# CS 540 Introduction to Artificial Intelligence
## Neural Networks (III)

University of Wisconsin-Madison

**Spring 2025**
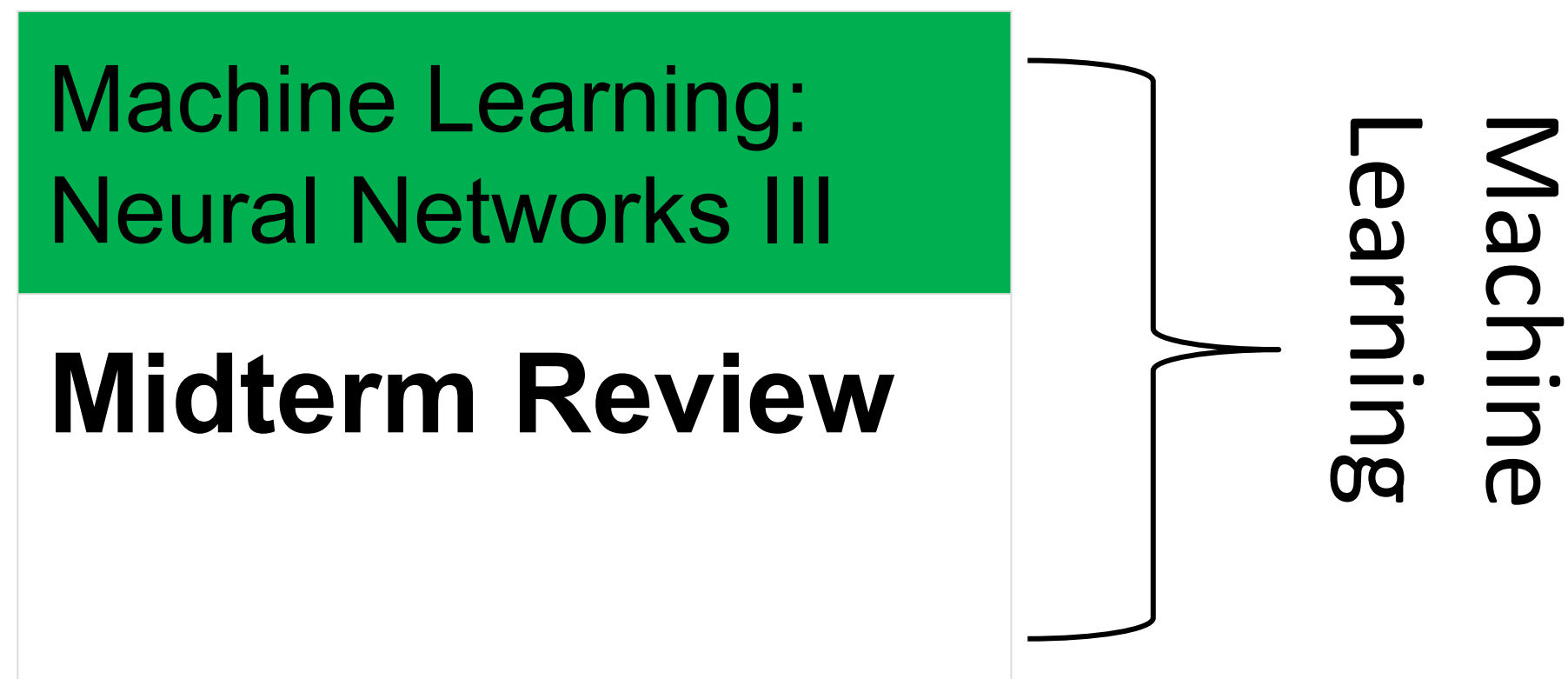
# Announcements

- **Homeworks**:
  - HW6 online, deadline on Monday **March. 17<sup>th</sup> at 11:59 PM**

- Midterm March 13<sup>th</sup> . More on next slide.

- Class roadmap:

| Machine Learning: Neural Networks III |
|---|
| **Midterm Review** |

Machine Learning

# Midterm Information

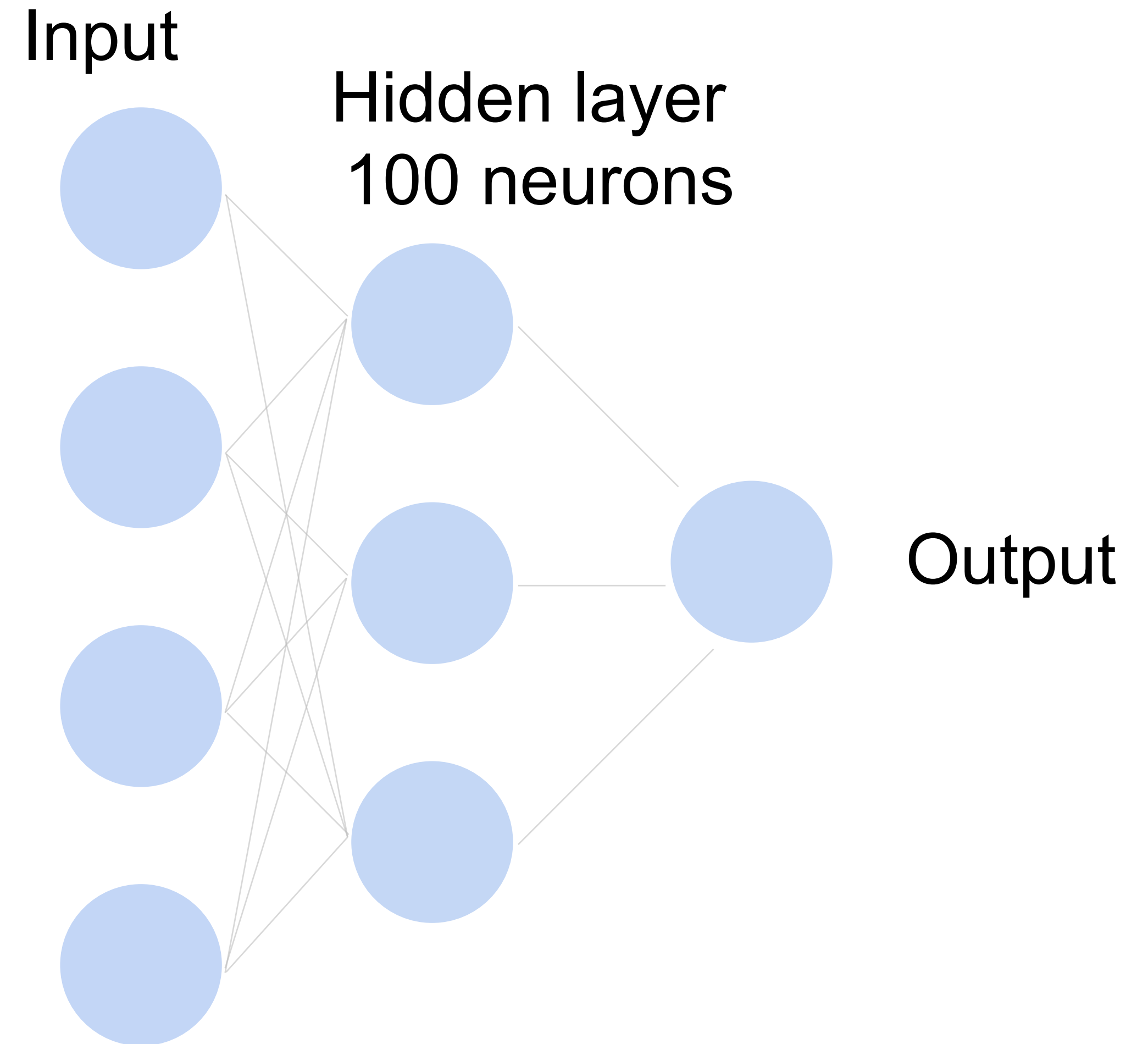- **Time:** March 13th 7:30-9 PM
- **Location:**
  - Section 001 : Ingraham Hall B10
  - Section 002 : Psychology  105
  - **Section 003: split in two locations according to the last name:**
    - **Chamberlin Hall 2103 ( last name starting with A-L)**
    - **Sterling Hall 1310 ( last name starting with M-Z)**
- McBurney students and students requesting alternate: reach out to your instructor if you have not received any email!
- Format: multiple choice
- Cheat sheet: single piece of paper, front and back
- Calculator: fine if it doesn't have an Internet connection
- Detailed topic list + practice on Piazza and Canvas

# How to train a neural network?

Update the weights W to minimize the loss function

$$L = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} \ell(\mathbf{x}, y)$$

**Use gradient descent!**
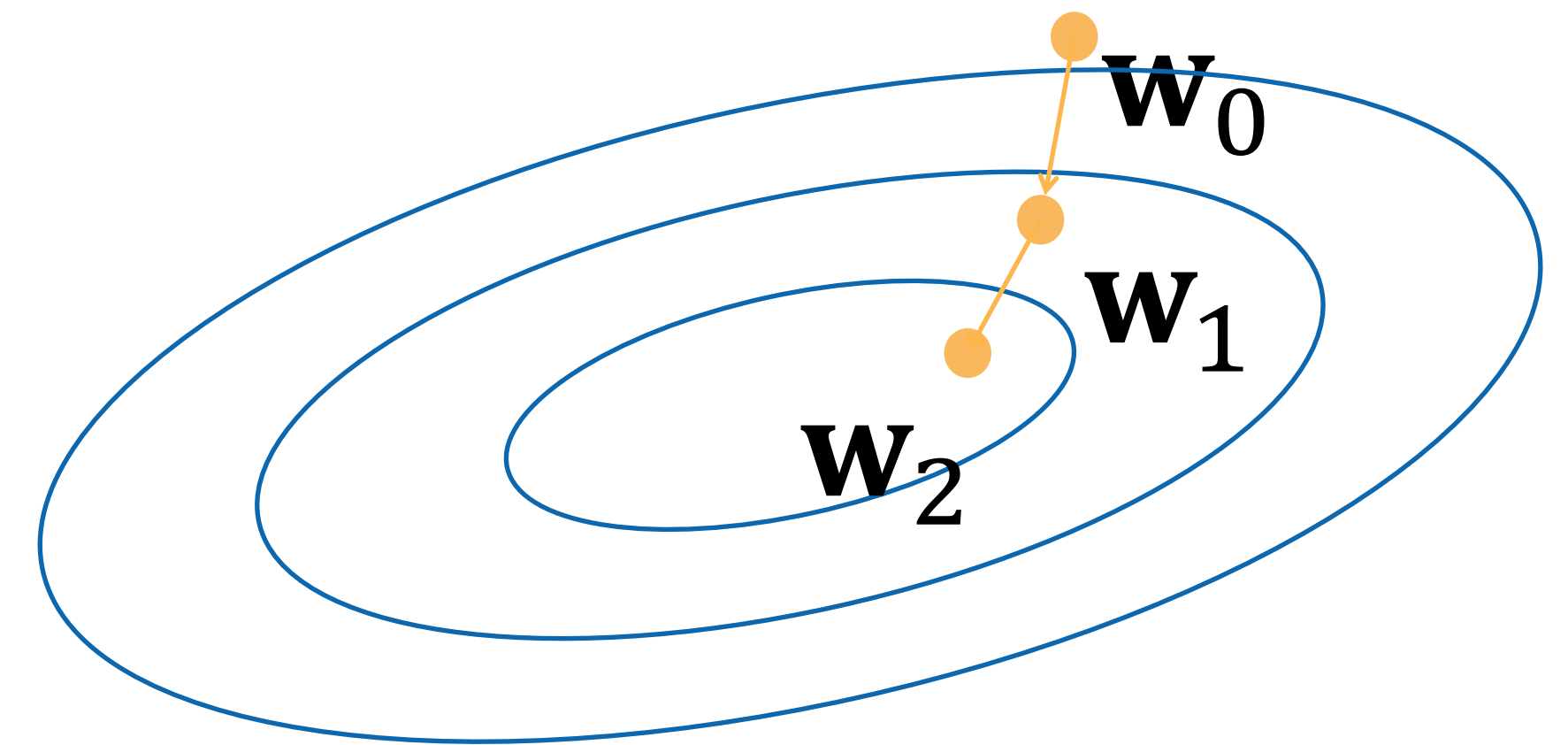
Input

Hidden layer
100 neurons

Output

# Gradient Descent

- Choose a learning rate $\alpha > 0$
- Initialize the model parameters $w_0$
- For t =1,2,…

  - Update parameters:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha \frac{\partial L}{\partial \mathbf{w}_{t-1}}$$

D can be very large. Expensive per iteration

$$= \mathbf{w}_{t-1} - \alpha \frac{1}{|D|} \sum_{(\mathbf{x},y) \in D} \frac{\partial \ell(\mathbf{x},y)}{\partial \mathbf{w}_{t-1}}$$
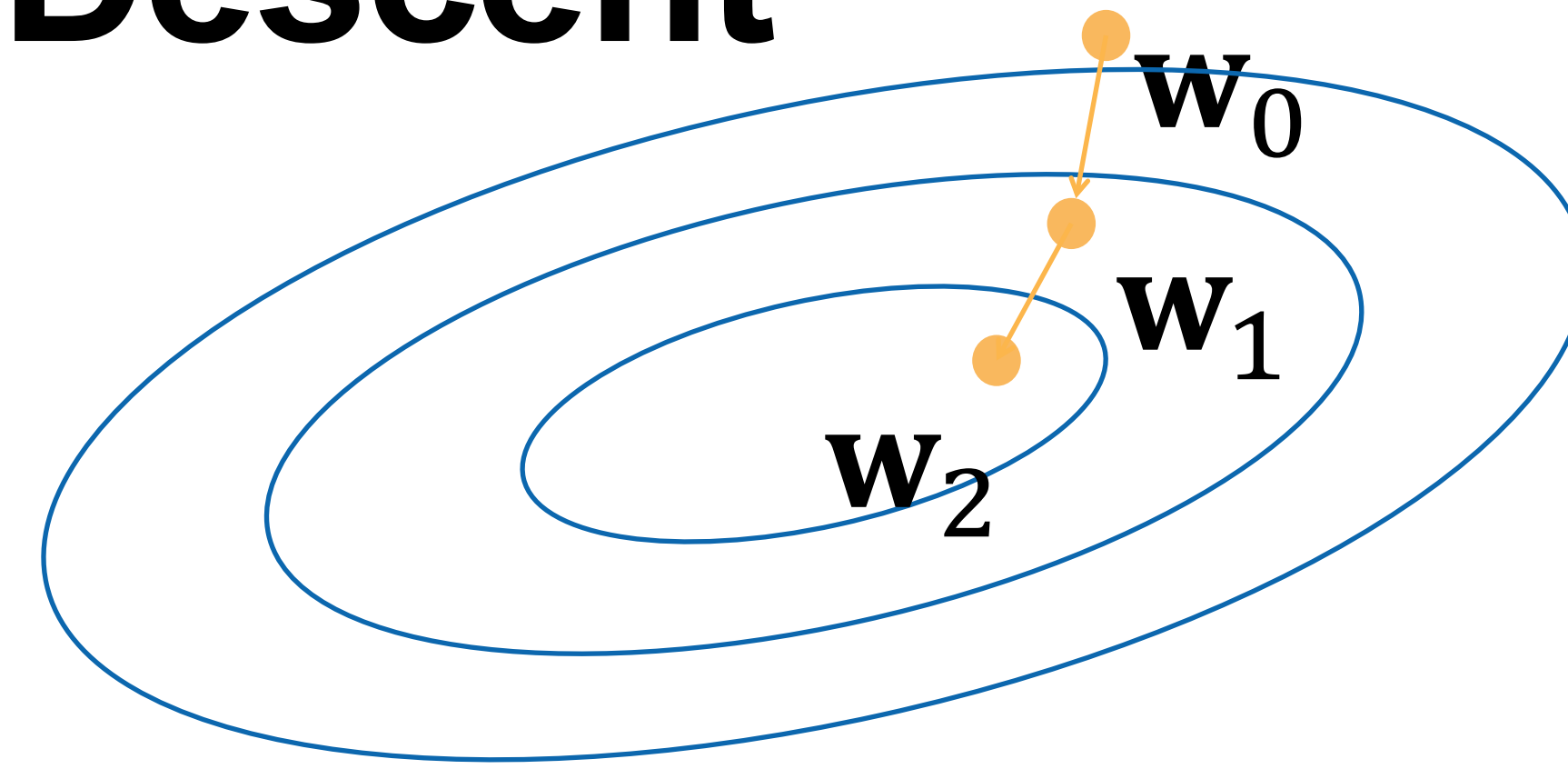
The gradient w.r.t. all parameters is obtained by concatenating the partial derivatives w.r.t. each parameter

- Repeat until converges

$\mathbf{w}_0$

$\mathbf{w}_1$

$\mathbf{w}_2$

# Minibatch Stochastic Gradient Descent



- Choose a learning rate $\alpha > 0$
- Initialize the model parameters $w_0$
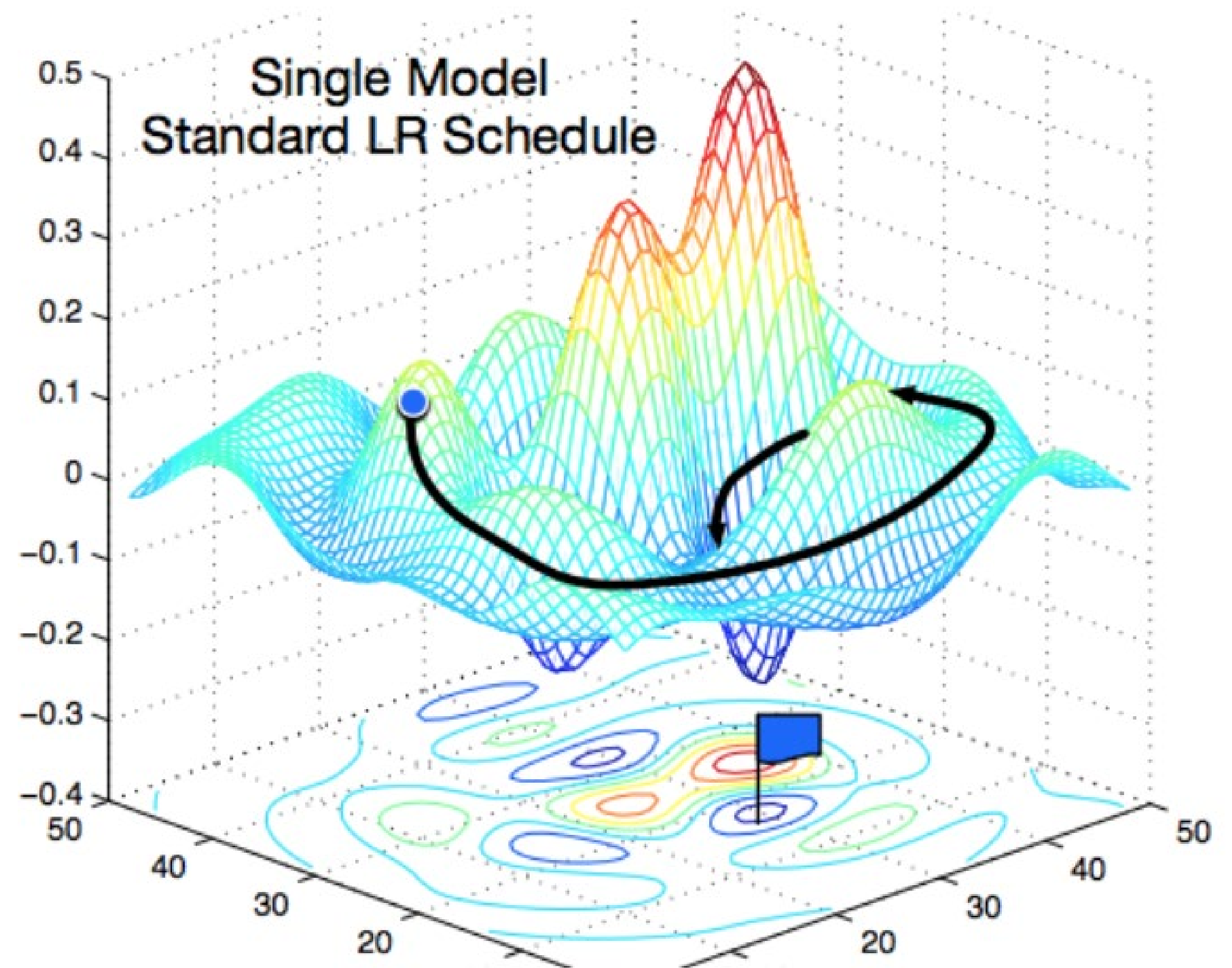- For t =1,2,…
  - **Randomly sample a subset (mini-batch)** $B$ $\subset D$ Update parameters:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha \frac{1}{|B|} \sum_{(\mathbf{x},y) \in B} \frac{\partial \ell(\mathbf{x}, y)}{\partial \mathbf{w}_{t-1}}$$

- Repeat until converges

# Non-convex Optimization



Single Model
Standard LR Schedule

[Gao and Li et al., 2018]

# Quiz Break

- What is the partial derivative $\frac{\partial f}{\partial w_1}$ of: $f(x_1, x_2, w_1, w_2, y) = y\log\sigma(w_1 x_1$

$+ w_2 x_2) + (1 - y)\log(1 - \sigma(w_1 x_1 + w_2 x_2))$ when $y = 1$ and $\sigma(z)$

$= \frac{1}{1+e^{-z}}$. **Hint**: $\frac{\partial\sigma}{\partial z} = \sigma(z)(1 - \sigma(z))$.

# Quiz Break

- What is the partial derivative $\frac{\partial f}{\partial w_1}$ of: $f(x_1, x_2, w_1, w_2, y) = y\log\sigma(w_1 x_1$

$+ w_2 x_2) + (1 - y)\log(1 - \sigma(w_1 x_1 + w_2 x_2))$ when $y = 1$ and $\sigma(z)$

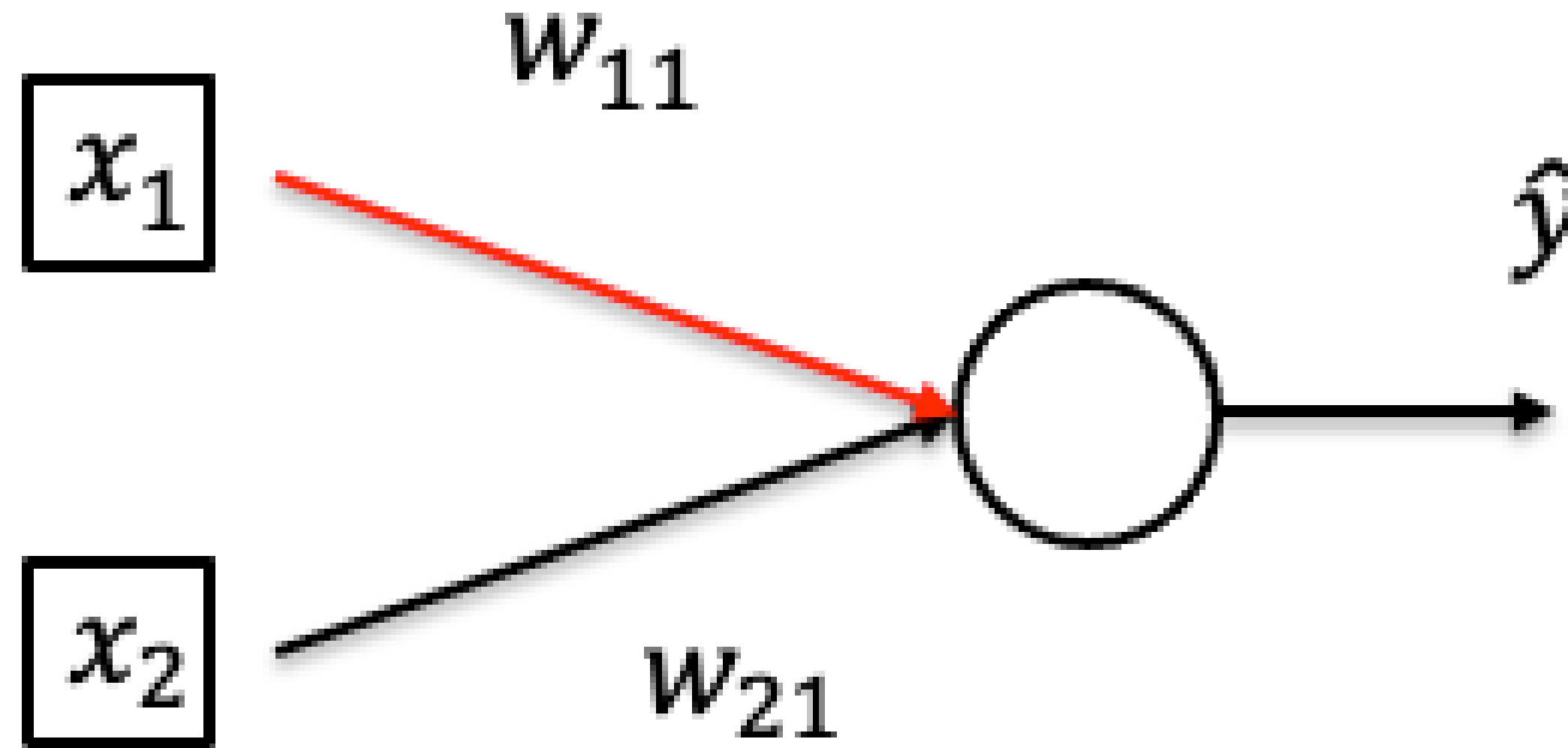$= \frac{1}{1 + e^{-z}}$. **Hint**: $\frac{\partial \sigma}{\partial z} = \sigma(z)(1 - \sigma(z))$.

Let $a = \sigma(z)$      Let $z = w_1 x_1 + w_2 x_2$      $\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial a}\frac{\partial a}{\partial z}\frac{\partial z}{w_1}$

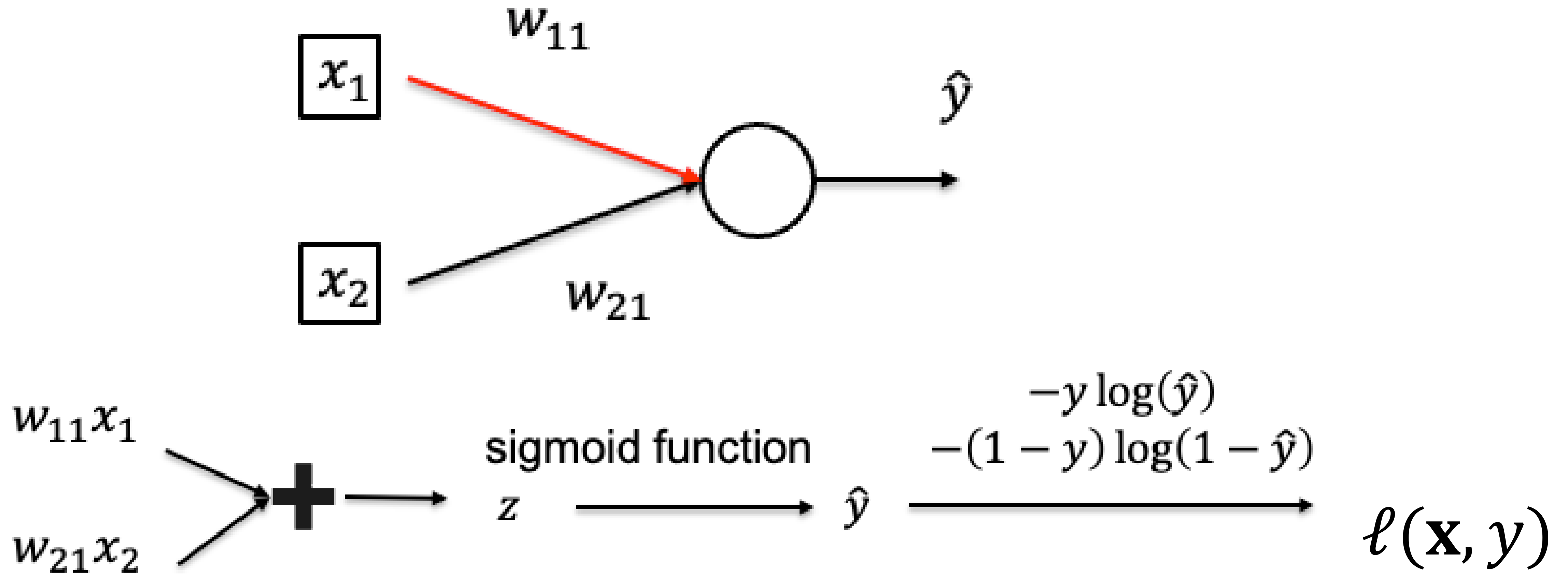$$\frac{\partial f}{\partial w_1} = \frac{y}{a}\sigma(z)(1 - \sigma(z))x_1 = (1 - \sigma(w_1 x_1 + w_2 x_2))x_1$$
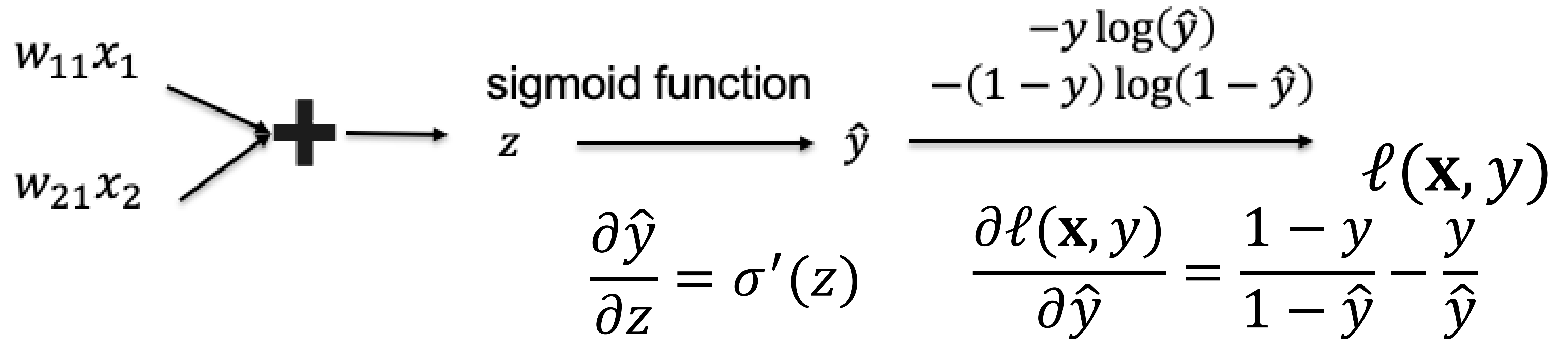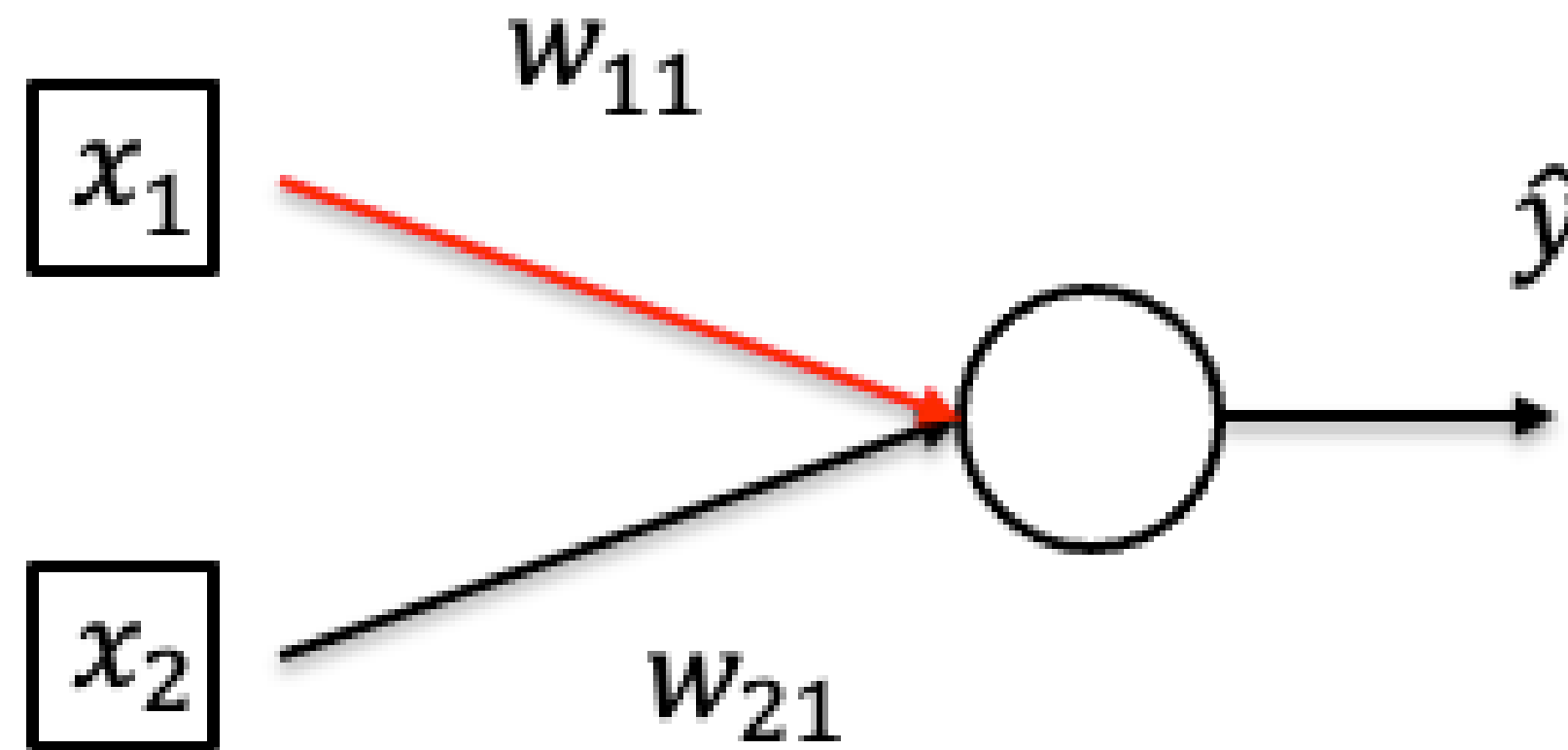
# Calculate Gradient (on one data point)



- Want to compute $\dfrac{\partial \ell(\mathbf{x}, y)}{\partial w_{11}}$

- Data point: $((x_1, x_2), y)$

# Calculate Gradient (on one data point)



$$w_{11}$$
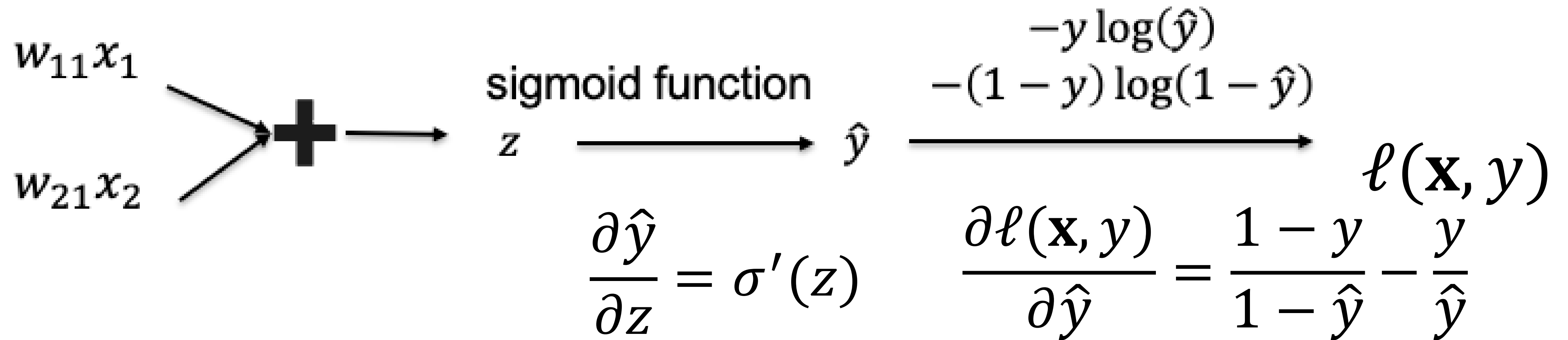
$$x_1 \xrightarrow{w_{11}} \hat{y}$$

$$x_2 \quad w_{21}$$

$$w_{11}x_1$$

$$w_{21}x_2$$

$$+ \longrightarrow \text{sigmoid function} \quad z \longrightarrow \hat{y} \quad \begin{array}{c} -y\log(\hat{y}) \\ -(1-y)\log(1-\hat{y}) \end{array} \longrightarrow \ell(\mathbf{x}, y)$$

Use chain rule!

# Calculate Gradient (on one data point)



$$w_{11}x_1$$
$$w_{21}x_2$$

$+$

sigmoid function

$z \longrightarrow \hat{y}$

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

$\ell(\mathbf{x}, y)$

$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z)$$

$$\frac{\partial \ell(\mathbf{x}, y)}{\partial \hat{y}} = \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}}$$
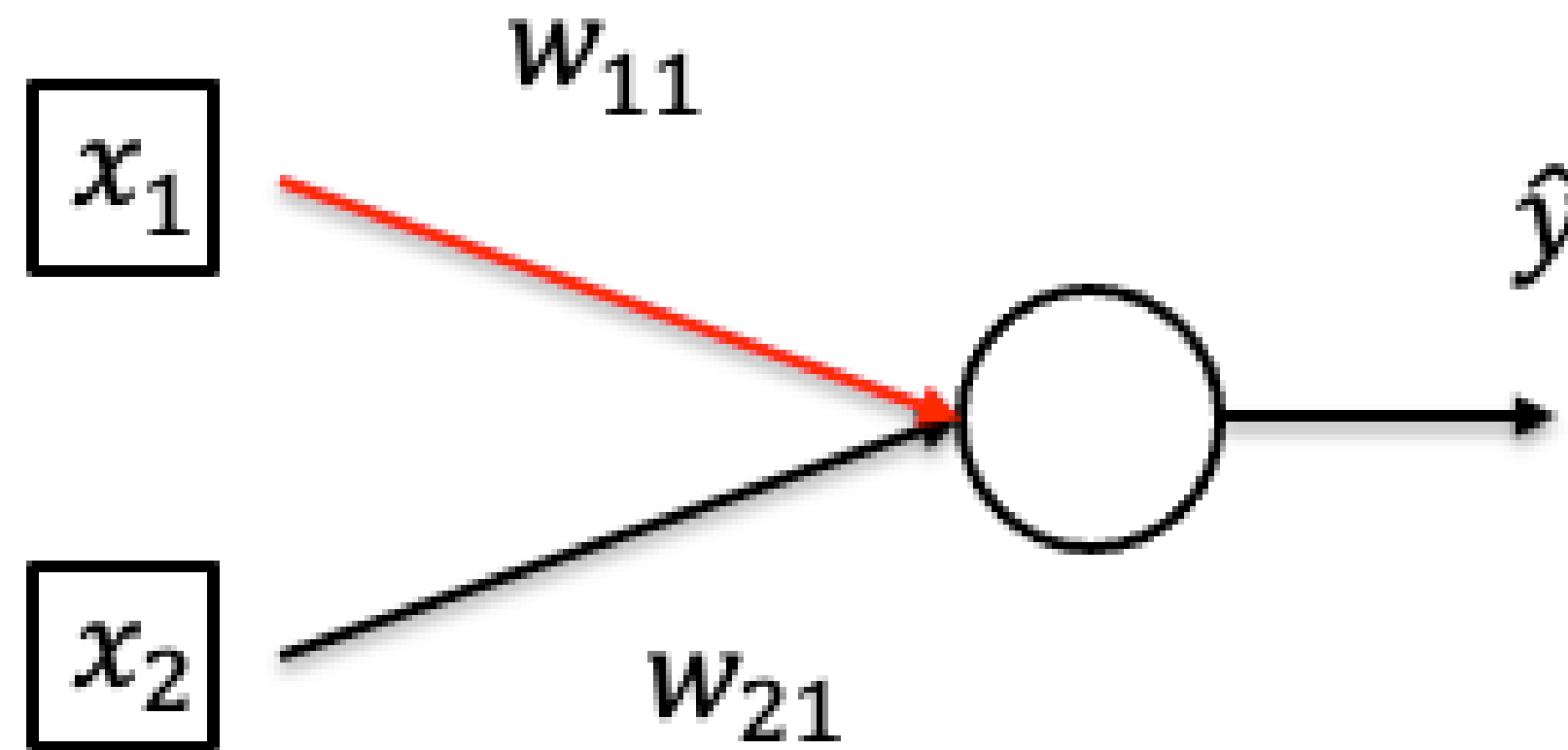
- By chain rule:

$$\frac{\partial l}{\partial w_{11}} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{11}}$$
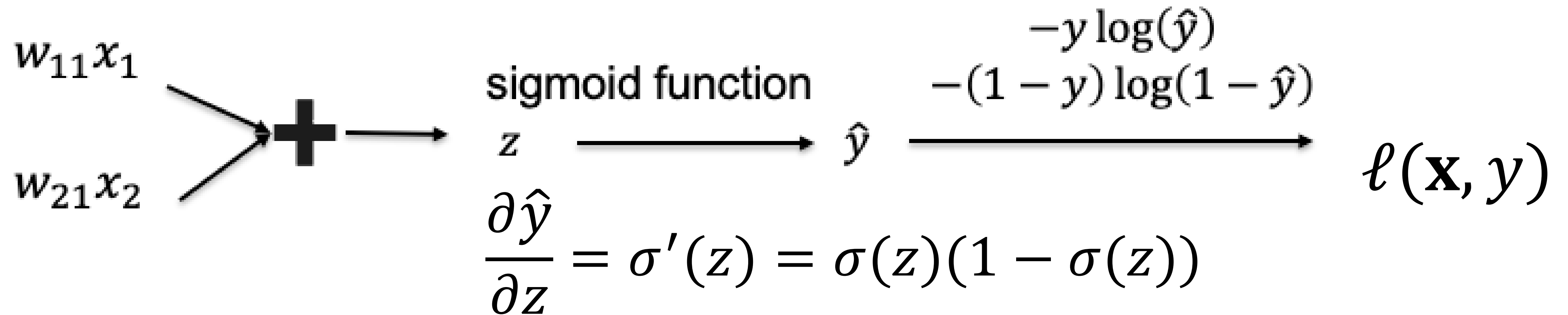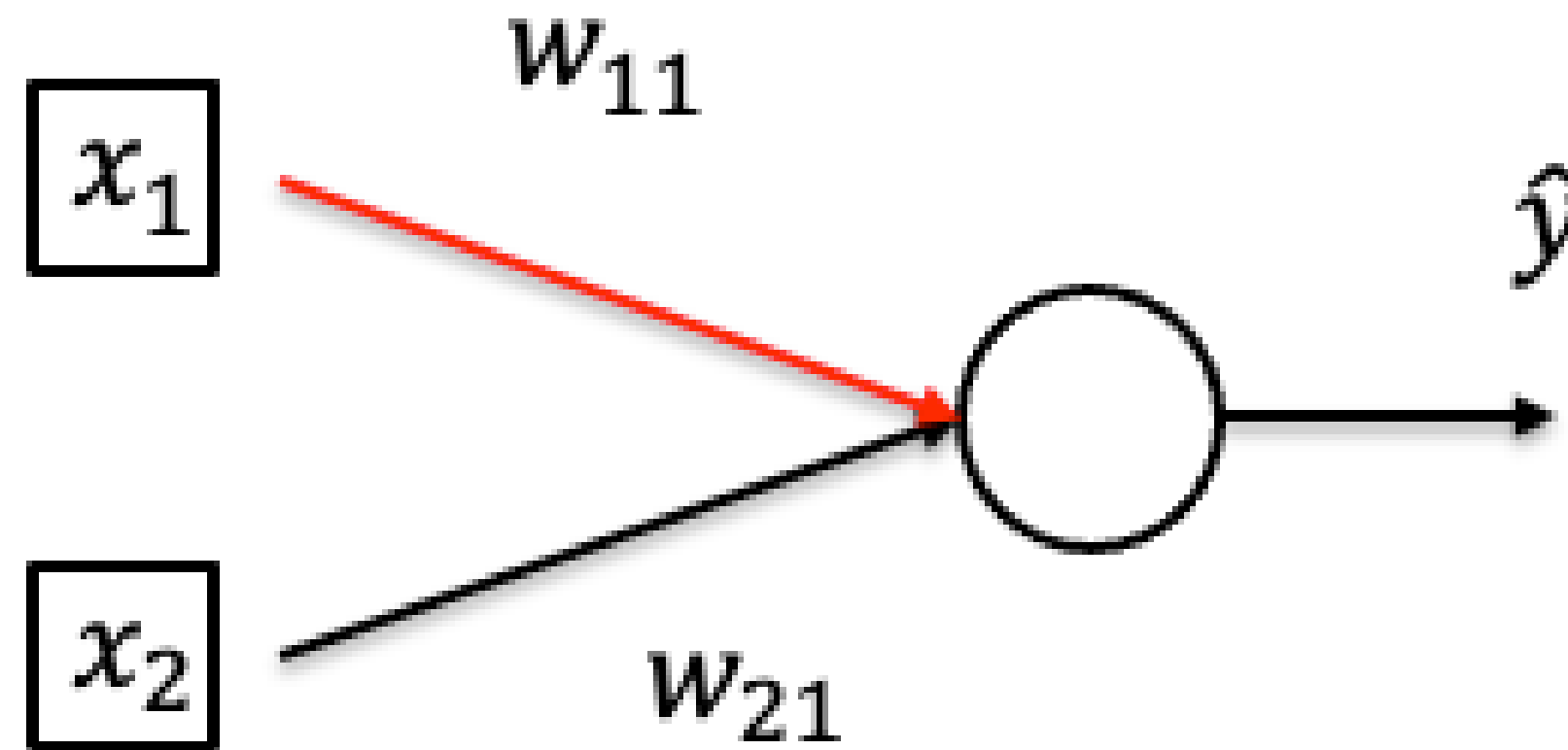
# Calculate Gradient (on one data point)

$x_1$ $\xrightarrow{w_{11}}$ $\hat{y}$

$x_2$ $w_{21}$

$w_{11}x_1$

$w_{21}x_2$ $\mathbf{+}$ $\xrightarrow{}$ $z$ $\xrightarrow{\text{sigmoid function}}$ $\hat{y}$ $\xrightarrow{\substack{-y\log(\hat{y}) \\ -(1-y)\log(1-\hat{y})}}$ $\ell(\mathbf{x}, y)$

$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) \qquad \frac{\partial \ell(\mathbf{x}, y)}{\partial \hat{y}} = \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}}$$
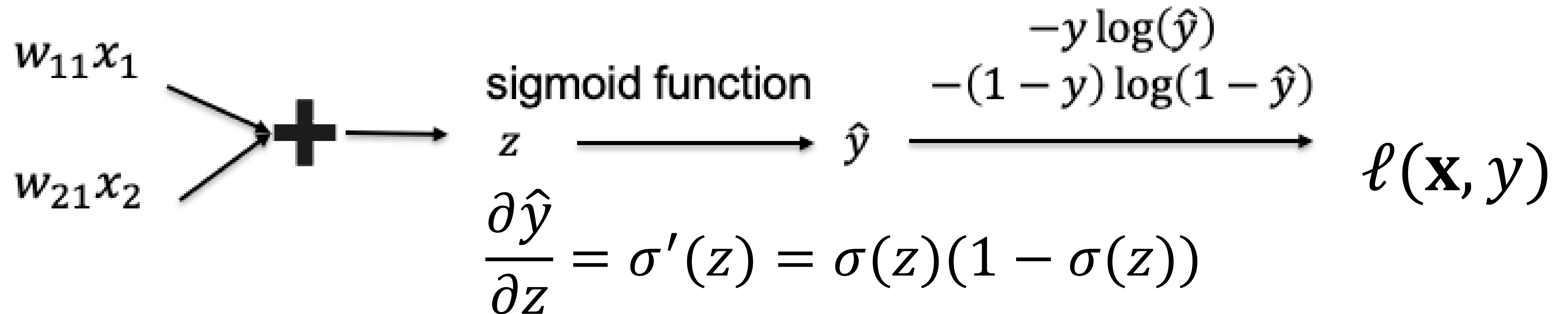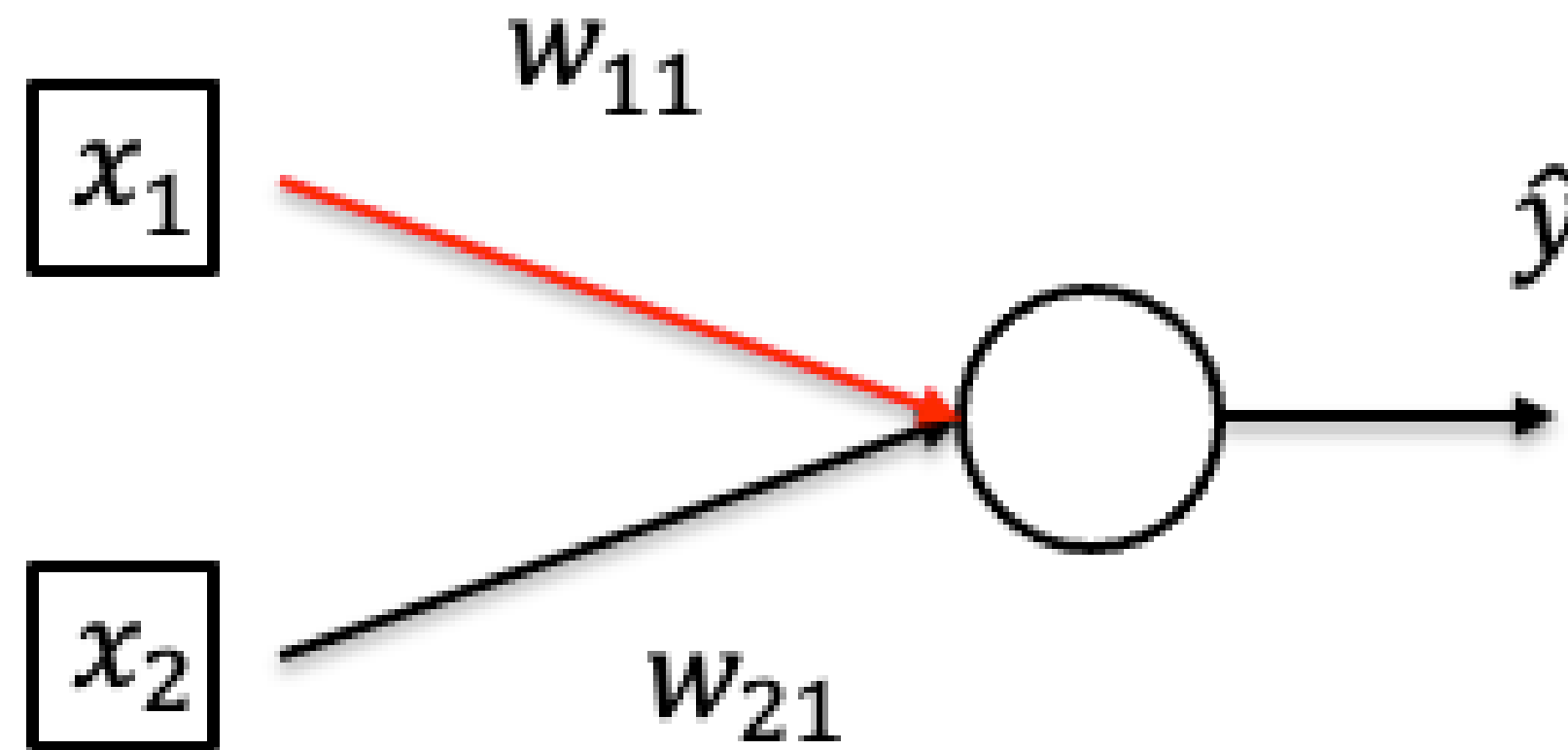
- By chain rule:

$$\frac{\partial l}{\partial w_{11}} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} x_1$$
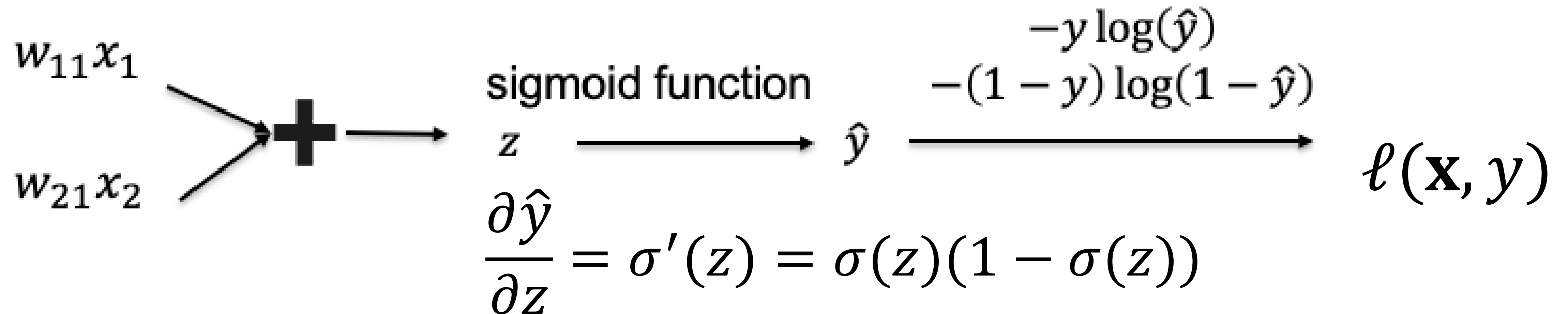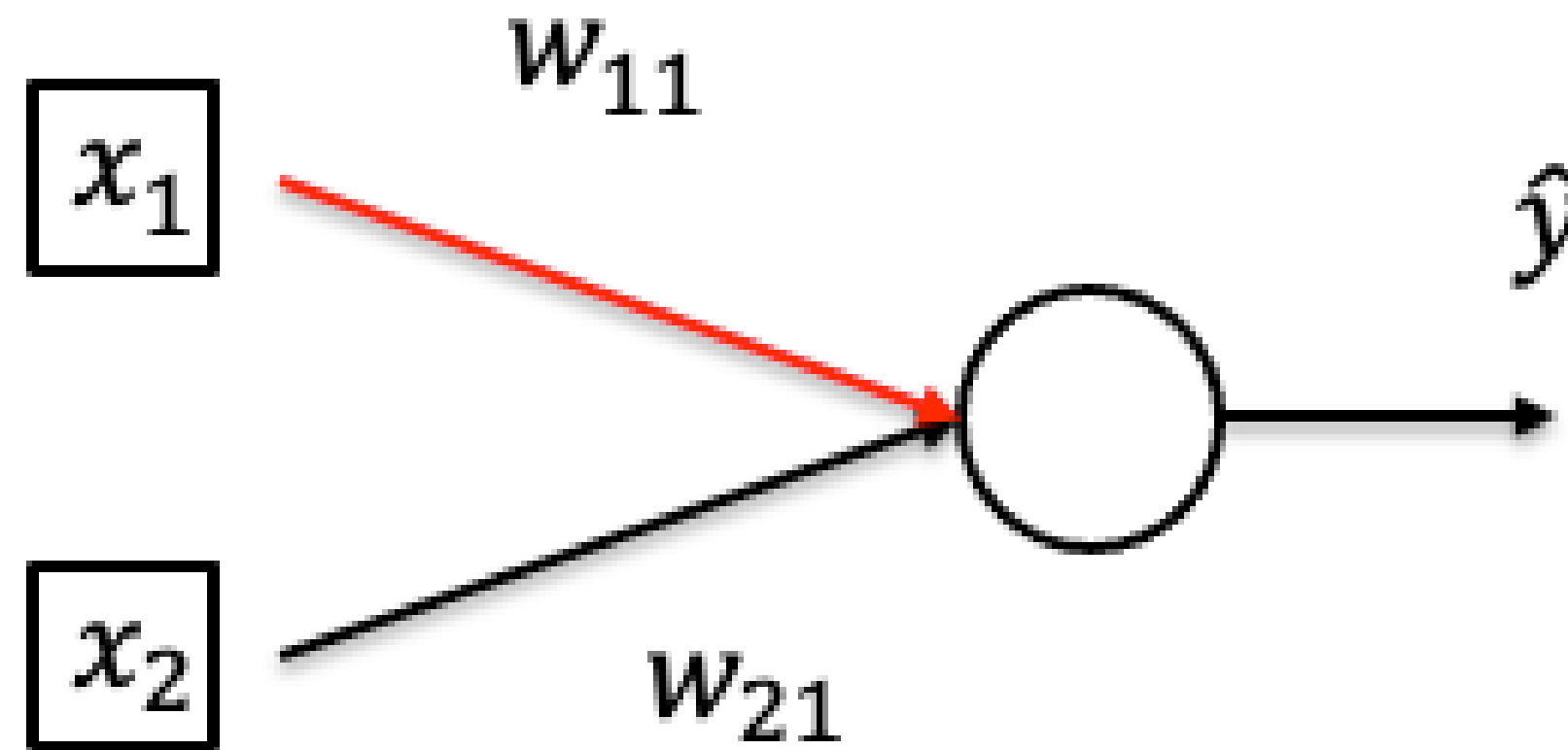
# Calculate Gradient (on one data point)



$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

- By chain rule:

$$\frac{\partial l}{\partial w_{11}} = \frac{\partial l}{\partial \hat{y}} \; \hat{y}(1 - \hat{y})x_1$$

# Calculate Gradient (on one data point)



$w_{11}$

$x_1$

$x_2$

$w_{21}$

$\hat{y}$

$w_{11}x_1$

$w_{21}x_2$

$+$

sigmoid function

$z \longrightarrow \hat{y}$

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

$\ell(\mathbf{x}, y)$

$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

- By chain rule: $\qquad \dfrac{\partial l}{\partial w_{11}} = \left(\dfrac{1-y}{1-\hat{y}} - \dfrac{y}{\hat{y}}\right)\hat{y}(1-\hat{y})x_1$
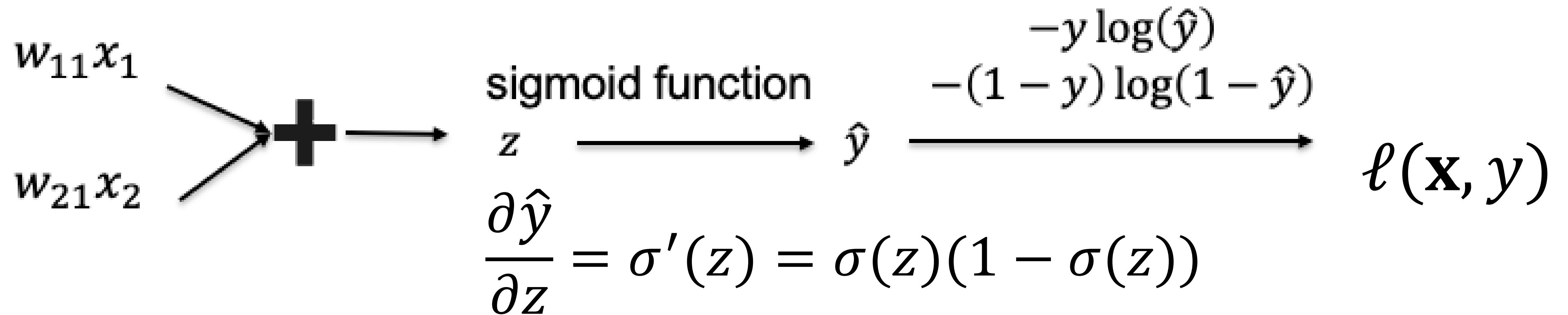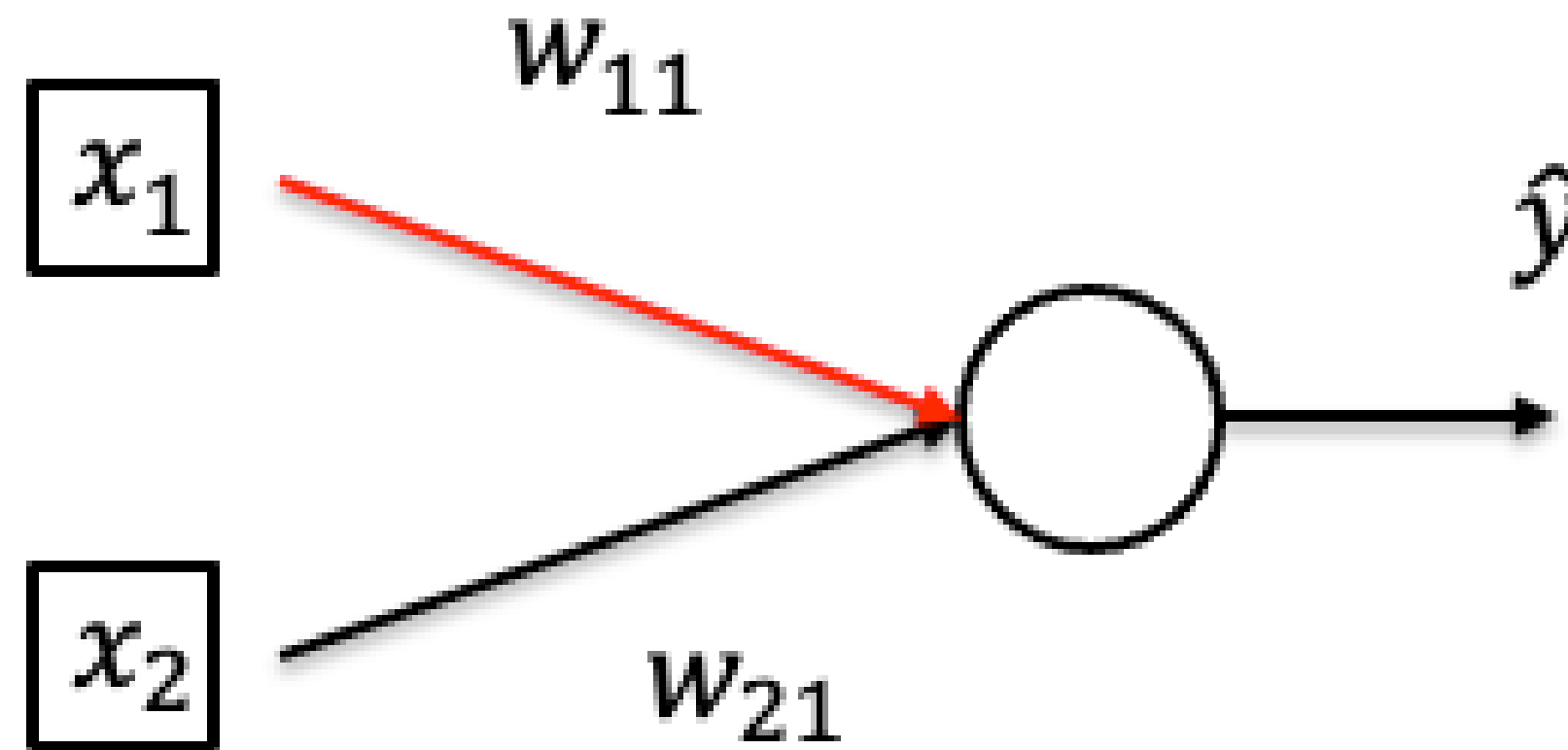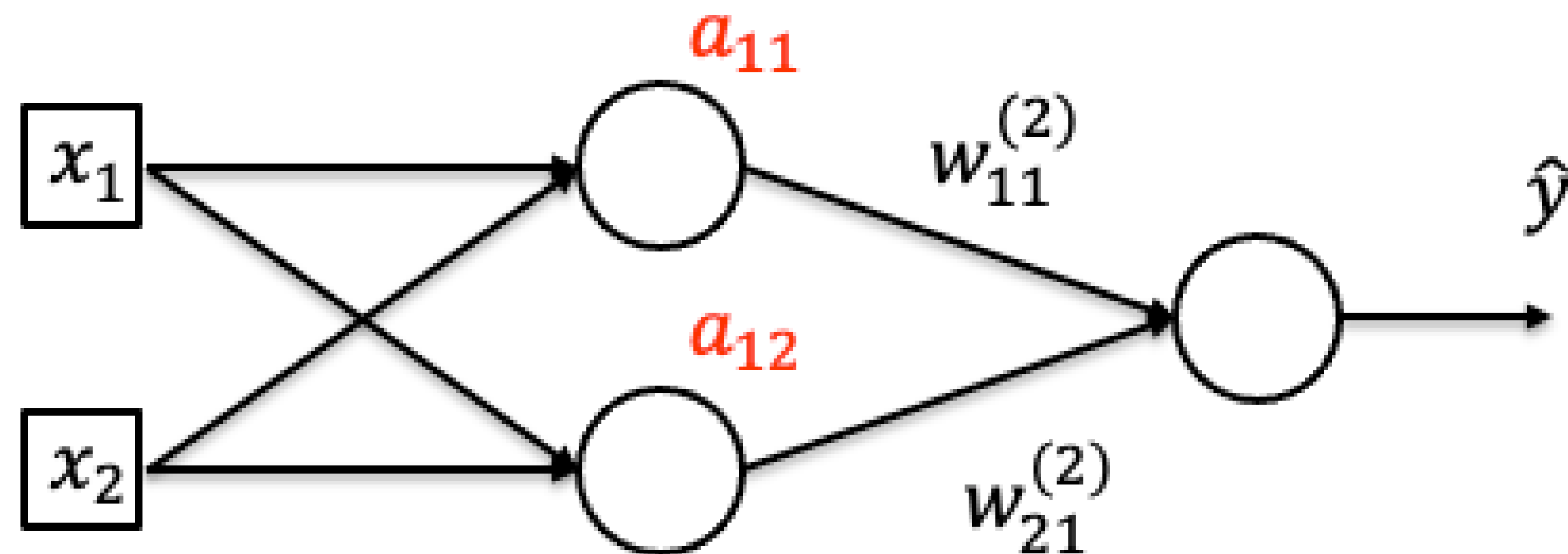
# Calculate Gradient (on one data point)



$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

- By chain rule:   $\dfrac{\partial l}{\partial w_{11}} = (\hat{y} - y)x_1$

# Calculate Gradient (on one data point)



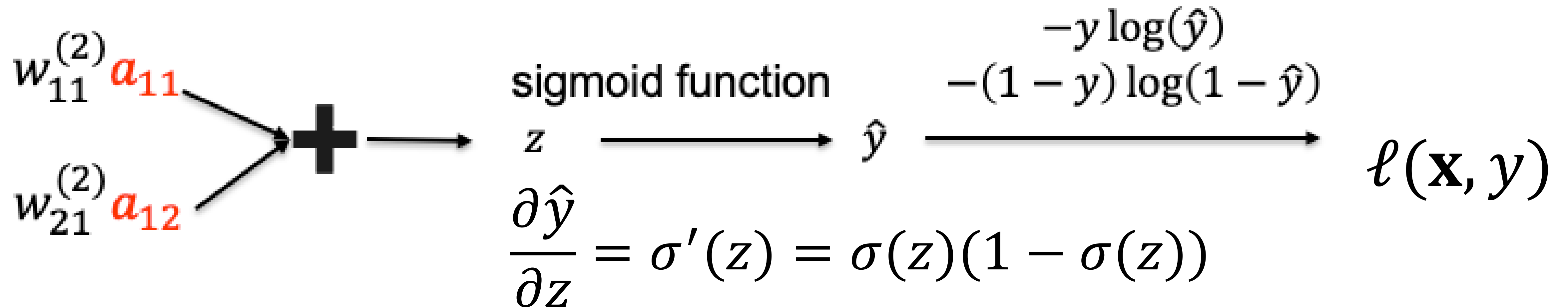$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

- By chain rule: $\dfrac{\partial l}{\partial x_1} = \dfrac{\partial l}{\partial \hat{y}} \dfrac{\partial \hat{y}}{\partial z} w_{11} = (\hat{y} - y) w_{11}$
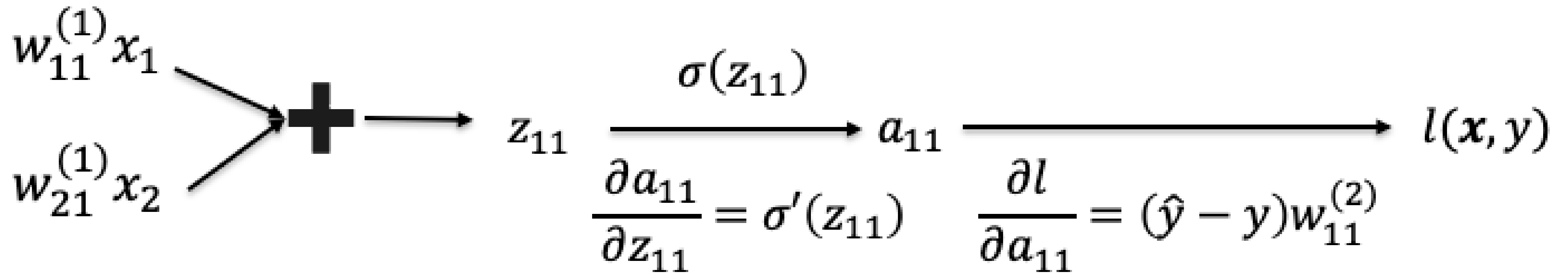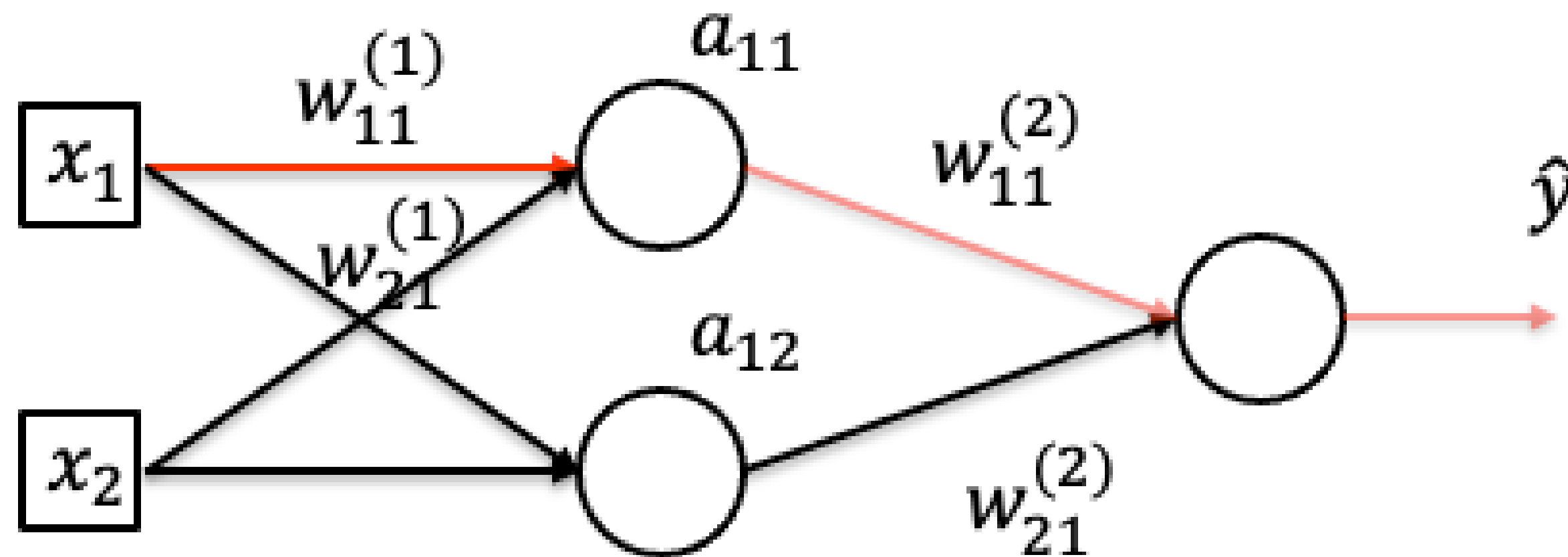
# Calculate Gradient (on one data point)



$w_{11}^{(2)} a_{11}$
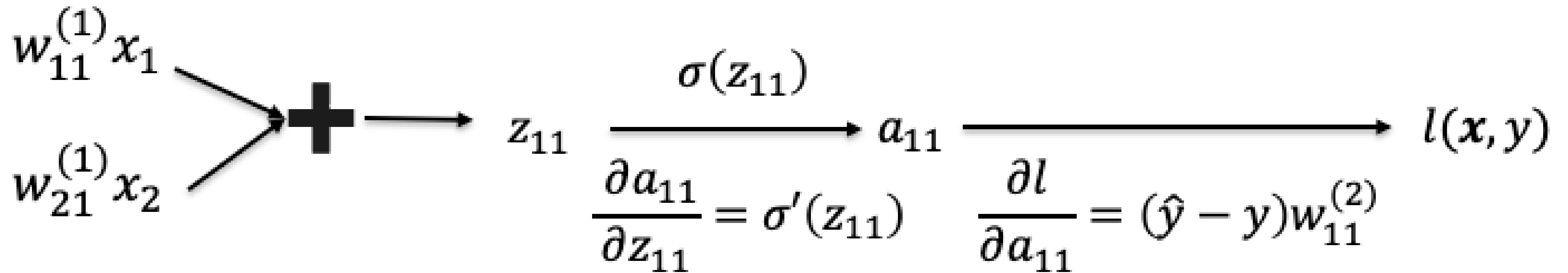
$w_{21}^{(2)} a_{12}$

$\mathbf{+}$

sigmoid function

$z \longrightarrow \hat{y}$

$\dfrac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$

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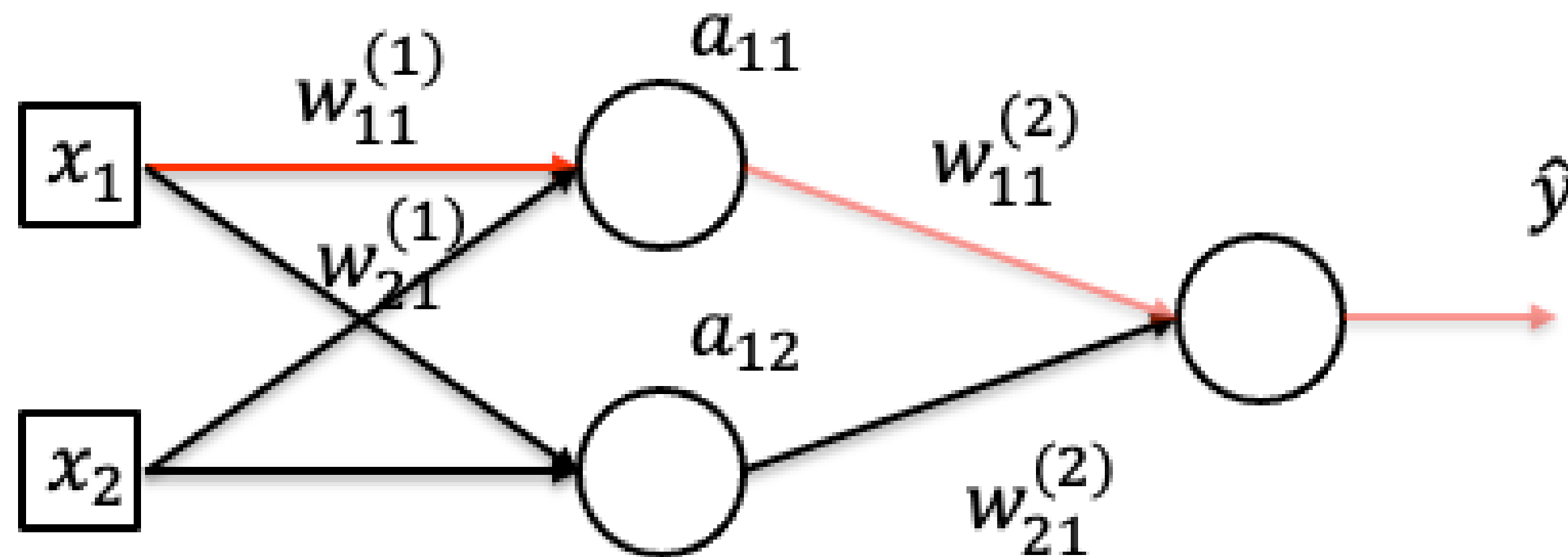$-(1-y)\log(1-\hat{y})$

$\ell(\mathbf{x}, y)$

Make it deeper

- By chain rule: $\quad \dfrac{\partial l}{\partial a_{11}} = (\hat{y} - y)w_{11}^{(2)}, \quad \dfrac{\partial l}{\partial a_{12}} = (\hat{y} - y)w_{21}^{(2)}$

# Calculate Gradient (on one data point)



$$\frac{\partial l}{\partial w_{11}^{(1)}} = \frac{\partial l}{\partial a_{11}}\frac{\partial a_{11}}{\partial w_{11}^{(1)}} = (\hat{y} - y)w_{11}^{(2)}\frac{\partial a_{11}}{\partial w_{11}^{(1)}}$$

- By chain rule:
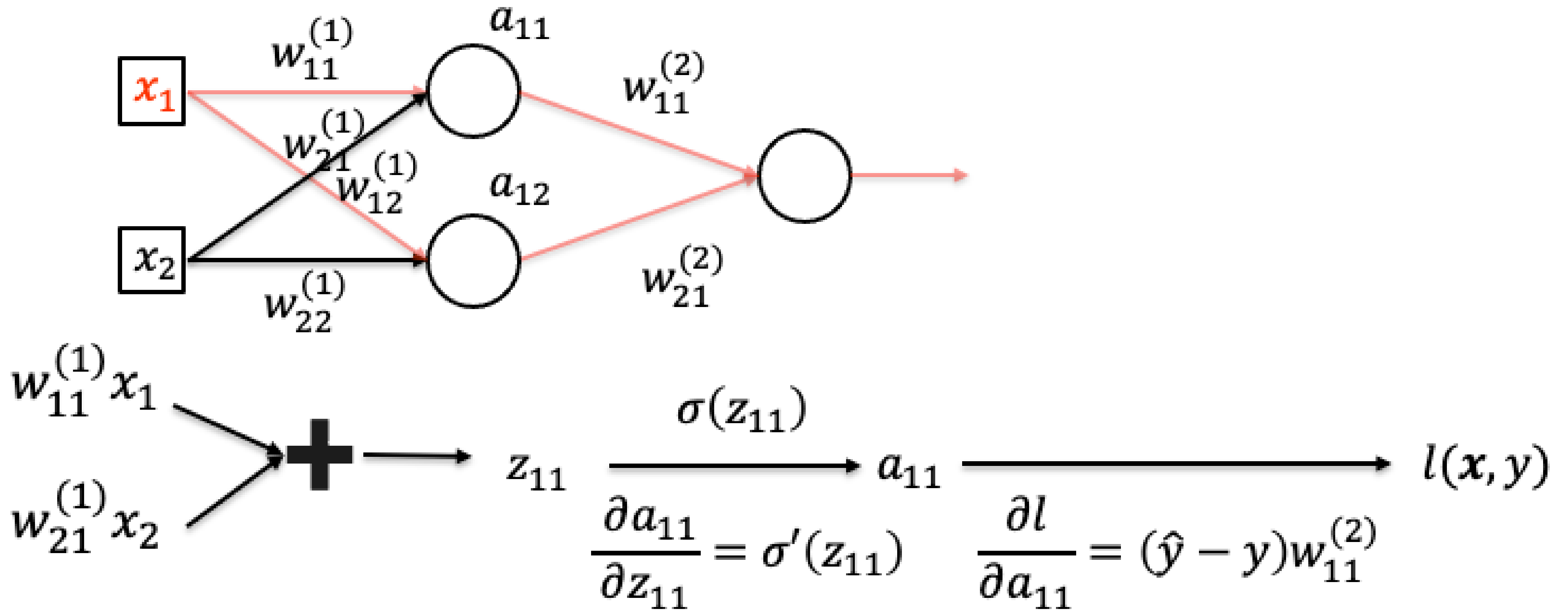
# Calculate Gradient (on one data point)



- By chain rule: $\dfrac{\partial l}{\partial w_{11}^{(1)}} = \dfrac{\partial l}{\partial a_{11}} \dfrac{\partial a_{11}}{\partial w_{11}^{(1)}} = (\hat{y} - y) w_{11}^{(2)} a_{11} (1 - a_{11}) x_1$

# Calculate Gradient (on one data point)



- By chain rule: $\dfrac{\partial l}{\partial x_1} = \dfrac{\partial l}{\partial a_{11}}\dfrac{\partial a_{11}}{\partial x_1} + \dfrac{\partial l}{\partial a_{12}}\dfrac{\partial a_{12}}{\partial x_1}$

# Quiz Break

Gradient Descent in neural network training computes the _____ of a loss function with respect to the model _____ until convergence.

A  gradients, parameters

B  parameters, gradients

C  loss, parameters

D parameters, loss

# Quiz Break

Gradient Descent in neural network training computes the _____ of a loss function with respect to the model _____ until convergence.

A  gradients, parameters

B  parameters, gradients

C  loss, parameters

D parameters, loss

# Quiz Break

Suppose you are given a dataset with 1,000,000 images to train with. Which of the following methods is more desirable if training resources are limit but enough accuracy is needed?

A  Gradient Descent

B  Stochastic Gradient Descent

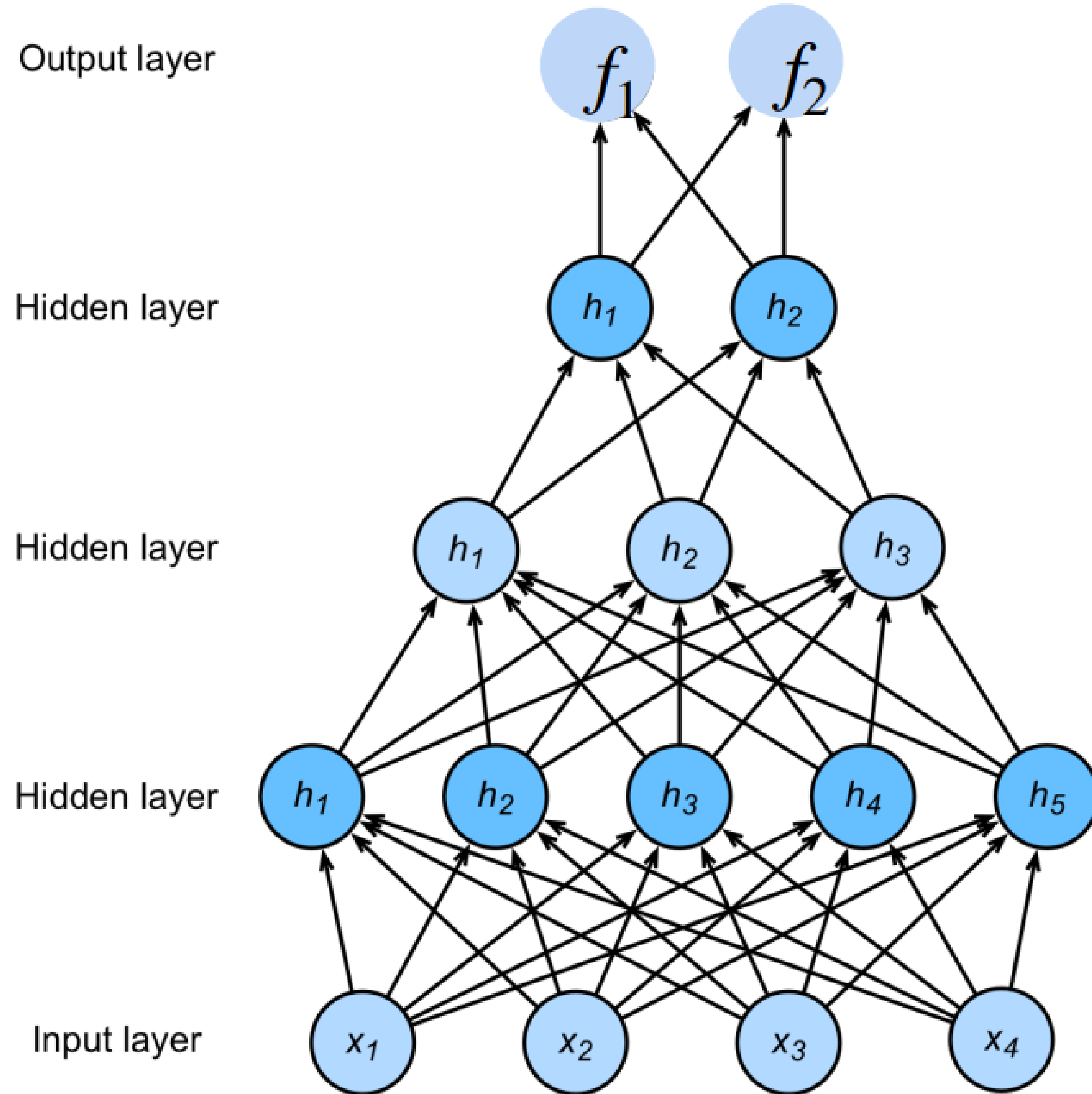C  Minibatch Stochastic Gradient Descent

D  Computation Graph

# Quiz Break

Suppose you are given a dataset with 1,000,000 images to train with. Which of the following methods is more desirable if training resources are limit but enough accuracy is needed?

A  Gradient Descent

B  Stochastic Gradient Descent

C  Minibatch Stochastic Gradient Descent

D  Computation Graph

# Neural Networks as a Computational Graph

# Deep neural networks (DNNs)



Output layer

Hidden layer

Hidden layer

Hidden layer

Input layer

$$\mathbf{h}_1 = \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

$$\mathbf{h}_2 = \sigma(\mathbf{W}^{(2)}\mathbf{h}_1 + \mathbf{b}^{(2)})$$

$$\mathbf{h}_3 = \sigma(\mathbf{W}^{(3)}\mathbf{h}_2 + \mathbf{b}^{(3)})$$

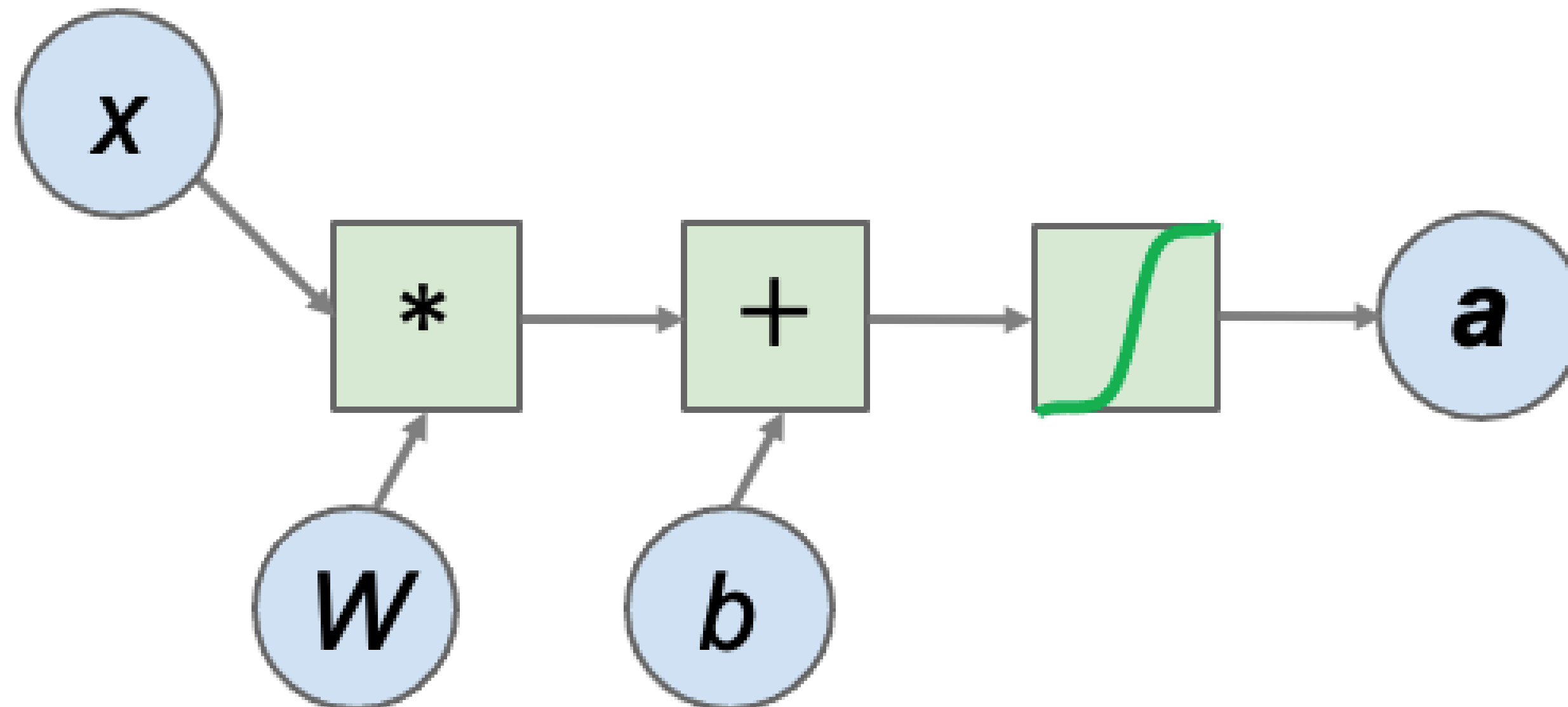$$\mathbf{f} = \mathbf{W}^{(4)}\mathbf{h}_3 + \mathbf{b}^{(4)}$$

$$\mathbf{p} = \text{softmax}(\mathbf{f})$$

NNs are composition of nonlinear functions
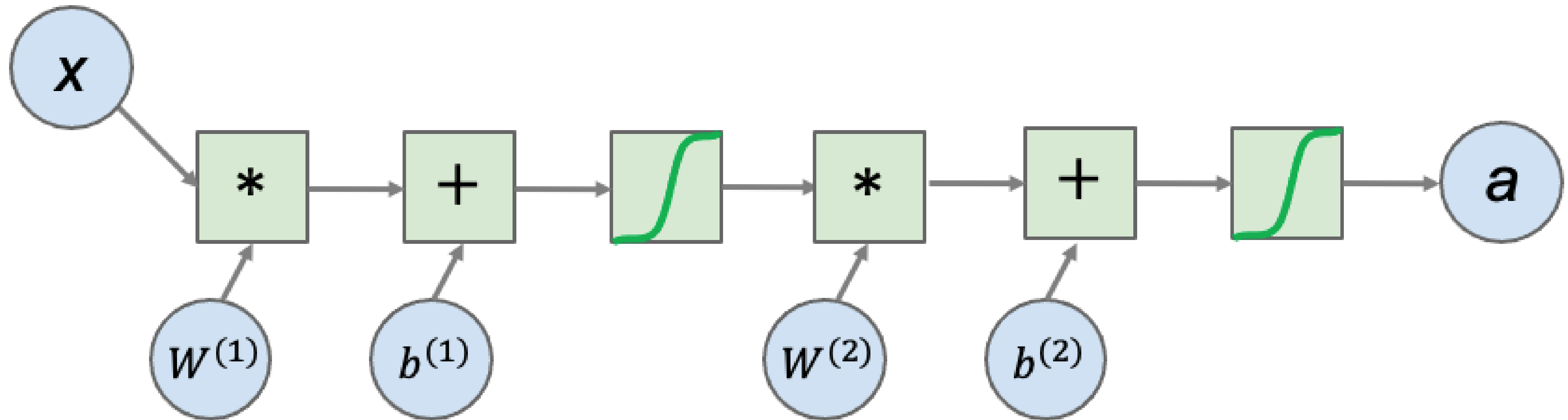
# Neural networks as variables + operations

$$\mathbf{a} = sigmoid(\mathbf{Wx} + \mathbf{b})$$

- Can describe with a **computational graph**

- Decompose functions into atomic operations

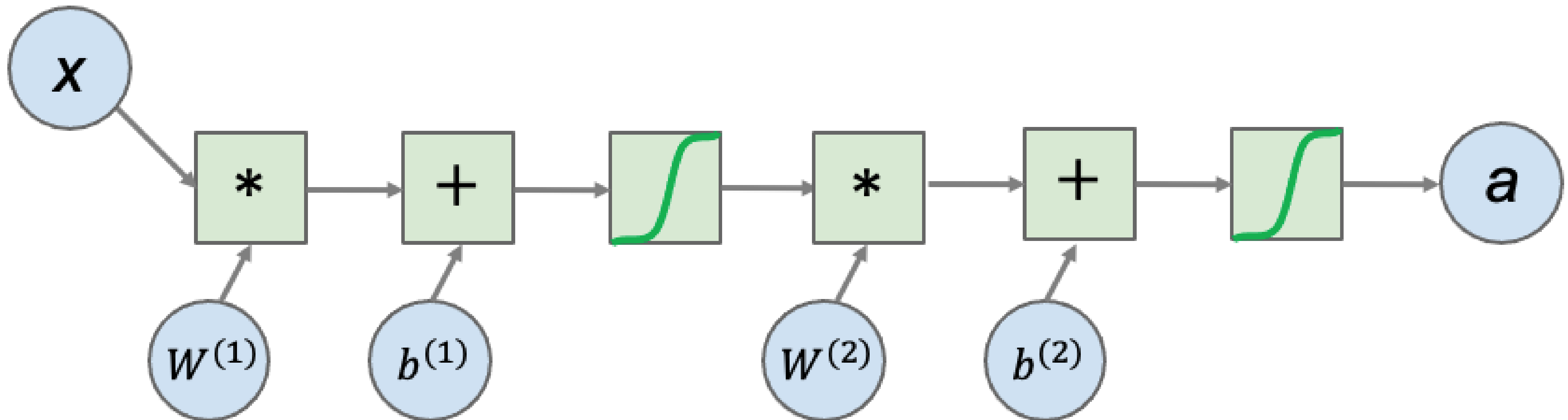- Separate data (variables) and computing (operations)

# Neural networks as a computational graph
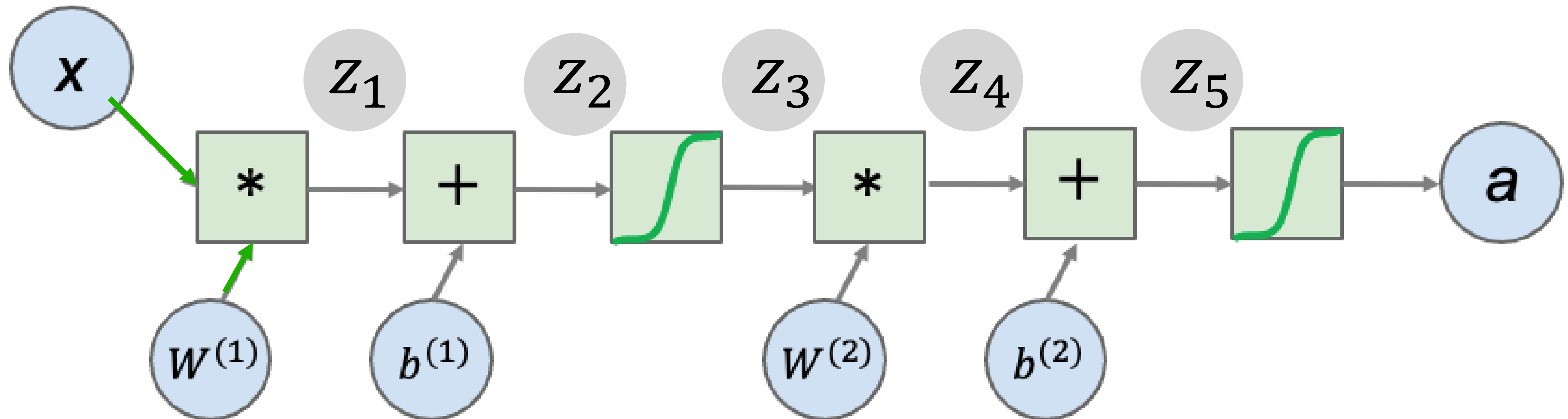
- A two-layer neural network

# Neural networks as a computational graph

- A two-layer neural network
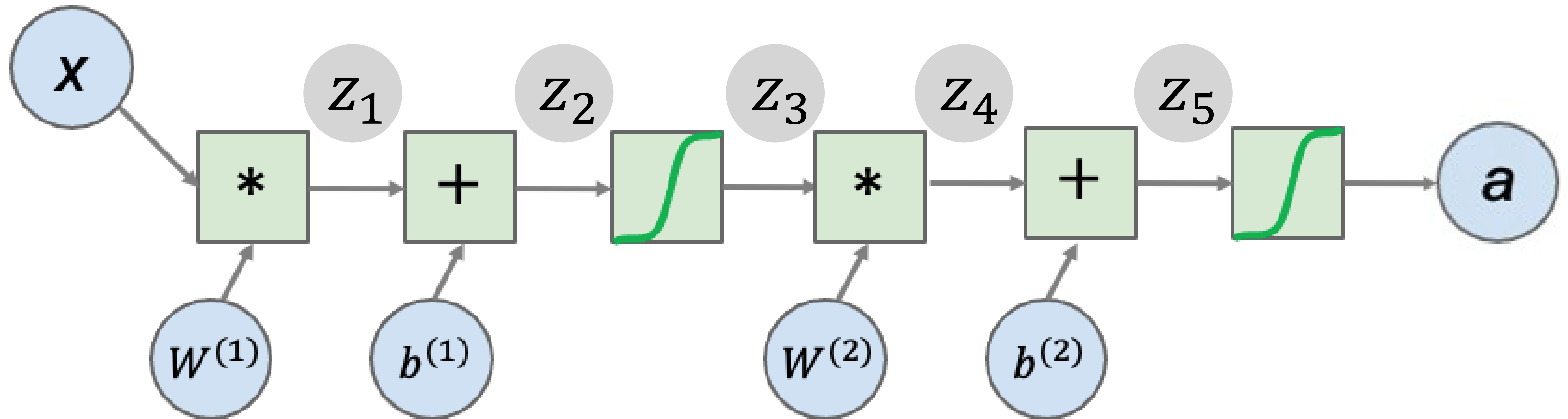- Forward propagation vs. backward propagation

# Neural networks: forward propagation

- A two-layer neural network
- Intermediate variables Z
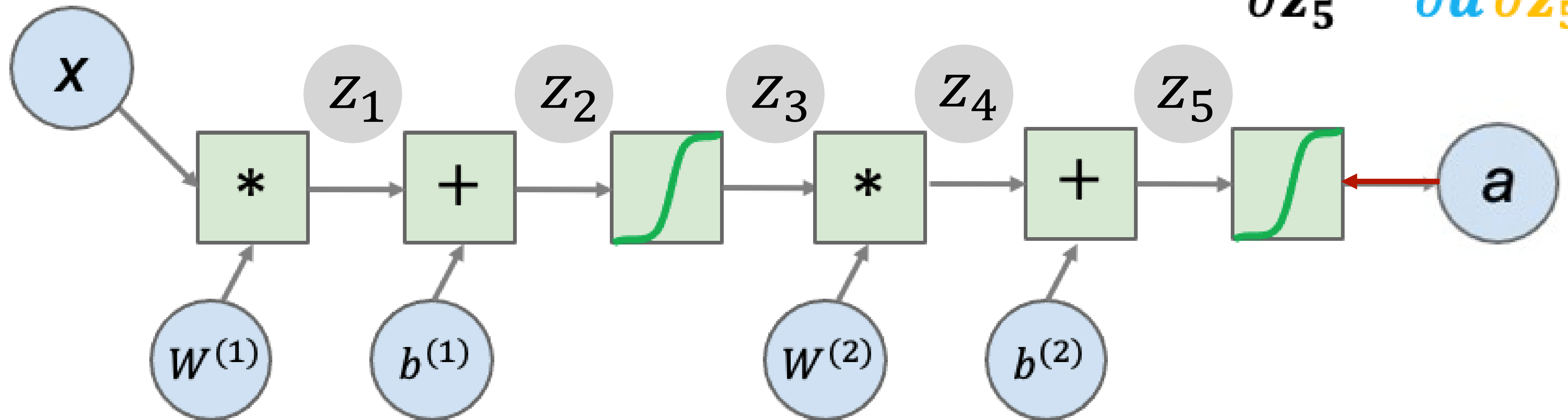
# Neural networks: backward propagation

- A two-layer neural network
- Assuming forward propagation is done
- Minimize a **loss function** L

# Neural networks: backward propagation

- A two-layer neural network
- Assuming forward propagation is done
- Minimize a **loss function** L

$$\frac{\partial L}{\partial z_5} = \frac{\partial L}{\partial a}\frac{\partial a}{\partial z_5}$$
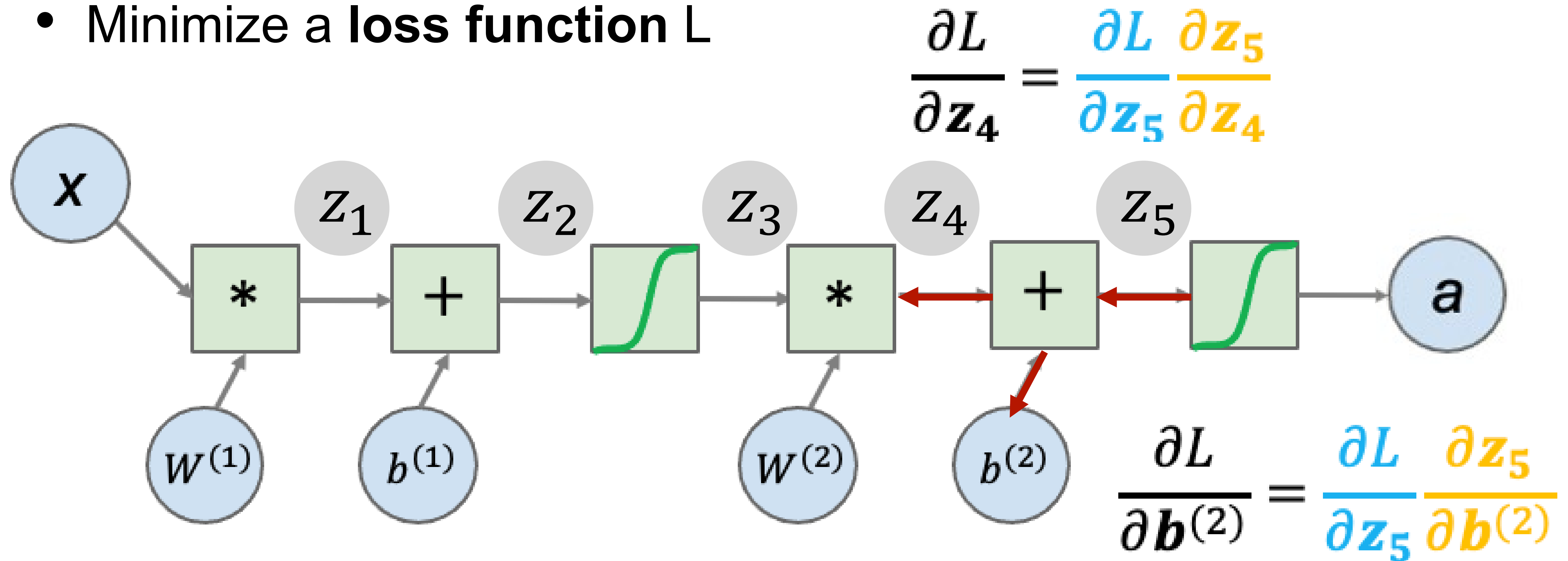
# Neural networks: backward propagation

- A two-layer neural network
- Assuming forward propagation is done
- Minimize a **loss function** L

$$\frac{\partial L}{\partial z_4} = \frac{\partial L}{\partial z_5} \frac{\partial z_5}{\partial z_4}$$



$$\frac{\partial L}{\partial b^{(2)}} = \frac{\partial L}{\partial z_5} \frac{\partial z_5}{\partial b^{(2)}}$$

# Neural networks: backward propagation

- A two-layer neural network
- Assuming forward propagation is done

$$\frac{\partial L}{\partial z_3} = \frac{\partial L}{\partial z_4}\frac{\partial z_4}{\partial z_3}$$



$$\frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial z_4}\frac{\partial z_4}{\partial W^{(2)}}$$

# Backward propagation: A modern treatment

- First, define a neural network as a computational graph
  - Nodes are variables and operations.
- Must be a directed graph
- All operations must be <span style="color:red">differentiable</span>.
- Backpropagation computes partial derivatives starting from the loss and then working backwards through the graph.

# Backward propagation: PyTorch

```python
for t in range(2000):

    # Forward pass: compute predicted y by passing x to the
    # override the __call__ operator so you can call them
    # doing so you pass a Tensor of input data to the Modul
    # a Tensor of output data.
    y_pred = model(xx)

    # Compute and print loss. We pass Tensors containing th
    # values of y, and the loss function returns a Tensor
    # loss.
    loss = loss_fn(y_pred, y)
    if t % 100 == 99:
        print(t, loss.item())

    # Zero the gradients before running the backward pass.
    model.zero_grad()

    # Backward pass: compute gradient of the loss with resp
    # parameters of the model. Internally, the parameters
    # in Tensors with requires_grad=True, so this call wil
    # all learnable parameters in the model.
    loss.backward()

    # Update the weights using gradient descent. Each param
    # we can access its gradients like we did before.
    with torch.no_grad():
        for param in model.parameters():
            param -= learning_rate * param.grad
```

Forward propagation

Backward propagation

Gradient Descent

Q1.1 Suppose we want to solve the following k-class classification problem with cross entropy loss
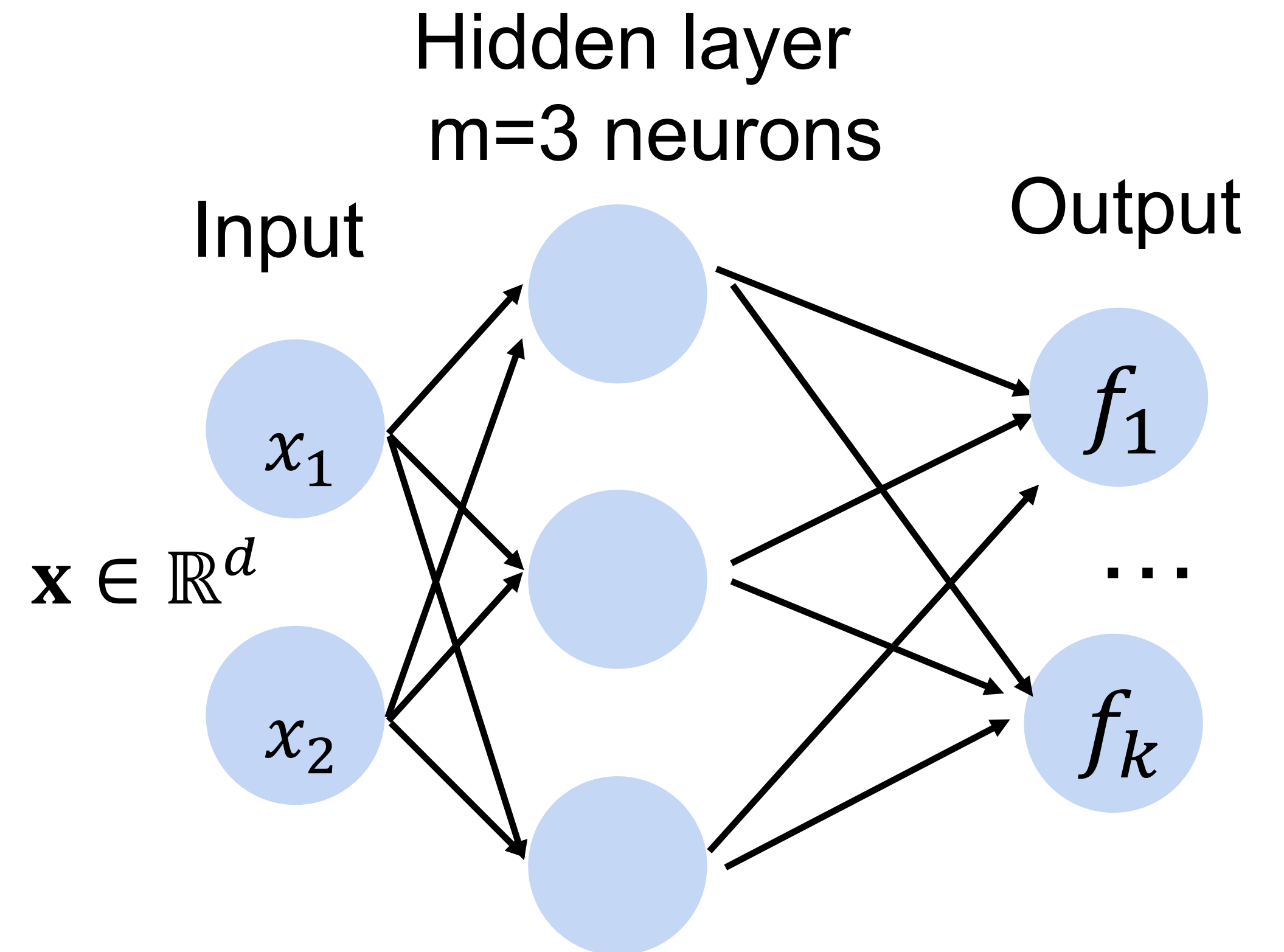
$\ell(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{j=1}^{k} y_j \log \hat{y}_j$ , where the ground truth and predicted probabilities $\mathbf{y}, \hat{\mathbf{y}} \in \mathbb{R}^k$. Recall that the

softmax function turns output into probabilities: $\hat{y}_j = \frac{\exp f_j(x)}{\sum_i^k \exp f_i(x)}$. What is the partial derivative $\partial_{f_j} \ell(\mathbf{y}, \hat{\mathbf{y}})$?

A. $\hat{y}_j - y_j$

B. $\exp(y_j) - y_j$

C. $y_j - \hat{y}_j$

Hidden layer
m=3 neurons

Input

Output



$\mathbf{x} \in \mathbb{R}^d$

$x_1$

$x_2$

$f_1$

$\cdots$

$f_k$

Q1.1 Suppose we want to solve the following k-class classification problem with cross entropy loss

$\ell(\mathbf{y}, \hat{\mathbf{y}}) = -\sum\limits_{j=1}^{k} y_j \log \hat{y}_j$ , where the ground truth and predicted probabilities $\mathbf{y}, \hat{\mathbf{y}} \in \mathbb{R}^k$. Recall that the

softmax function turns output into probabilities: $\hat{y}_j = \dfrac{\exp f_j(x)}{\sum_i^k \exp f_i(x)}$. What is the partial derivative $\partial_{f_j} \ell(\mathbf{y}, \hat{\mathbf{y}})$?

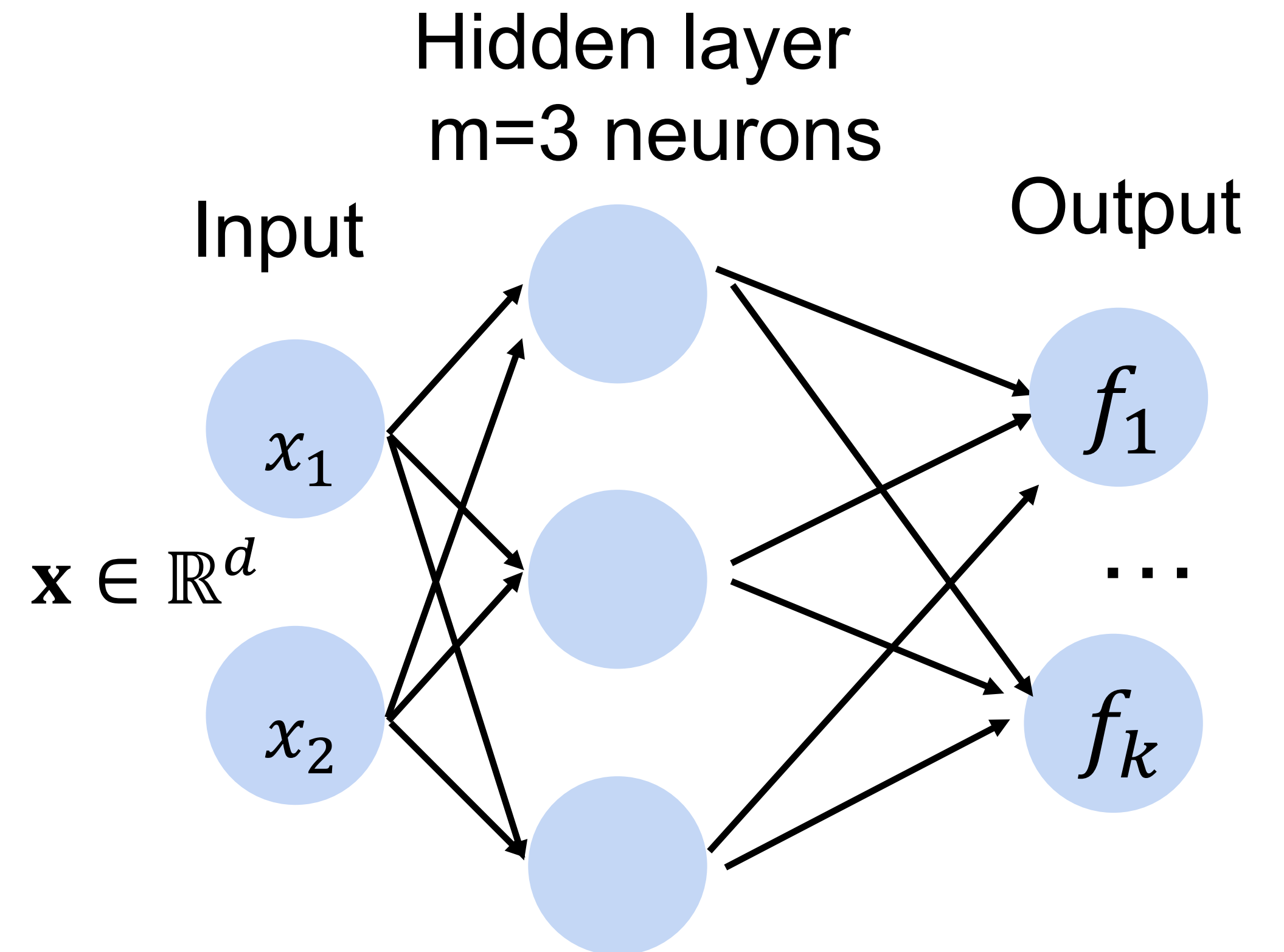A. $\hat{y}_j - y_j$

B. $\exp(y_j) - y_j$

C. $y_j - \hat{y}_j$

Rewrite $\quad \ell(\mathbf{y}, \hat{\mathbf{y}}) = -\sum\limits_{j=1}^{k} y_j \log \dfrac{\exp(f_j)}{\sum_{i=1}^{k} \exp(f_i)}$

$= \sum\limits_{j=1}^{k} y_j \log \sum\limits_{i=1}^{k} \exp(f_i) - \sum\limits_{j=1}^{k} y_j f_j$

$= \log \sum\limits_{i=1}^{k} \exp(f_i) - \sum\limits_{j=1}^{k} y_j f_j .$

We have $\quad \partial_{f_j} \ell(\mathbf{y}, \hat{\mathbf{y}}) = \dfrac{\exp(f_j)}{\sum_{i=1}^{k} \exp(f_k)} - y_j = \hat{y}_j - y_j$

Hidden layer
m=3 neurons

Input

Output

$x_1$

$\mathbf{x} \in \mathbb{R}^d$

$x_2$

$f_1$

...

$f_k$

# Numerical Stability

# Gradients for Neural Networks

- Compute the gradient of the loss $\ell$ w.r.t. $\mathbf{W}_t$

$$\frac{\partial \ell}{\partial \mathbf{W}^t} = \frac{\partial \ell}{\partial \mathbf{h}^d} \frac{\partial \mathbf{h}^d}{\partial \mathbf{h}^{d-1}} \ldots \frac{\partial \mathbf{h}^{t+1}}{\partial \mathbf{h}^t} \frac{\partial \mathbf{h}^t}{\partial \mathbf{W}^t}$$

Multiplication of *many* matrices

Wikipedia

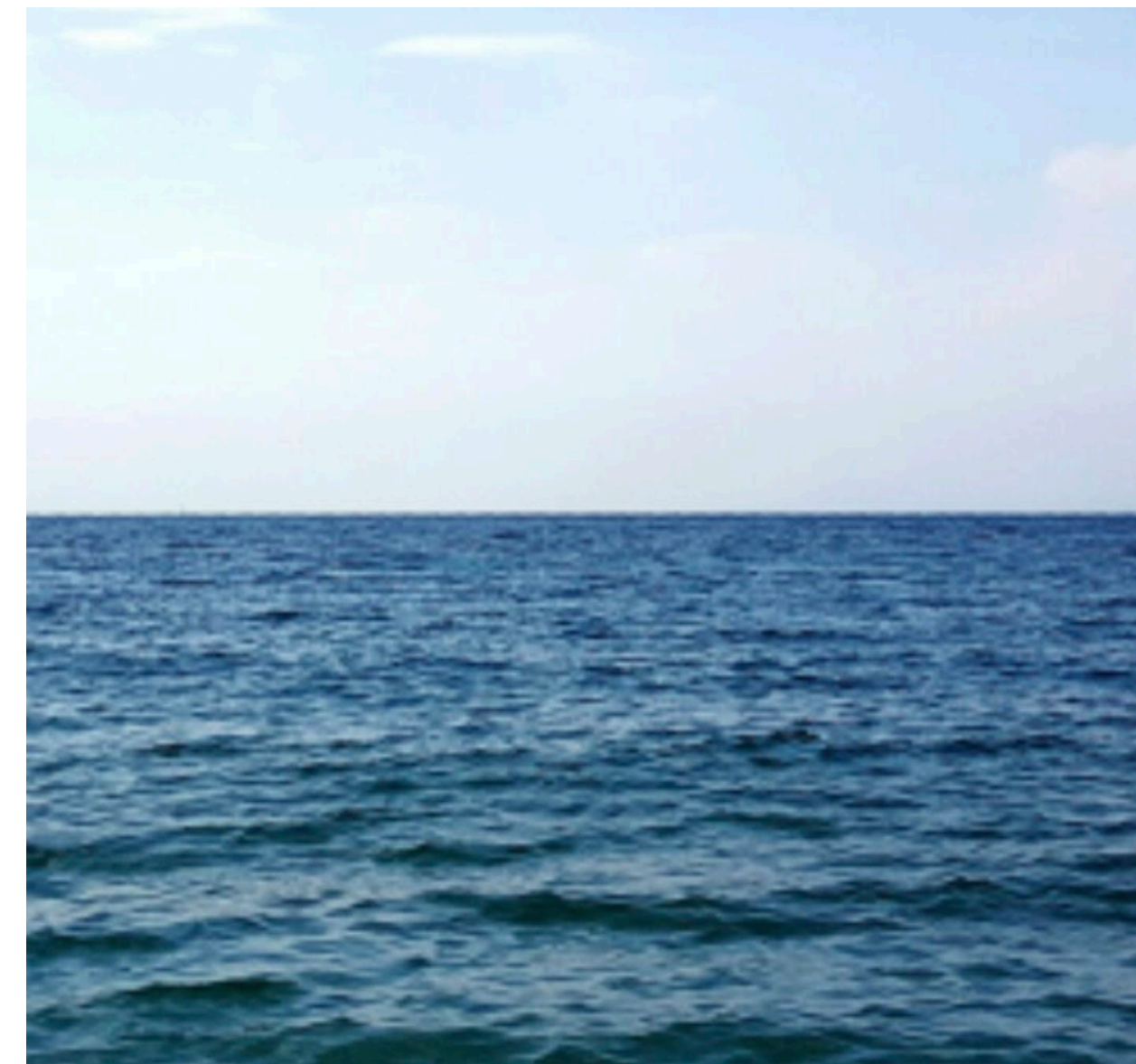# Two Issues for Deep Neural Networks

$$\prod_{i=t}^{d-1} \frac{\partial \mathbf{h}^{i+1}}{\partial \mathbf{h}^i}$$

Gradient Exploding

Gradient Vanishing





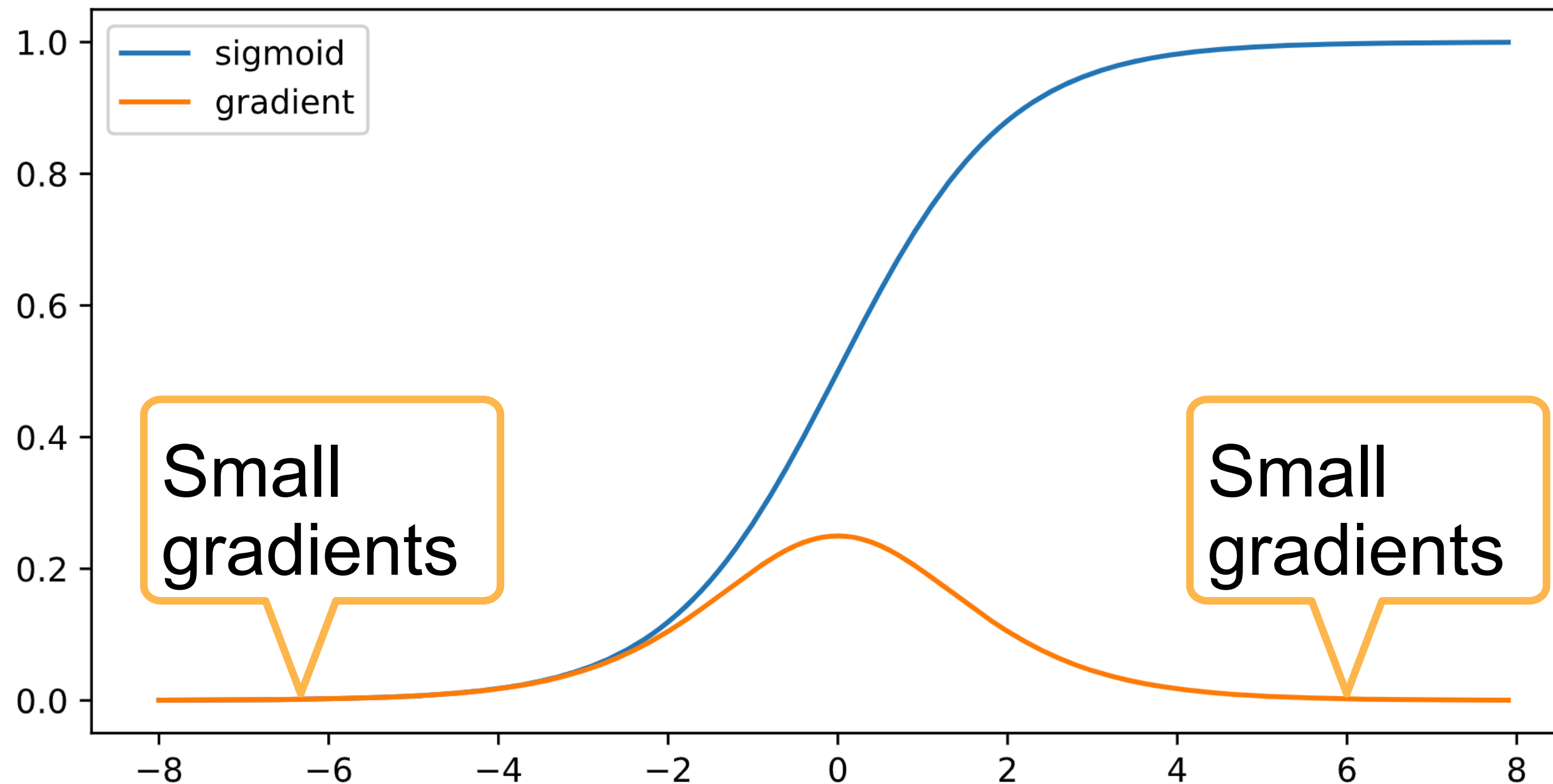$$1.5^{100} \approx 4 \times 10^{17}$$

$$0.8^{100} \approx 2 \times 10^{-10}$$

# Issues with Gradient Exploding

- Value out of range: infinity value (NaN)
- Sensitive to learning rate (LR)
  - Not small enough LR → larger gradients
  - Too small LR → No progress
  - May need to change LR dramatically during training

# Gradient Vanishing

- Use sigmoid as the activation function

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \sigma'(x) = \sigma(x)(1 - \sigma(x))$$

# Issues with Gradient Vanishing

- Gradients with value 0
- No progress in training
  - No matter how to choose learning rate
- Severe with bottom layers (those near the input)
  - Only top layers (near output) are well trained
  - No benefit to make networks deeper

# How to stabilize training?

# Stabilize Training: Practical Considerations

- Goal: make sure gradient values are in a proper range
  - E.g. in [1e-6, 1e3]
- Multiplication → plus
  - Architecture change (e.g., ResNet)
- Normalize
  - Batch Normalization, Gradient clipping
- Proper activation functions

Quiz. Which of the following are TRUE about the vanishing gradient problem in neural networks? Multiple answers are possible.

A. Deeper neural networks tend to be more susceptible to vanishing gradients.

B. Using the ReLU function can reduce this problem.

C. If a network has the vanishing gradient problem for one training point due to the sigmoid function, it will also have a vanishing gradient for every other training point.

D. Networks with sigmoid functions don't suffer from the vanishing gradient problem if trained with the cross-entropy loss.

Quiz. Which of the following are TRUE about the vanishing gradient problem in neural networks? Multiple answers are possible?

A. Deeper neural networks tend to be more susceptible to vanishing gradients.

B. Using the ReLU function can reduce this problem.

C. If a network has the vanishing gradient problem for one training point due to the sigmoid function, it will also have a vanishing gradient for every other training point.

D. Networks with sigmoid functions don't suffer from the vanishing gradient problem if trained with the cross-entropy loss.

Quiz. Let's compare sigmoid with rectified linear unit (ReLU). Which of the following statement is NOT true?

A. Sigmoid function is more expensive to compute

B. ReLU has non-zero gradient everywhere

C. The gradient of Sigmoid is always less than 0.3

D. The gradient of ReLU is constant for positive input

Quiz. Let's compare sigmoid with rectified linear unit (ReLU). Which of the following statement is NOT true?

A. Sigmoid function is more expensive to compute

B. ReLU has non-zero gradient everywhere

C. The gradient of Sigmoid is always less than 0.3

D. The gradient of ReLU is constant for positive input

Q5. A Leaky ReLU is defined as *f(x)=max(0.1x, x).* Let f'(0)=1. Does it have non-zero gradient everywhere??
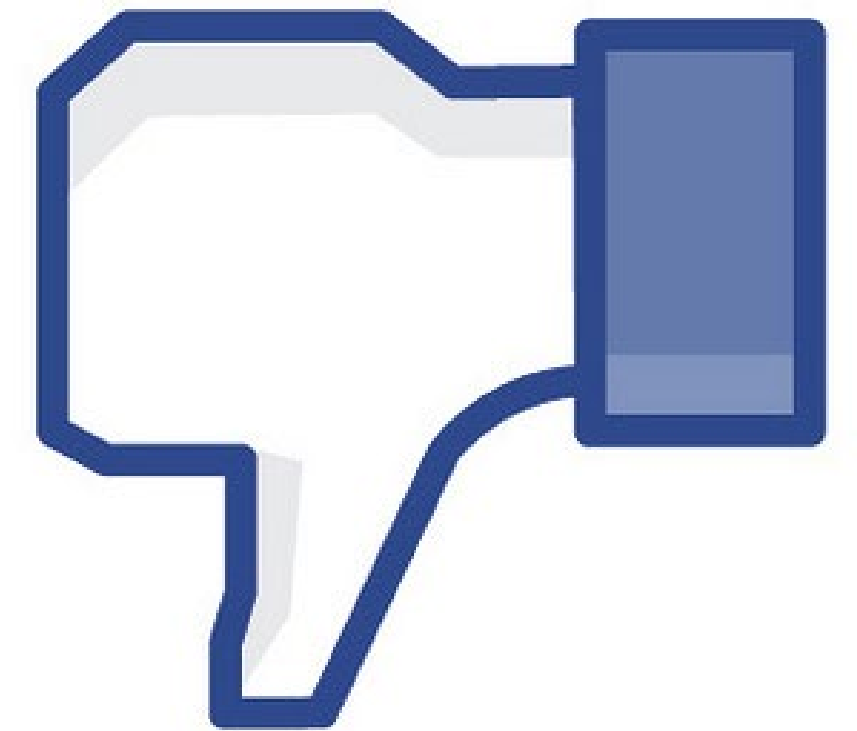
A.Yes

B. No

Q5. A Leaky ReLU is defined as *f(x)=max(0.1x, x).* Let f'(0)=1. Does it have non-zero gradient everywhere??

A.Yes

B. No

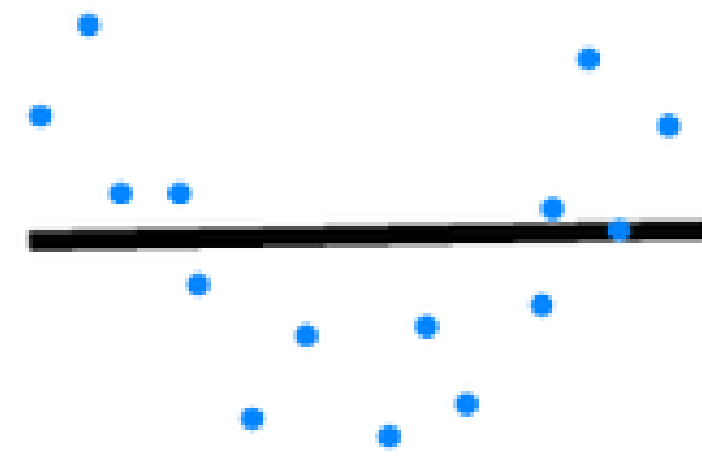# Generalization & Regularization

# How good are the models?

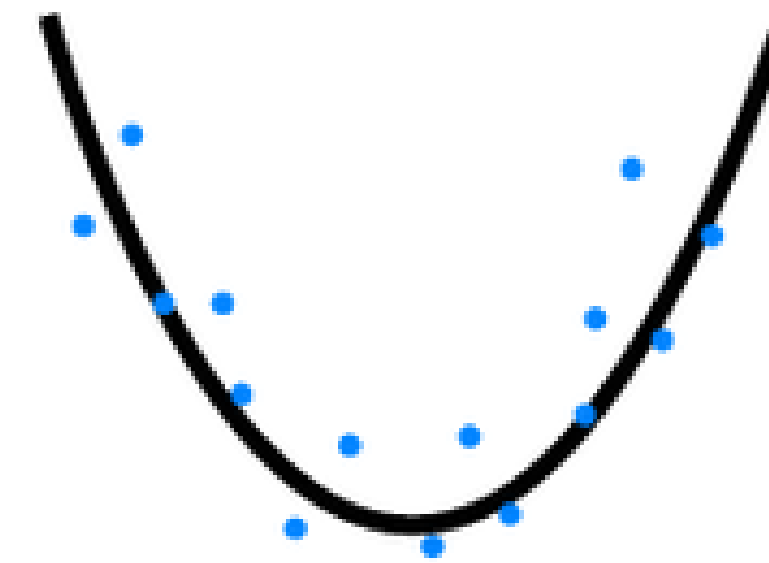# Training Error and Generalization Error

- Training error: model error on the training data
- **Generalization error**: model error on new data
- Example: practice a future exam with past exams
  - Doing well on past exams (training error) doesn't guarantee a good score on the future exam (generalization error)

# Underfitting
# Overfitting



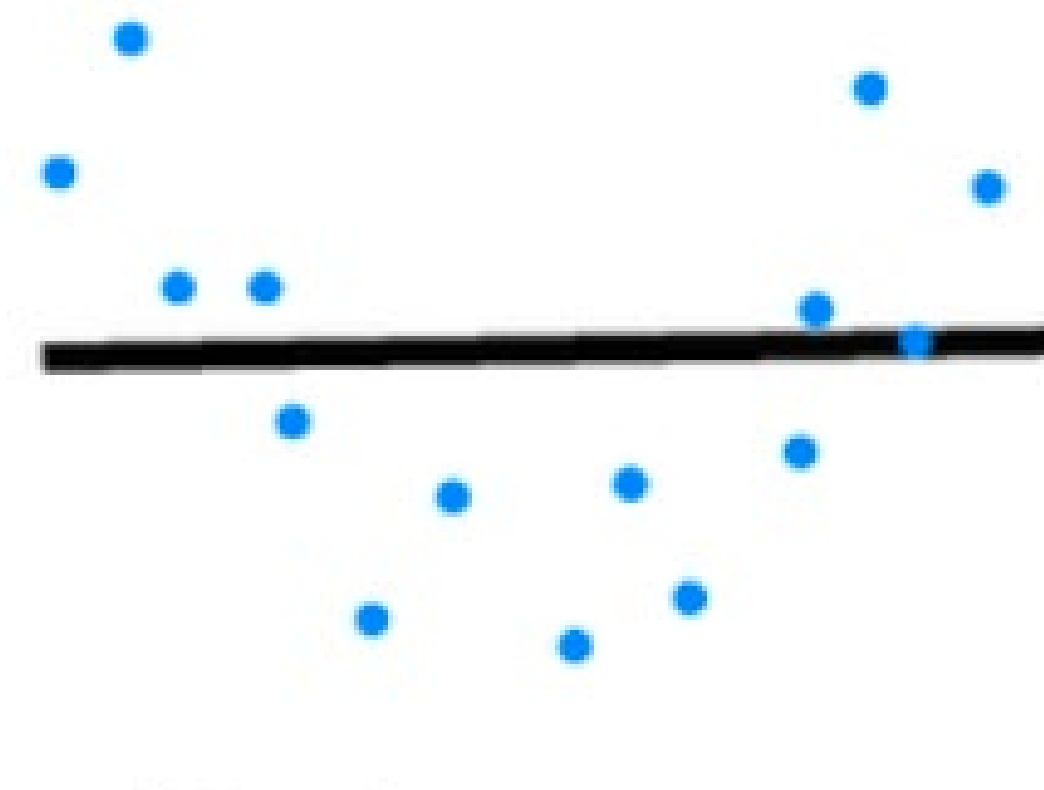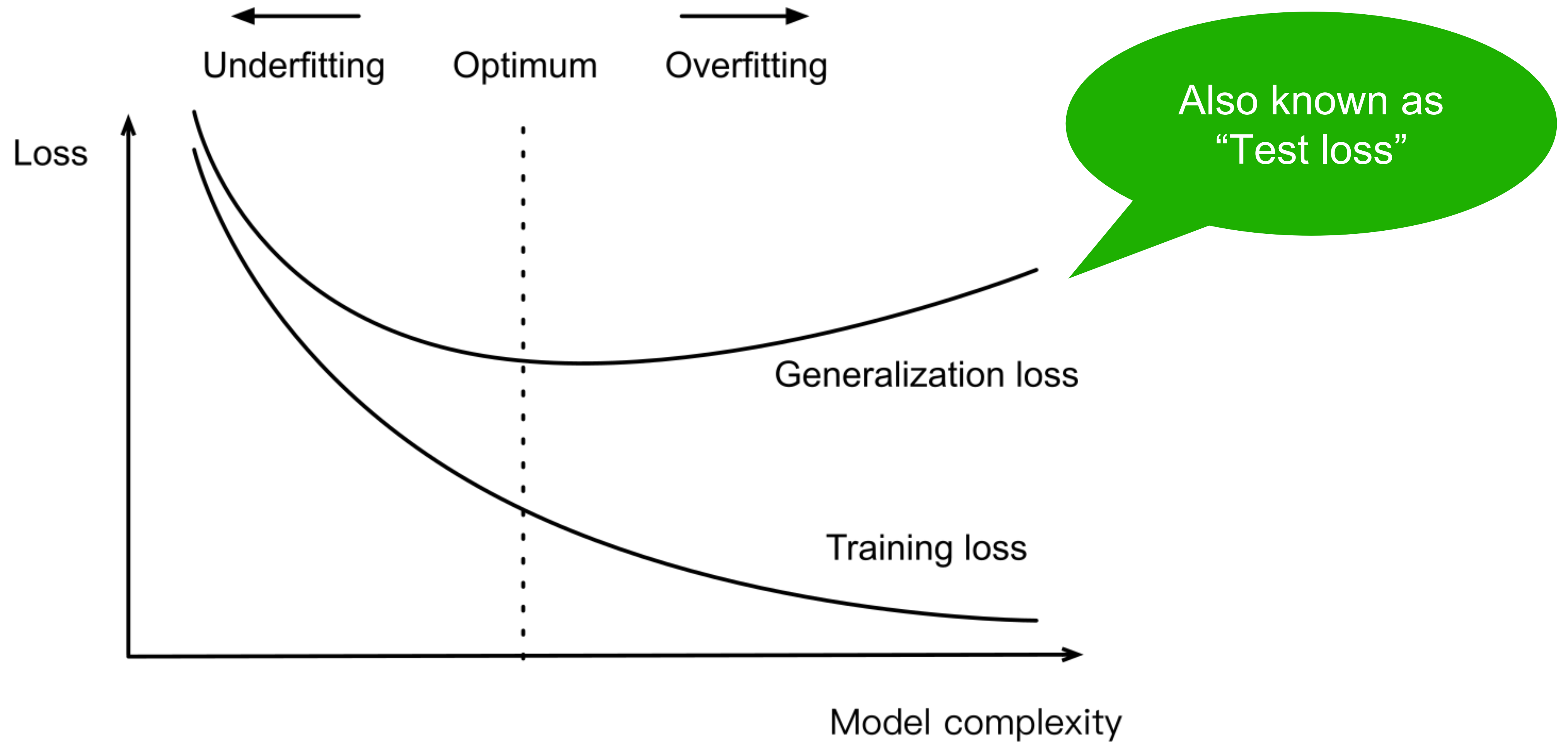Underfitting      Desired      Overfitting

Image credit: hackernoon.com

# Model Capacity

- The ability to fit variety of functions
- Low capacity models struggles to fit training set
  - Underfitting
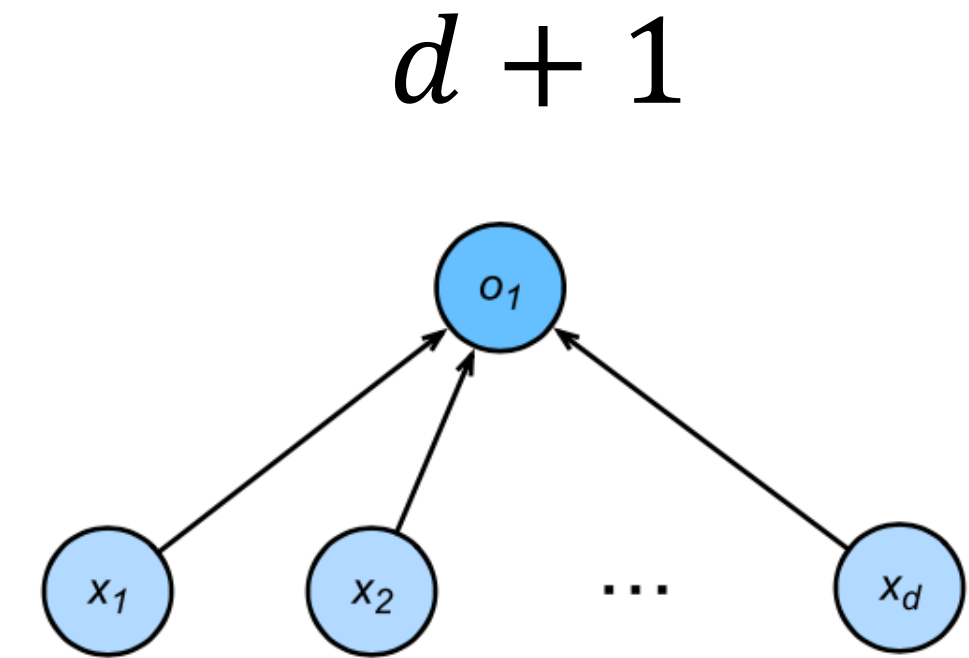- High capacity models can memorize the training set
  - Overfitting
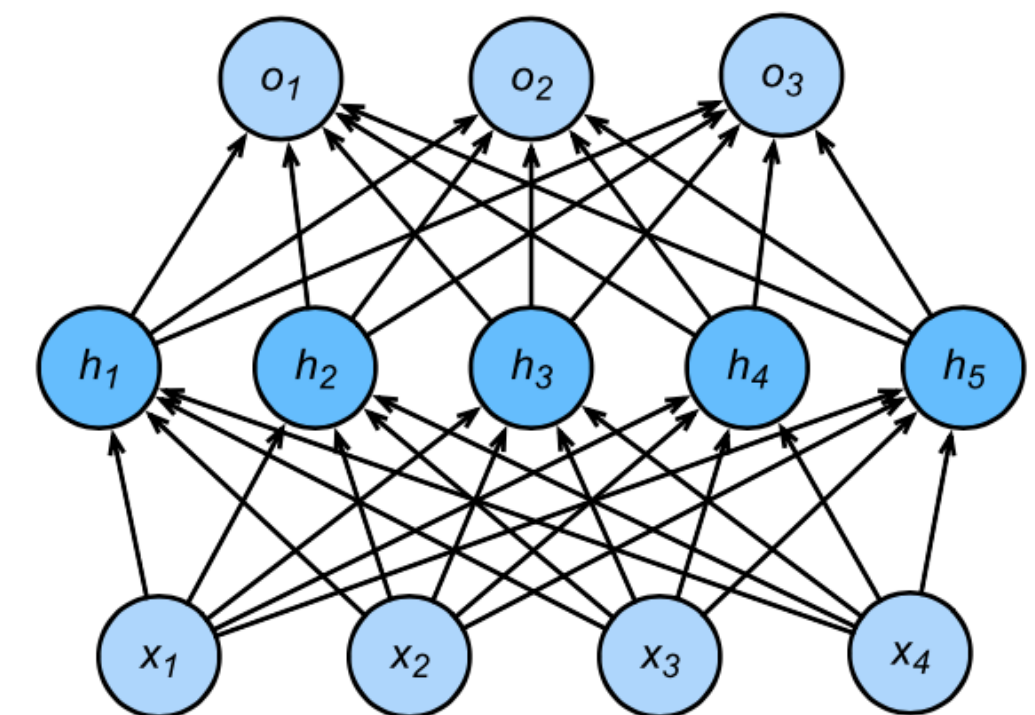
# Influence of Model Complexity



* Recent research has challenged this view for some types of models.

# Estimate Neural Network Capacity

- It's hard to compare complexity between different families of models.

  - e.g. K-NN vs neural networks

- Given a model family, two main factors matter:

  - The number of parameters

  - The values taken by each parameter

$$d + 1$$



$$(d + 1)m + (m + 1)k$$

# Data Complexity

- Multiple factors matters
  - # of examples
  - # of features in each example
  - time/space structure
  - # of labels

# Quiz Break: When training a neural network, which one below indicates that the network has overfit the training data?

A. Training loss is low and generalization loss is high.

B. Training loss is low and generalization loss is low.

C. Training loss is high and generalization loss is high.

D. Training loss is high and generalization loss is low.

E. None of these.

# Quiz Break: When training a neural network, which one below indicates that the network has overfit the training data?

A. Training loss is low and generalization loss is high.

B. Training loss is low and generalization loss is low.

C. Training loss is high and generalization loss is high.

D. Training loss is high and generalization loss is low.

E. None of these.

# Quiz Break: Adding more layers to a multi-layer perceptron may cause _____.

A. Vanishing gradients during back propagation.

B. A more complex decision boundary.

C. Underfitting.

D. Higher test loss.

E. None of these.

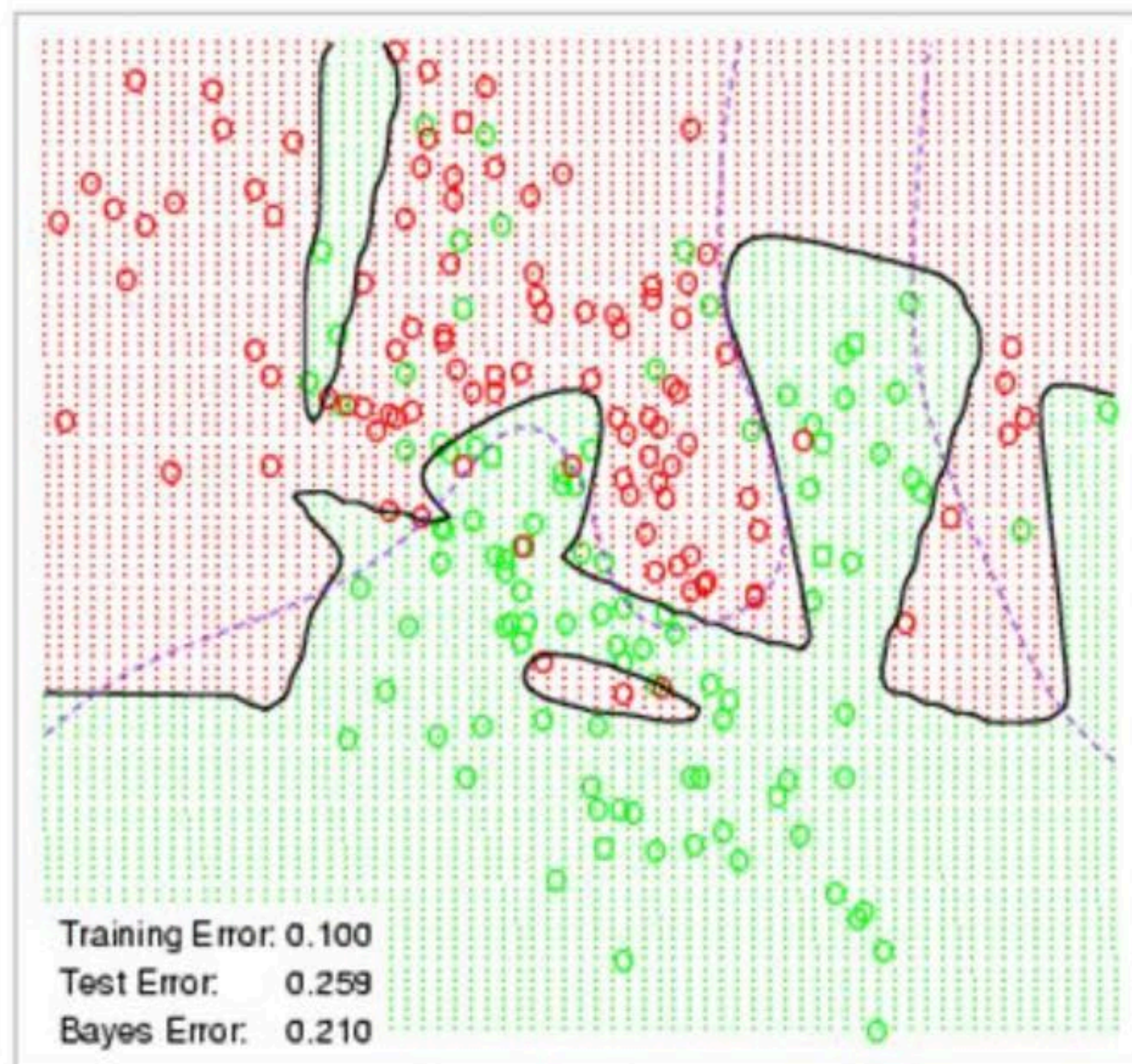# Quiz Break: Adding more layers to a multi-layer perceptron may cause _____. (Multiple answers)

A. Vanishing gradients during back propagation.

B. A more complex decision boundary.
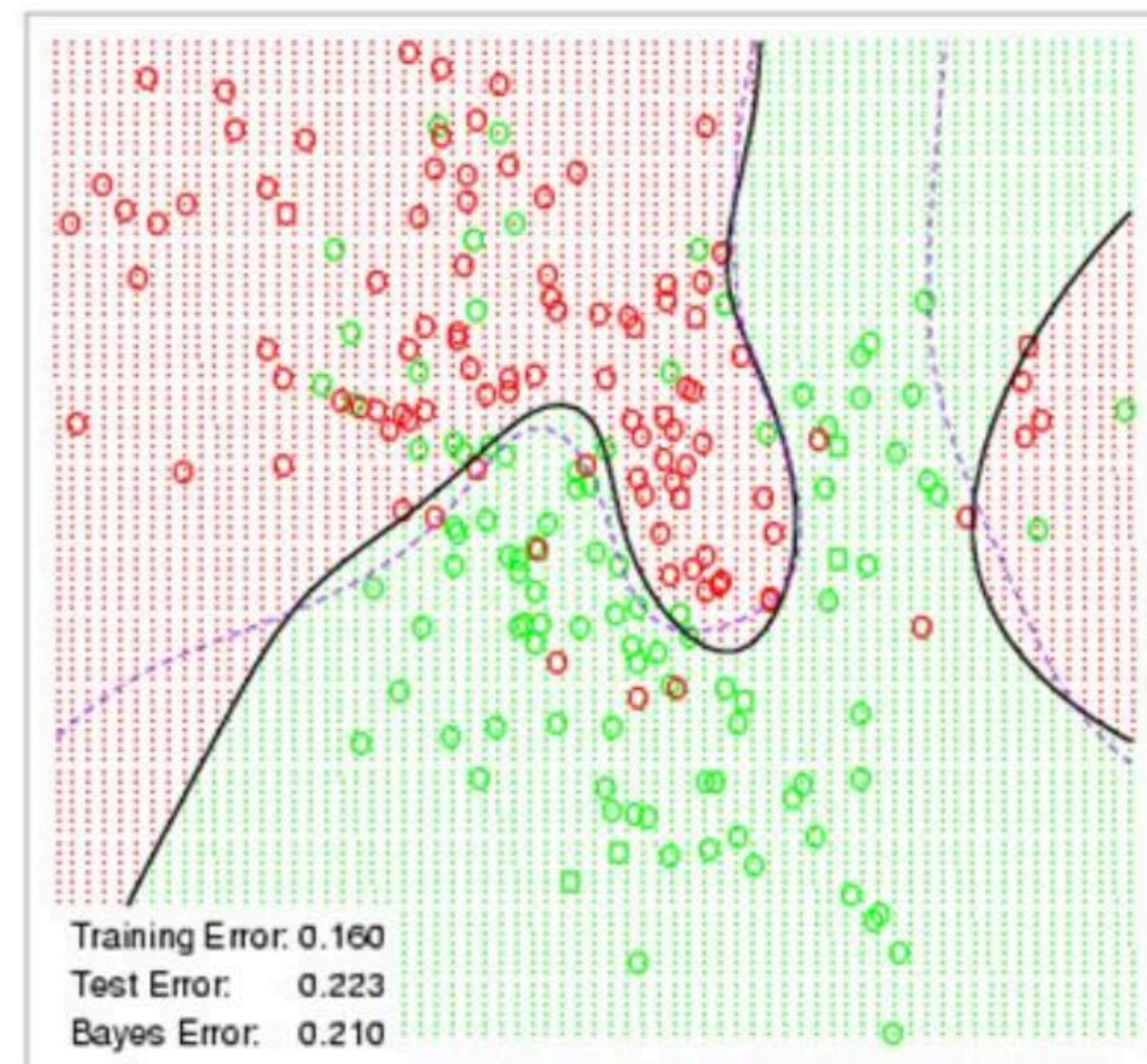
C. Underfitting.

D. Higher test loss.

E. None of these.

# How to regularize the model for better generalization?

# Weight Decay



Neural Network - 10 Units, No Weight Decay

Training Error: 0.100
Test Error: 0.259
Bayes Error: 0.210

Neural Network - 10 Units, Weight Decay=0.02

Training Error: 0.160
Test Error: 0.223
Bayes Error: 0.210

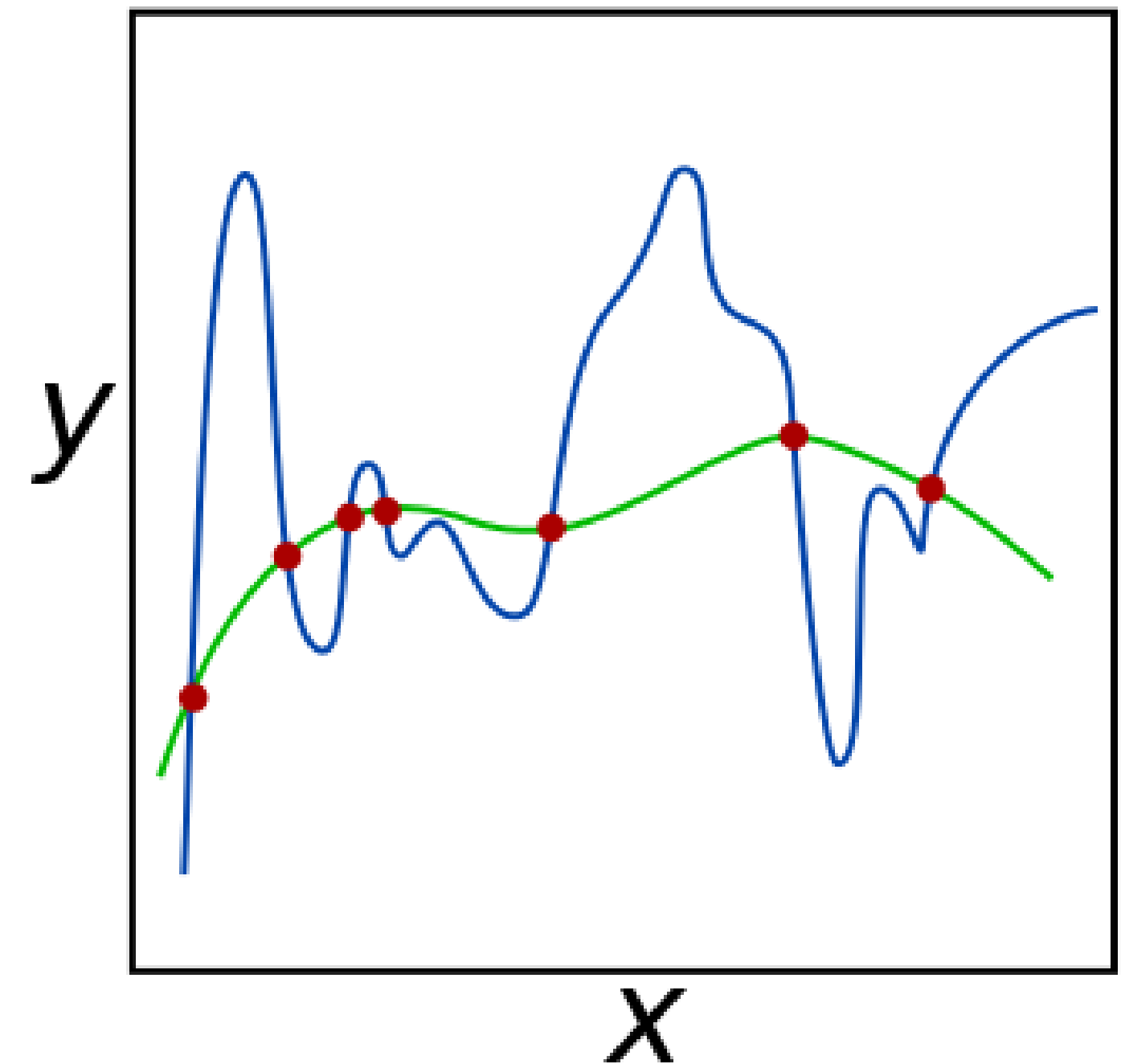# Squared Norm Regularization as Hard Constraint

- Reduce model complexity by limiting value range

$$minL(\mathbf{w}, b) \; subject \; to \parallel \mathbf{w} \parallel^2 \leq B$$

- Often do not regularize bias $b$

  - Doing or not doing has little difference in practice

- A small $B$ means more regularization

# Squared Norm Regularization as Soft Constraint

- We can rewrite the hard constraint version as

$$minL(\mathbf{w}, b) + \frac{\lambda}{2} \parallel \mathbf{w} \parallel^2$$
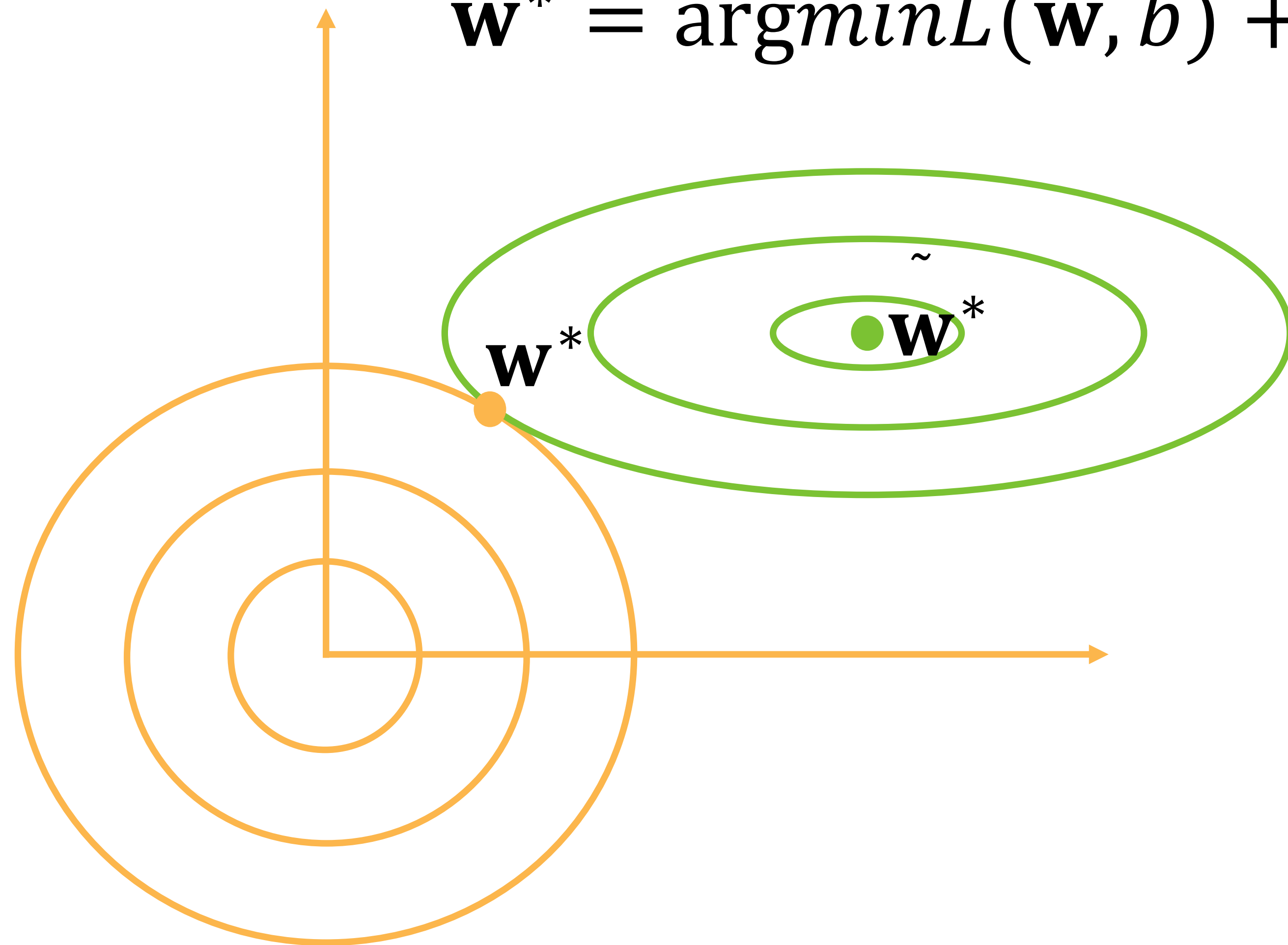
# Squared Norm Regularization as Soft Constraint

- We can rewrite the hard constraint version as

$$minL(\mathbf{w}, b) + \frac{\lambda}{2} \parallel \mathbf{w} \parallel^2$$

- Hyper-parameter $\lambda$ controls regularization importance
- $\lambda = 0$ : no effect

$\lambda \to \infty, \mathbf{w}^* \to \mathbf{0}$

# Illustrate the Effect on Optimal Solutions

$$\mathbf{w}^* = \text{arg}minL(\mathbf{w}, b) + \frac{\lambda}{2} \parallel \mathbf{w} \parallel^2$$



$$\tilde{\mathbf{w}}^* = \text{arg}minL(\mathbf{w}, b)$$

# Dropout

Hinton et al.

# Apply Dropout

- Often apply dropout on the output of hidden fully-connected layers
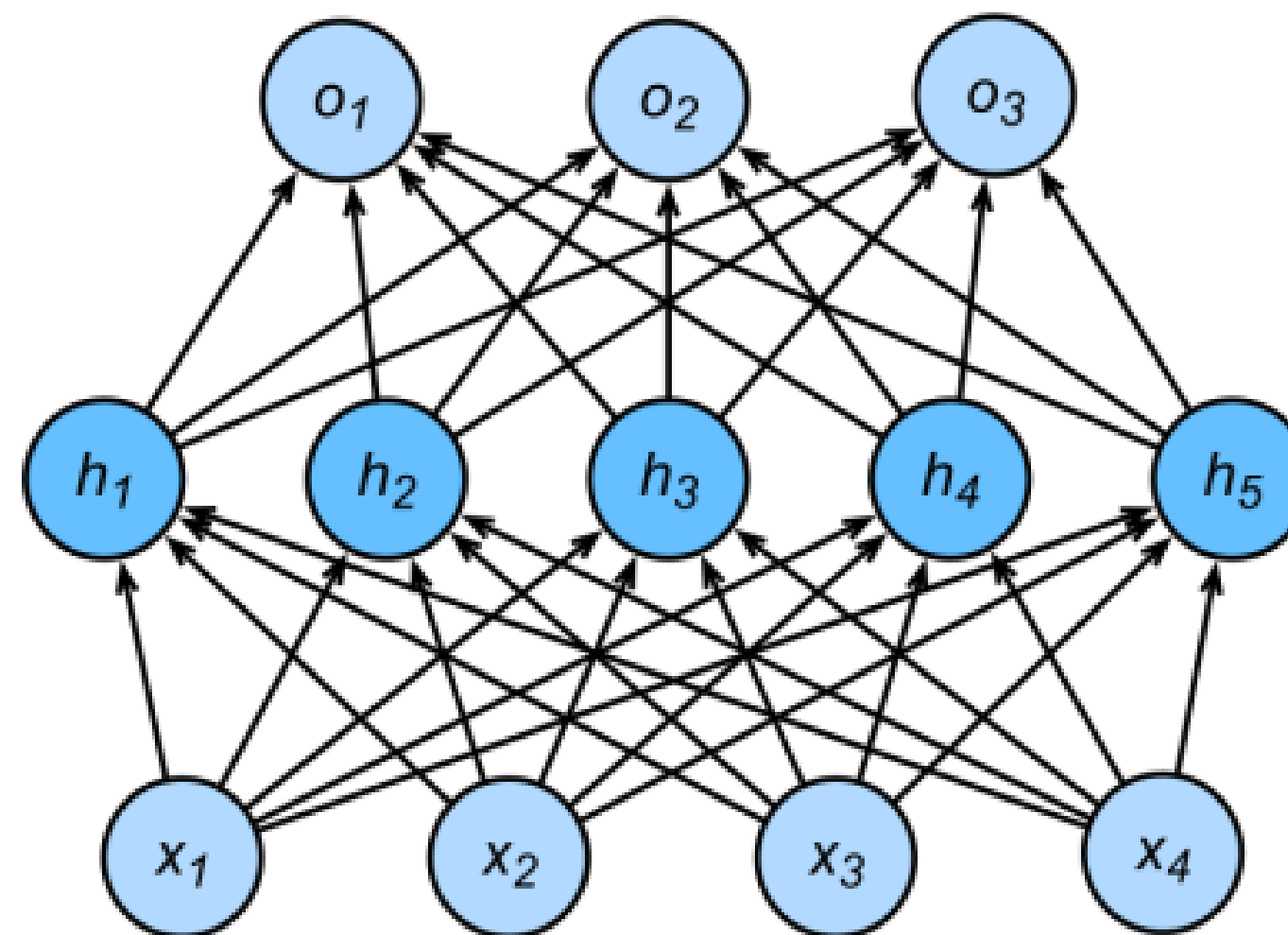
$$\mathbf{h} = \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

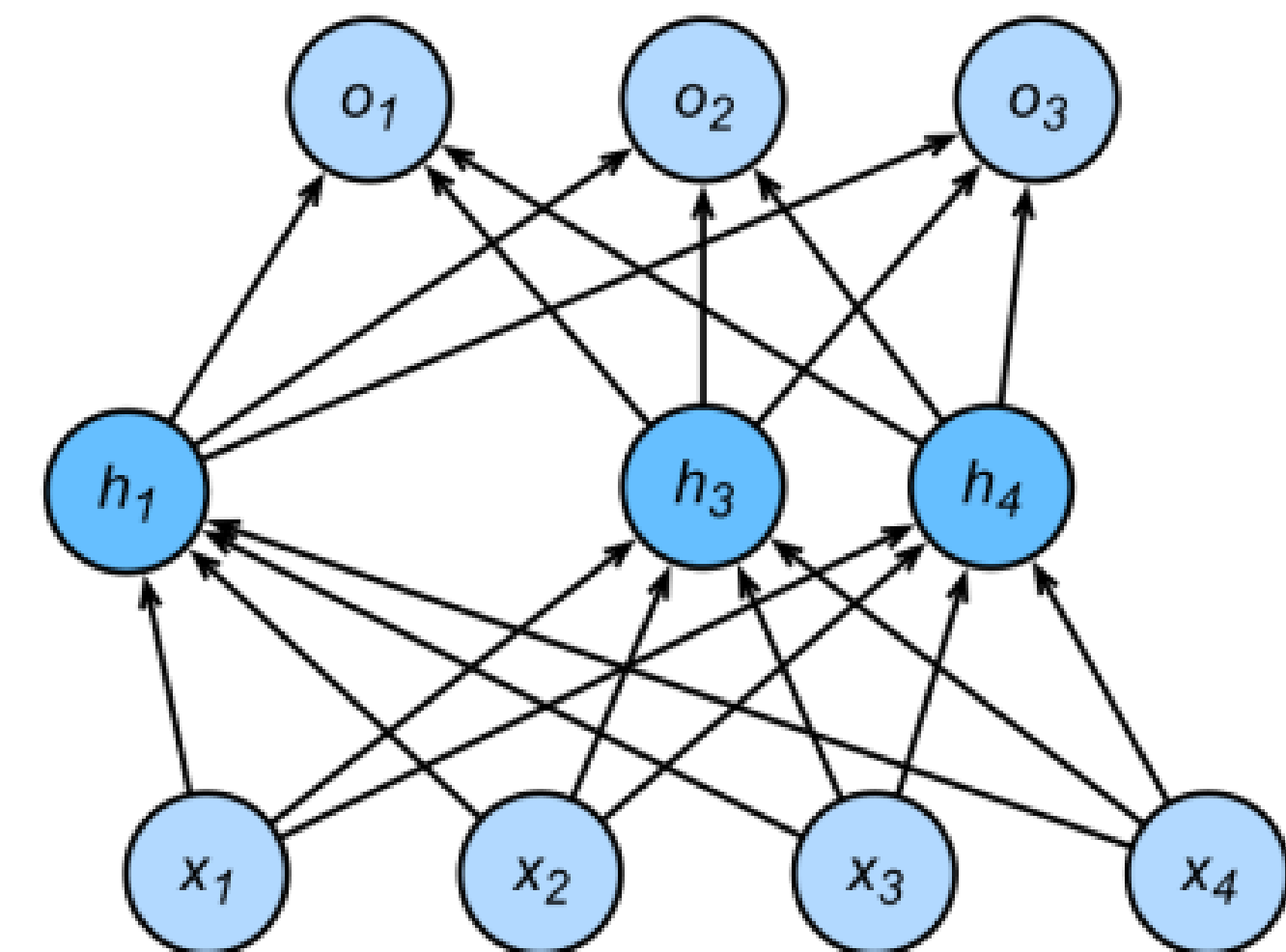$$\mathbf{h}' = \text{dropout}(\mathbf{h})$$

$$\mathbf{o} = \mathbf{W}^{(2)}\mathbf{h}' + \mathbf{b}^{(2)}$$

$$\mathbf{p} = \text{softmax}(\mathbf{o})$$



MLP with one hidden layer
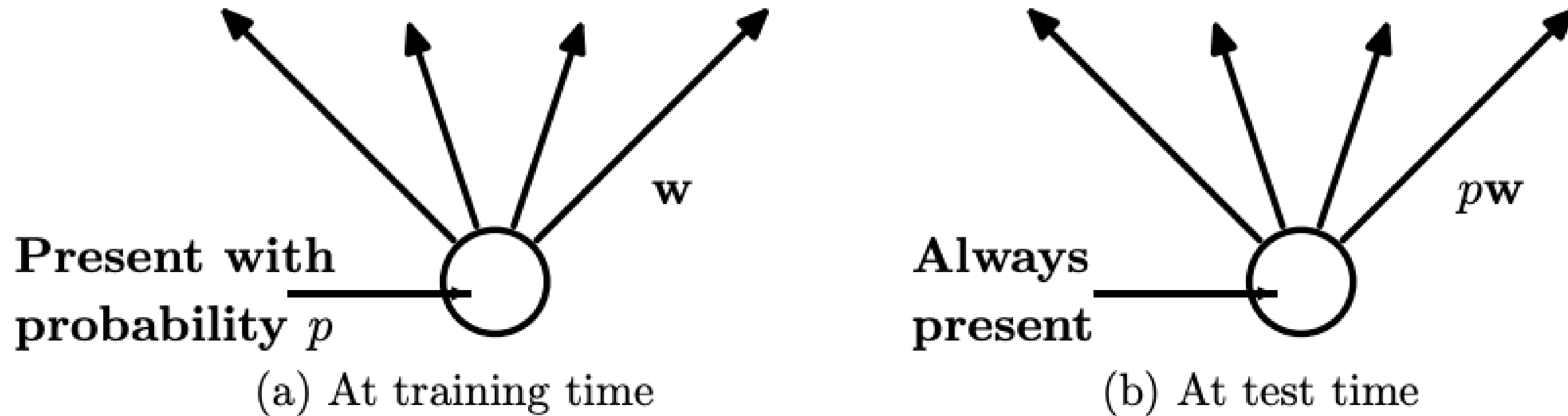
Hidden layer after dropout

# Dropout



(a) At training time       (b) At test time

Figure 2: **Left**: A unit at training time that is present with probability $p$ and is connected to units in the next layer with weights $\mathbf{w}$. **Right**: At test time, the unit is always present and the weights are multiplied by $p$. The output at test time is same as the expected output at training time.
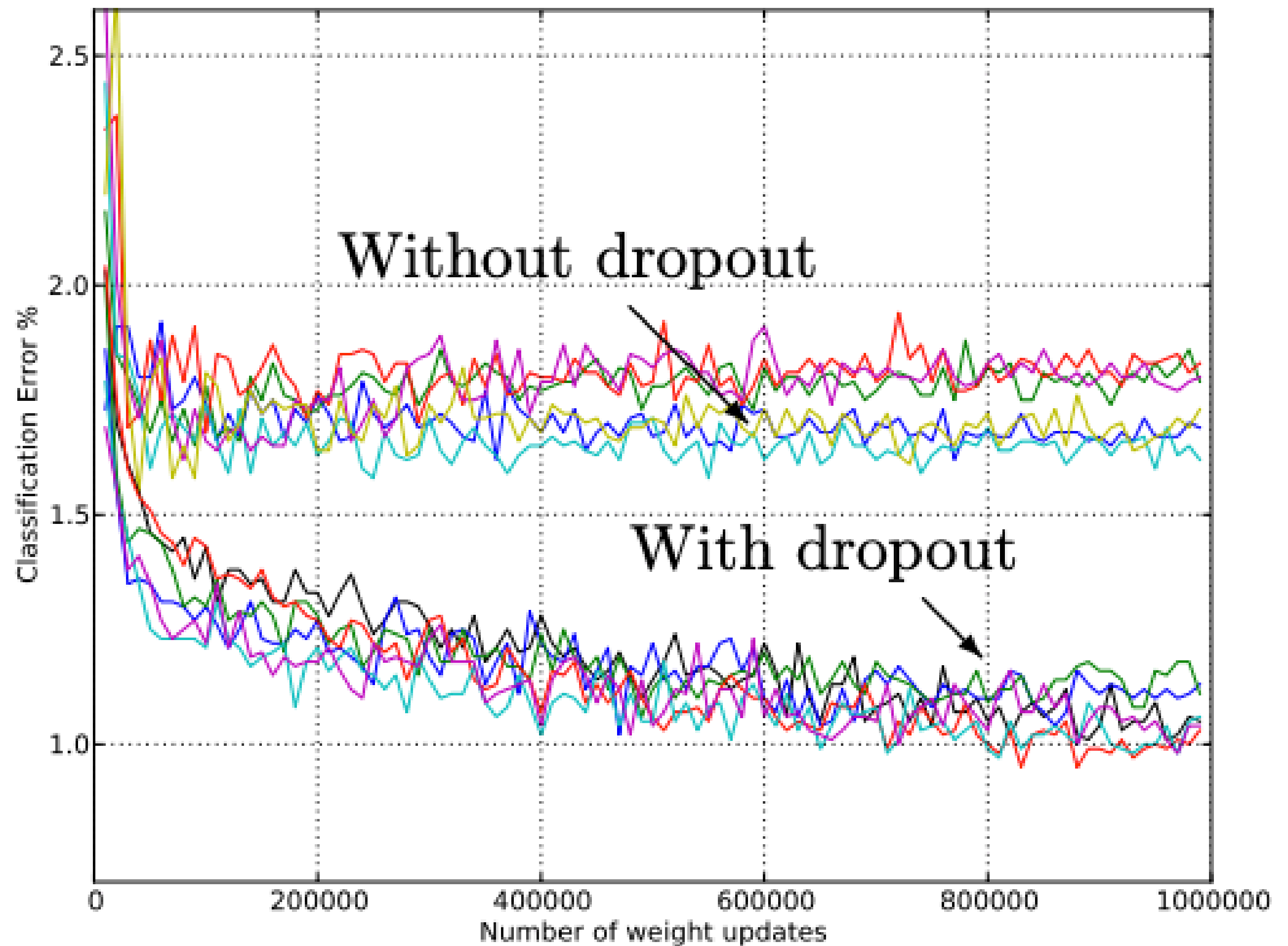
# Dropout

Hinton et al.



Figure 4: Test error for different architectures with and without dropout. The networks have 2 to 4 hidden layers each with 1024 to 2048 units.
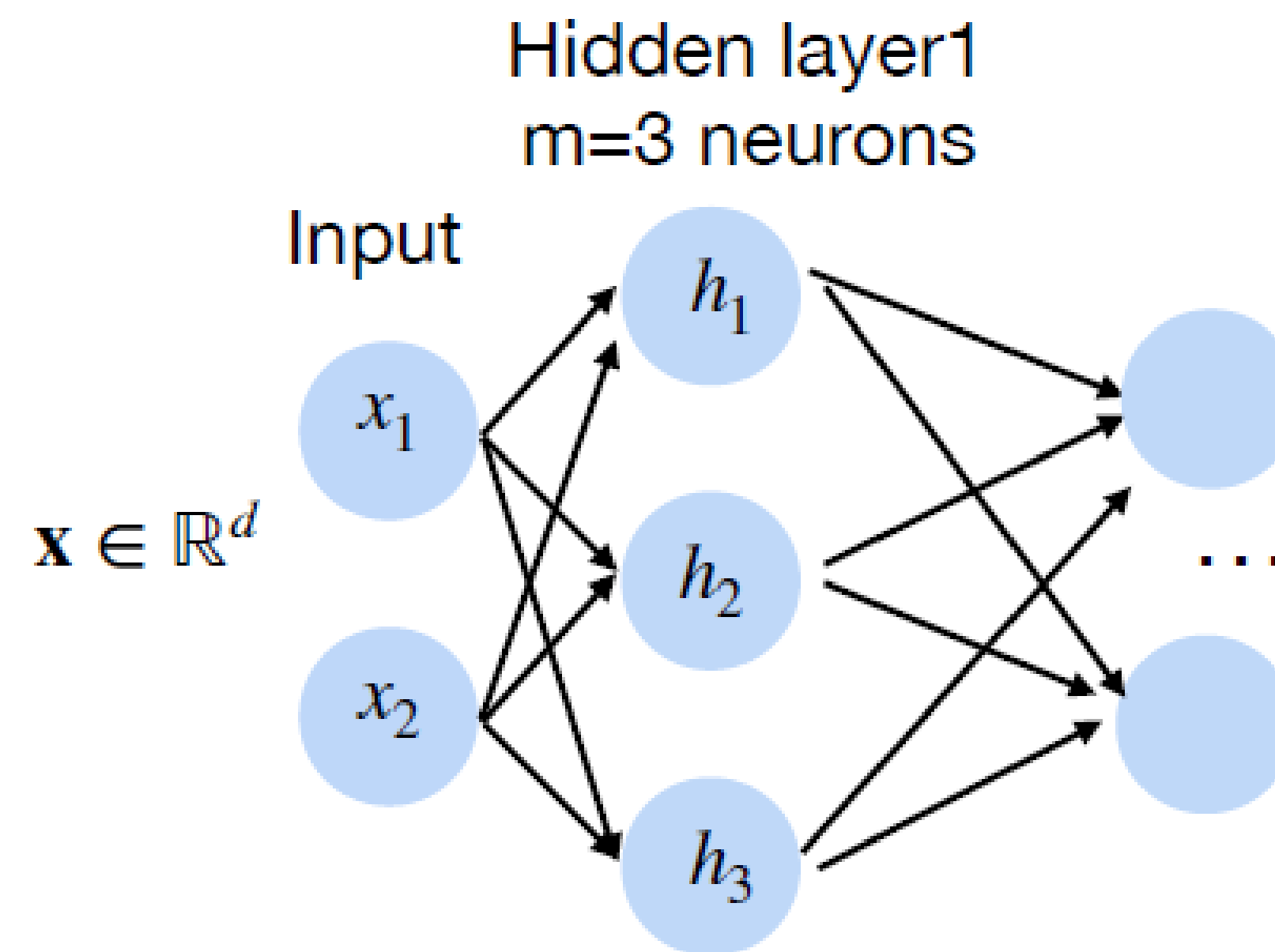
Quiz Break, Q4.1:

In standard dropout regularization, with dropout probability p, each

intermediate activation h is replaced by a random variable h' as: $h' = \begin{cases} 0 \text{ with probability } p \\ ? \text{ otherwise} \end{cases}$.

To make $E[h'] = h$. What is "?" ?

A.  h

B.  h/p

C. h/(1-p)

D. h(1-p)



Input

$\mathbf{x} \in \mathbb{R}^d$

$x_1$
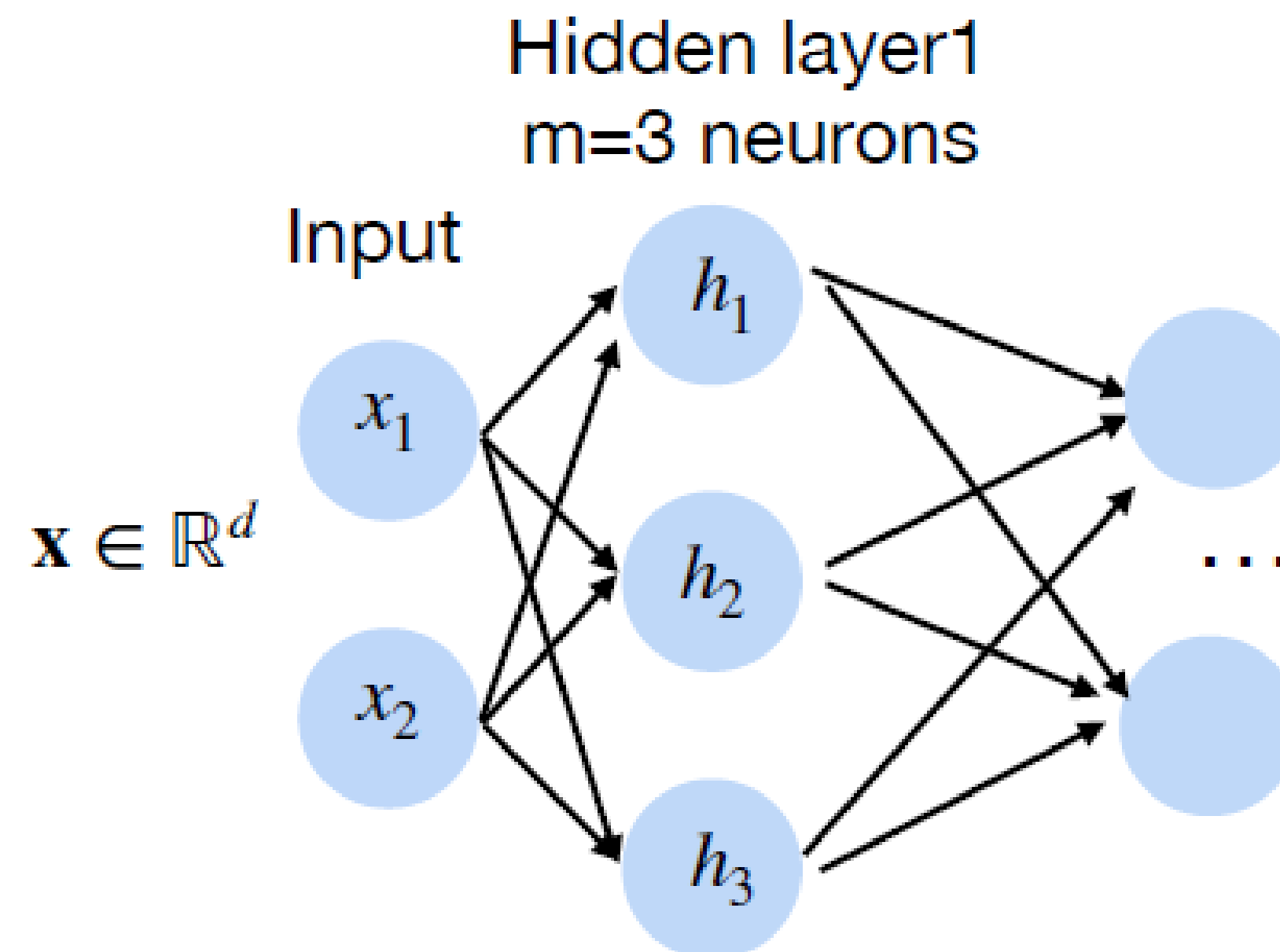
$x_2$

Hidden layer1
m=3 neurons

$h_1$

$h_2$

$h_3$

Quiz Break, Q4.1:

In standard dropout regularization, with dropout probability p, each

intermediate activation h is replaced by a random variable h' as: $h' = \begin{cases} 0 \text{ with probability } p \\ ? \text{ otherwise} \end{cases}$.

To make  $E[h'] = h$. What is "?" ?

A.  h

B.  h/p

C. h/(1-p)

D. h(1-p)

Hidden layer1
m=3 neurons

Input

$h_1$

$x_1$

$\mathbf{x} \in \mathbb{R}^d$

$h_2$

$x_2$

$h_3$

...

# What we've learned today…

- Deep neural networks

  - Computational graph (forward and backward propagation)

- Numerical stability in training

  - Gradient vanishing/exploding

- Generalization and regularization

  - Overfitting, underfitting

  - Weight decay and dropout