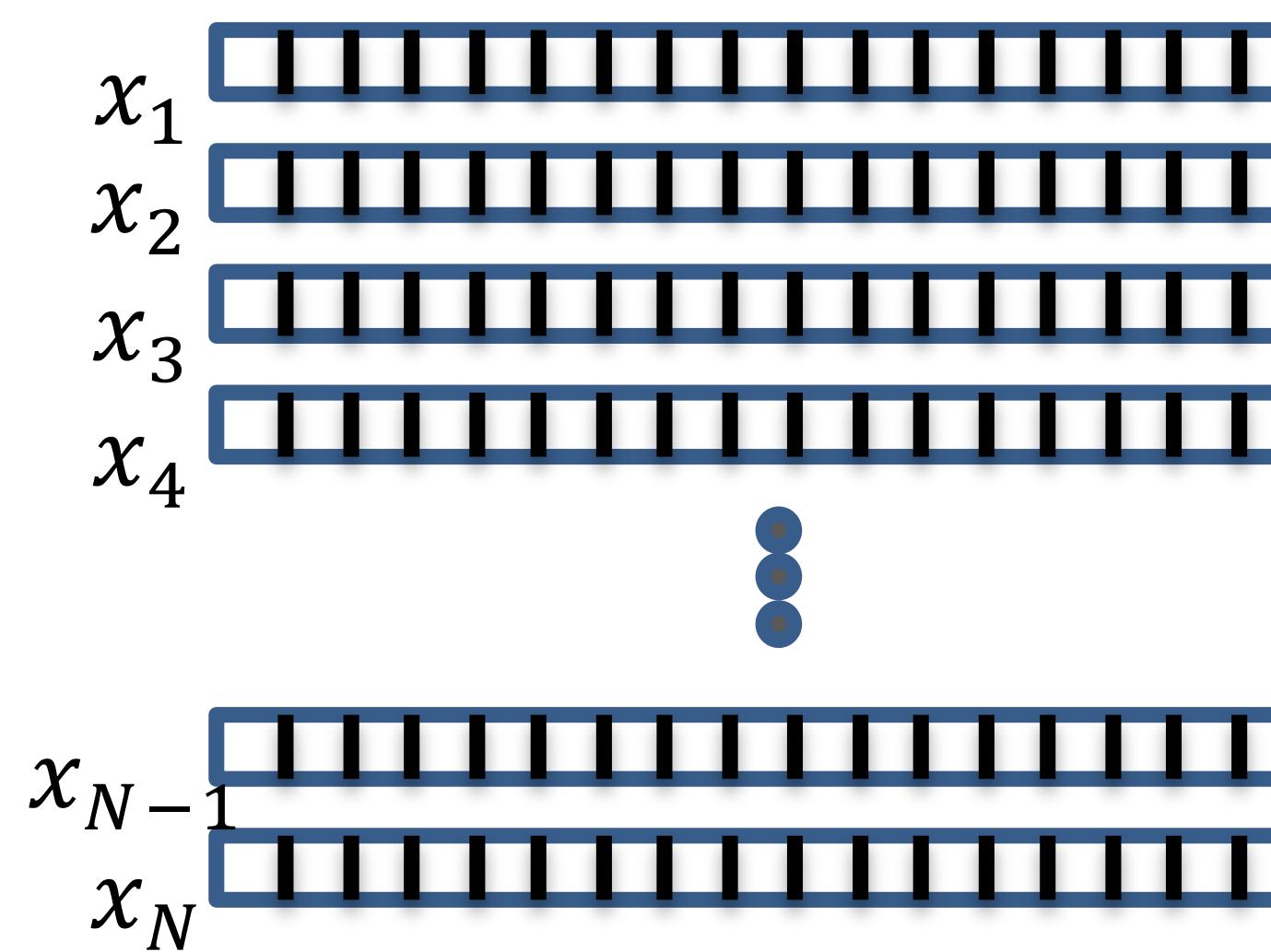


# Revisiting our Visualization

- When the starting point is “far away” from the optimal, then we typically move in the right direction
- When the current point is near the optimal value, it is very random whether we move in the correct direction or not

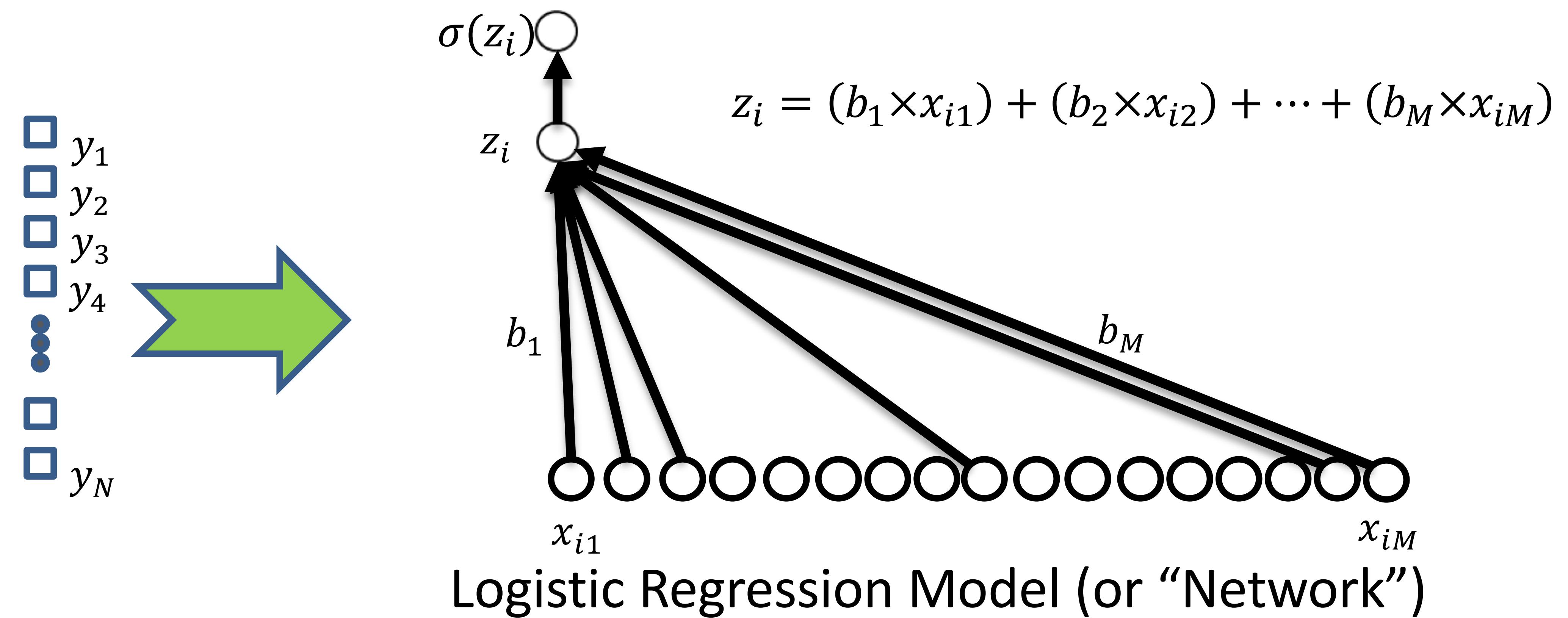


# Learning Model Parameters (Recap)



Training Set

Stochastic Gradient Descent



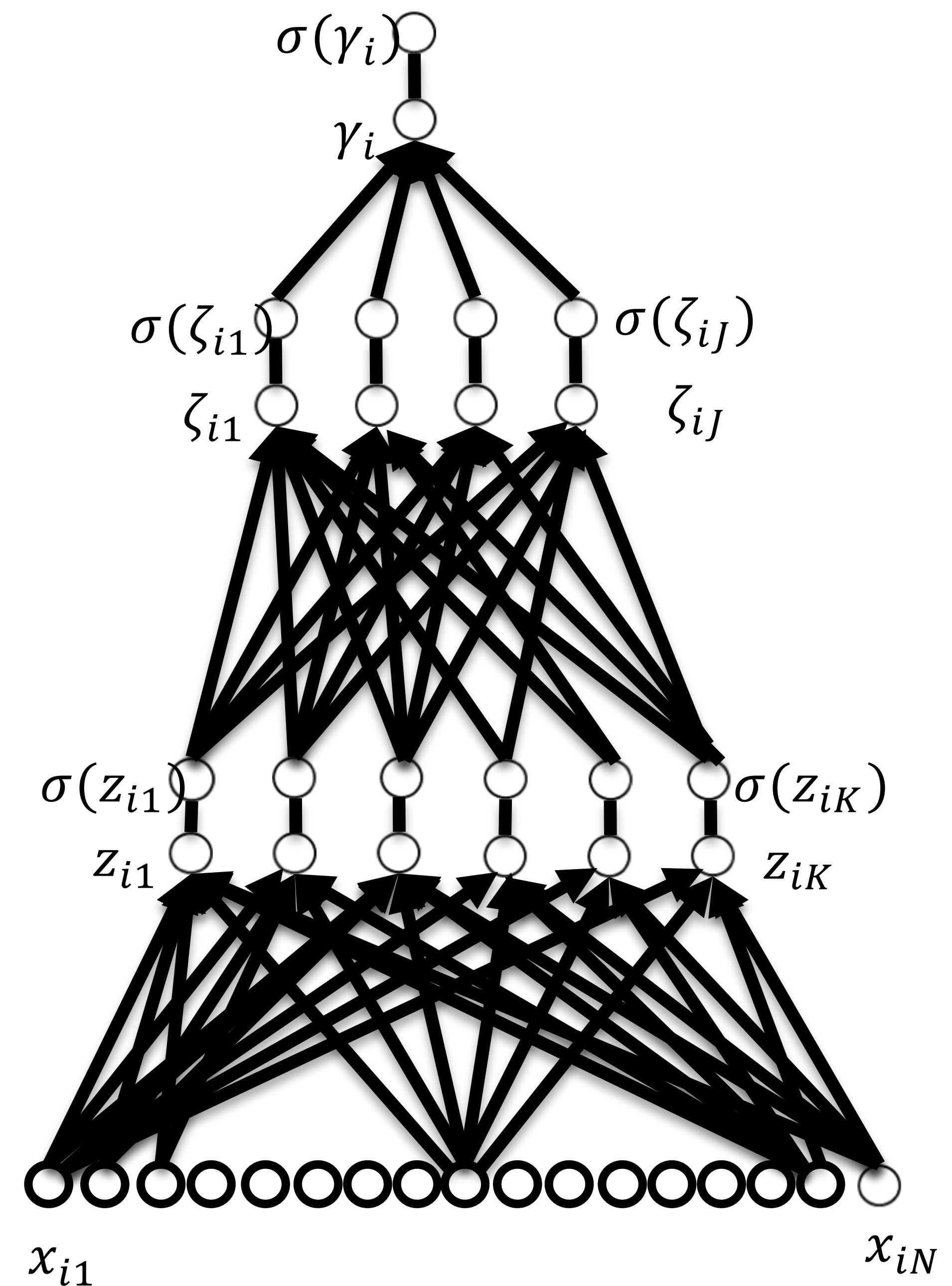
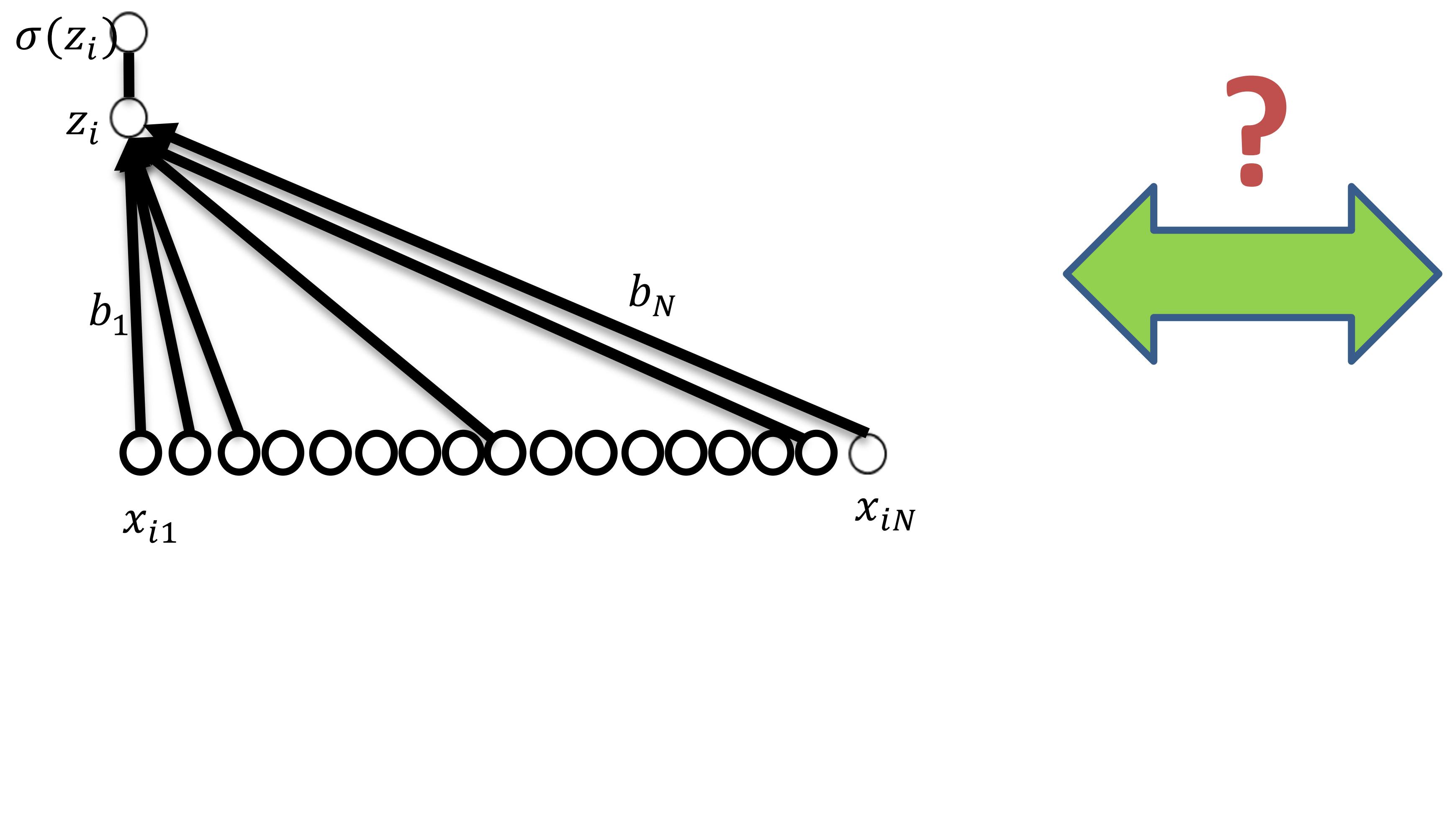
Logistic Regression Model (or “Network”)

Learned  
Parameters  $(b_0, b_1, \dots, b_N)$

An Introduction to Model Validation

# **ESTIMATING PERFORMANCE**

# Do we always want to increase complexity?

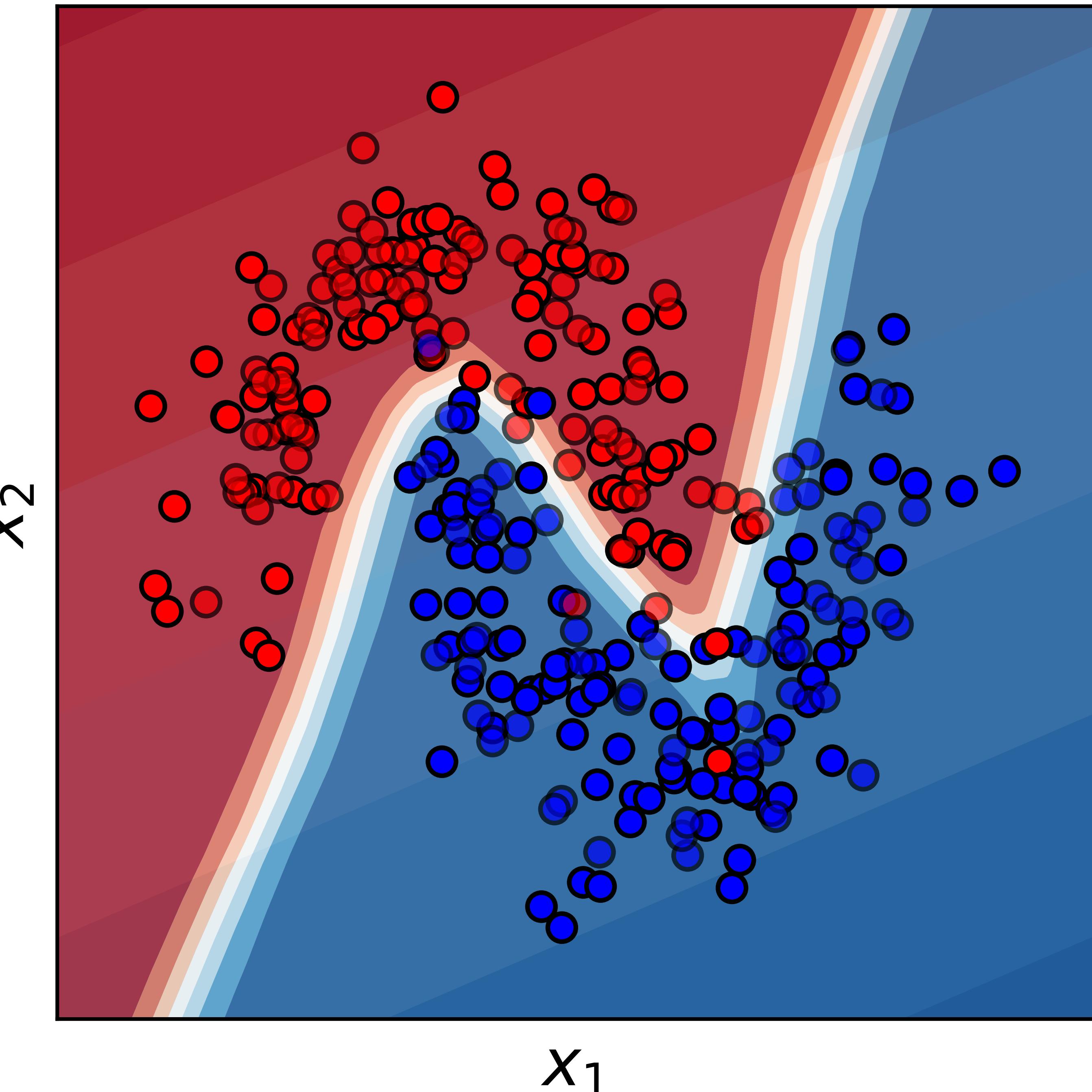


# Creating deep models can help us learn complex relationships

My learning multiple layers and transformation, it is possible to have a non-linear classifier capable of more accurately capturing the data.

However, a deep model can also give “perfect” performance on the training dataset and **fail completely** in the real world

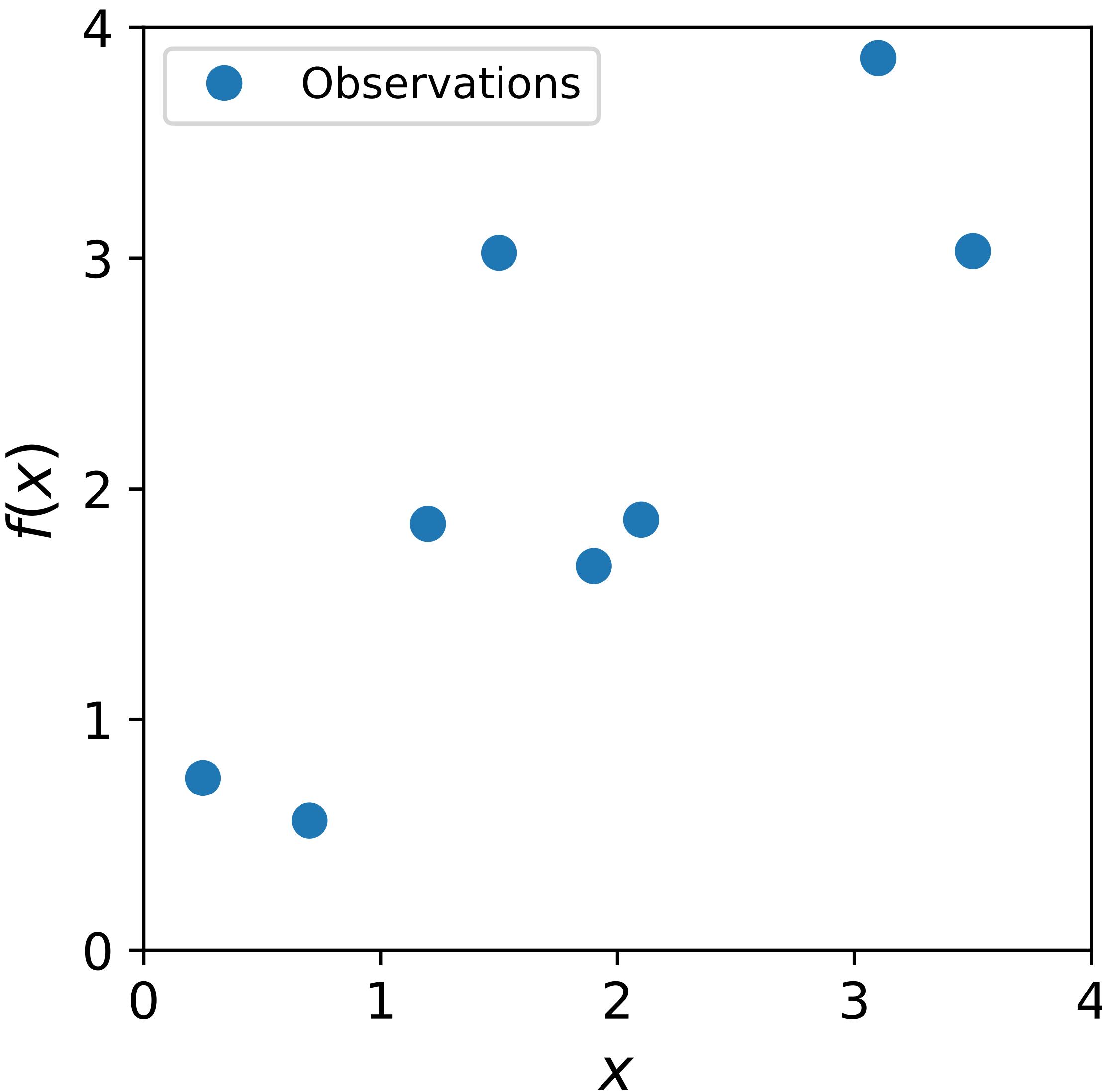
We need to *validate* the performance



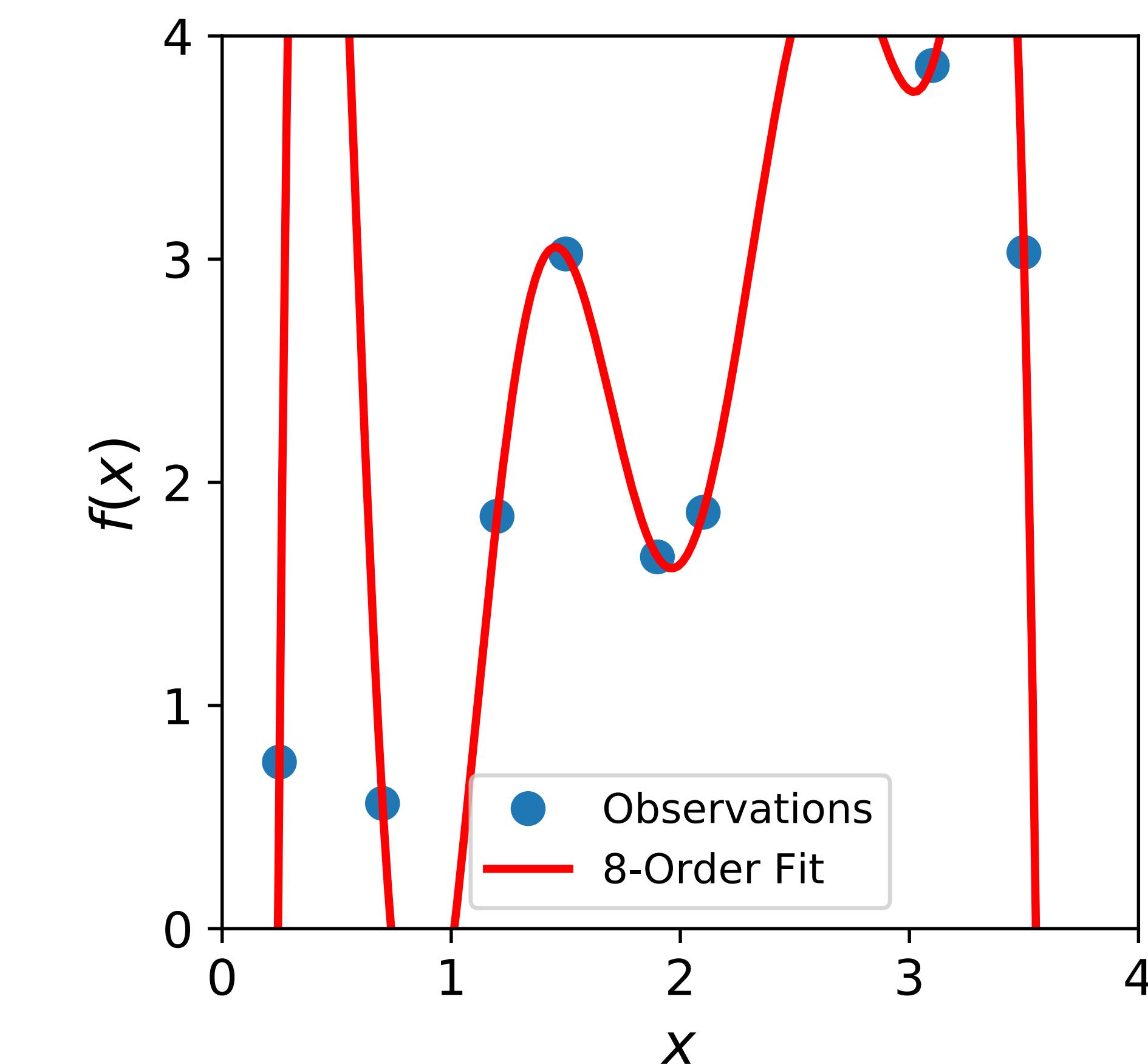
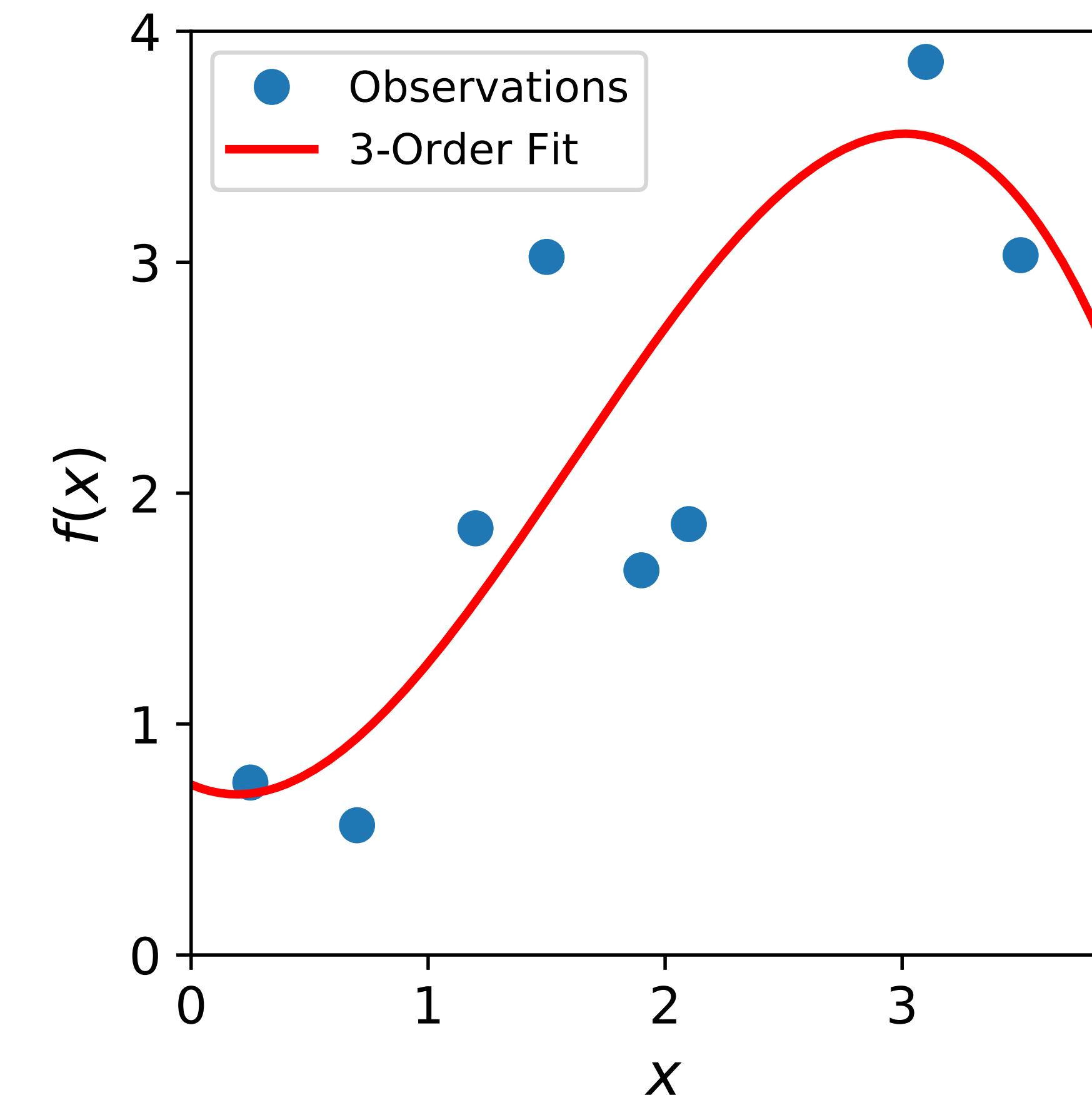
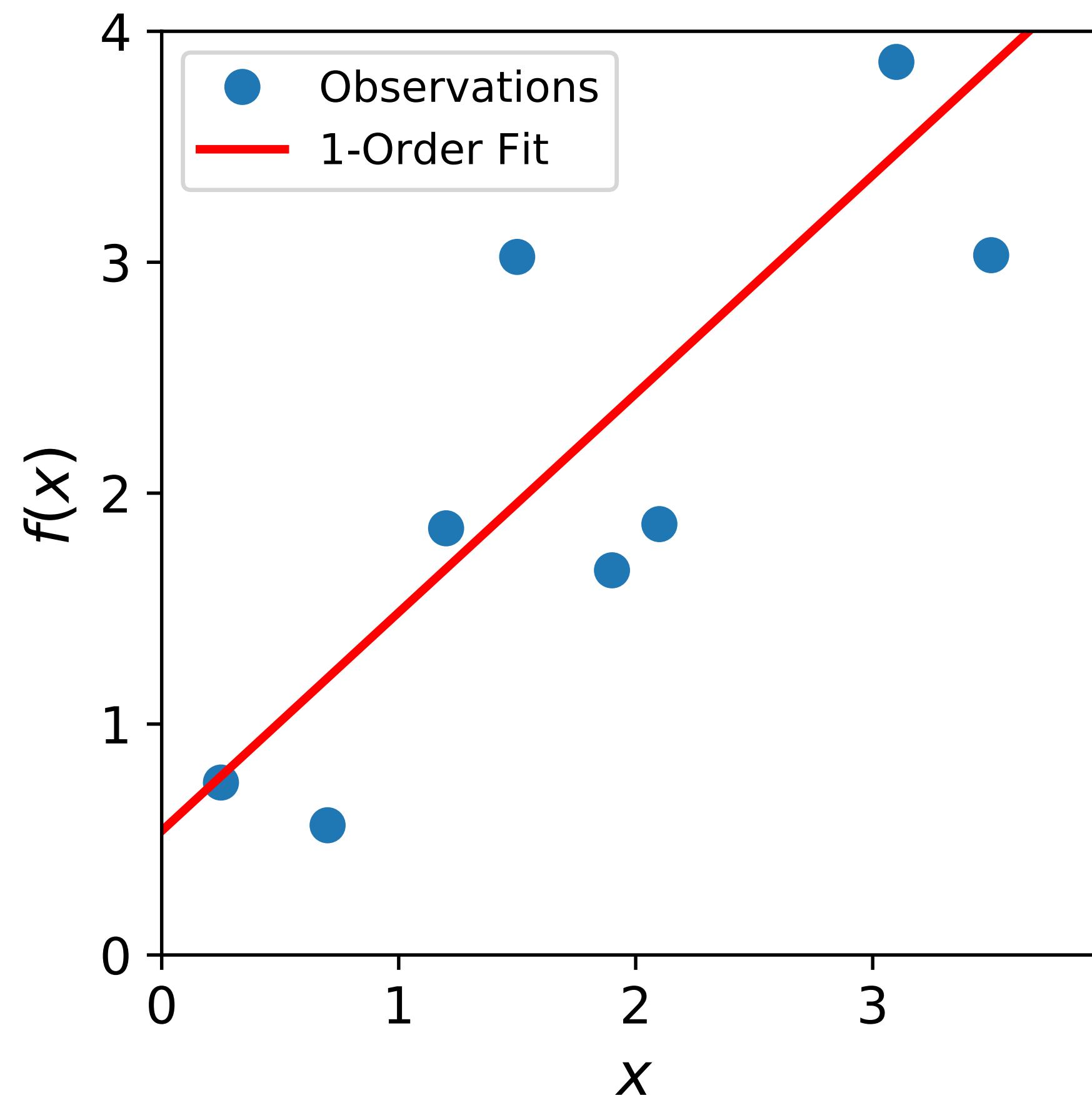
# Overfitting

“Overfitting” is when the learned model increases complexity to fit the observed training data *too well*  
– will not work to predict future data!

What would we want to use to fit these example data points?



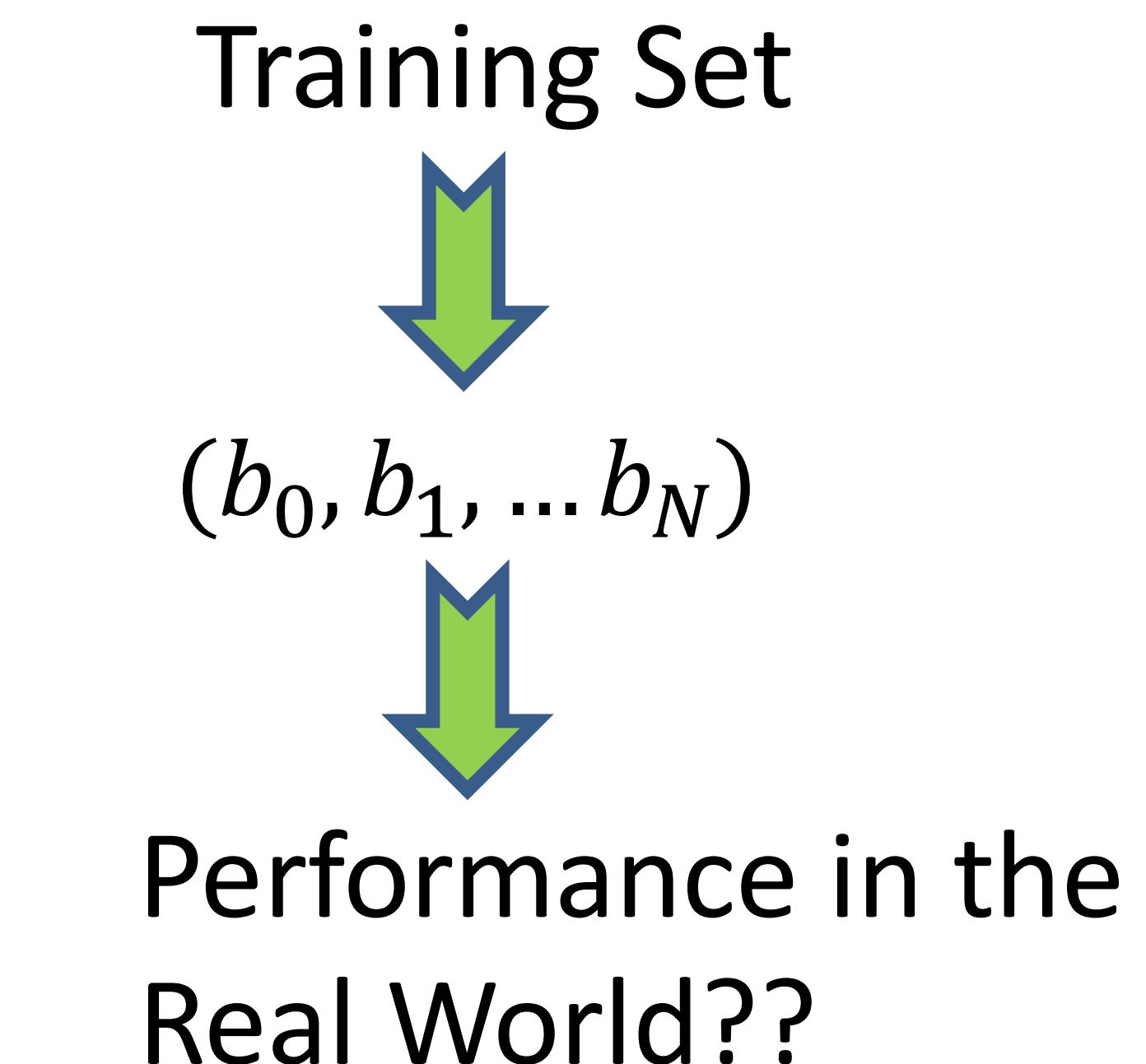
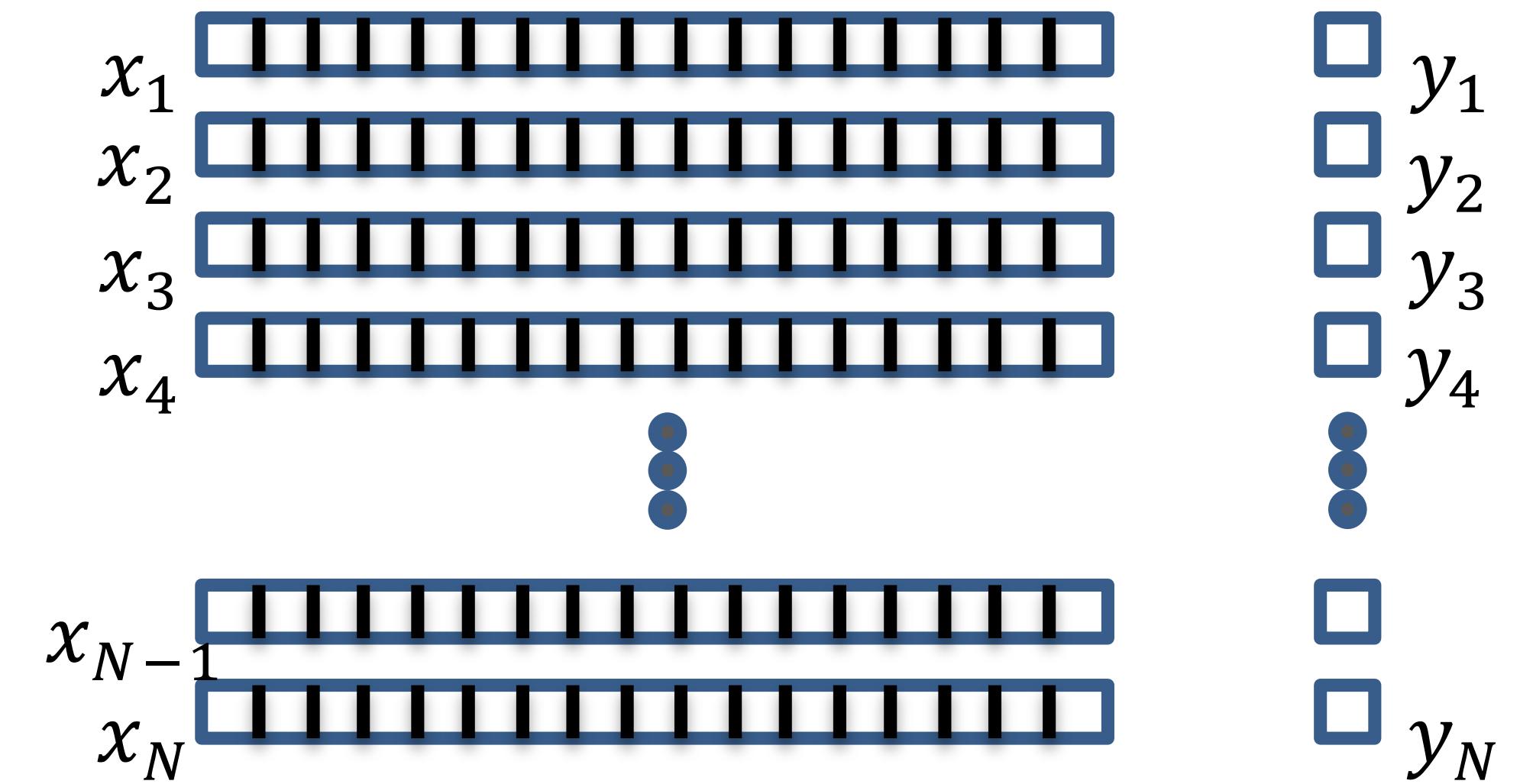
# Classic Example: Increasing Polynomial Order



Increasing complexity visually looks ridiculous...

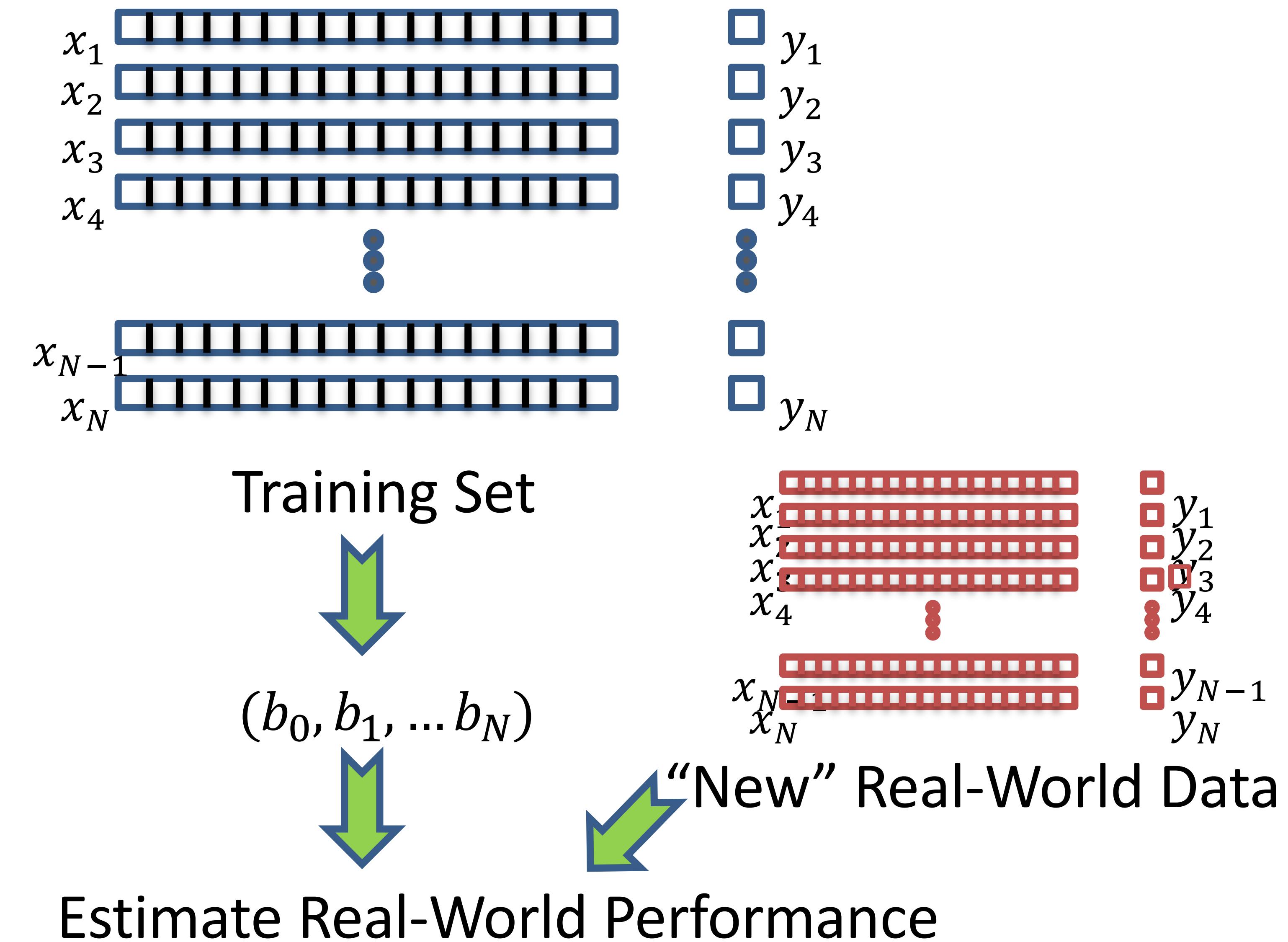
# What happens in overfitting?

- We are increasing the number of parameters in the model, which means:
  - More parameters to estimate (all of their errors add up)
  - Can learn *complex* relationships, maybe too complex for reality
- When we *overfit*, this means that we will not *generalize*
  - We want our models and analysis to generalize, or provide accurate predictions on newly collected data

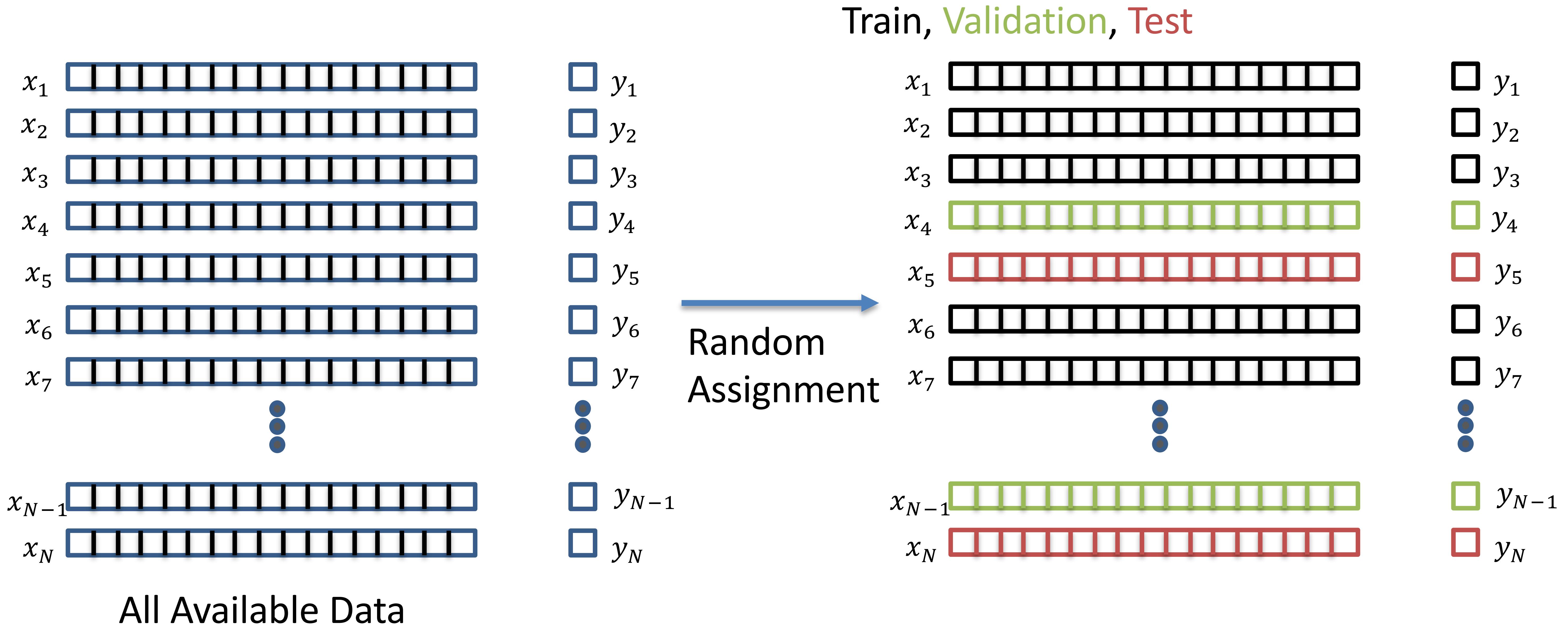


# Standard Validation Strategy

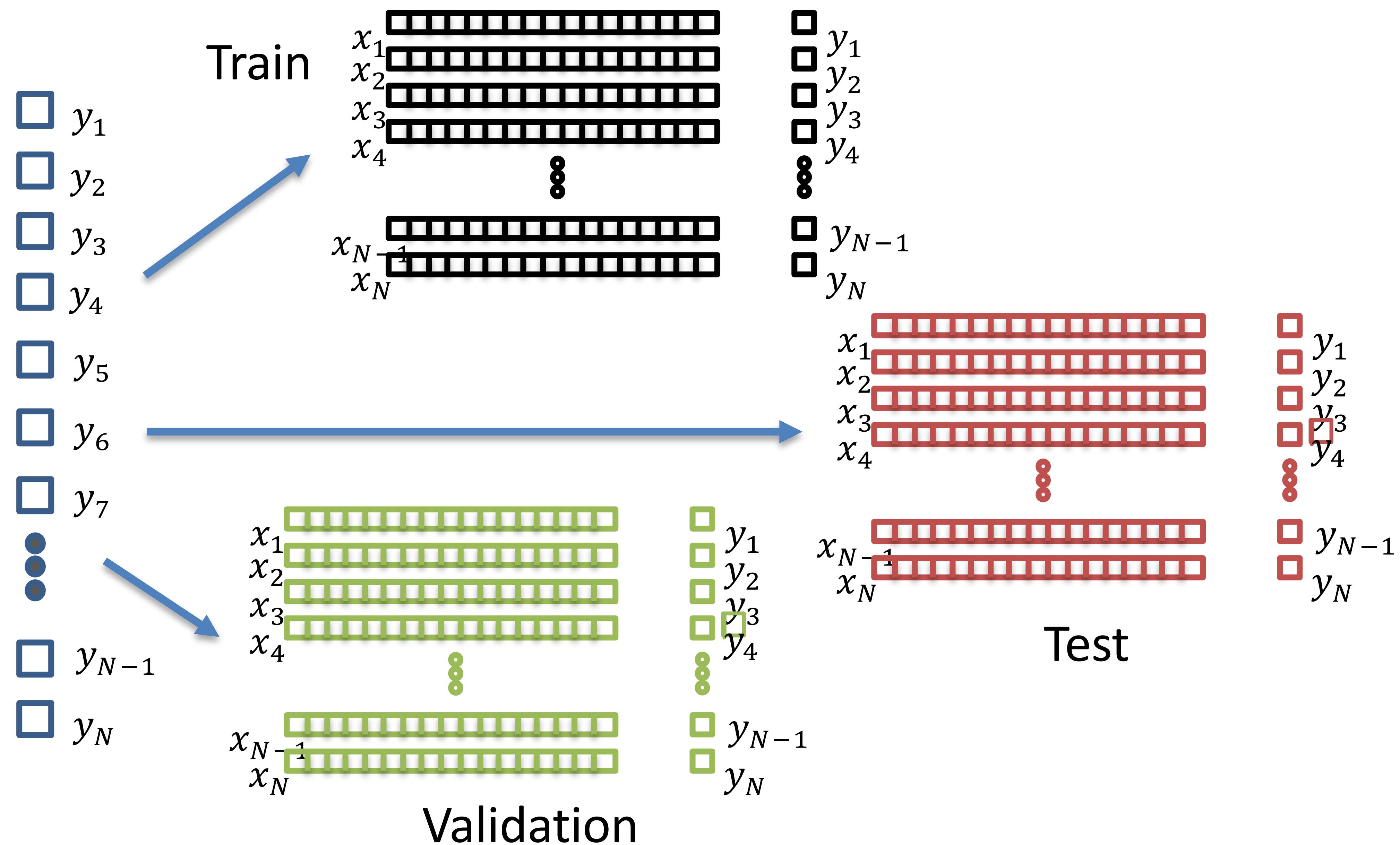
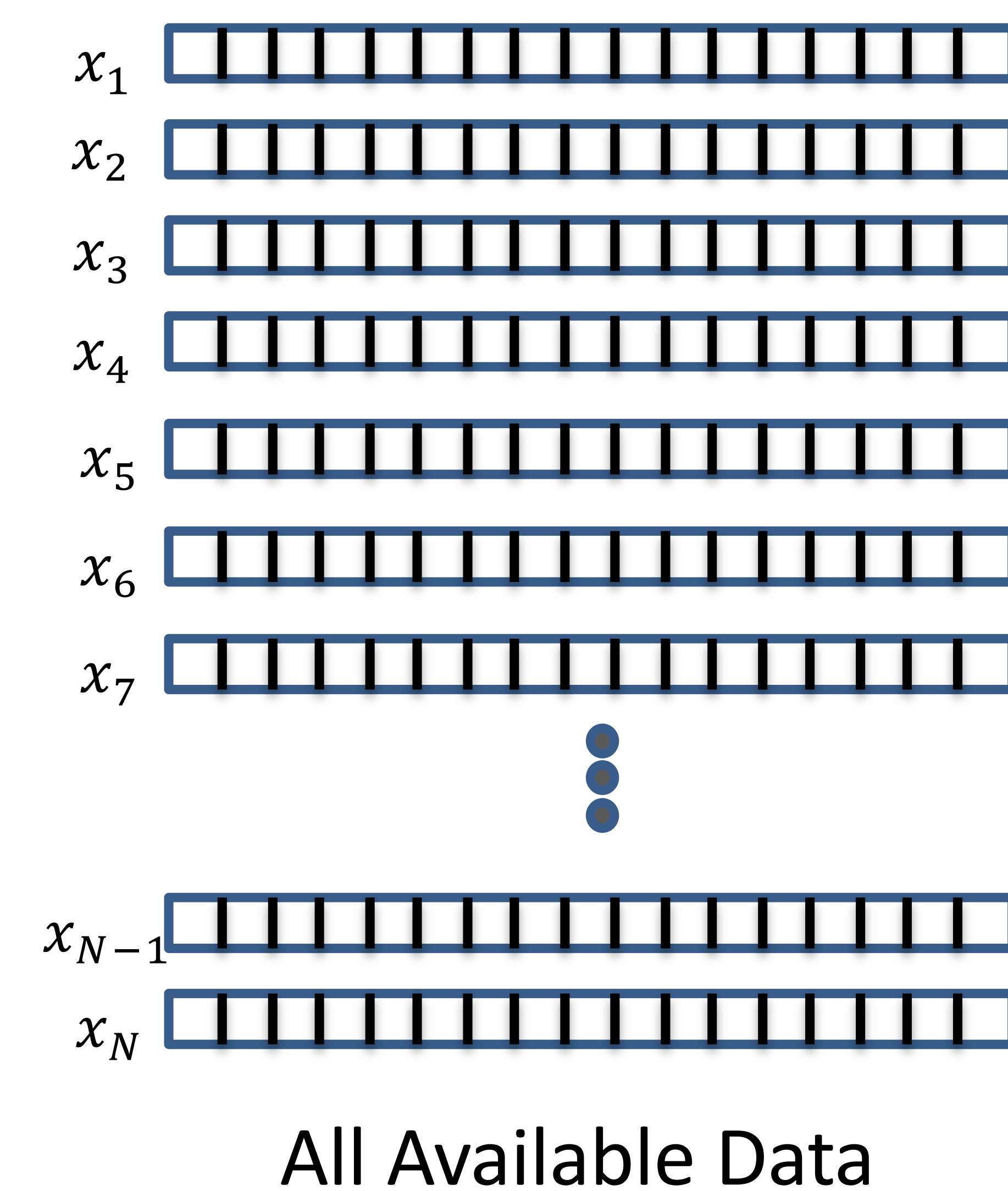
- In the end, we want to know how the network will perform *in the real world*
- Standard approach: actually try it in the real world
- This is costly; instead, can we use existing data to estimate performance?



# Split Data into Separate Groups



# Split Data into Separate Groups

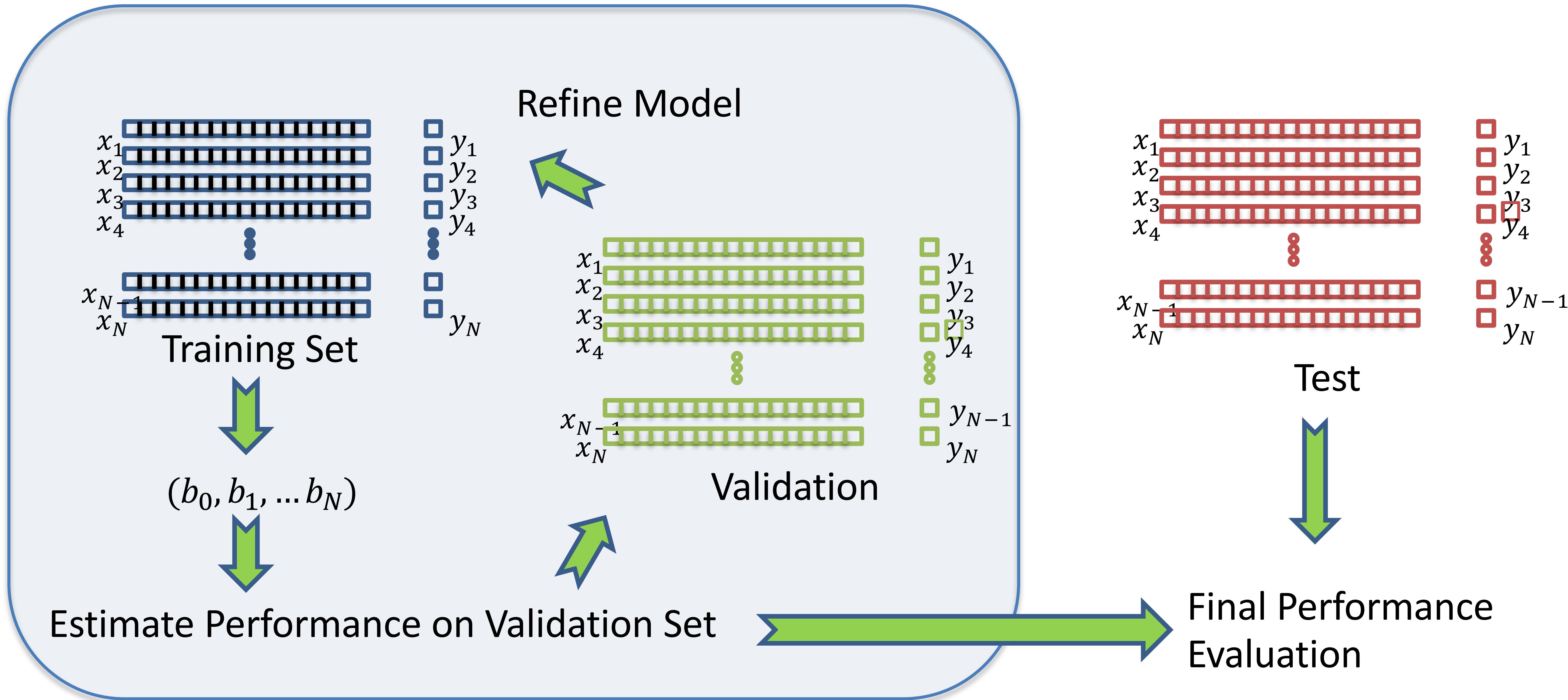


# Test set

- Standard practice to create prior to any analysis—**will not be used to learn or fit any parameters**
- After learning the network, can evaluate its performance on the test set
  - This data was not included in the training/fitting, so it is analogous to running a new synthetic experiment
- Ideally, the test set will be used *once*.
  - Reusing the test set leads to bias; performance estimates will be optimistic

# Validation Set

- Want to be able to compare which approach is best
  - Problematic if we only want to use a test set once
  - Can create a second held-out dataset
- The validation data is not used for learning parameters, but can be used repeatedly to estimate performance of a model
- We can pick the model with the best performance on the validation set, and run a final evaluation on the test data

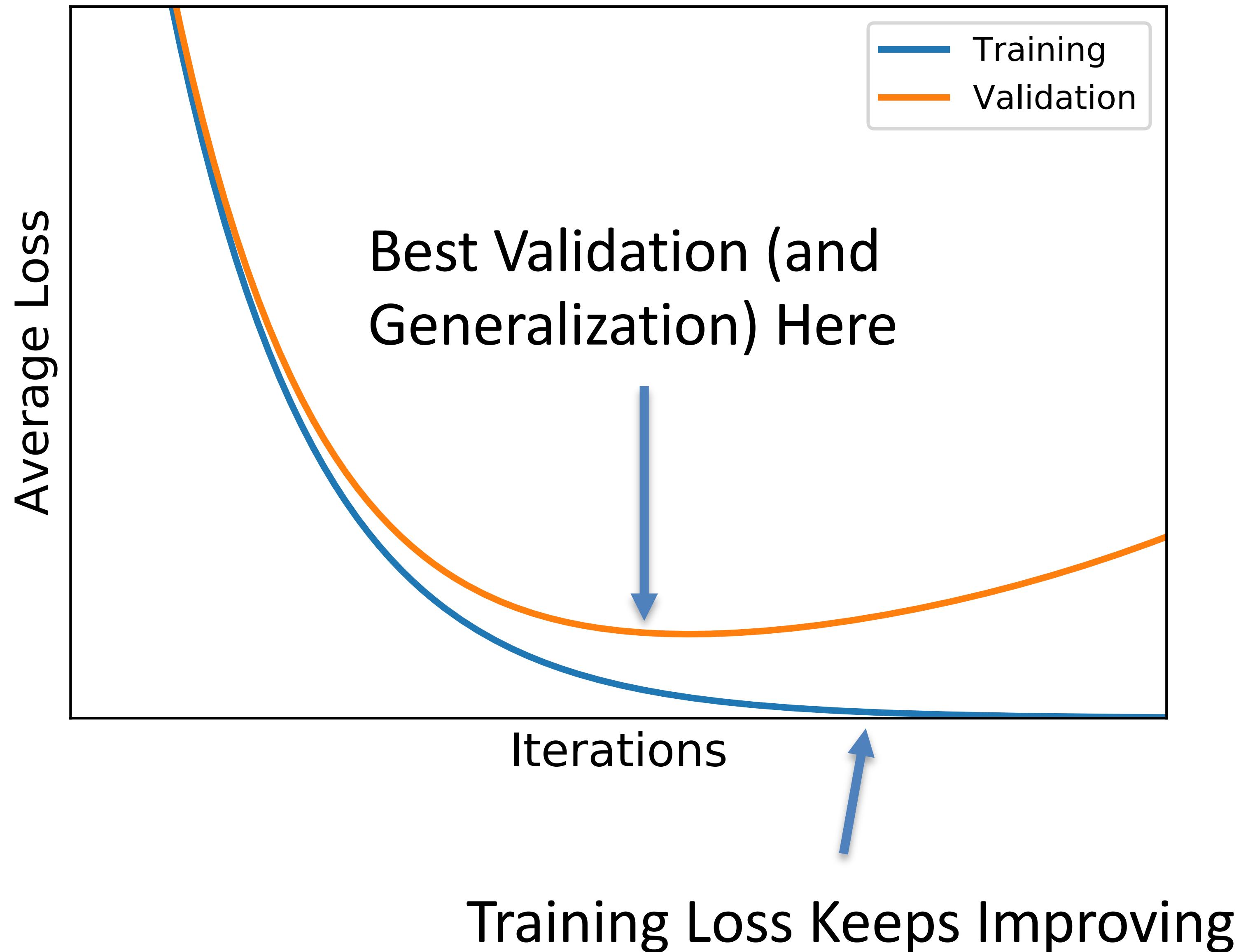


# Validation During Optimization

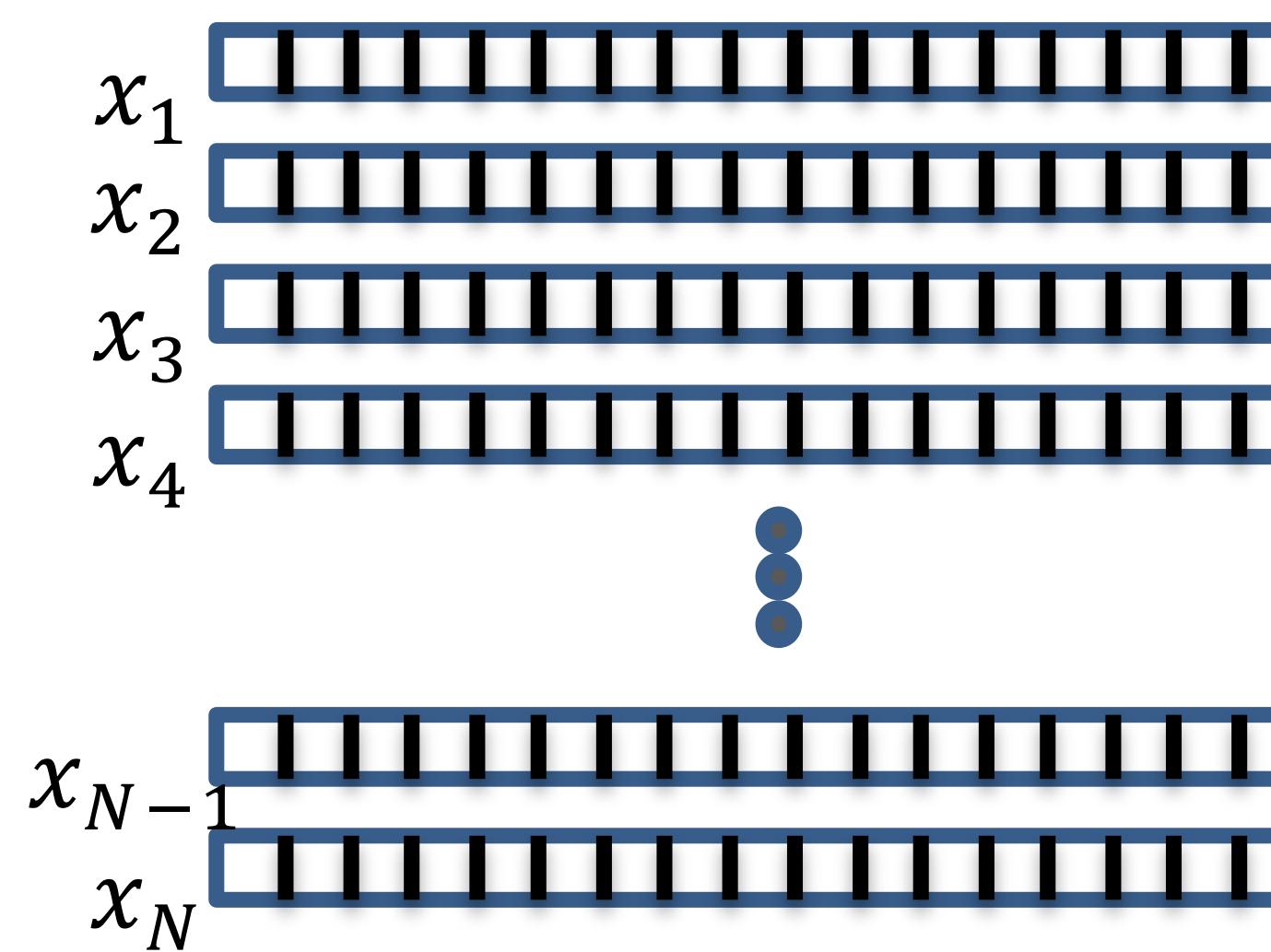
- We want to maximize the *generalization* of the network, not our previous optimization goal
- In practice, can we validate the model while running the optimization loop?

# Early Stopping

- During optimization, we can check the validation loss as we go.
- Instead of optimizing to convergence, we can optimize until the *validation* loss stops improving
  - Saves computational cost
  - Performs better on validation (and test) sets
- Widely used technique in the field

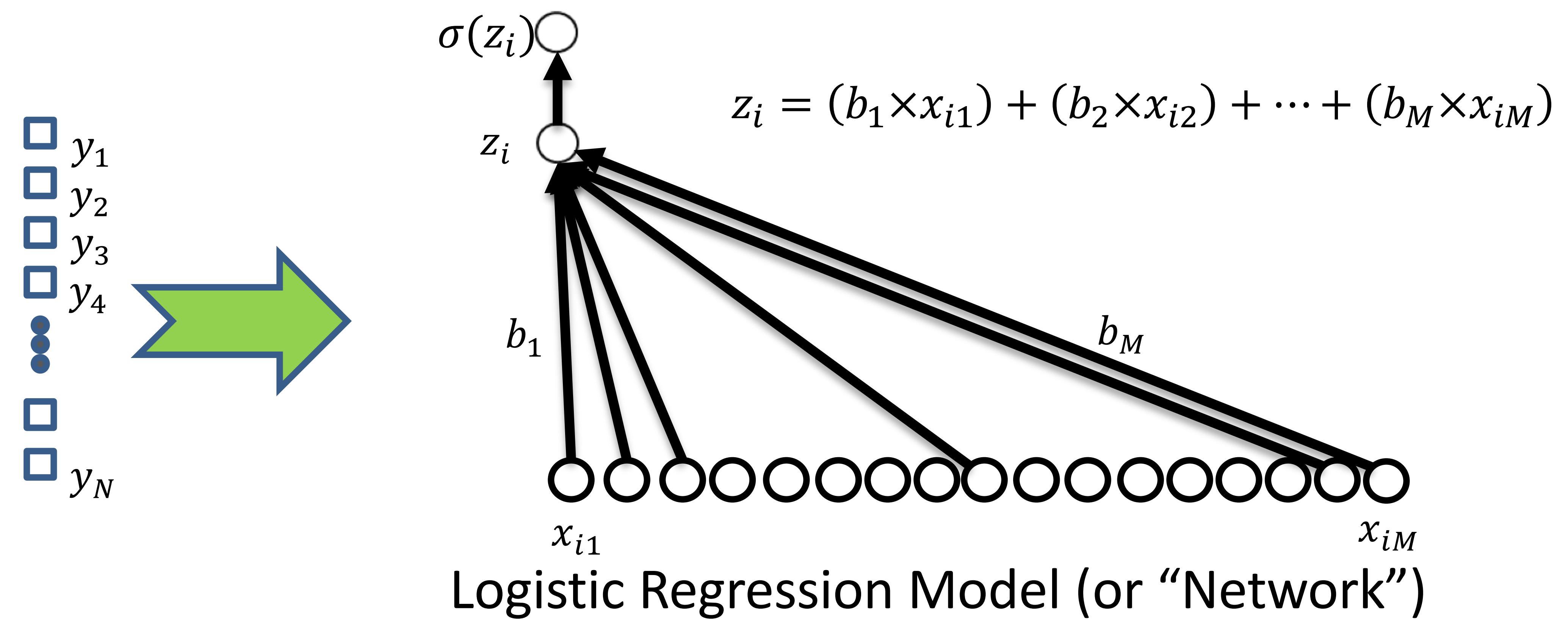


# Learning Model Parameters (Recap)



Training Set

Stochastic Gradient Descent



Learned  
Parameters  $(b_0, b_1, \dots, b_N)$

# Conclusions/Next Steps

- Stochastic Gradient methods are the *modus operandi* of learning deep networks
  - Many extensions such as momentum, preconditioning (will use in hands-on sessions)
- Proper model validation is critical to estimate real-world performance
- Hands on section you will implement stochastic gradient descent and a multi-layer perceptron

Thanks for your attention!

# WRAPPING UP