CMSC 475/675 Neural Networks
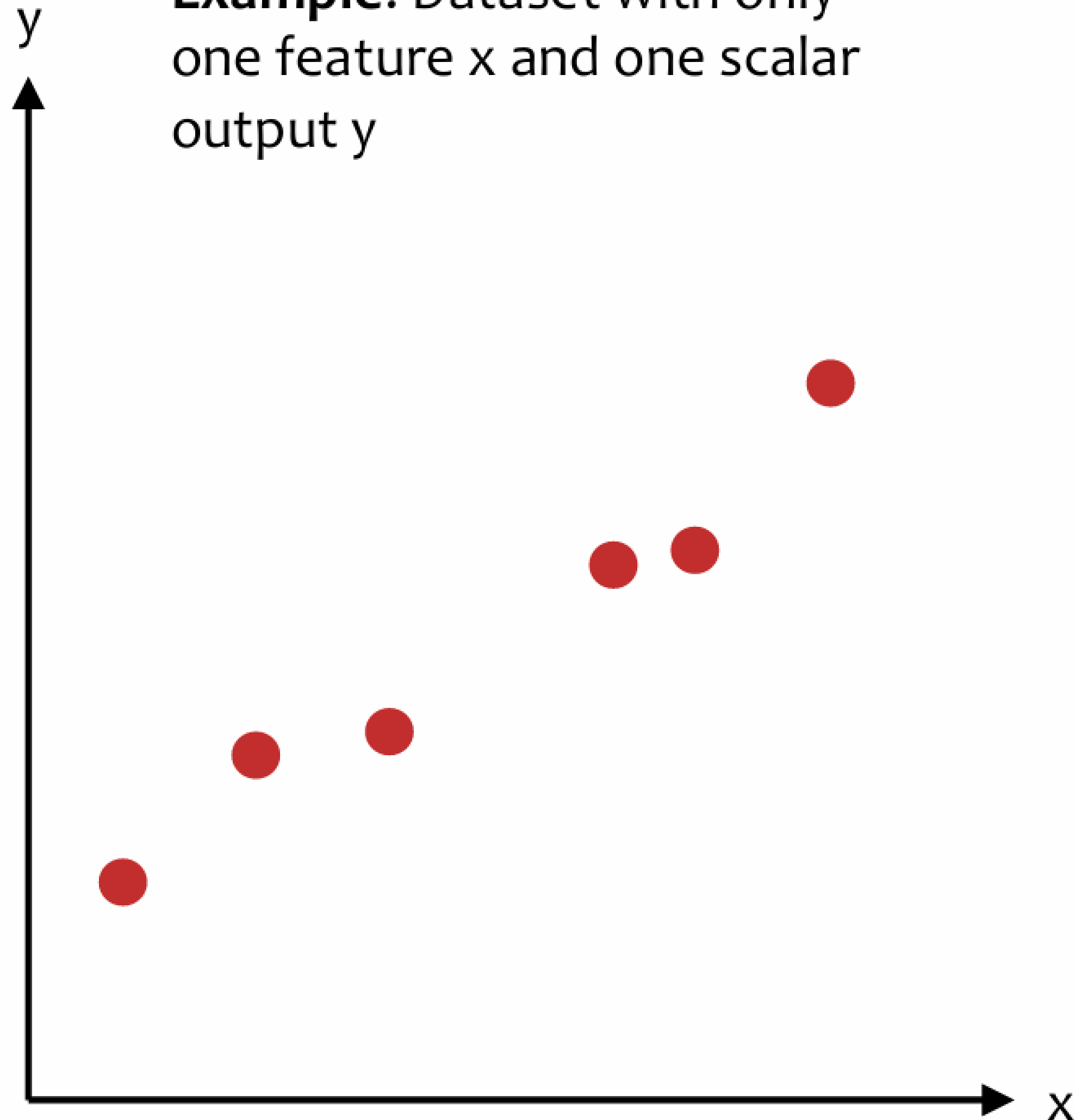
# Lecture 3: Gradient Descent



Some slides adapted from Suren Jayasuriya (ASU), Matt Gormley (CMU)
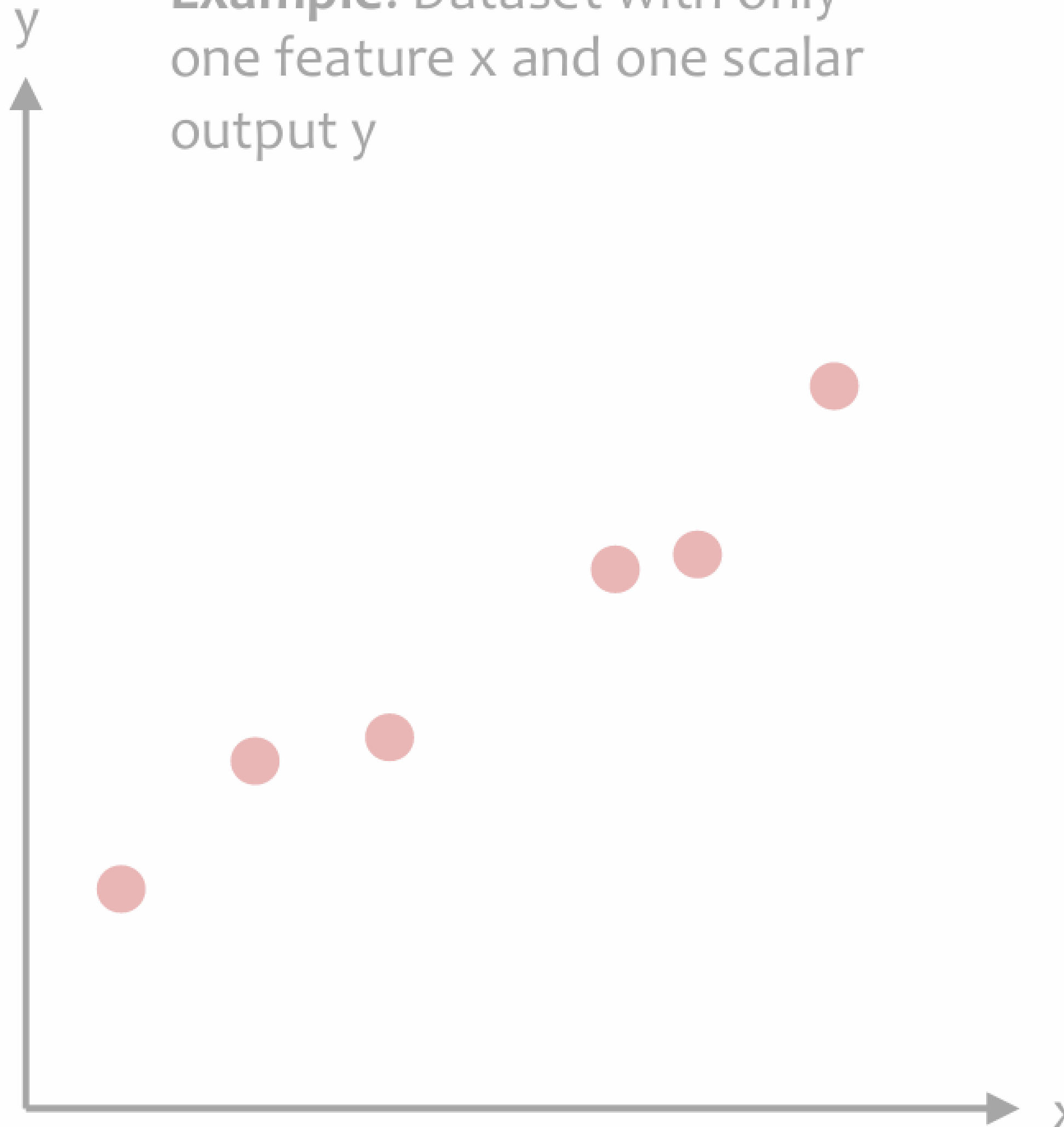
# Recap: Linear Regression

**Example:** Dataset with only one feature x and one scalar output y
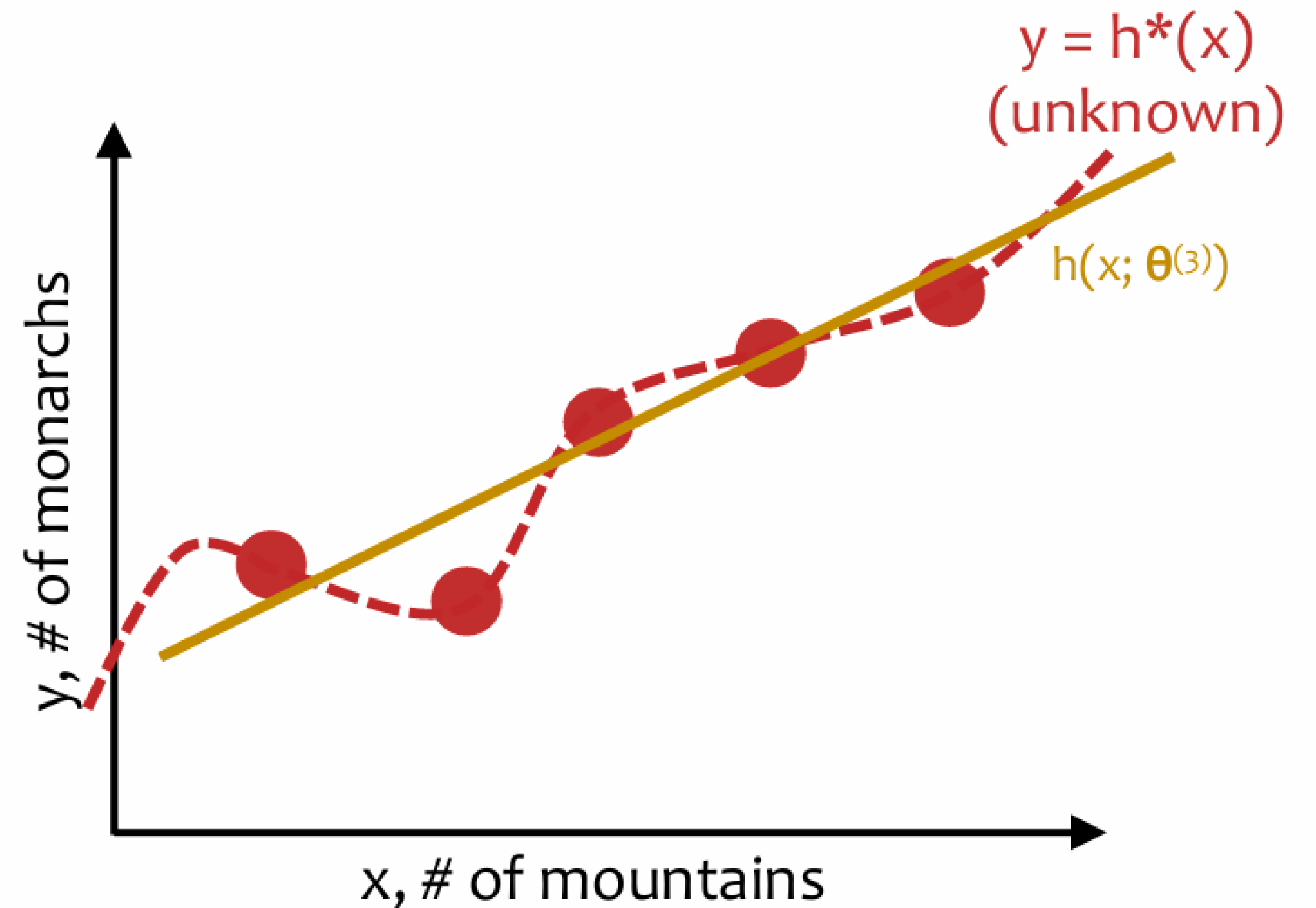
**Q: What is the function that best fits these points?**

# Recap: Linear Regression



**Example:** Dataset with only one feature x and one scalar output y
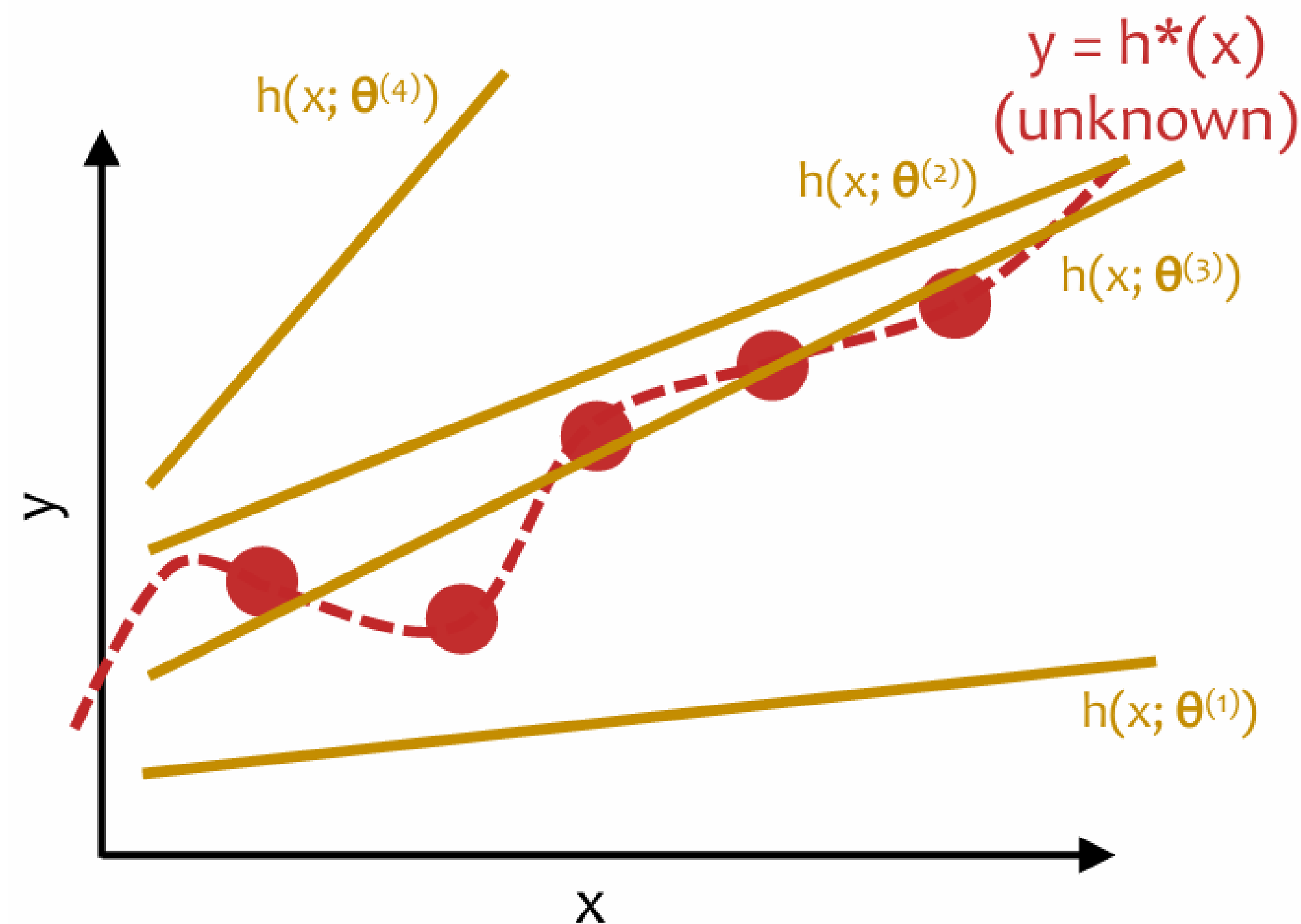
y

**Q: What is the function that best fits these points?**

$y = h^*(x)$
(unknown)

$h(x; \theta^{(3)})$

y, # of monarchs

x, # of mountains

Naïve Idea:

Linear Regression by Random Guessing

!!!

# Linear Regression by Rand. Guessing

**Optimization Method #0: Random Guessing**

1. Pick a random $\theta$

2. Evaluate $J(\theta)$

3. Repeat steps 1 and 2 many times

4. Return $\theta$ that gives smallest $J(\theta)$



$h(x; \theta^{(4)})$

$y = h^*(x)$ (unknown)

$h(x; \theta^{(2)})$
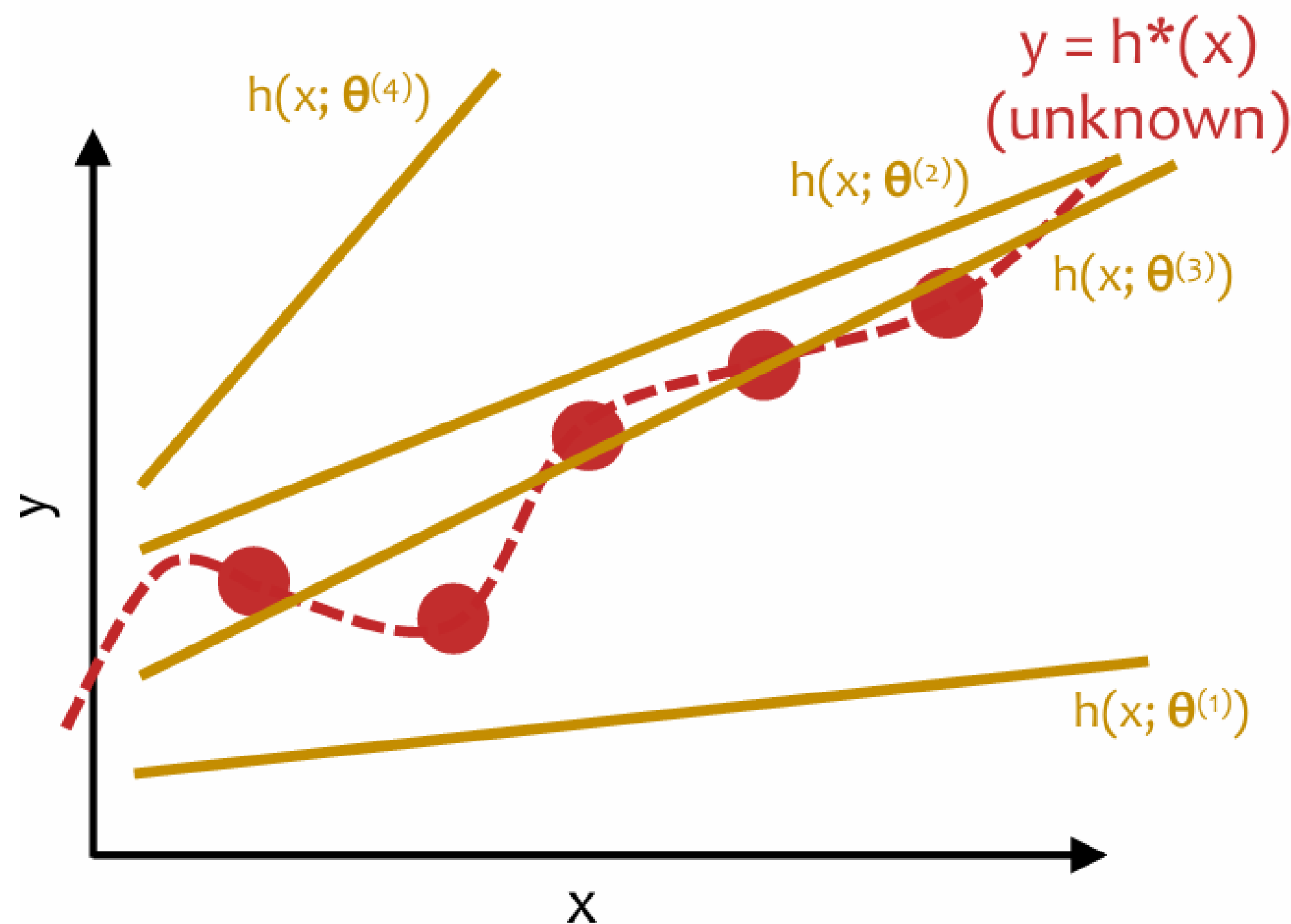
$h(x; \theta^{(3)})$

$h(x; \theta^{(1)})$

y

x

**For Linear Regression:**

- target function $h^*(x)$ is **unknown**

- only have access to $h^*(x)$ through **training examples** $(x^{(i)}, y^{(i)})$

- want $h(x; \theta^{(t)})$ that **best approximates** $h^*(x)$

- **enable generalization** w/inductive bias that restricts hypothesis class to **linear functions**
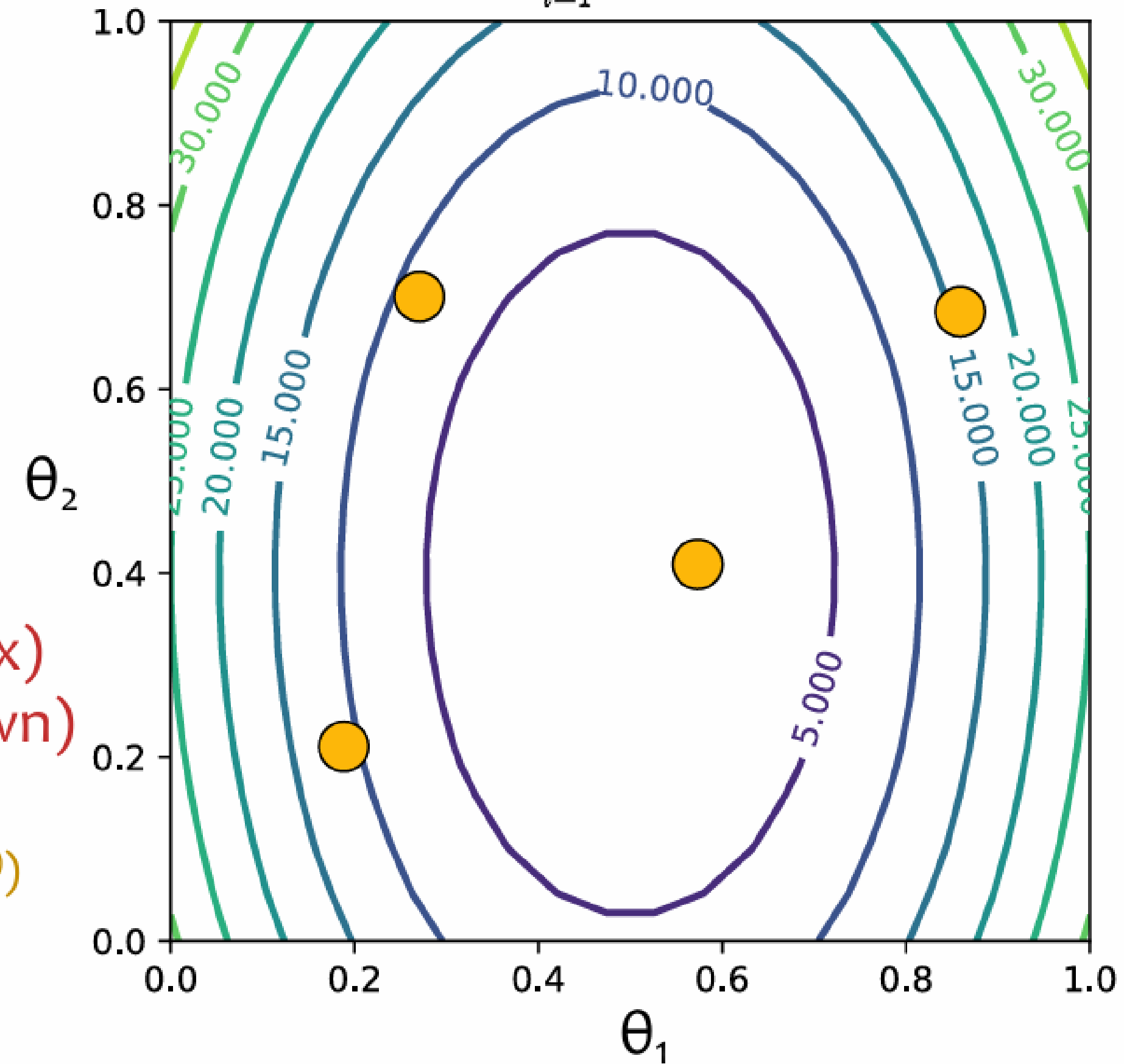
# Linear Regression by Rand. Guessing

$$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^{N} \left( y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)} \right)^2$$

**Optimization Method #0: Random Guessing**

1. Pick a random $\boldsymbol{\theta}$

2. Evaluate $J(\boldsymbol{\theta})$

3. Repeat steps 1 and 2 many times

4. Return $\boldsymbol{\theta}$ that gives smallest $J(\boldsymbol{\theta})$



| t | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|---|---|---|
| 1 | 0.2 | 0.2 | 10.4 |
| 2 | 0.3 | 0.7 | 7.2 |
| 3 | 0.6 | 0.4 | 1.0 |
| 4 | 0.9 | 0.7 | 16.2 |

Better Idea:
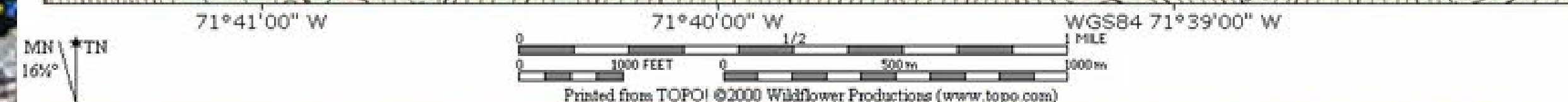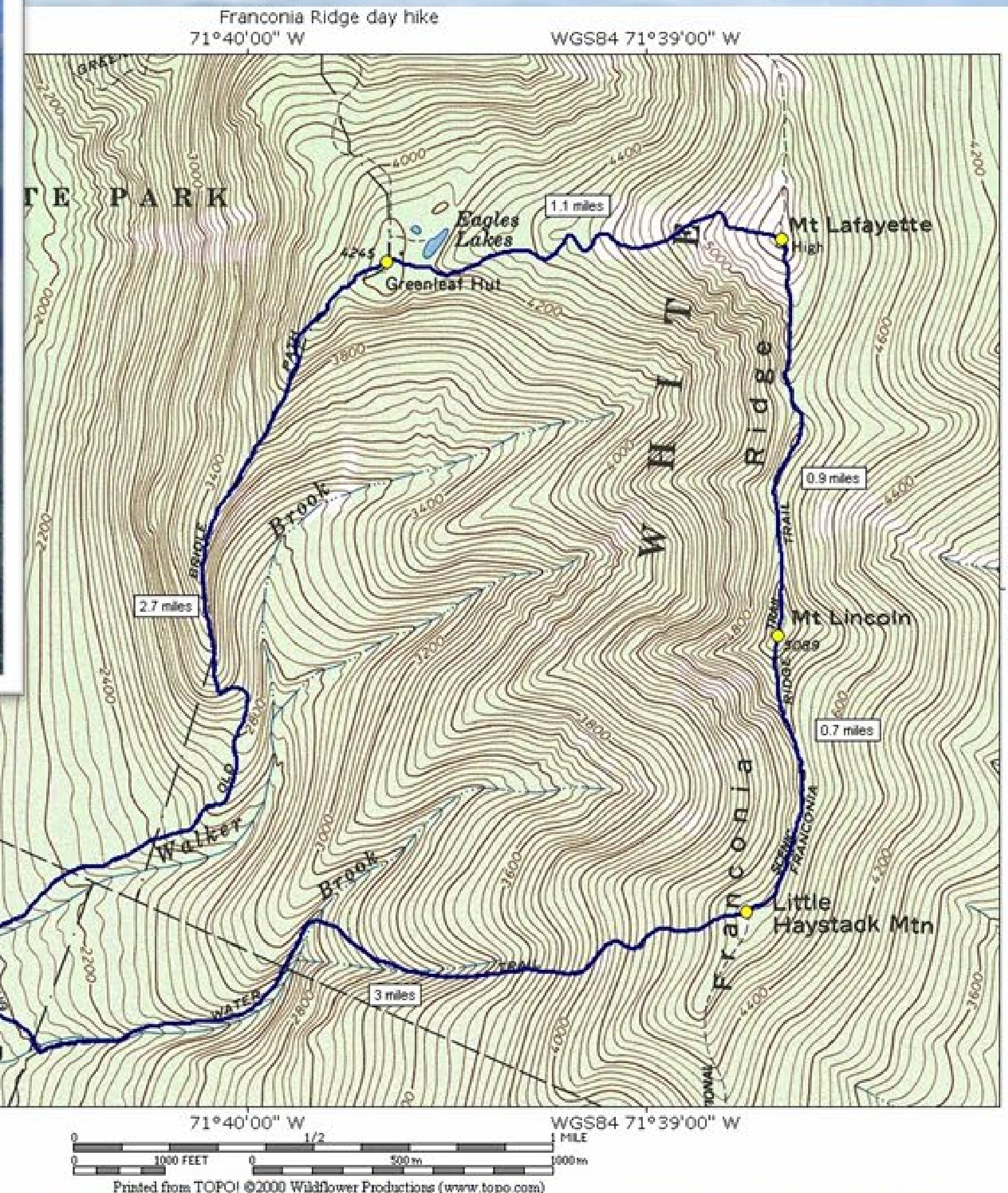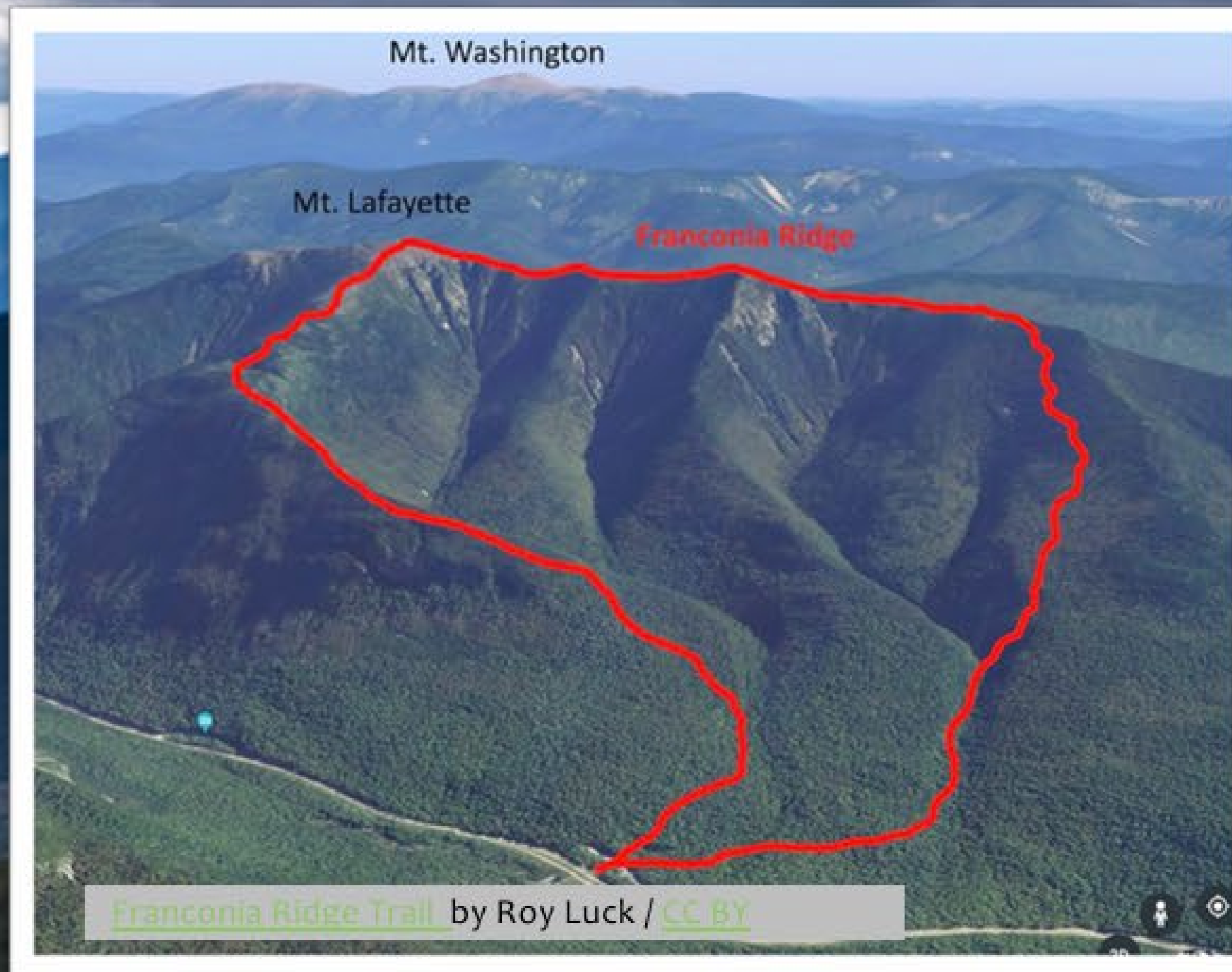
An optimization algorithm called

"Gradient Descent"

# Recall: Gradient of a vector

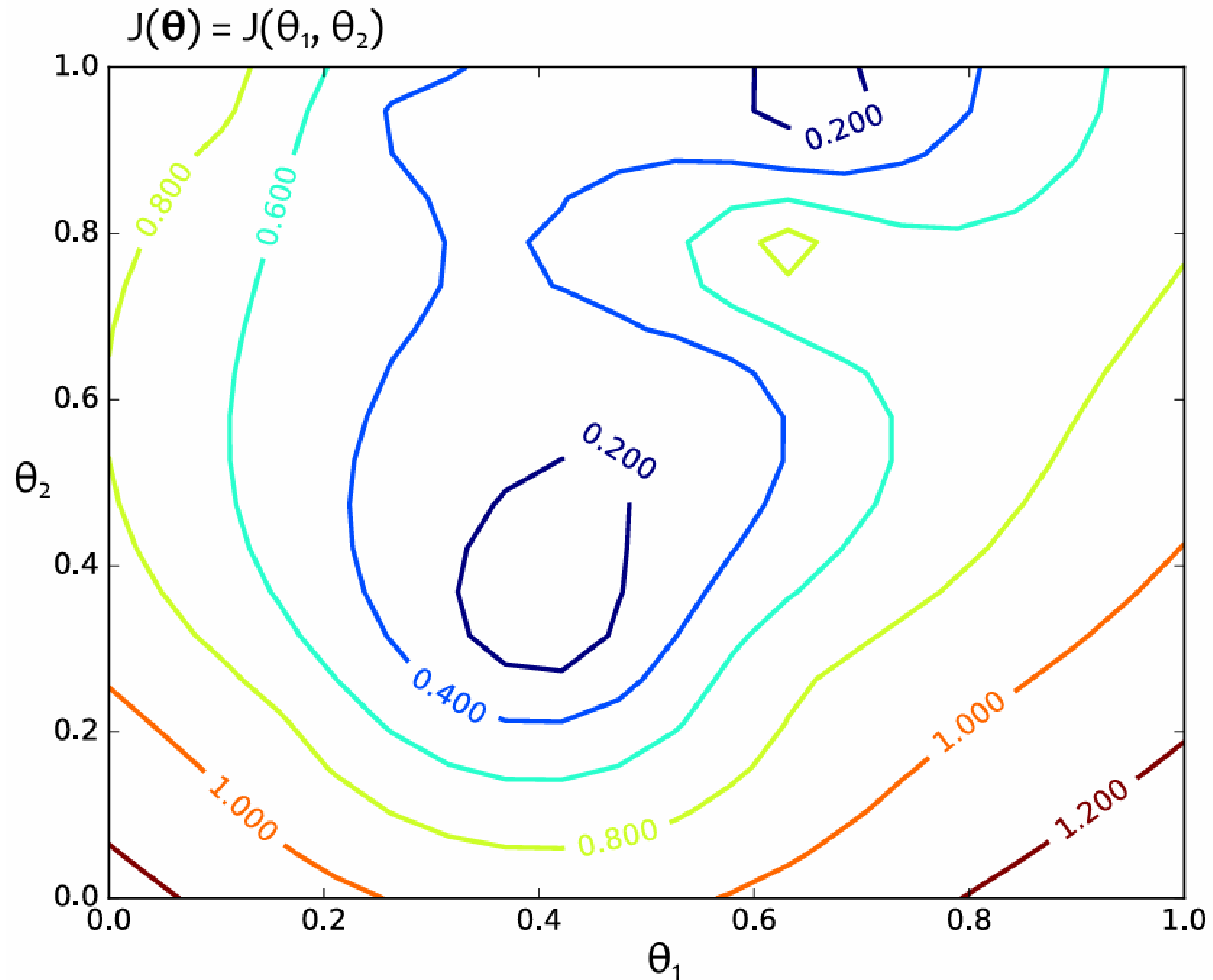**Def:** The gradient of $J : \mathbb{R}^M \to \mathbb{R}$ is

$$\nabla J(\boldsymbol{\theta}) = \begin{bmatrix} \dfrac{\partial J(\boldsymbol{\theta})}{\partial \theta_1} \\ \dfrac{\partial J(\boldsymbol{\theta})}{\partial \theta_2} \\ \vdots \\ \dfrac{\partial J(\boldsymbol{\theta})}{\partial \theta_M} \end{bmatrix}$$

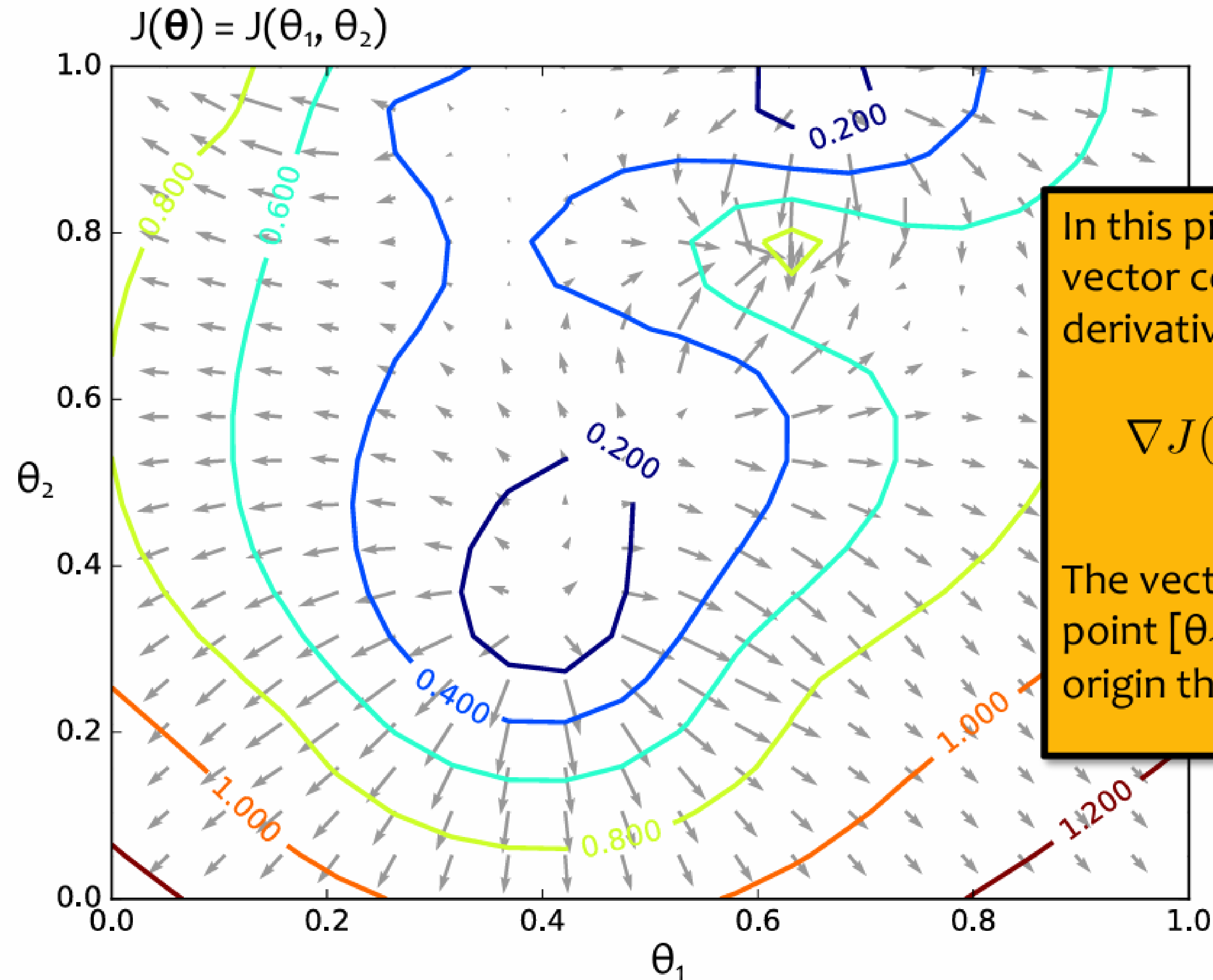Each entry is a first-order partial derivative

# Topographical Maps



Mt. Washington

Mt. Lafayette

Franconia Ridge

Franconia Ridge Trail by Roy Luck / CC BY

Franconia Ridge day hike
71°40'00" W          WGS84 71°39'00" W

Mt Lafayette
High

Eagles Lakes

Greenleaf Hut

1.1 miles

0.9 miles

2.7 miles

Mt. Lincoln

0.7 miles

Brook

Walker

Brook

Little Haystack Mtn

3 miles

Lafayette Campground

BM 1770
Gravel Pit

0.2 miles

Pemigewasset

Dry

71°41'00" W          71°40'00" W          WGS84 71°39'00" W

MN   TN
16%

1000 FEET        500 m        1000 m

Printed from TOPO! ©2000 Wildflower Productions (www.topo.com)

# Gradients



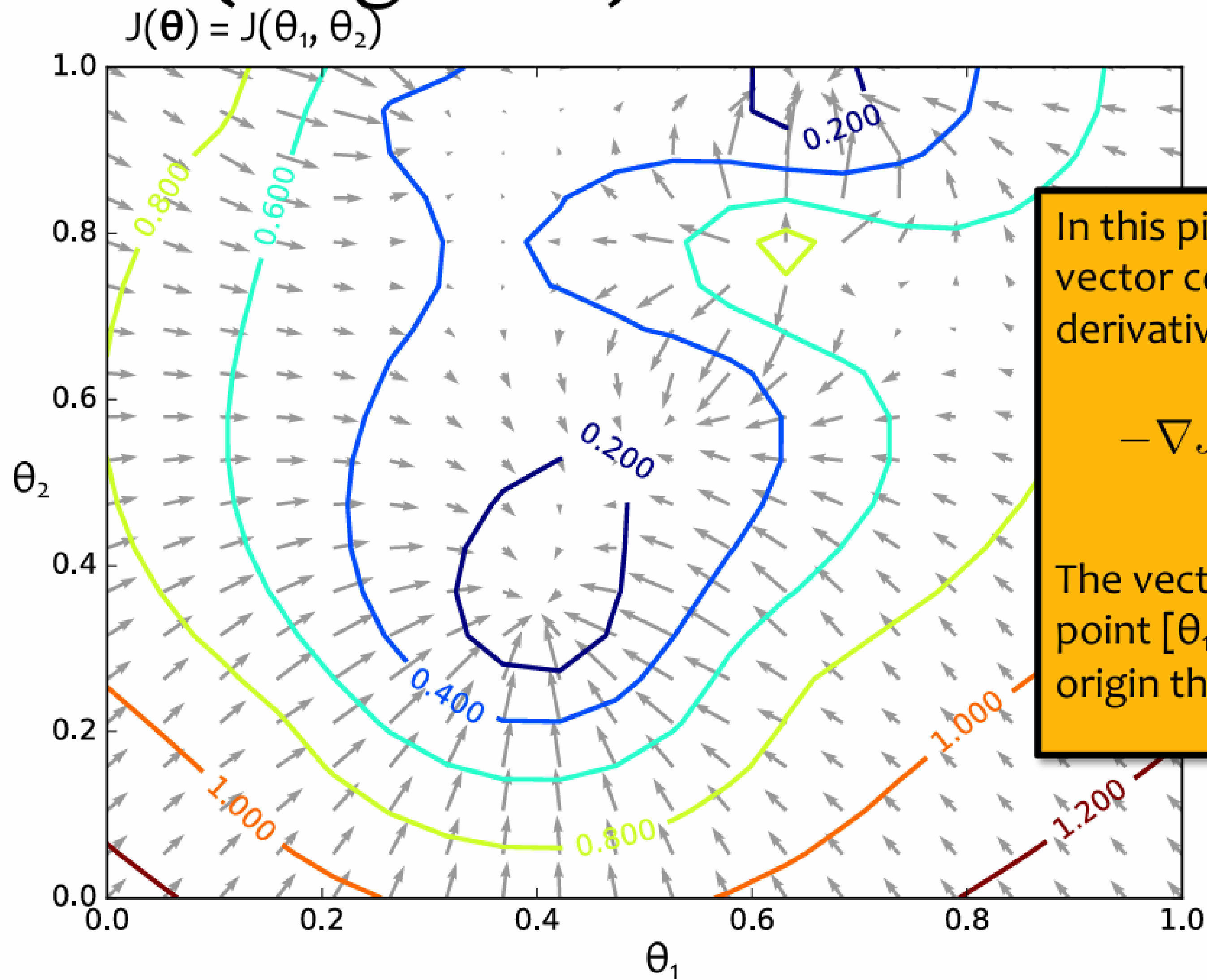$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2)$

In this picture, each arrow is a 2D vector consisting of two partial derivatives.
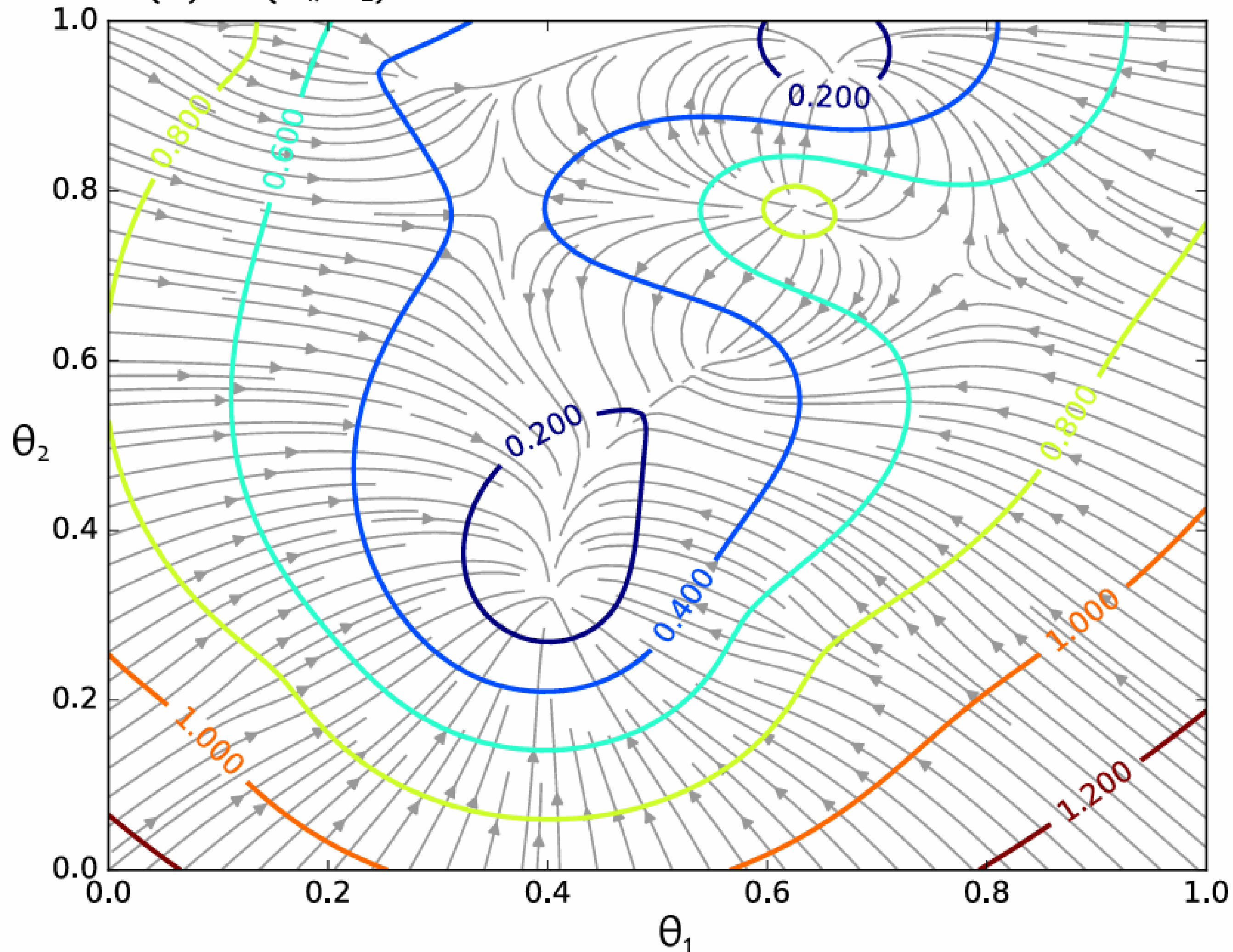
$$\nabla J(\theta_1, \theta_2) = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\[1em] \frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

The vector is evaluated at the point $[\theta_1, \theta_2]^\mathsf{T}$ and plotted with its origin there as well.

These are the **gradients** that
Gradient **Ascent** would follow.

# (Negative) Gradients

$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2)$



In this picture, each arrow is a 2D vector consisting of two partial derivatives.

$$-\nabla J(\theta_1, \theta_2) = \begin{bmatrix} -\frac{\partial J}{\partial \theta_1} \\[2mm] -\frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

The vector is evaluated at the point $[\theta_1, \theta_2]^{\mathsf{T}}$ and plotted with its origin there as well.

These are the **negative** gradients that
Gradient **Descent** would follow.
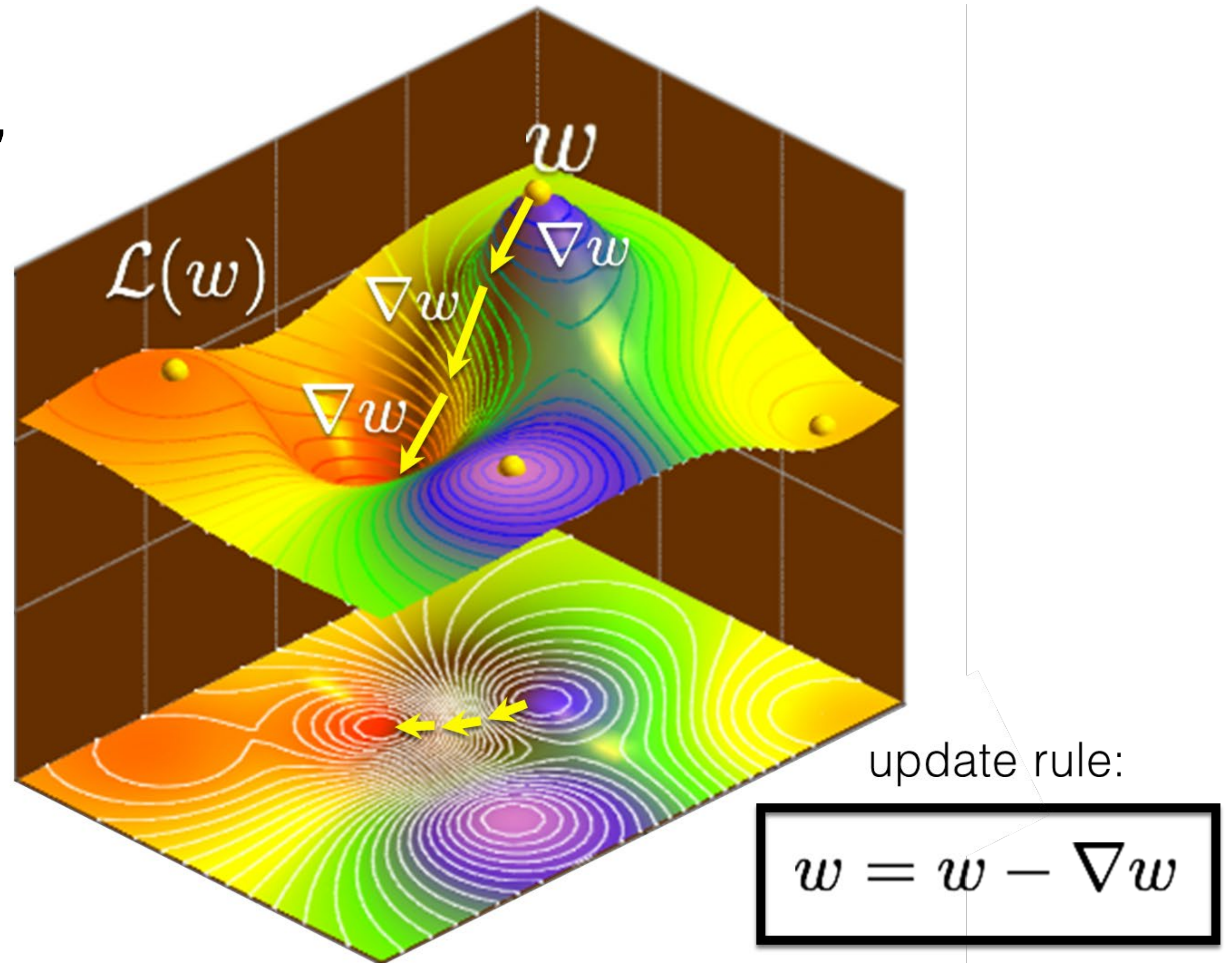
# (Negative) Gradient *Paths*

$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2)$



Shown are the **paths** that Gradient Descent
would follow if it were making **infinitesimally
small steps.**

# Gradient Descent
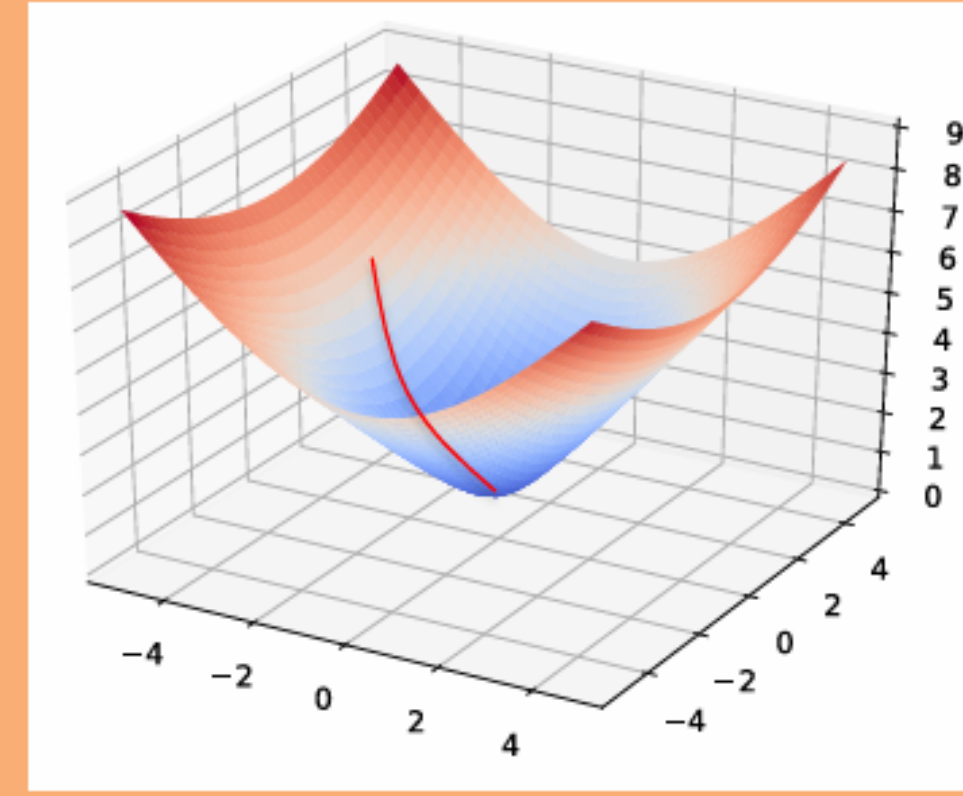
- Go down the path of "steepest descent"



update rule:

$$w = w - \nabla w$$

# Gradient Descent

- Go down the path of steepest descent



**Algorithm 1** Gradient Descent

1: **procedure** $\text{GD}(\mathcal{D}, \boldsymbol{\theta}^{(0)})$
2:      $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$
3:      **while** not converged **do**
4:          $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
5:      **return** $\boldsymbol{\theta}$

In order to apply GD to Linear Regression all we need is the **gradient** of the objective function (i.e. vector of partial derivatives).
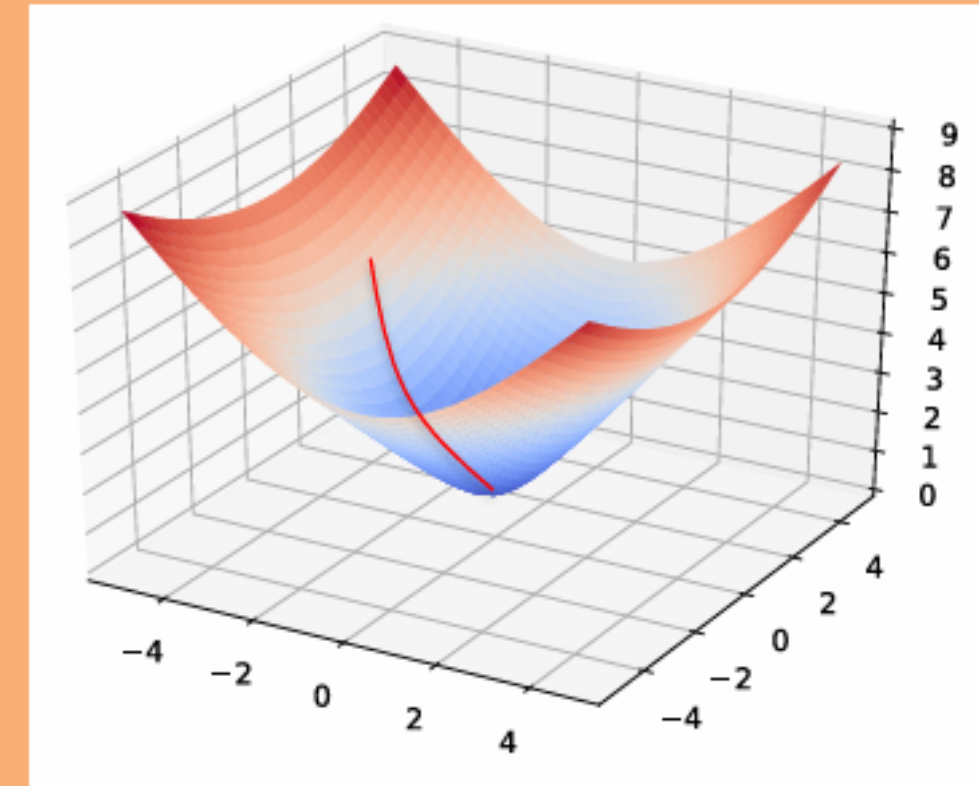
$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \begin{bmatrix} \frac{d}{d\theta_1} J(\boldsymbol{\theta}) \\ \frac{d}{d\theta_2} J(\boldsymbol{\theta}) \\ \vdots \\ \frac{d}{d\theta_M} J(\boldsymbol{\theta}) \end{bmatrix}$$

# Gradient Descent

- Go down the path of steepest descent

**Algorithm 1** Gradient Descent

1: **procedure** $\text{GD}(\mathcal{D}, \boldsymbol{\theta}^{(0)})$
2: $\quad \boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{(0)}$
3: $\quad$ **while** not converged **do**
4: $\quad\quad \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
5: $\quad$ **return** $\boldsymbol{\theta}$

There are many possible ways to detect **convergence**. For example, we could check whether the L2 norm of the gradient is below some small tolerance.

$$||\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})||_2 \leq \epsilon$$

Alternatively we could check that the reduction in the objective function from one iteration to the next is small.

# Linear Regression by Gradient Desc.

**Optimization Method #1:
Gradient Descent**

1. Pick a random **θ**

2. Repeat:
   a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
   b. Step opposite gradient

3. Return **θ** that gives smallest $J(\boldsymbol{\theta})$

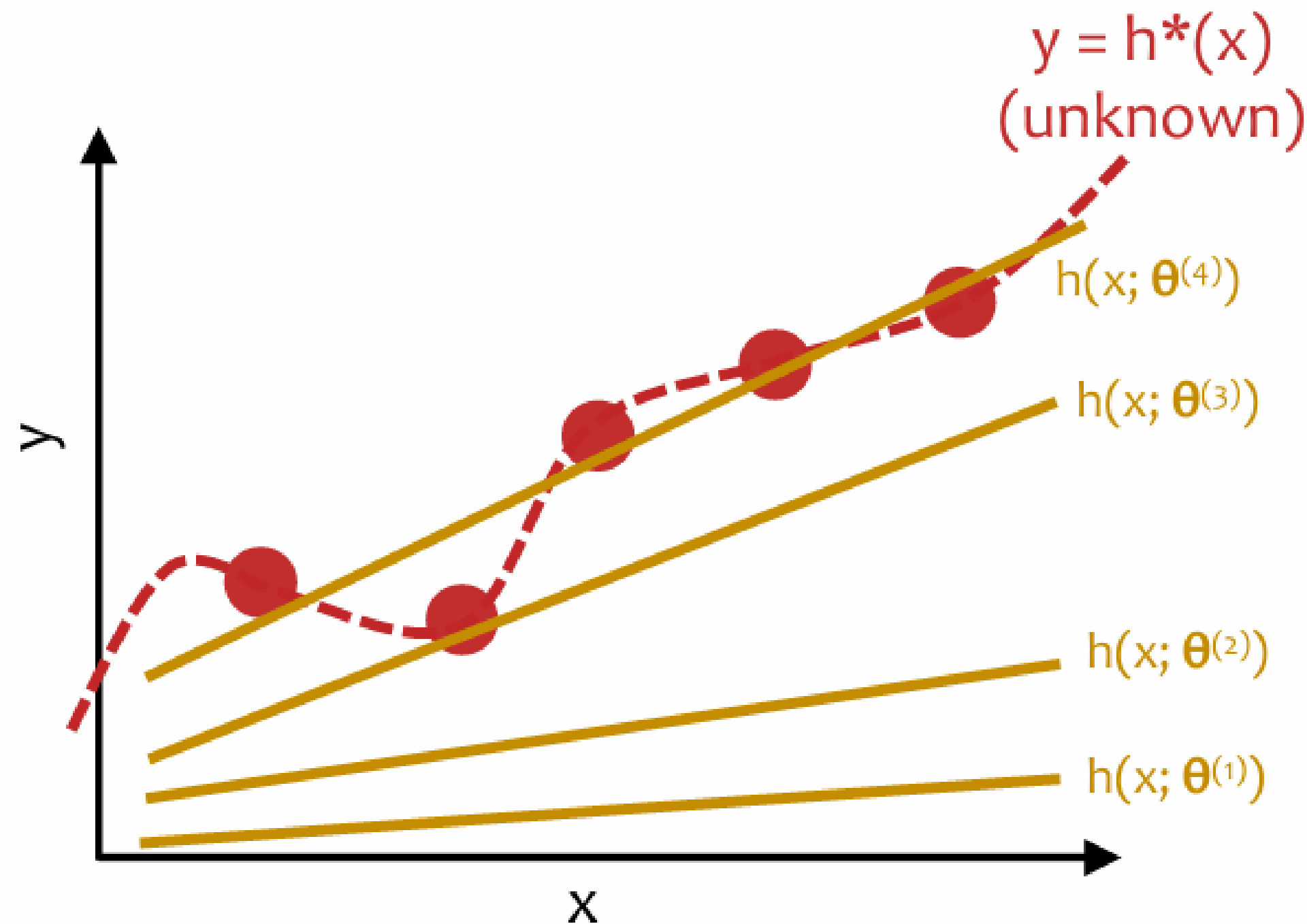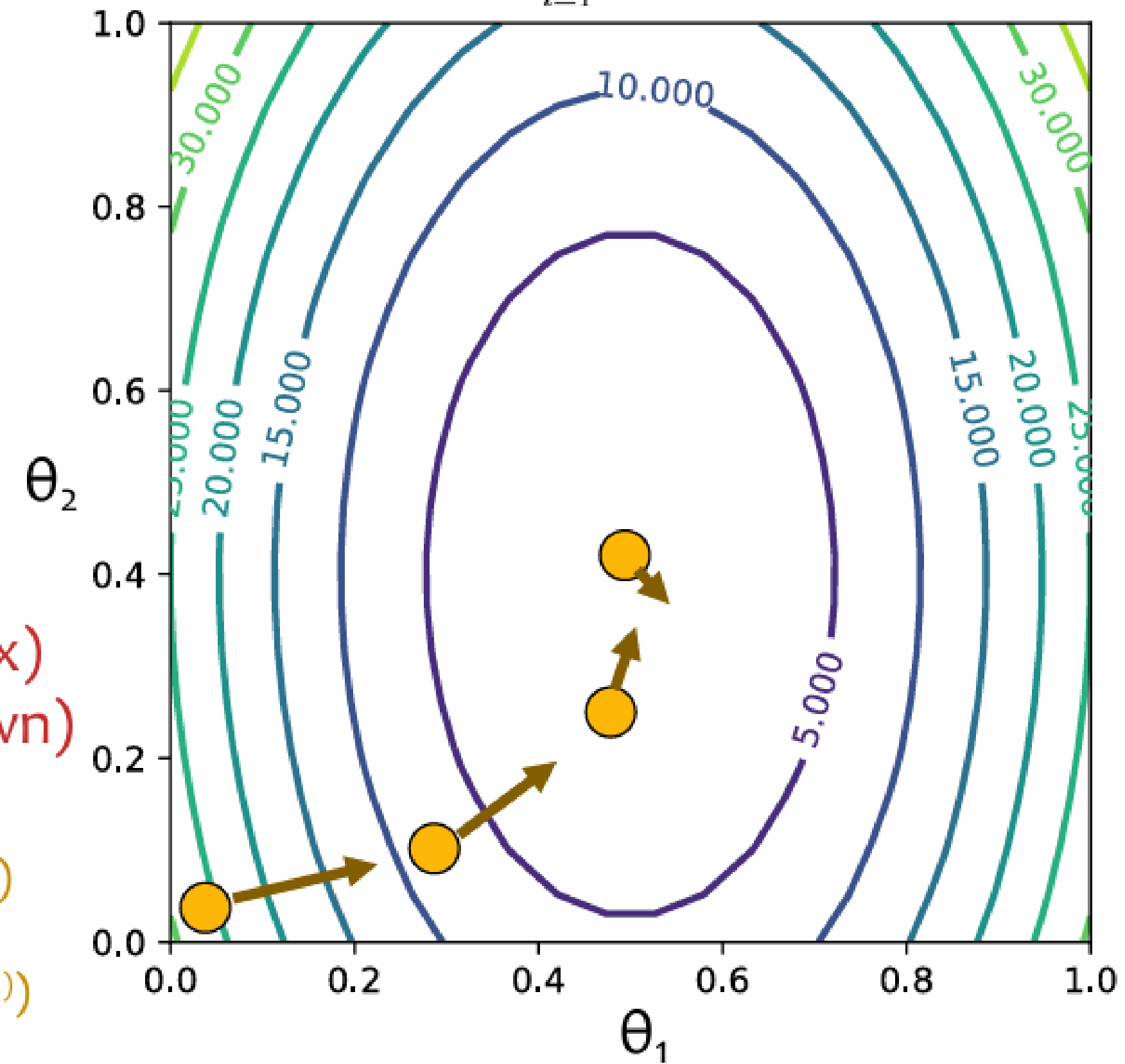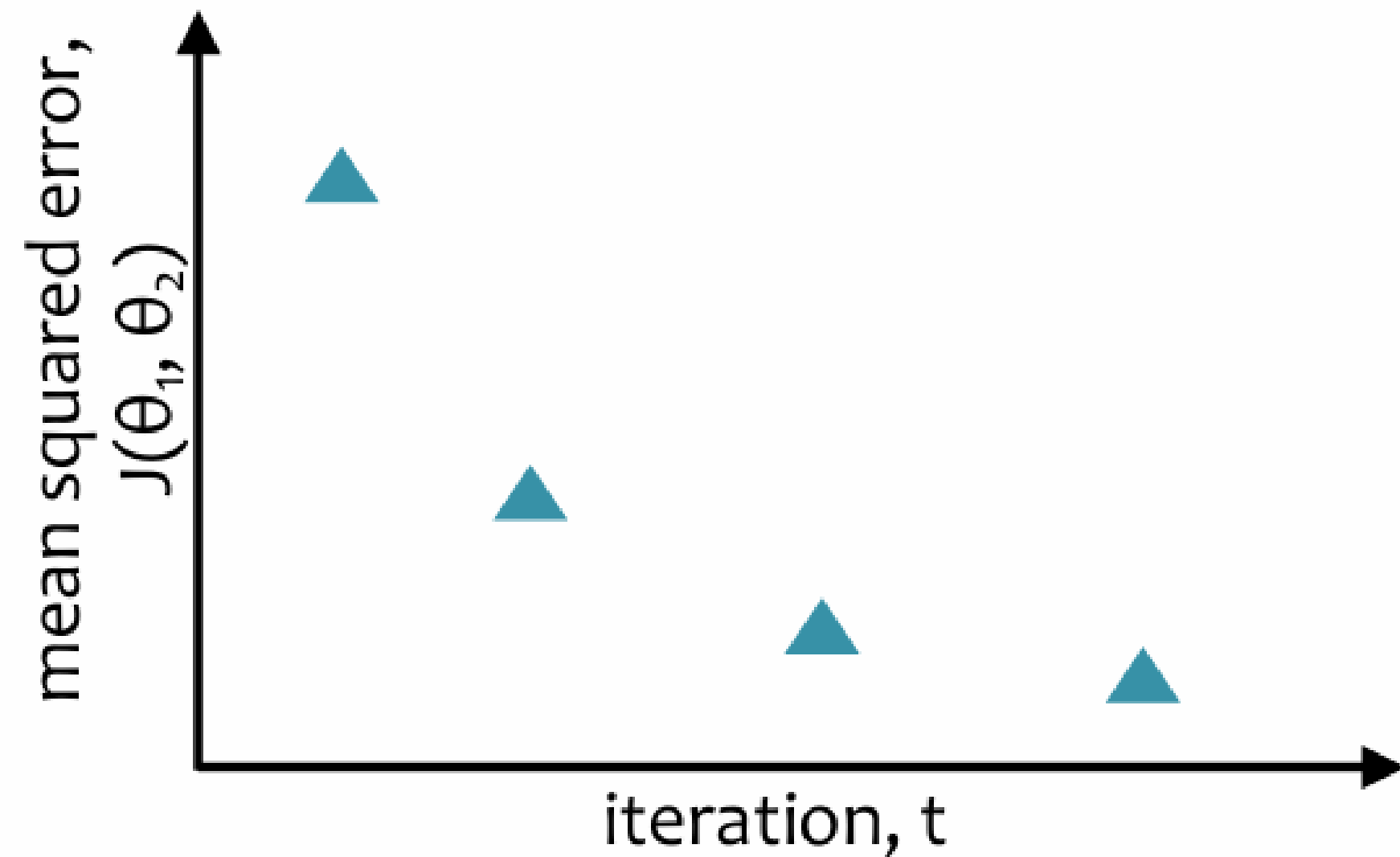$$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2) = \frac{1}{N} \sum_{i=1}^{N} \left( y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)} \right)^2$$



| t | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|------|------|------|
| 1 | 0.01 | 0.02 | 25.2 |
| 2 | 0.30 | 0.12 | 8.7 |
| 3 | 0.51 | 0.30 | 1.5 |
| 4 | 0.59 | 0.43 | 0.2 |

# Linear Regression by Gradient Desc.

$$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2) = \frac{1}{N}\sum_{i=1}^{N}\left(y^{(i)} - \boldsymbol{\theta}^T\mathbf{x}^{(i)}\right)^2$$

**Optimization Method #1:**
**Gradient Descent**

1. Pick a random $\boldsymbol{\theta}$

2. Repeat:
   a. Evaluate gradient $\nabla J(\boldsymbol{\theta})$
   b. Step opposite gradient

3. Return $\boldsymbol{\theta}$ that gives smallest $J(\boldsymbol{\theta})$

y = h*(x)
(unknown)

h(x; $\boldsymbol{\theta}^{(4)}$)

h(x; $\boldsymbol{\theta}^{(3)}$)

h(x; $\boldsymbol{\theta}^{(2)}$)

h(x; $\boldsymbol{\theta}^{(1)}$)

| t | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|---|---|---|
| 1 | 0.01 | 0.02 | 25.2 |
| 2 | 0.30 | 0.12 | 8.7 |
| 3 | 0.51 | 0.30 | 1.5 |
| 4 | 0.59 | 0.43 | 0.2 |

# Linear Regression by Gradient Desc.



$$J(\boldsymbol{\theta}) = J(\theta_1, \theta_2) = \frac{1}{N}\sum_{i=1}^{N}\left(y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)}\right)^2$$

| t | $\theta_1$ | $\theta_2$ | $J(\theta_1, \theta_2)$ |
|---|---|---|---|
| 1 | 0.01 | 0.02 | 25.2 |
| 2 | 0.30 | 0.12 | 8.7 |
| 3 | 0.51 | 0.30 | 1.5 |
| 4 | 0.59 | 0.43 | 0.2 |

**Gradient Descent**

For each example sample $\{x_i, y_i\}$

1. Predict

    a. Forward pass $\qquad \hat{y} = f_{\mathrm{MLP}}(x_i; \theta)$

    b. Compute Loss $\qquad \mathcal{L}_i$

2. Update

    a. Back Propagation $\qquad \dfrac{\partial \mathcal{L}}{\partial \theta}$

vector of parameter partial derivatives

    b. Gradient update $\qquad \theta \leftarrow \theta - \eta \dfrac{\partial \mathcal{L}}{\partial \theta}$

vector of parameter update equations

# Gradient Descent

# GD for LR: Python Step-by-Step Example

- https://colab.research.google.com/drive/17dK6cynECzk2ObyCqDk5gKcUyN1kMjSR?usp=sharing

# Learning rates



original $\theta$

negative gradient direction

$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

Step size: learning rate
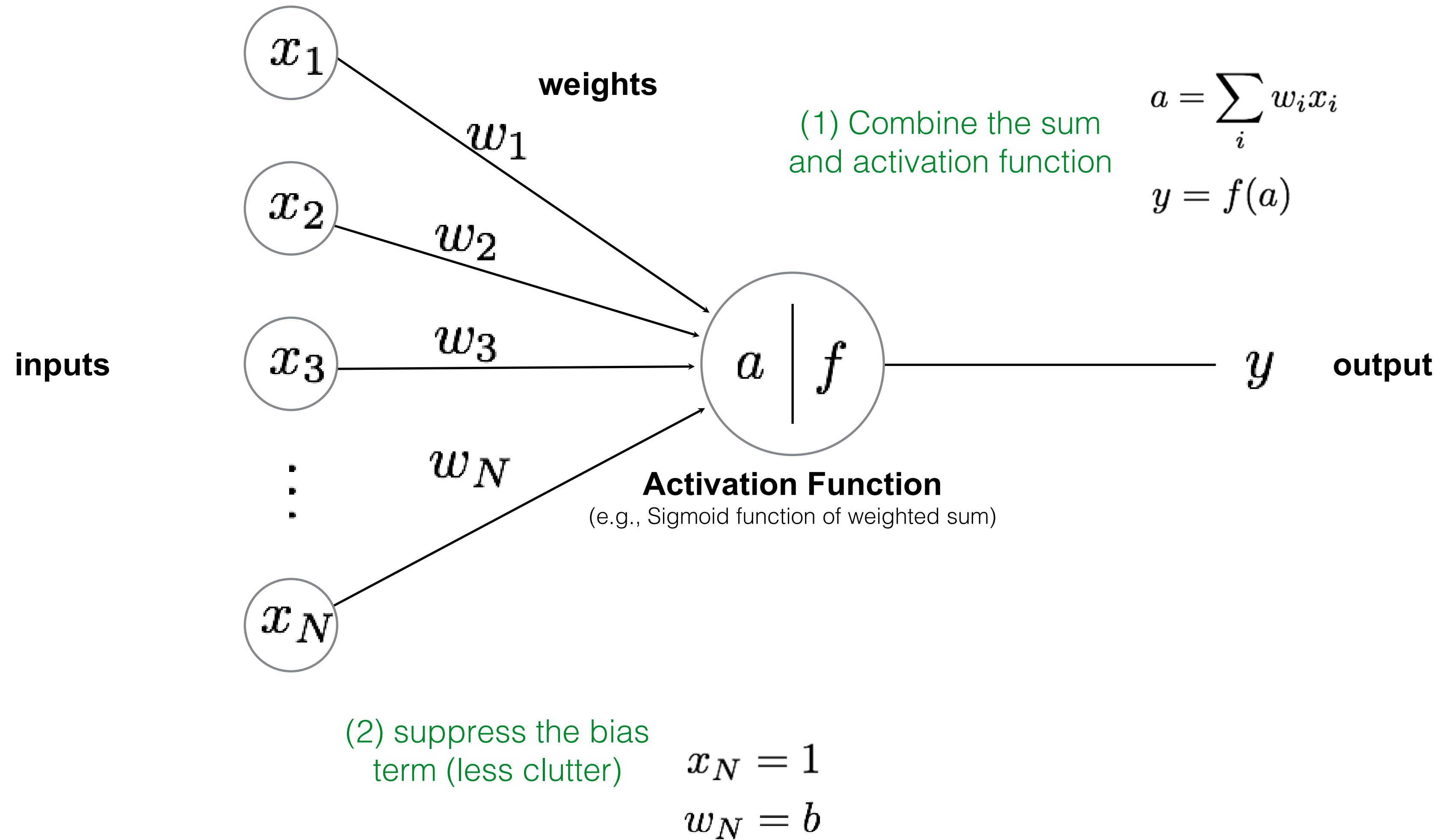Too big: will miss the minimum
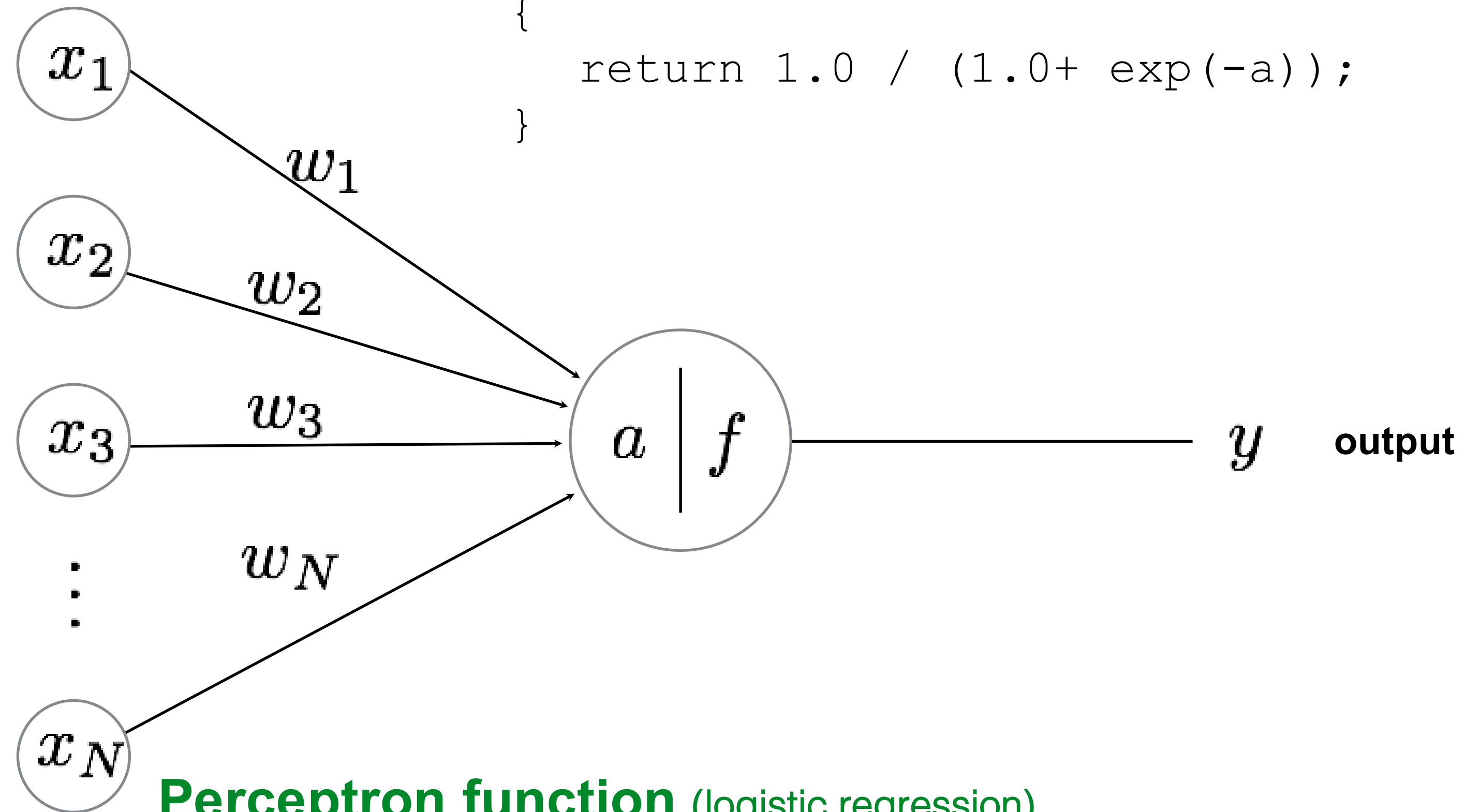Too small: slow convergence

# The Perceptron

# The Perceptron

# Another way to draw it…



inputs

weights

(1) Combine the sum
and activation function

$$a = \sum_i w_i x_i$$

$$y = f(a)$$

**Activation Function**
(e.g., Sigmoid function of weighted sum)

output

(2) suppress the bias
term (less clutter)

$$x_N = 1$$

$$w_N = b$$

# Programming the 'forward pass'

**Activation function** (sigmoid, logistic function)

```
float f(float a)
{
    return 1.0 / (1.0+ exp(-a));
}
```

$x_1$

$x_2$

$x_3$

$\vdots$

$x_N$

$w_1$

$w_2$

$w_3$

$w_N$

$a \mid f$

$y$ **output**

**Perceptron function** (logistic regression)

```
float perceptron(vector<float> x, vector<float> w)
{
    float a  = dot(x,w);
    return f(a);
}
```

# Neural networks

Connect a bunch of perceptrons together …

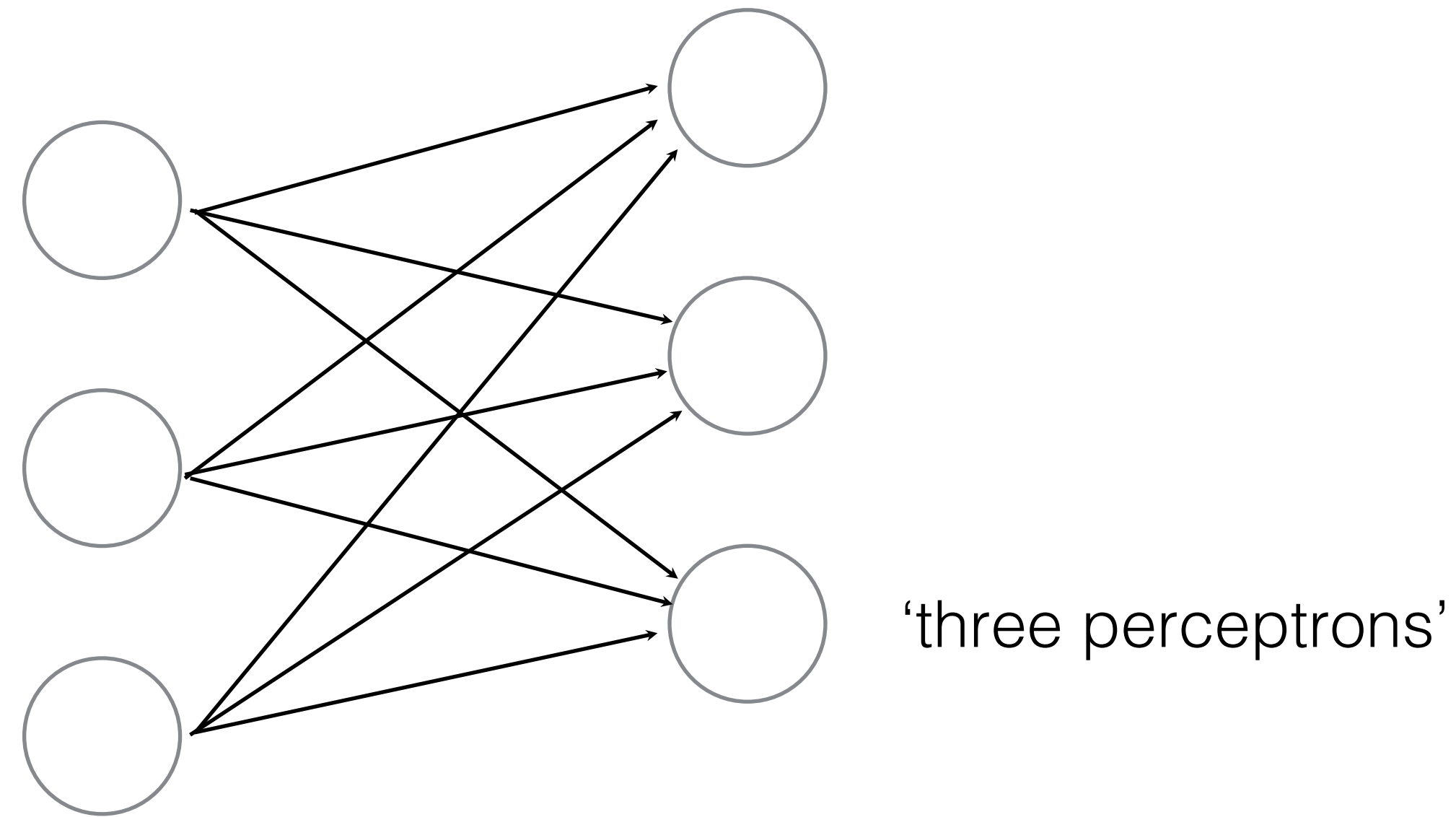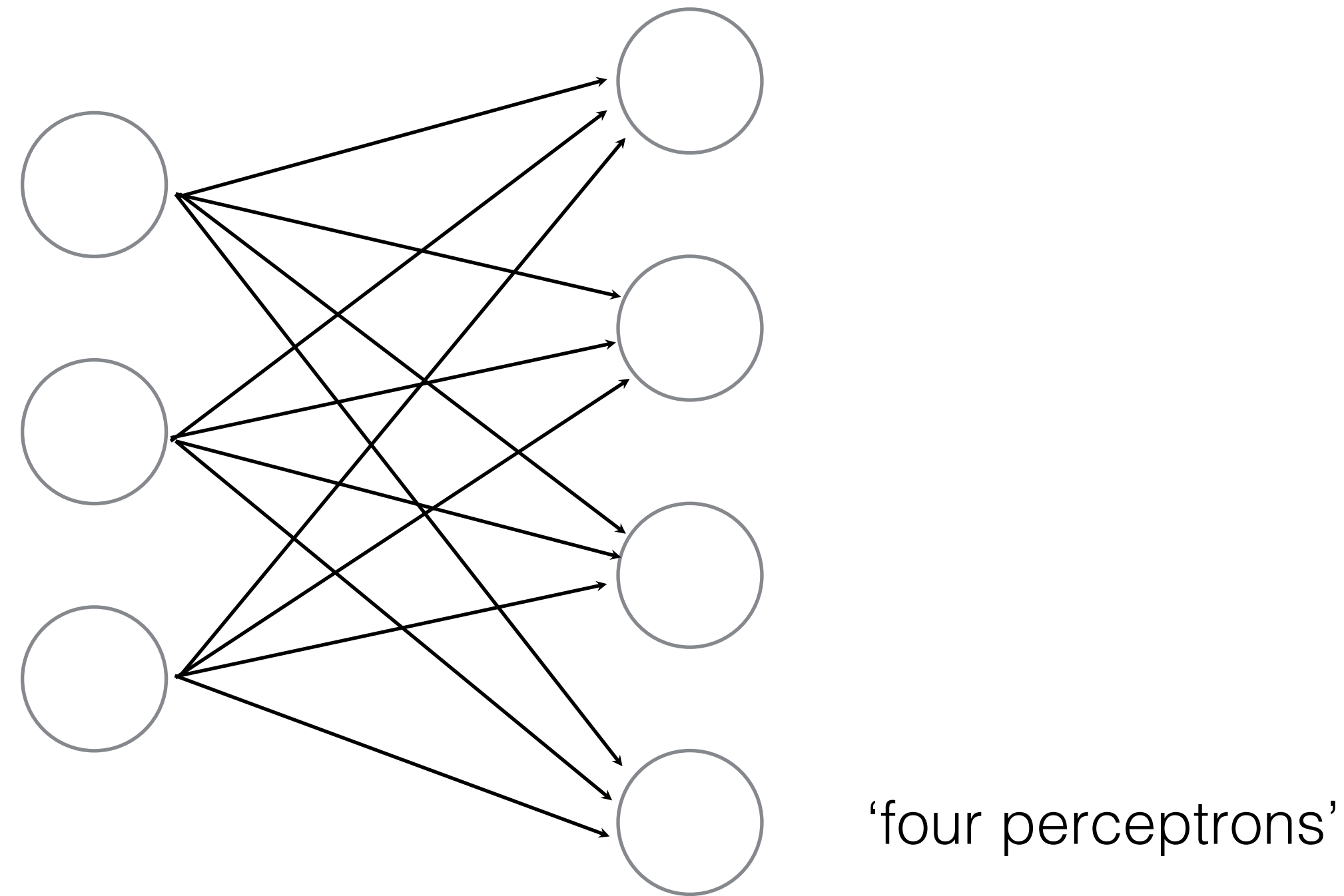Connect a bunch of perceptrons together …
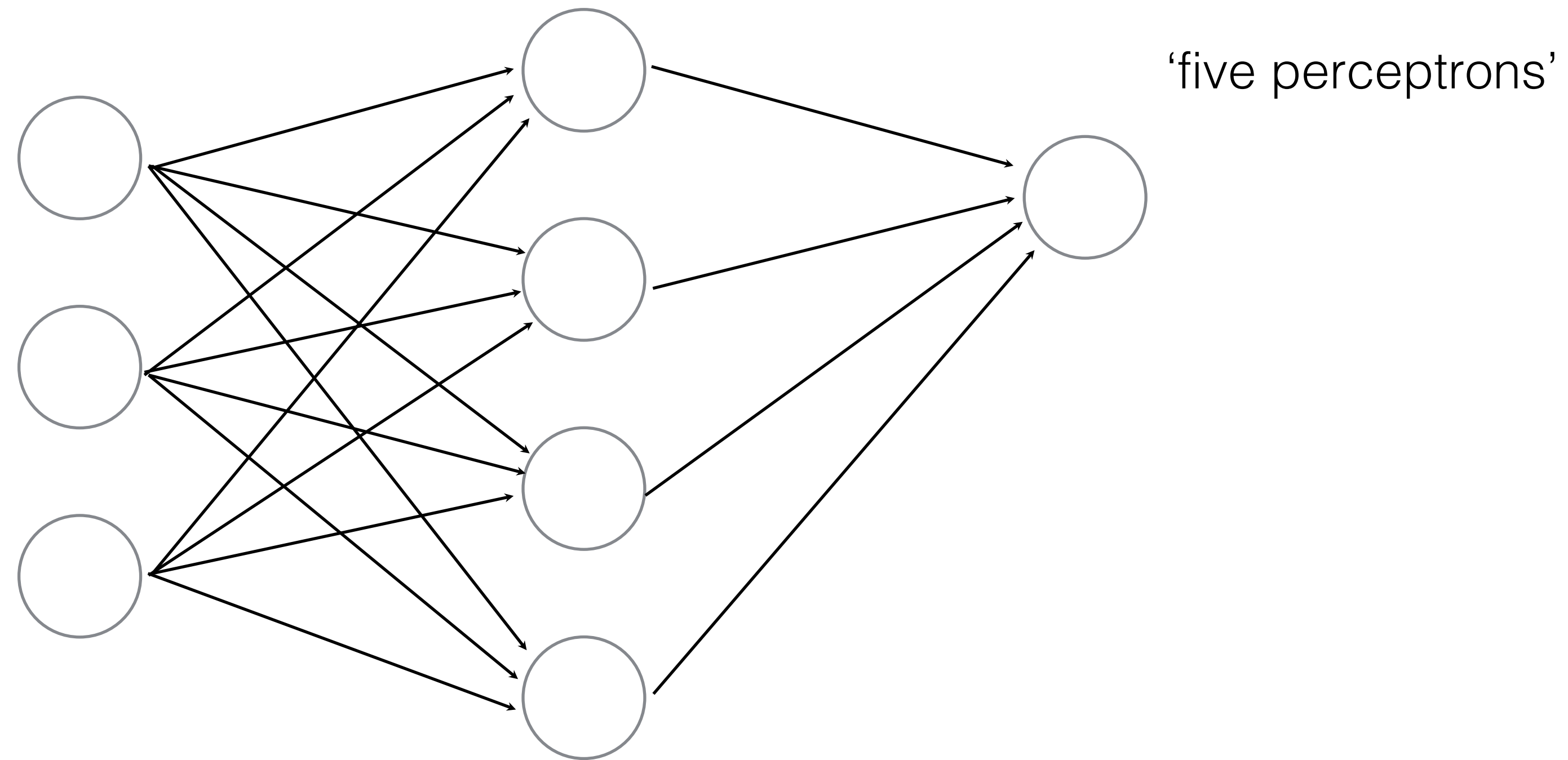
# Neural Network

a collection of connected perceptrons

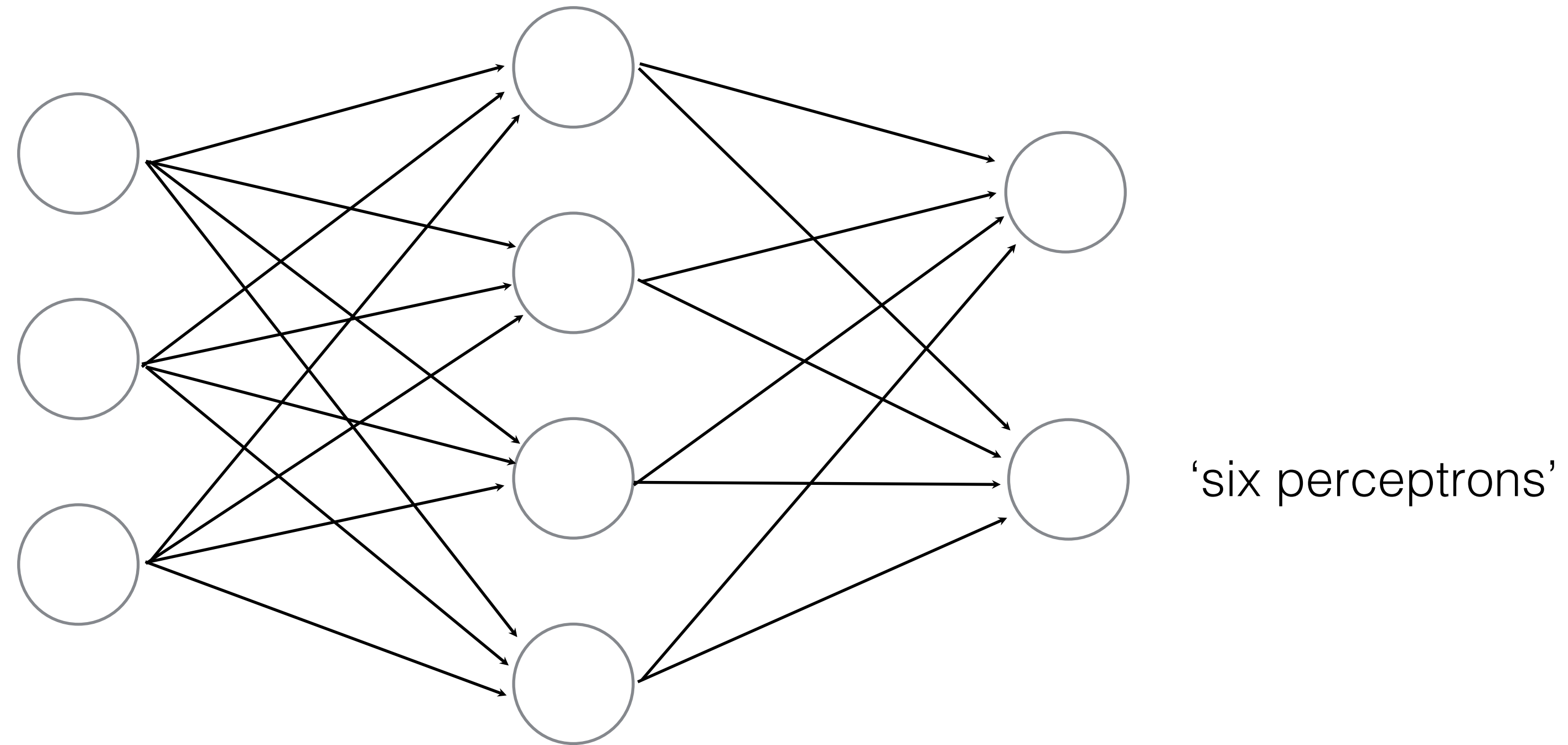# Connect a bunch of perceptrons together …

# Neural Network

a collection of connected perceptrons
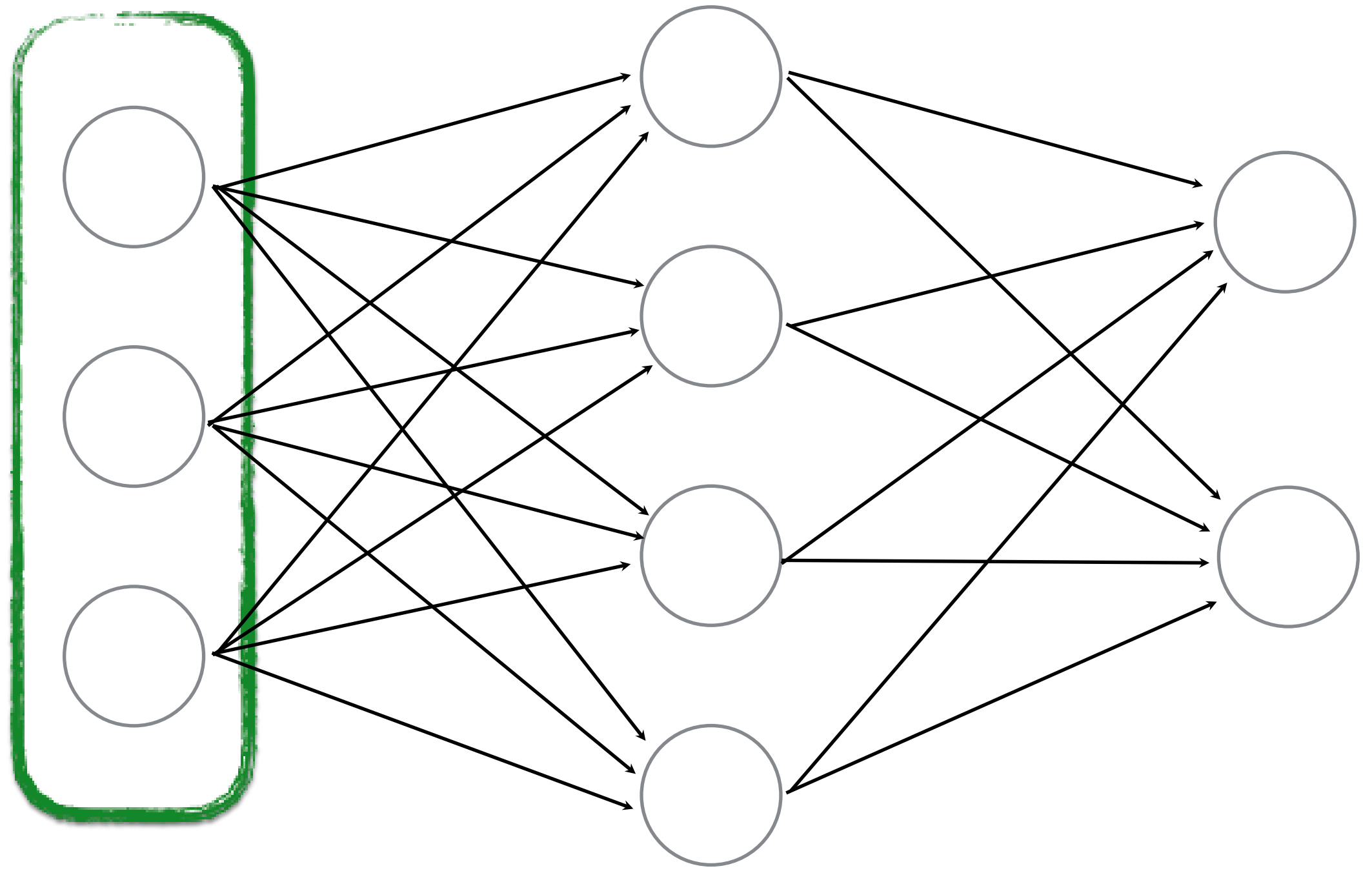
# Connect a bunch of perceptrons together …

# Neural Network

a collection of connected perceptrons

'one perceptron'
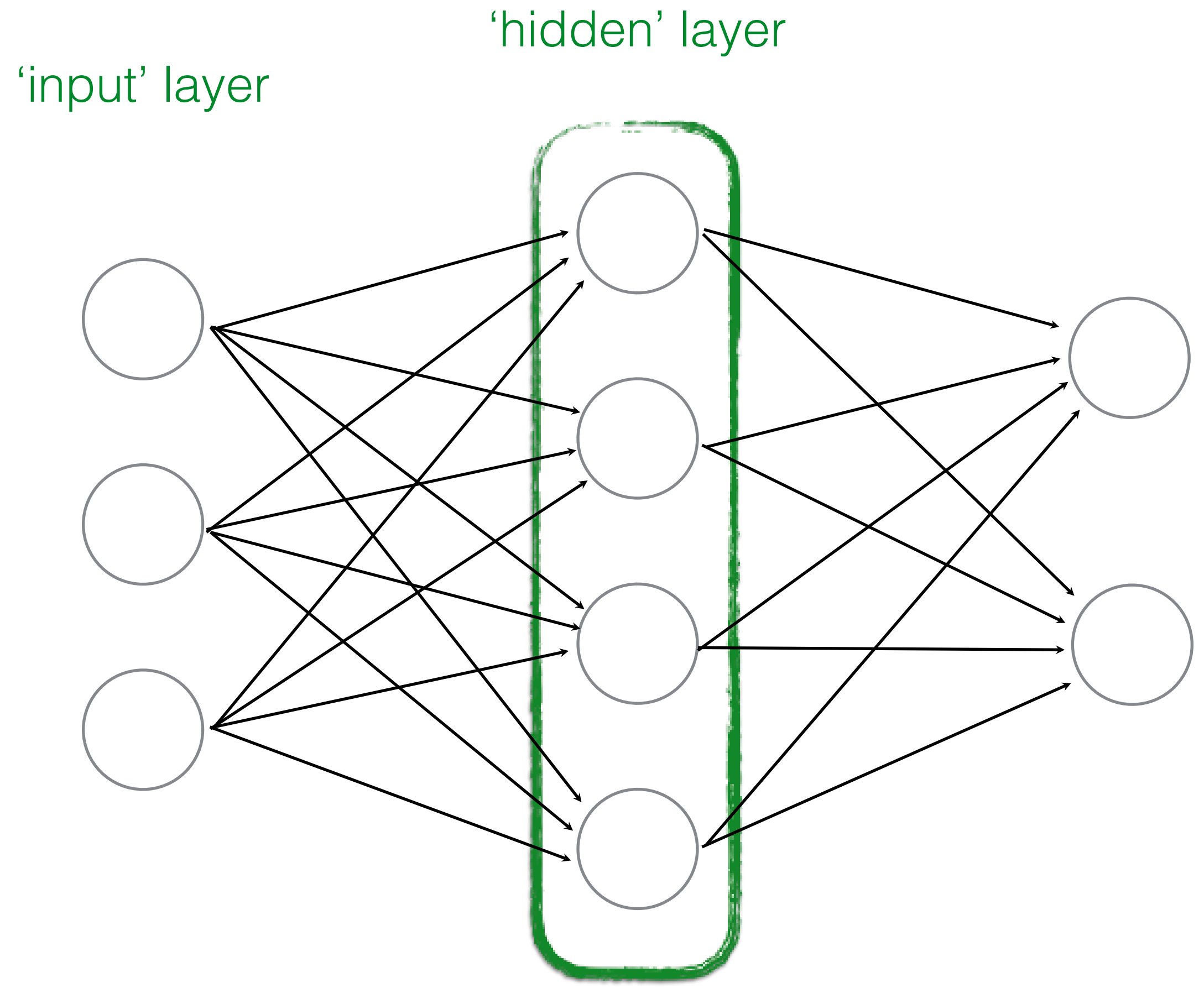
# Connect a bunch of perceptrons together …

# Neural Network

a collection of connected perceptrons

'two perceptrons'

# Connect a bunch of perceptrons together …

# Neural Network

a collection of connected perceptrons



'three perceptrons'

Connect a bunch of perceptrons together …

# Neural Network

a collection of connected perceptrons



'four perceptrons'

# Connect a bunch of perceptrons together …

# Neural Network

a collection of connected perceptrons



'five perceptrons'

# Connect a bunch of perceptrons together …

# Neural Network

a collection of connected perceptrons



'six perceptrons'

# Some terminology…

‘input’ layer



…also called a **Multi-layer Perceptron** (MLP)

# Some terminology…

'hidden' layer

'input' layer



…also called a **Multi-layer Perceptron** (MLP)

# Some terminology…
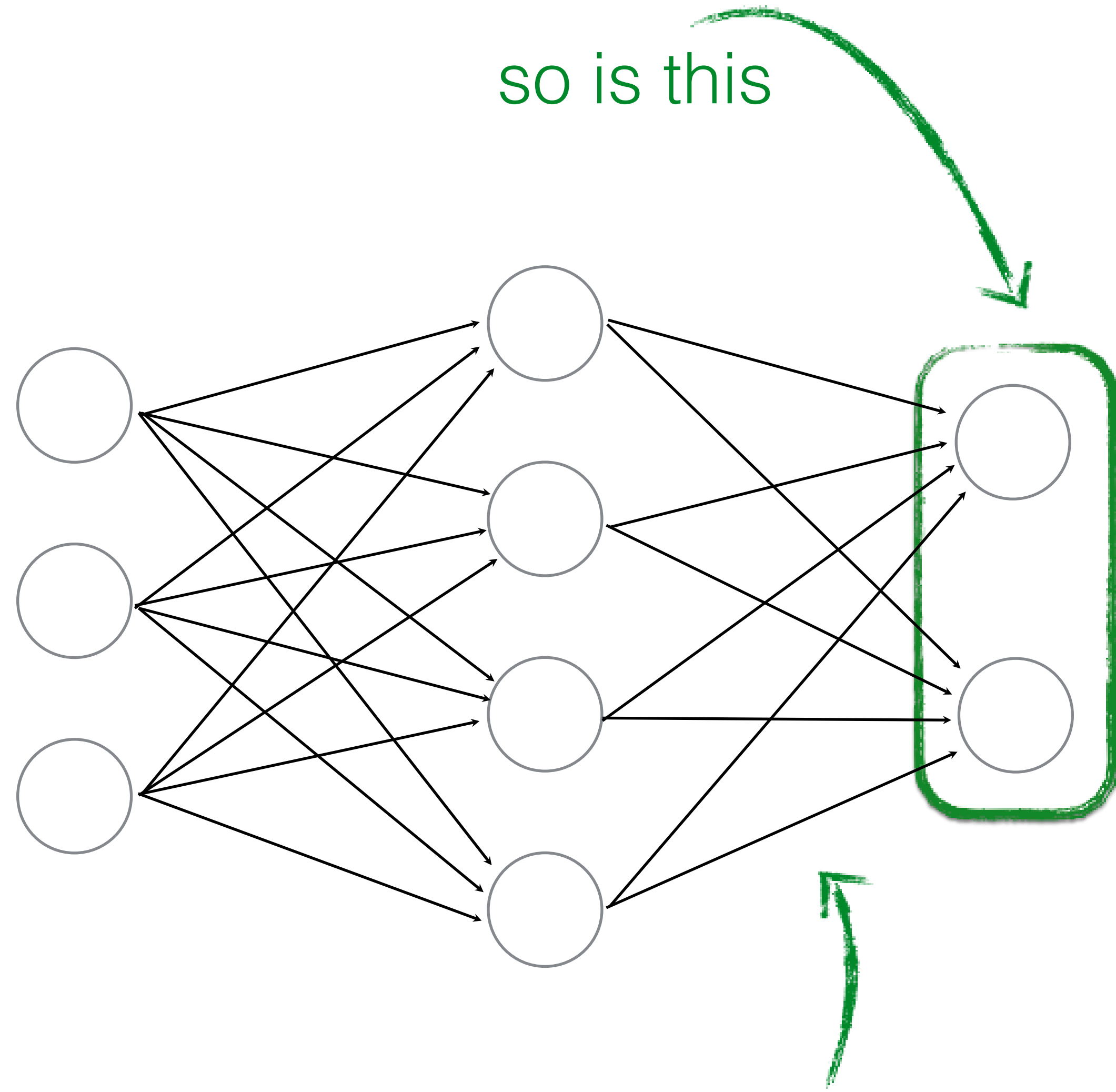
'hidden' layer

'input' layer

'output' layer

…also called a **Multi-layer Perceptron** (MLP)

this layer is a
'fully connected layer'

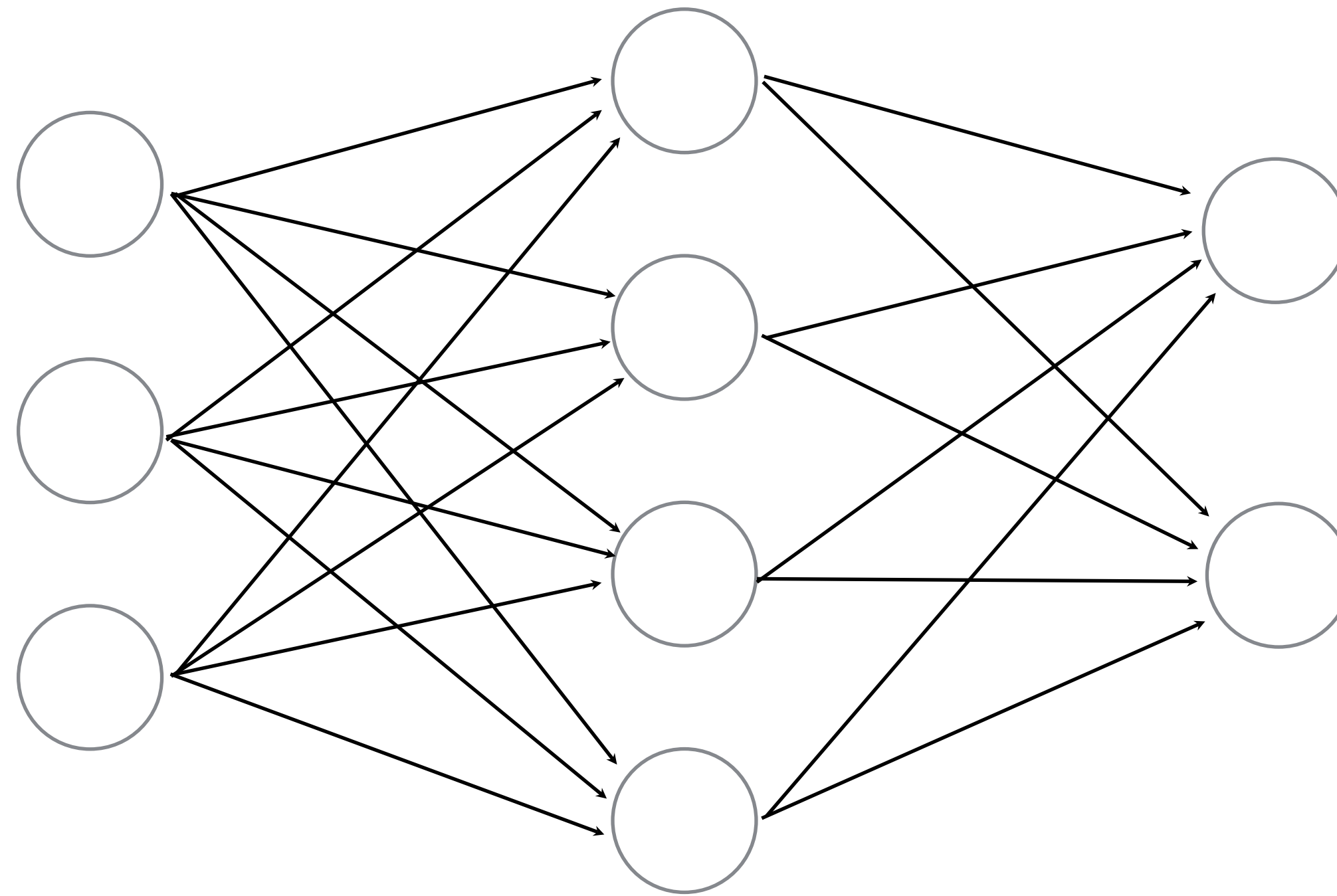all pairwise neurons <u>between</u> layers are connected

so is this

all pairwise neurons <u>between</u> layers are connected
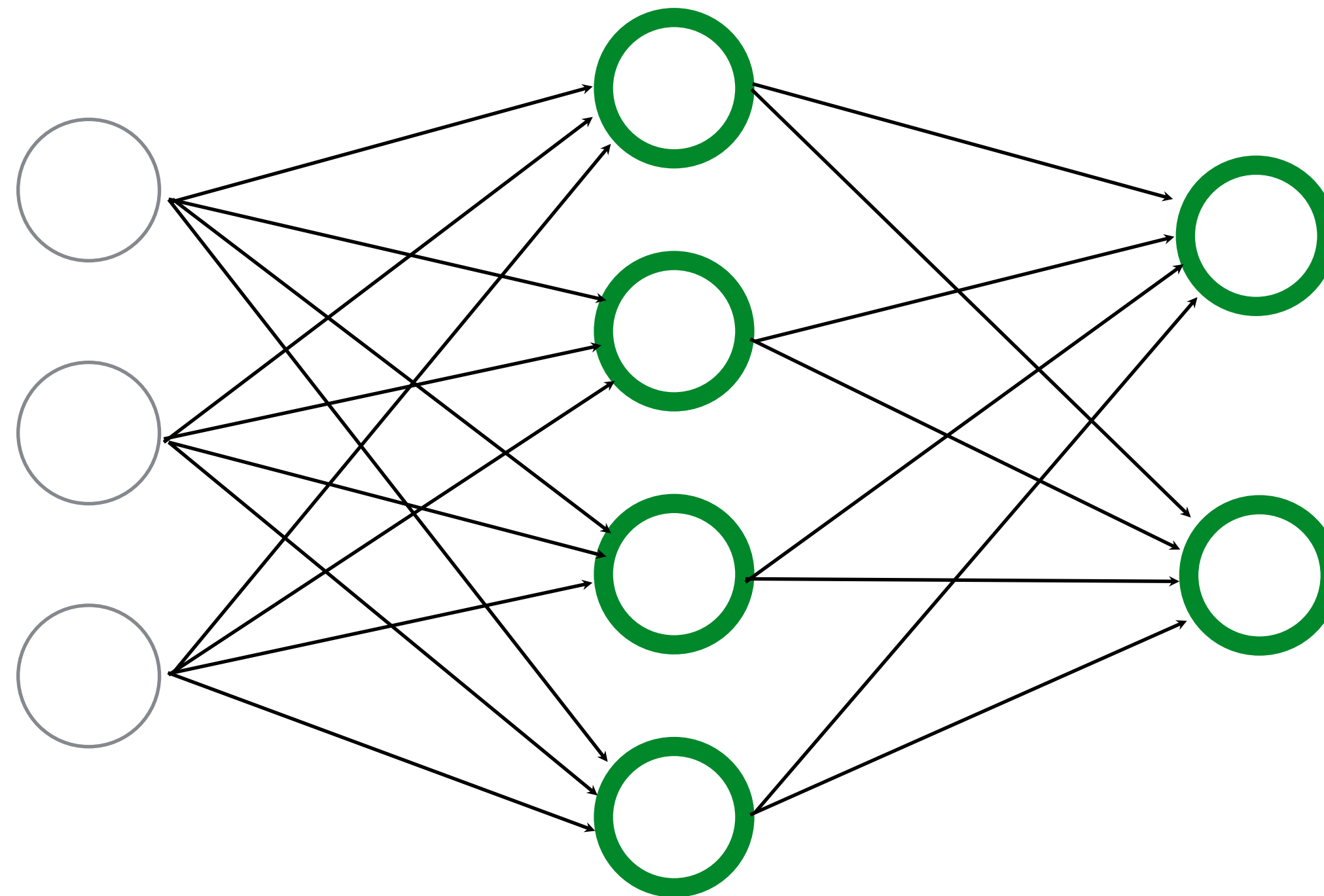
How many neurons (perceptrons)?

How many weights (edges)?

How many learnable parameters total?

*How many neurons (perceptrons)?*     4 + 2 = 6

*How many weights (edges)?*
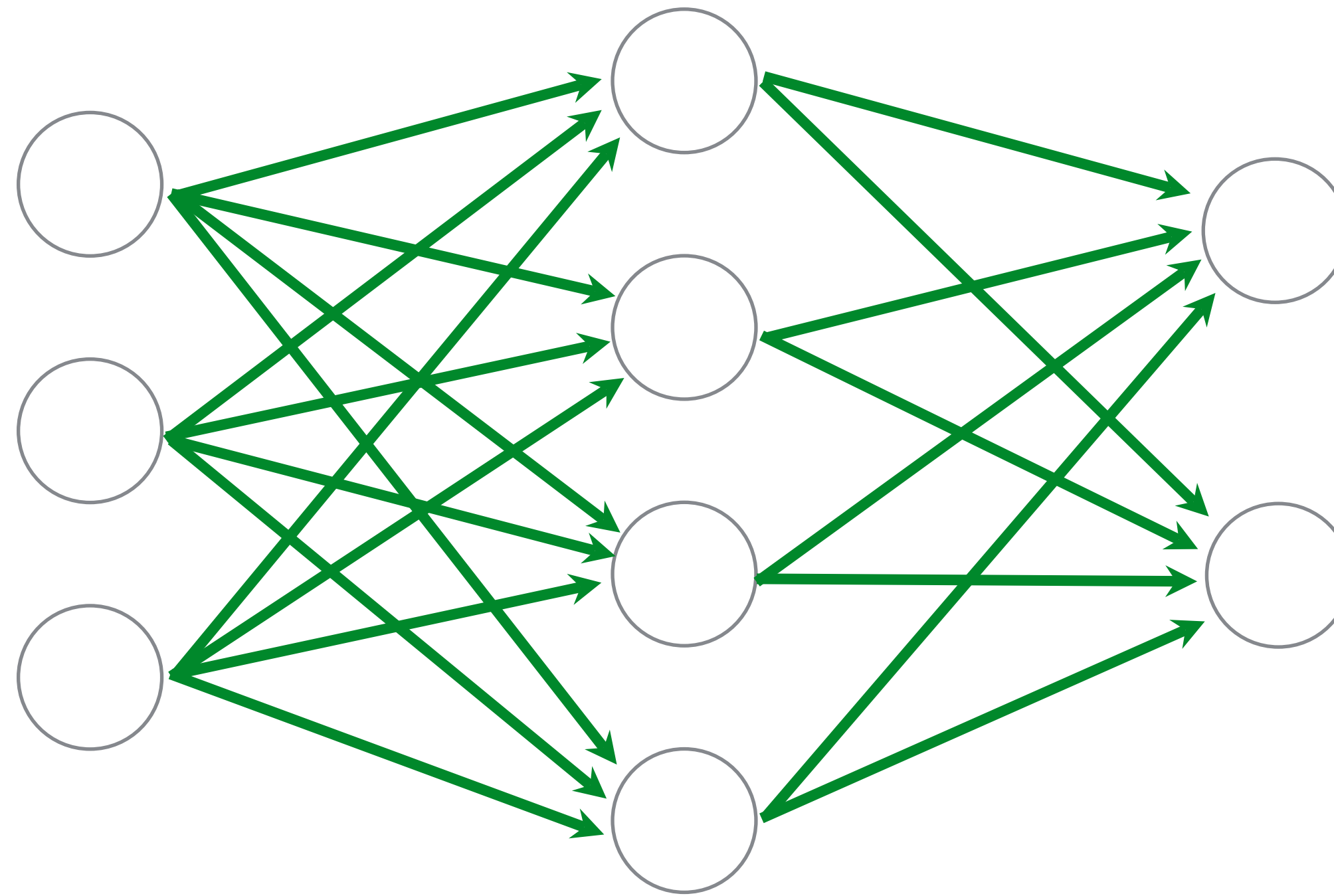
*How many learnable parameters total?*

How many neurons (perceptrons)?     4 + 2 = 6

How many weights (edges)?     (3 x 4) + (4 x 2) = 20
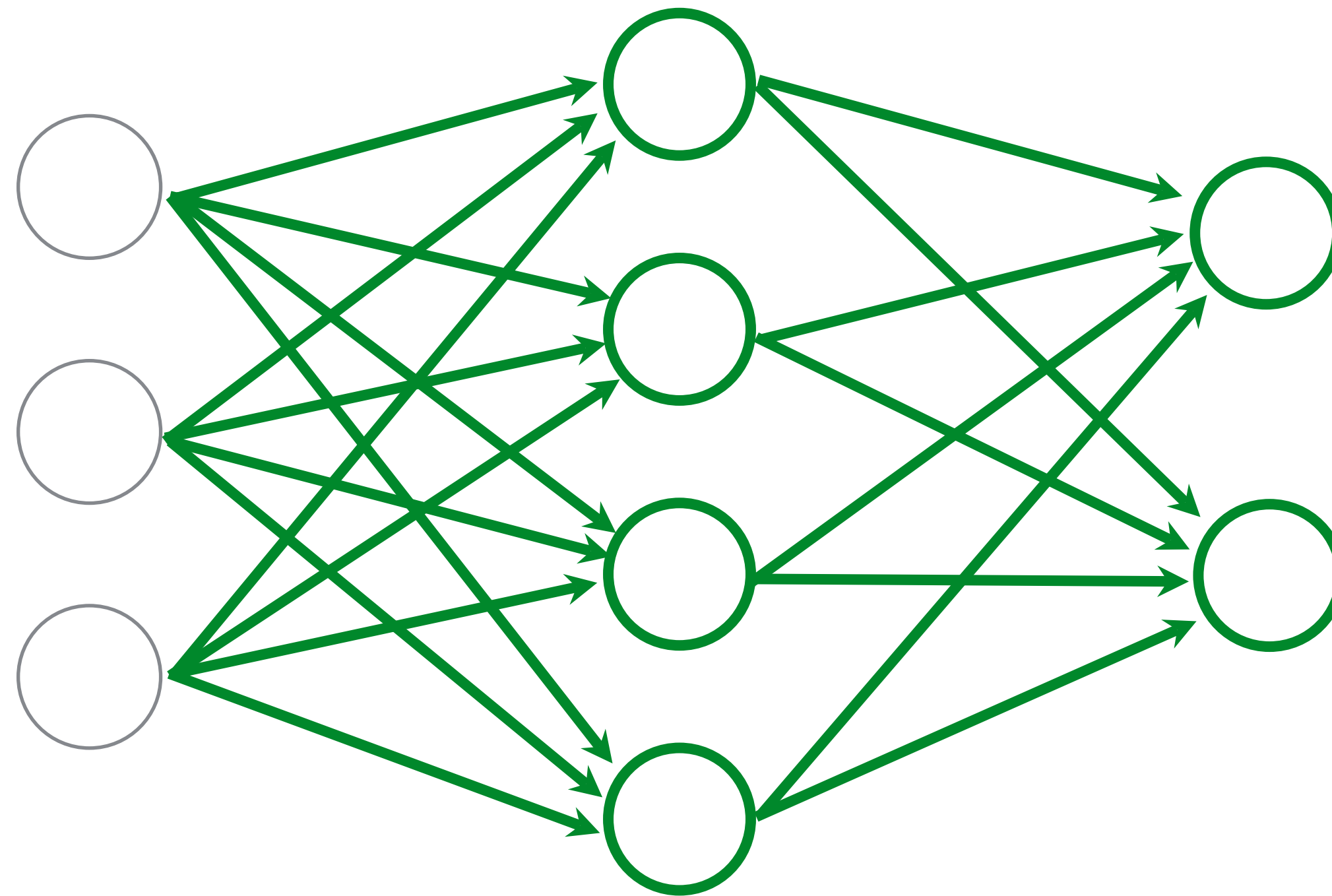


How many learnable parameters total?

*How many neurons (perceptrons)?*            $4 + 2 = 6$

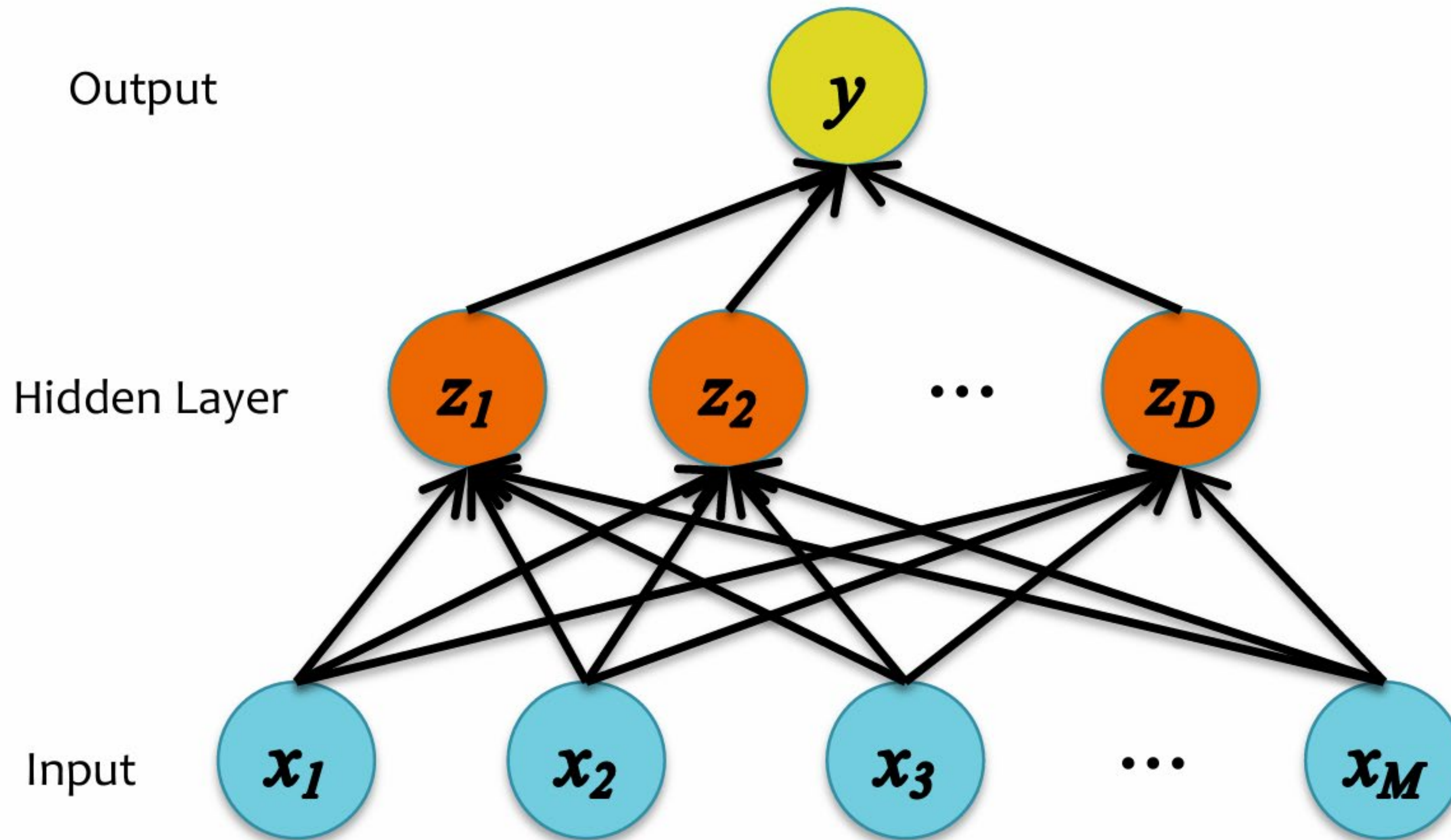*How many weights (edges)?*            $(3 \times 4) + (4 \times 2) = 20$

*How many learnable parameters total?*            $20 + 4 + 2 = 26$
bias terms

# Single Output Neural Networl
**Let's write the equation**

# Objective Functions for NNs

1. **Quadratic Loss:**
   - the same objective as Linear Regression
   - i.e. mean squared error

$$J = \ell_Q(y, y^{(i)}) = \frac{1}{2}(y - y^{(i)})^2$$

$$\frac{dJ}{dy} = y - y^{(i)}$$

2. **Binary Cross-Entropy:**
   - the same objective as Binary Logistic Regression
   - i.e. negative log likelihood
   - This requires our output y to be a probability in [0,1]

$$J = \ell_{CE}(y, y^{(i)}) = -(y^{(i)} \log(y) + (1 - y^{(i)}) \log(1 - y))$$

$$\frac{dJ}{dy} = -\left(y^{(i)} \frac{1}{y} + (1 - y^{(i)}) \frac{1}{y - 1}\right)$$

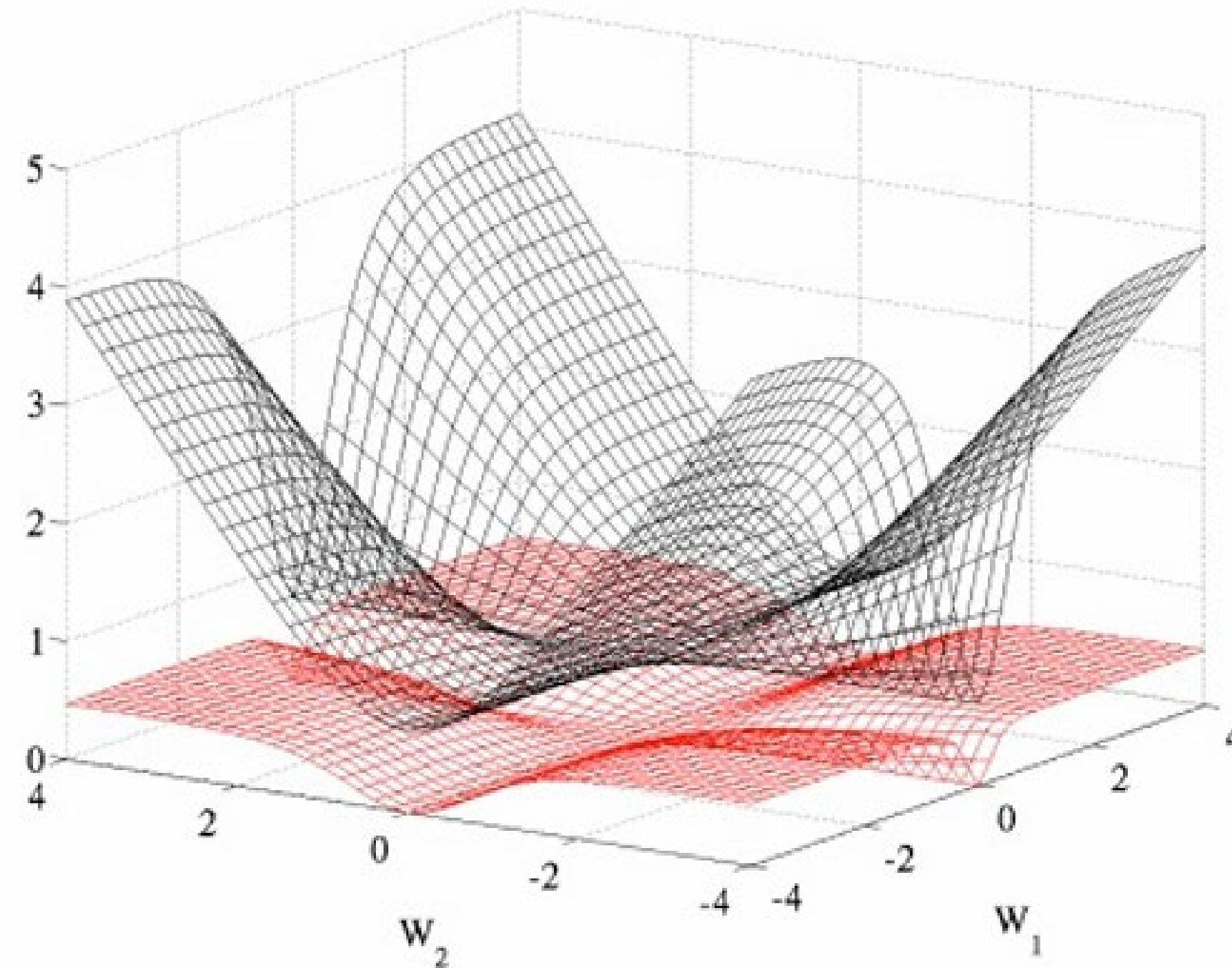# Objective Functions for NNs

**Cross-entropy vs. Quadratic loss**
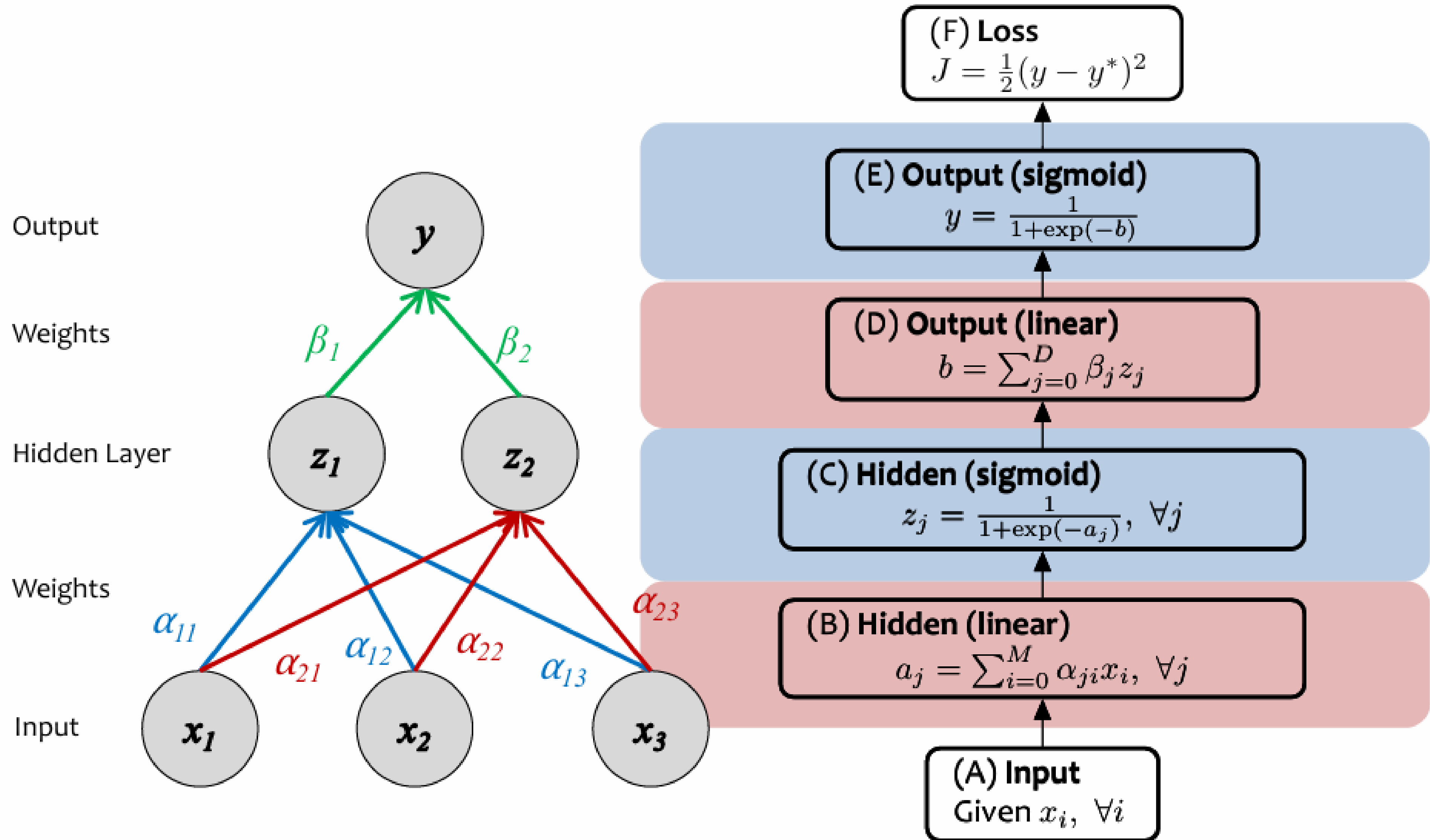


Figure 5: Cross entropy (black, surface on top) and quadratic (red, bottom surface) cost as a function of two weights (one at each layer) of a network with two layers, $W_1$ respectively on the first layer and $W_2$ on the second, output layer.
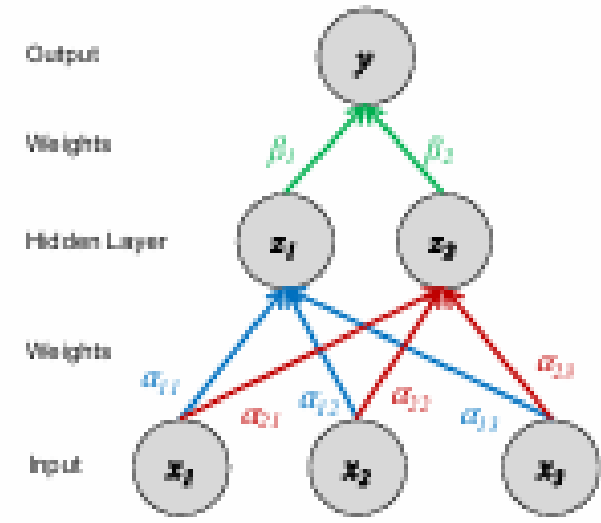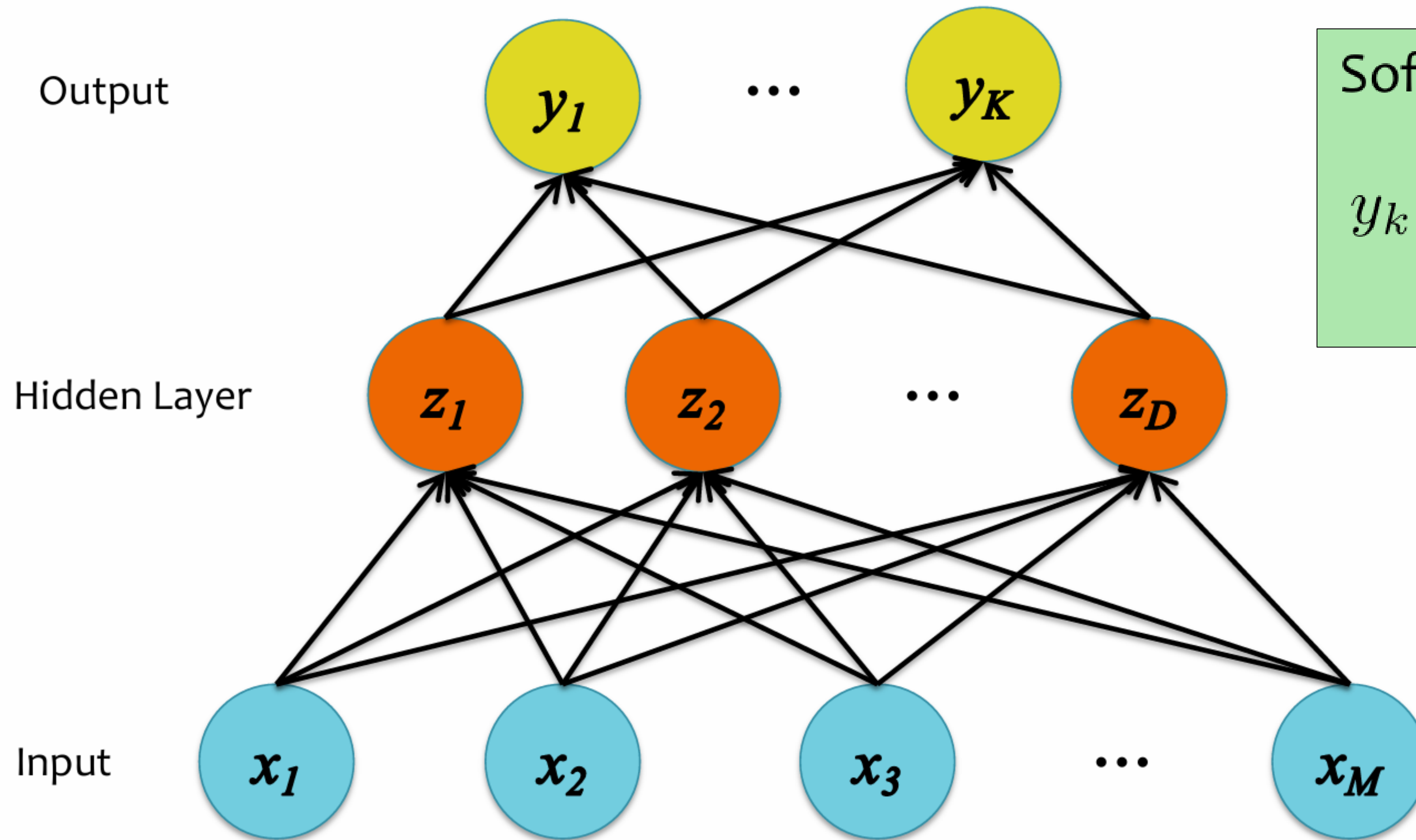
# Backpropagation



| | Forward | Backward |
|---|---|---|
| Loss | $J = y^* \log y + (1 - y^*) \log(1 - y)$ | $g_y = \dfrac{y^*}{y} + \dfrac{(1 - y^*)}{y - 1}$ |
| Sigmoid | $y = \dfrac{1}{1 + \exp(-b)}$ | $g_b = g_y \dfrac{\partial y}{\partial b}, \; \dfrac{\partial y}{\partial b} = \dfrac{\exp(-b)}{(\exp(-b) + 1)^2}$ |
| Linear | $b = \displaystyle\sum_{j=0}^{D} \beta_j z_j$ | $g_{\beta_j} = g_b \dfrac{\partial b}{\partial \beta_j}, \; \dfrac{\partial b}{\partial \beta_j} = z_j$ <br><br> $g_{z_j} = g_b \dfrac{\partial b}{\partial z_j}, \; \dfrac{\partial b}{\partial z_j} = \beta_j$ |
| Sigmoid | $z_j = \dfrac{1}{1 + \exp(-a_j)}$ | $g_{a_j} = g_{z_j} \dfrac{\partial z_j}{\partial a_j}, \; \dfrac{\partial z_j}{\partial a_j} = \dfrac{\exp(-a_j)}{(\exp(-a_j) + 1)^2}$ |
| Linear | $a_j = \displaystyle\sum_{i=0}^{M} \alpha_{ji} x_i$ | $g_{\alpha_{ji}} = g_{a_j} \dfrac{\partial a_j}{\partial \alpha_{ji}}, \; \dfrac{\partial a_j}{\partial \alpha_{ji}} = x_i$ <br><br> $g_{x_i} = \displaystyle\sum_{j=0}^{D} g_{a_j} \dfrac{\partial a_j}{\partial x_i}, \; \dfrac{\partial a_j}{\partial x_i} = \alpha_{ji}$ |

# Multi-class Output



Output

$y_1$ ... $y_K$

Hidden Layer

$z_1$ $z_2$ ... $z_D$

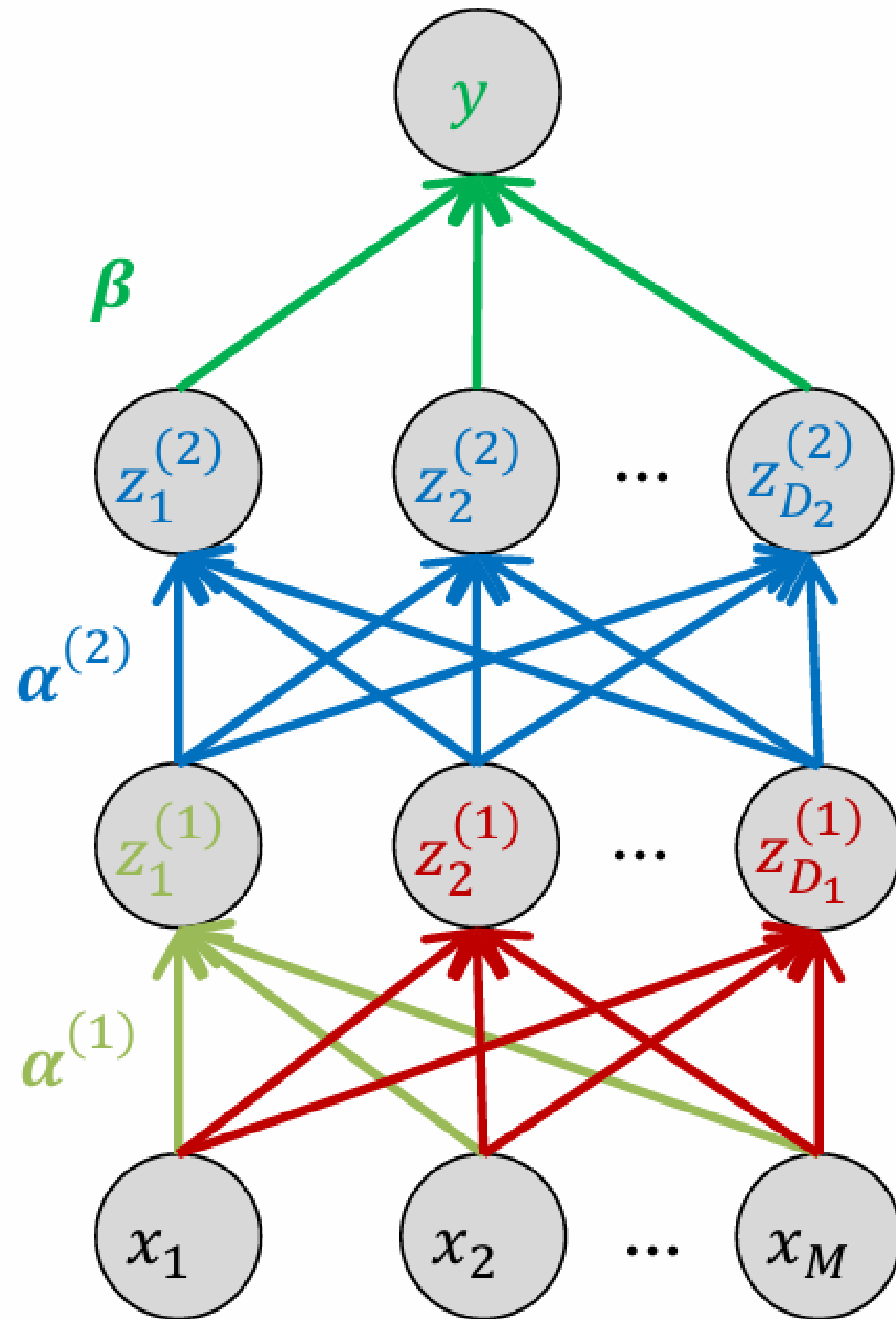Input

$x_1$ $x_2$ $x_3$ ... $x_M$

Softmax:

$$y_k = \frac{\exp(b_k)}{\sum_{l=1}^{K} \exp(b_l)}$$

# Two-Layer NN: How do we train this model?



$\boldsymbol{\beta} \in \mathbb{R}^{D_2}$

$\beta_0 \in \mathbb{R}$

$\boldsymbol{\alpha}^{(2)} \in \mathbb{R}^{M \times D_2}$

$\boldsymbol{b}^{(2)} \in \mathbb{R}^{D_2}$

$\boldsymbol{\alpha}^{(1)} \in \mathbb{R}^{M \times D_1}$

$\boldsymbol{b}^{(1)} \in \mathbb{R}^{D_1}$

$y = \sigma((\boldsymbol{\beta})^T \boldsymbol{z}^{(2)} + \beta_0)$

$\boldsymbol{z}^{(2)} = \sigma((\boldsymbol{\alpha}^{(2)})^T \boldsymbol{z}^{(1)} + \boldsymbol{b}^{(2)})$

$\boldsymbol{z}^{(1)} = \sigma((\boldsymbol{\alpha}^{(1)})^T \boldsymbol{x} + \boldsymbol{b}^{(1)})$