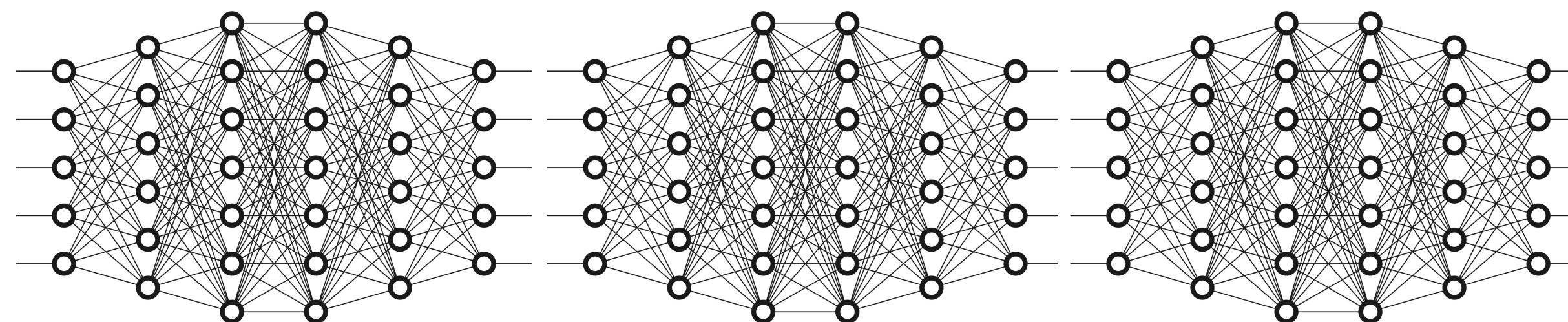




CPSC 425: Computer Vision



Lecture 19: Neural Networks 1

Menu for Today

Topics:

- **Neural Networks** introduction
- **Activation functions** softmax, relu
- **2-layer** fully connected net
- **Backprop** intro

Readings:

- **Today's Lecture:** Szeliski 5.1.3, 5.3-5.4, Justin Johnson Michigan EECS 498/598

Reminders:

- **Assignment 5:** due today
- **Assignment 6:** Deep Learning is now available

Recall: Linear Classifier

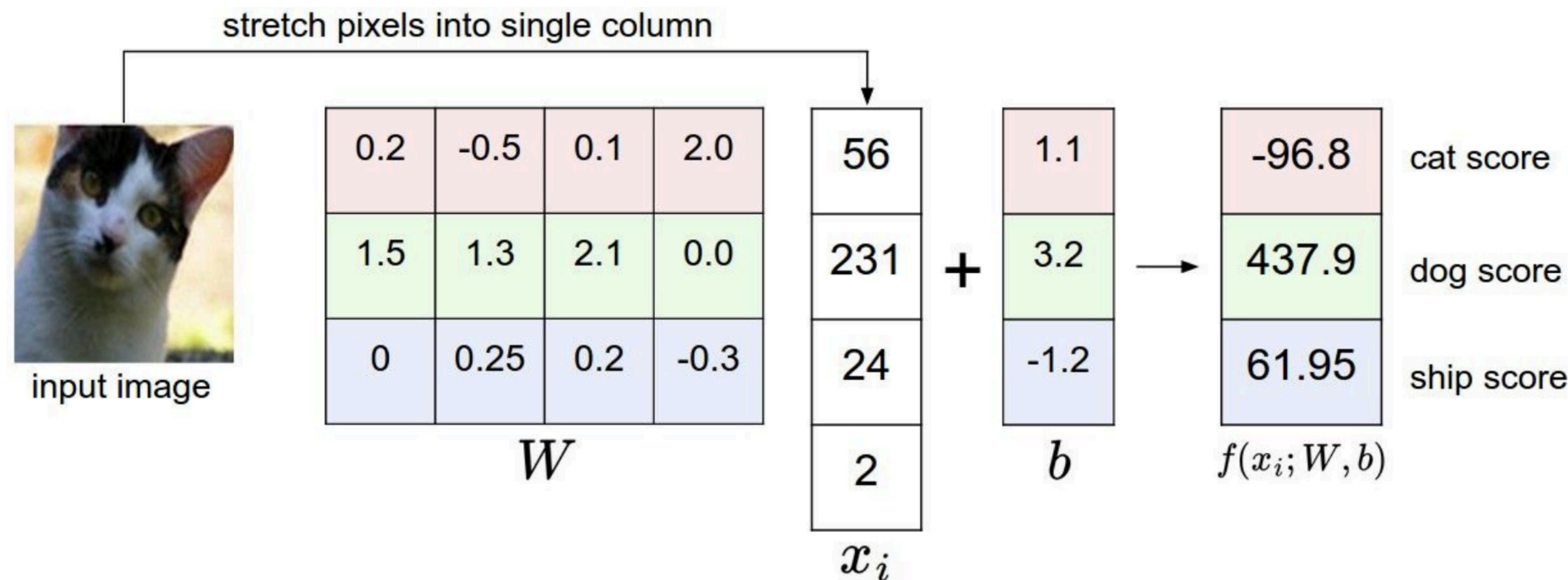
Defines a score function:

$$f(\mathbf{x}_i, \mathbf{W}, \mathbf{b}) = \mathbf{W}\mathbf{x}_i + \mathbf{b}$$

The diagram illustrates the components of the linear classifier equation. The input features \mathbf{x}_i are represented by a green box labeled "image features". The weights \mathbf{W} are represented by a purple box labeled "weights (parameters)". The bias vector \mathbf{b} is represented by a blue box labeled "bias vector (parameters)". The equation $f(\mathbf{x}_i, \mathbf{W}, \mathbf{b}) = \mathbf{W}\mathbf{x}_i + \mathbf{b}$ is shown below, with the terms corresponding to the colored boxes.

Recall: Linear Classifier

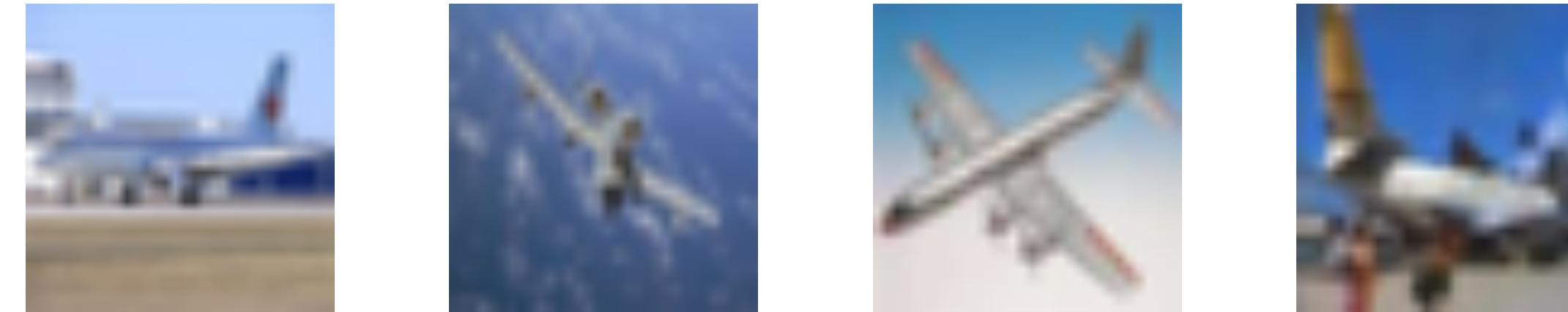
Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)



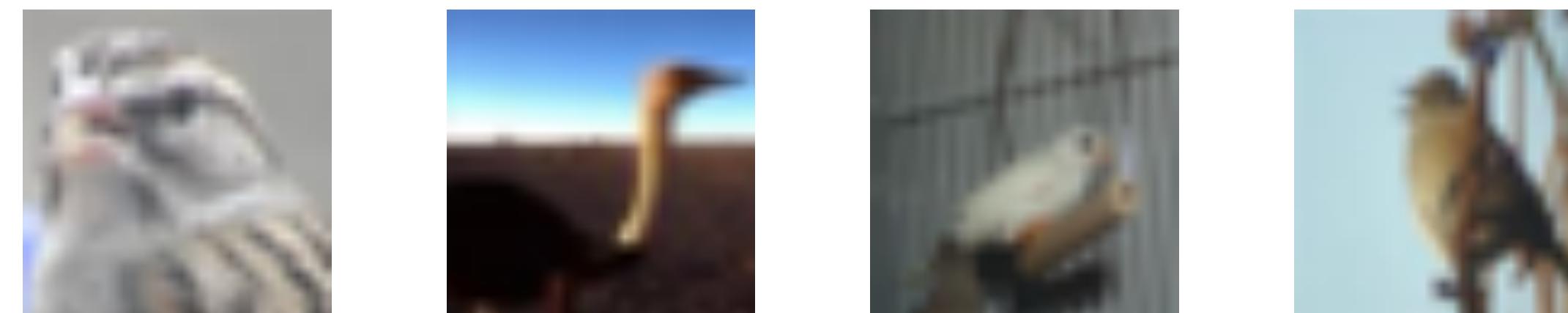
Linear Classification

- Let's start by using 2 classes, e.g., bird and plane
- Apply labels (y) to training set:

$y = +1$



$y = -1$

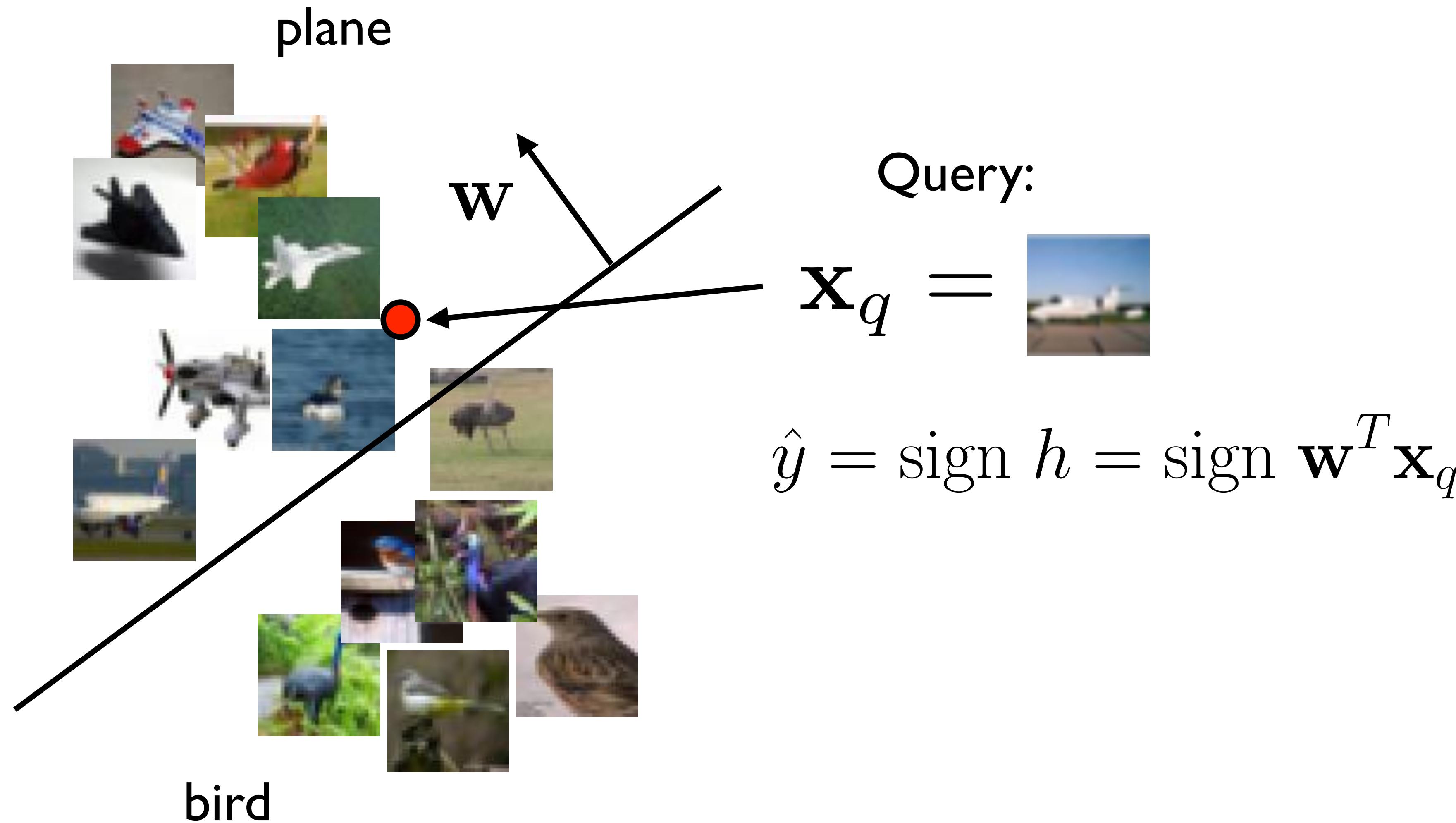


- Use a linear model to regress y from x

$$\hat{y} = \text{sign } h = \text{sign } \mathbf{w}^T \mathbf{x}_q$$

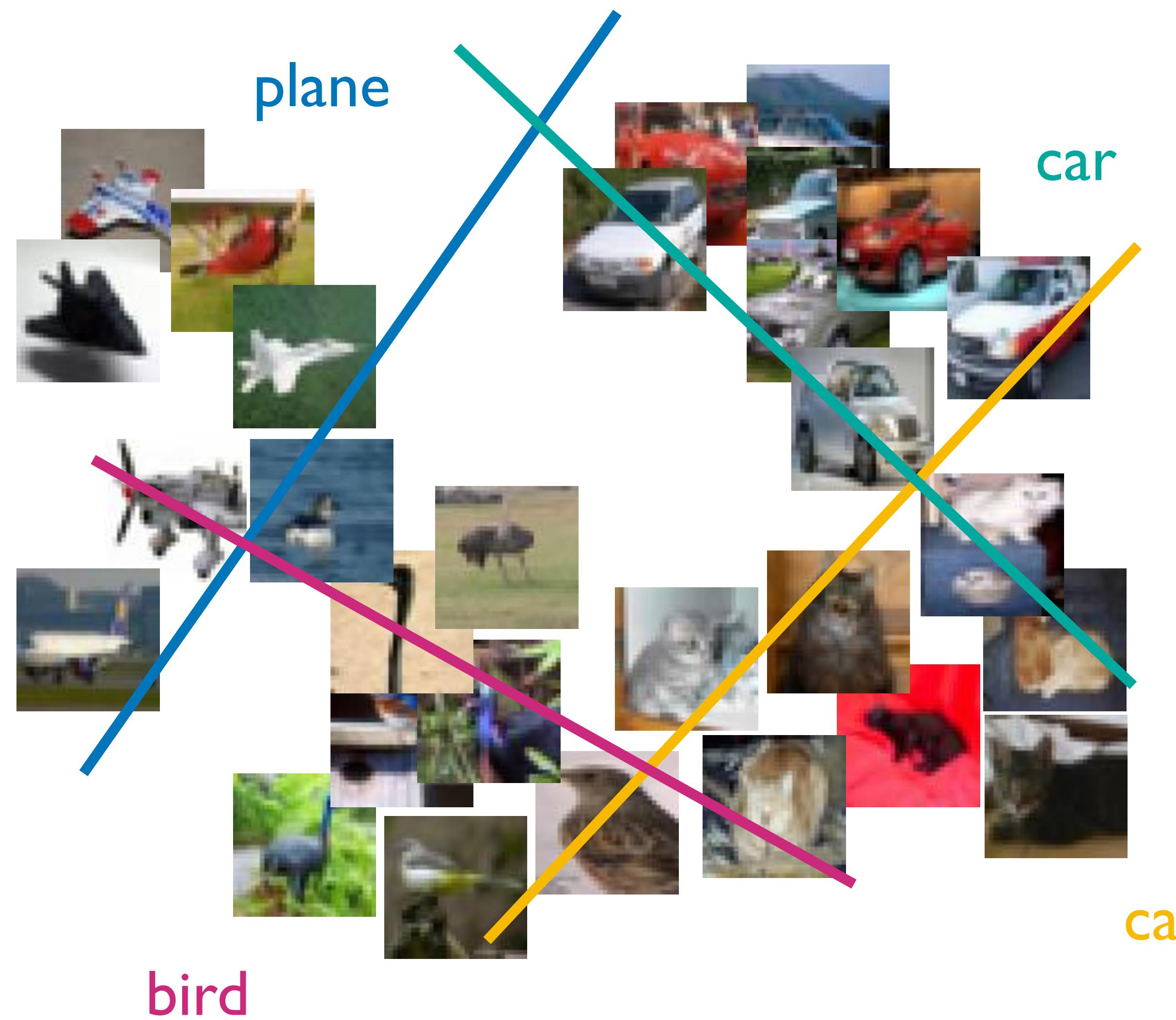
2-class Linear Classification

- Separating hyperplane, projection to a line defined by \mathbf{w}



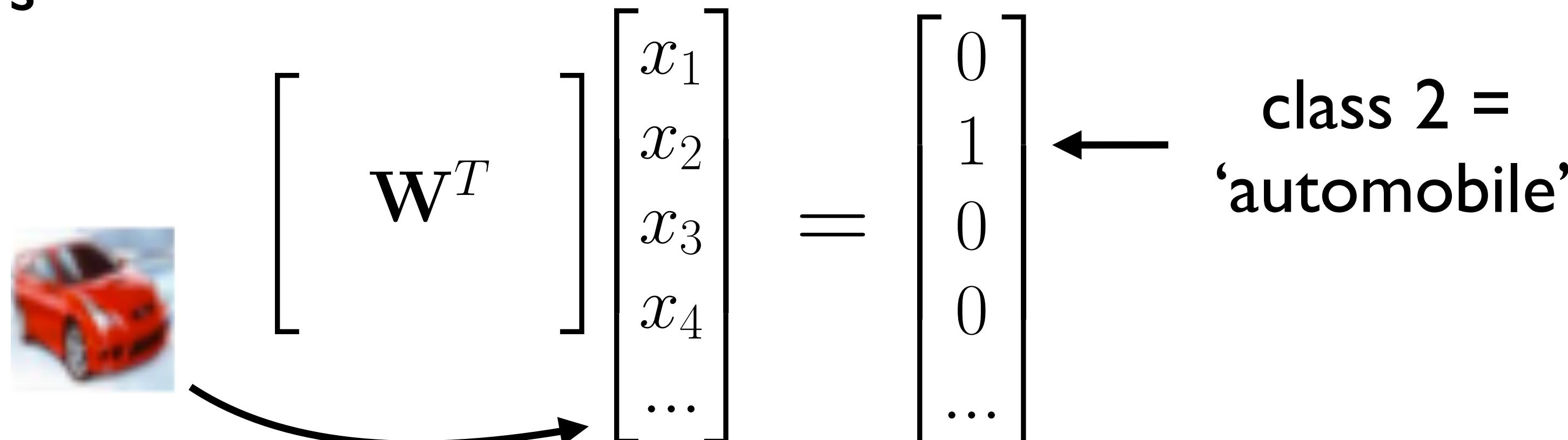
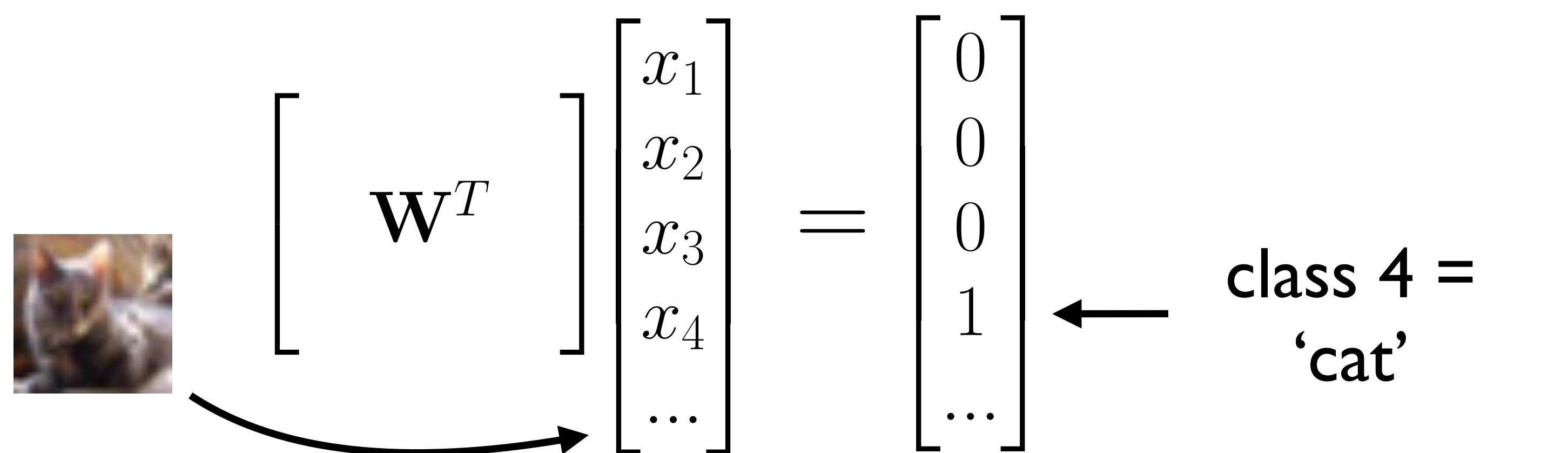
N-class Linear Classification

- One hot regression = I vs all classifiers



One-Hot Regression

- A better solution is to regress to one-hot targets = 1 vs all classifiers

$$\begin{bmatrix} \mathbf{W}^T \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \quad \leftarrow \text{class 2 = 'automobile'}$$

$$\begin{bmatrix} \mathbf{W}^T \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \vdots \end{bmatrix} \quad \leftarrow \text{class 4 = 'cat'}$$


One-Hot Regression

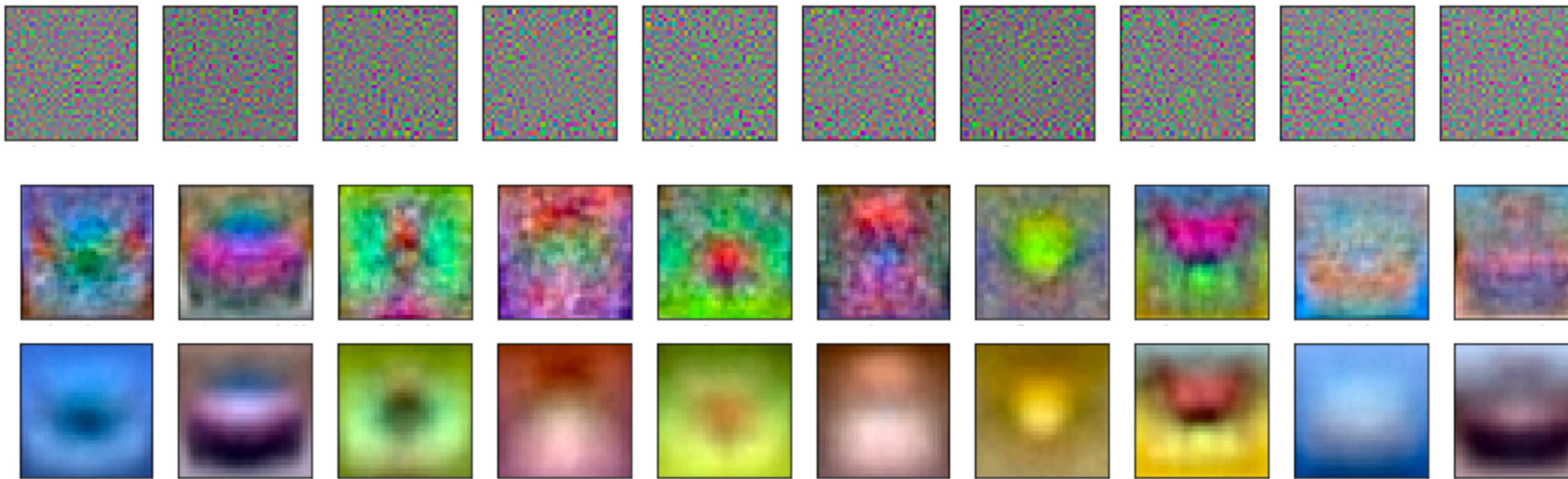
- Transpose (to match Project 3 notebook)


$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots \\ x_{21} & x_{22} & x_{23} & \dots \\ x_{31} & x_{32} & x_{33} & \dots \\ \dots \end{bmatrix} \begin{bmatrix} \mathbf{W} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 0 & 1 & \dots \\ \dots & \dots \end{bmatrix} \begin{matrix} \text{auto} \\ \text{cat} \end{matrix}$$
$$\mathbf{X}\mathbf{W} = \mathbf{T}$$

- Solve regression problem by Least Squares

Regularized Classification

- Add regularization to CIFAR10 linear classifier

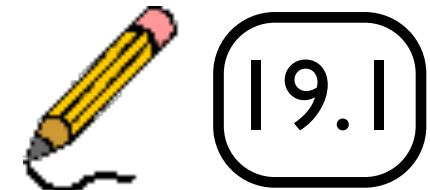


- Row 1 = overfitting, Row 3 = oversmoothing?

$$e = |\mathbf{X}\mathbf{W} - \mathbf{T}|^2 + \lambda|\mathbf{W}|^2$$

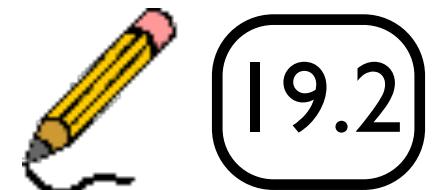
Softmax + Logistic Outputs

- Linear regression to one-hot targets is a bit strange..
- Output could be very large, and scores $>> 1$ are penalised even for the correct class, ditto scores $<< 1$ for incorrect
- How about restricting output scores to 0-1?



Softmax + Cross Entropy

- What is the gradient of the softmax linear classifier?
- We could use L2 loss, but we'll use cross entropy instead
- This has a sound motivation — it is a measure of the difference between probability distributions
- It also leads to a simple update rule



Linear + Softmax Regression

- We found the following gradient descent update rule

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \alpha(\mathbf{h} - \mathbf{t})\mathbf{x}^T$$

↑ ↑ ↗
prediction targets data

- This applies to:

Linear regression $\mathbf{h} = \mathbf{W}^T \mathbf{x}$ L2 loss

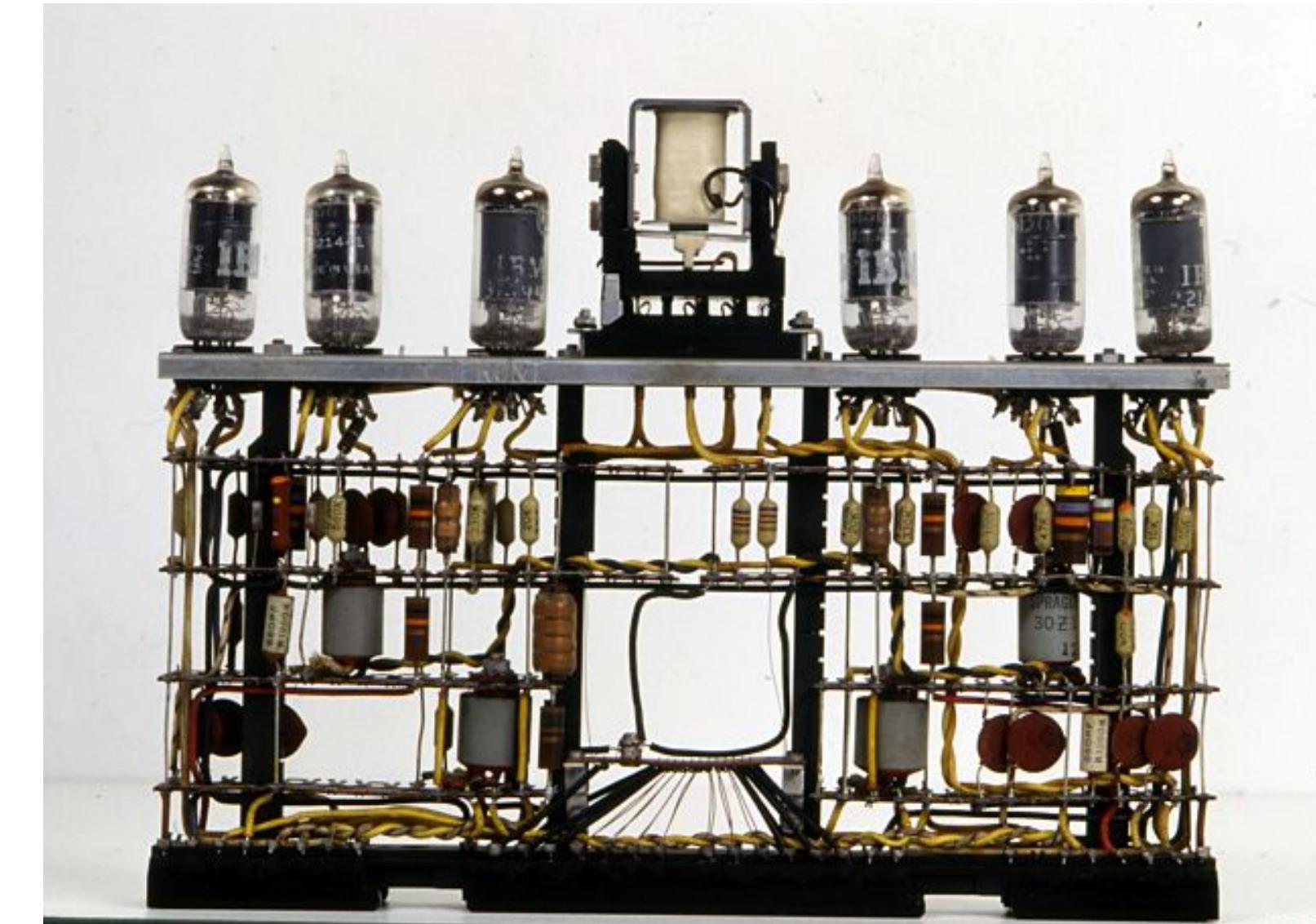
Softmax regression $\mathbf{h} = \sigma(\mathbf{W}^T \mathbf{x})$ cross-entropy loss

- The same update rule with a binary prediction function

$$\mathbf{h} = \mathbb{1}_{\max}(\mathbf{W}^T \mathbf{x})$$

implements the multiclass Perceptron learning rule

History of the Perceptron



[I.B.M. Italia]

- This machine (IBM 704) was used by Frank Rosenblatt to implement the perceptron in 1958
- Based on his statements, the New York Times reported it as: "the embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence."

2-class Perceptron Classifier

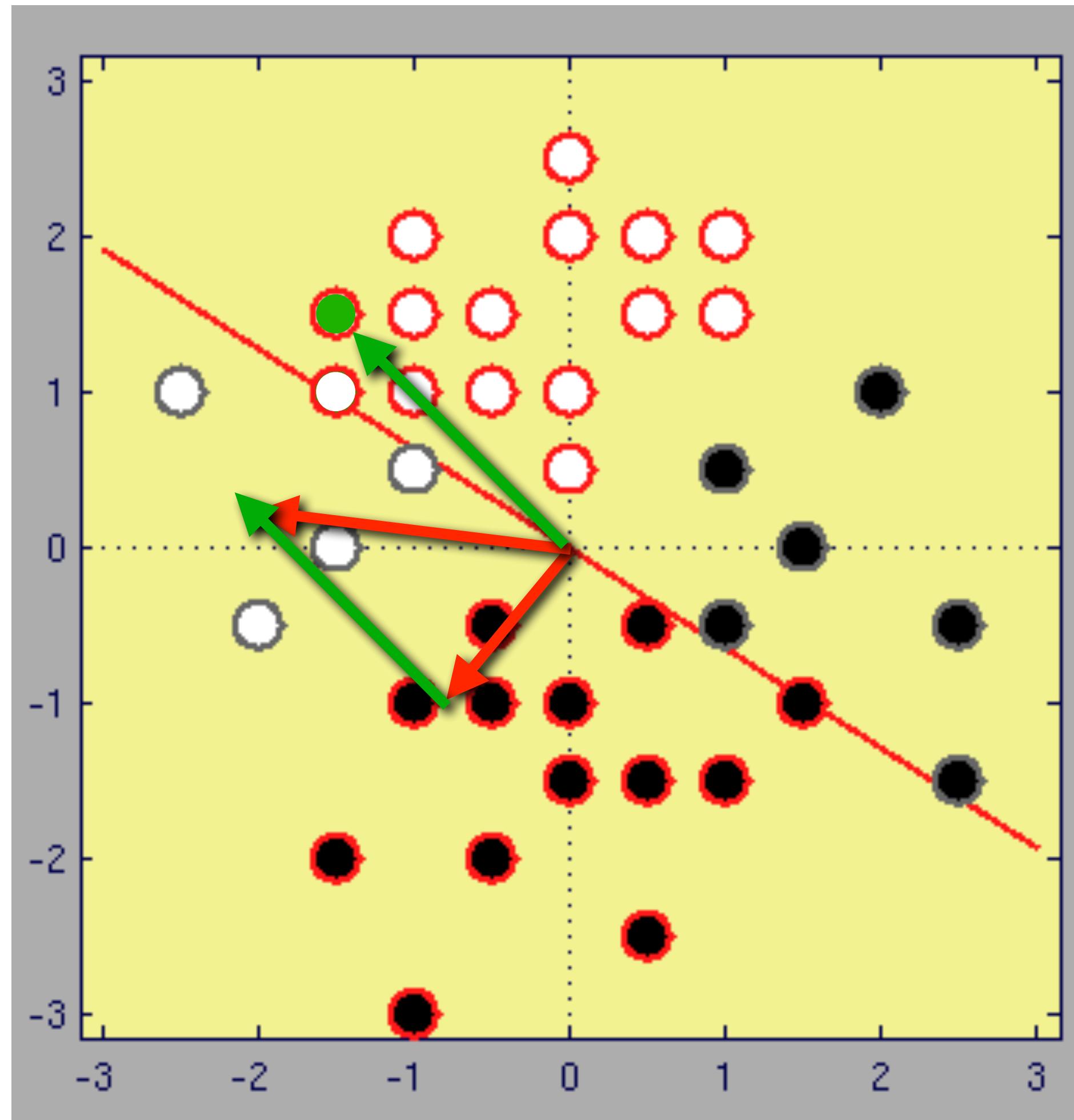
- Classification function is

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$

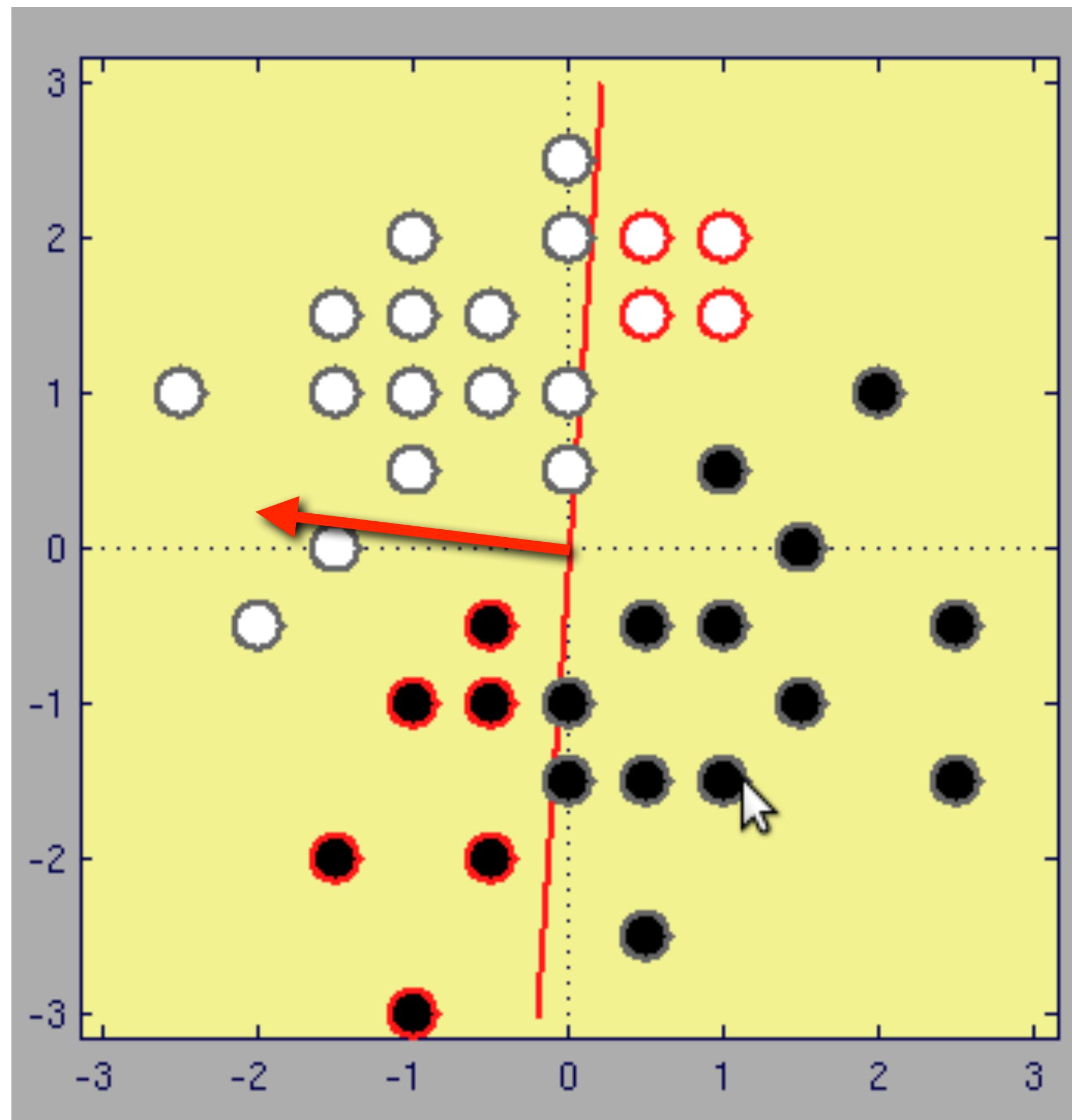
- Linear function of the data (\mathbf{x}) followed by 0/1 activation
- Update rule: present data \mathbf{x}
 - if correctly classified, do nothing
 - if incorrectly classified, update the weight vector

$$\mathbf{w}_{n+1} = \mathbf{w}_n + y_i \mathbf{x}_i$$

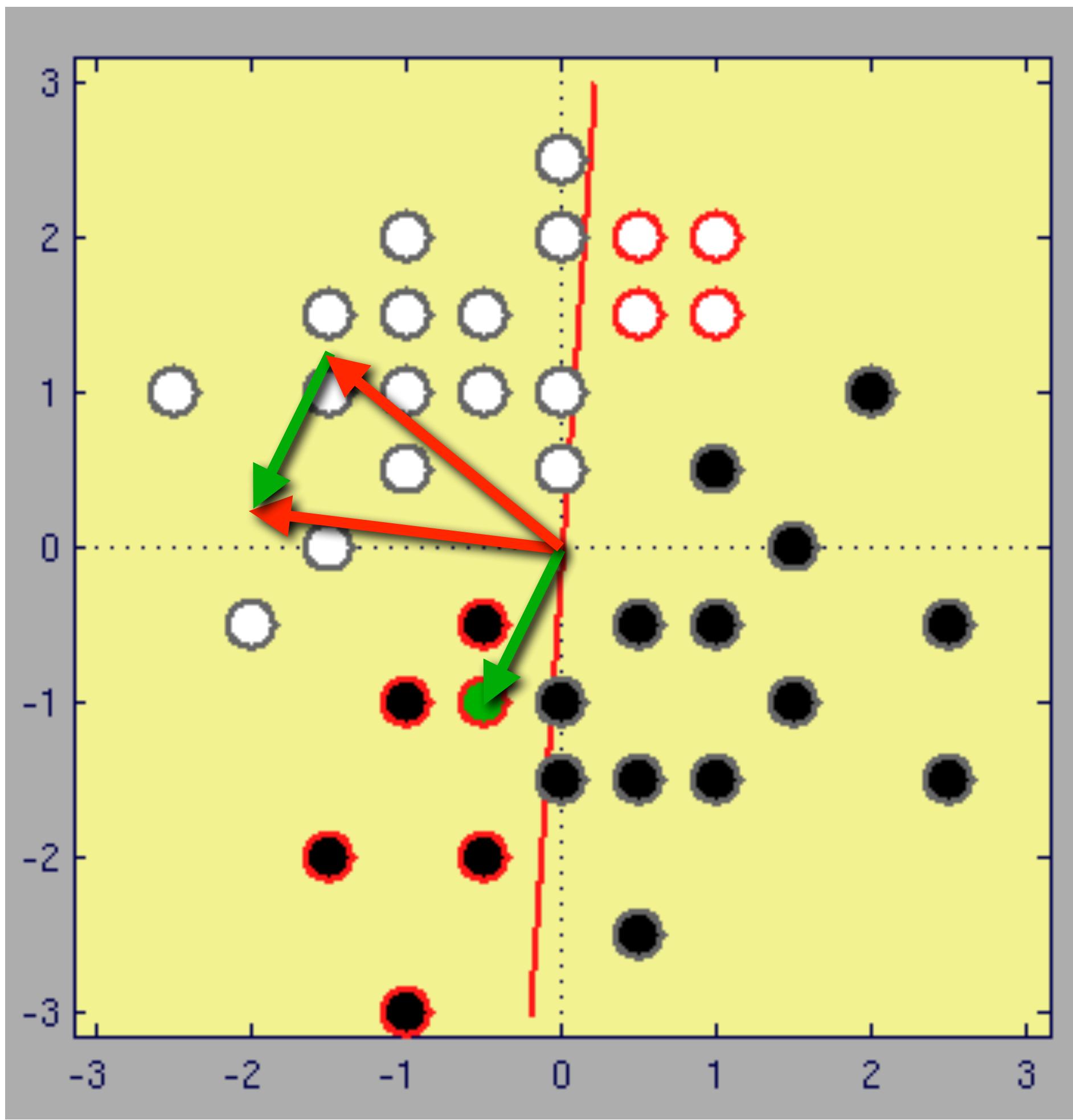
Example of Perceptron Learning



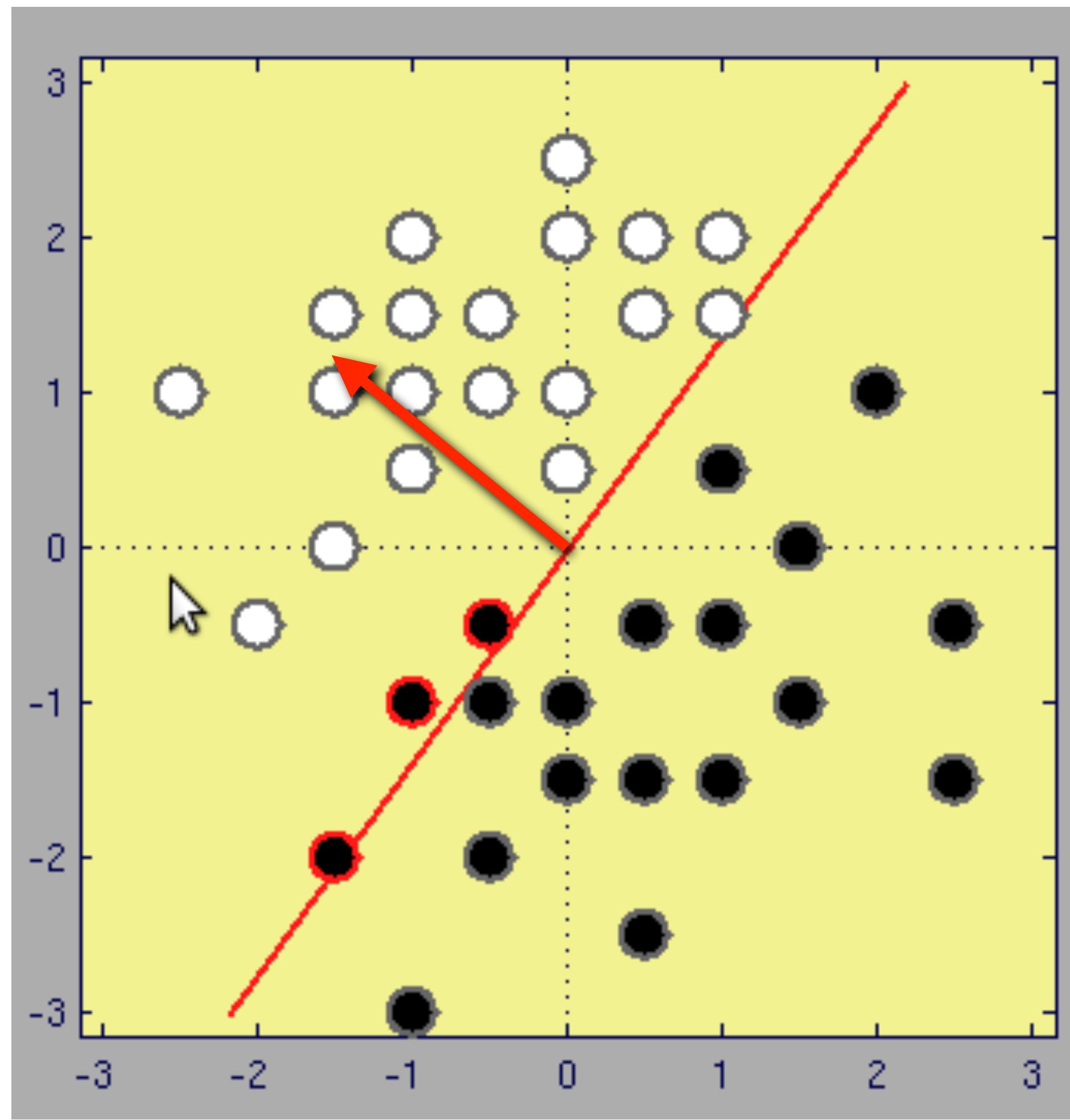
Example of Perceptron Learning



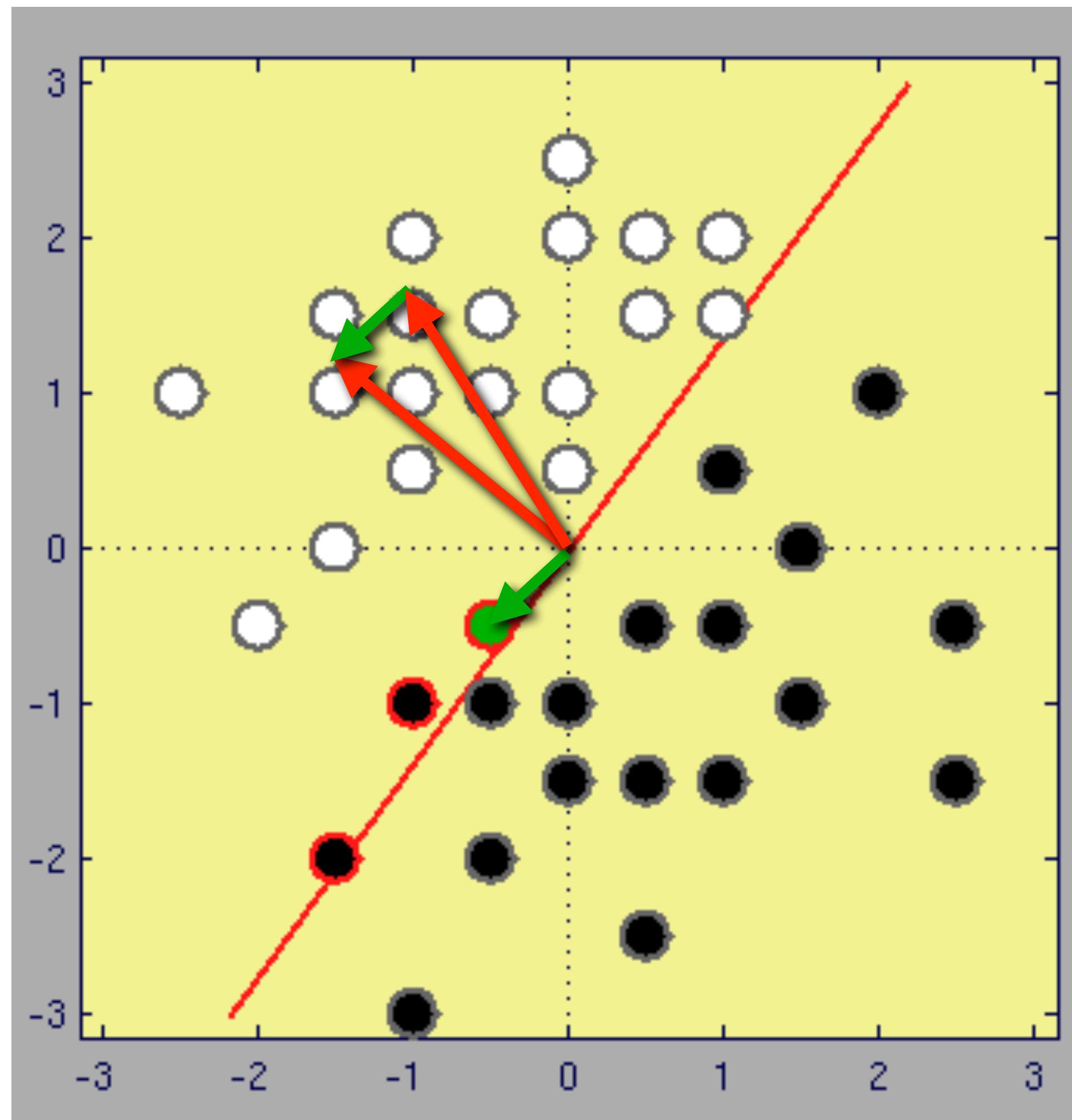
Example of Perceptron Learning



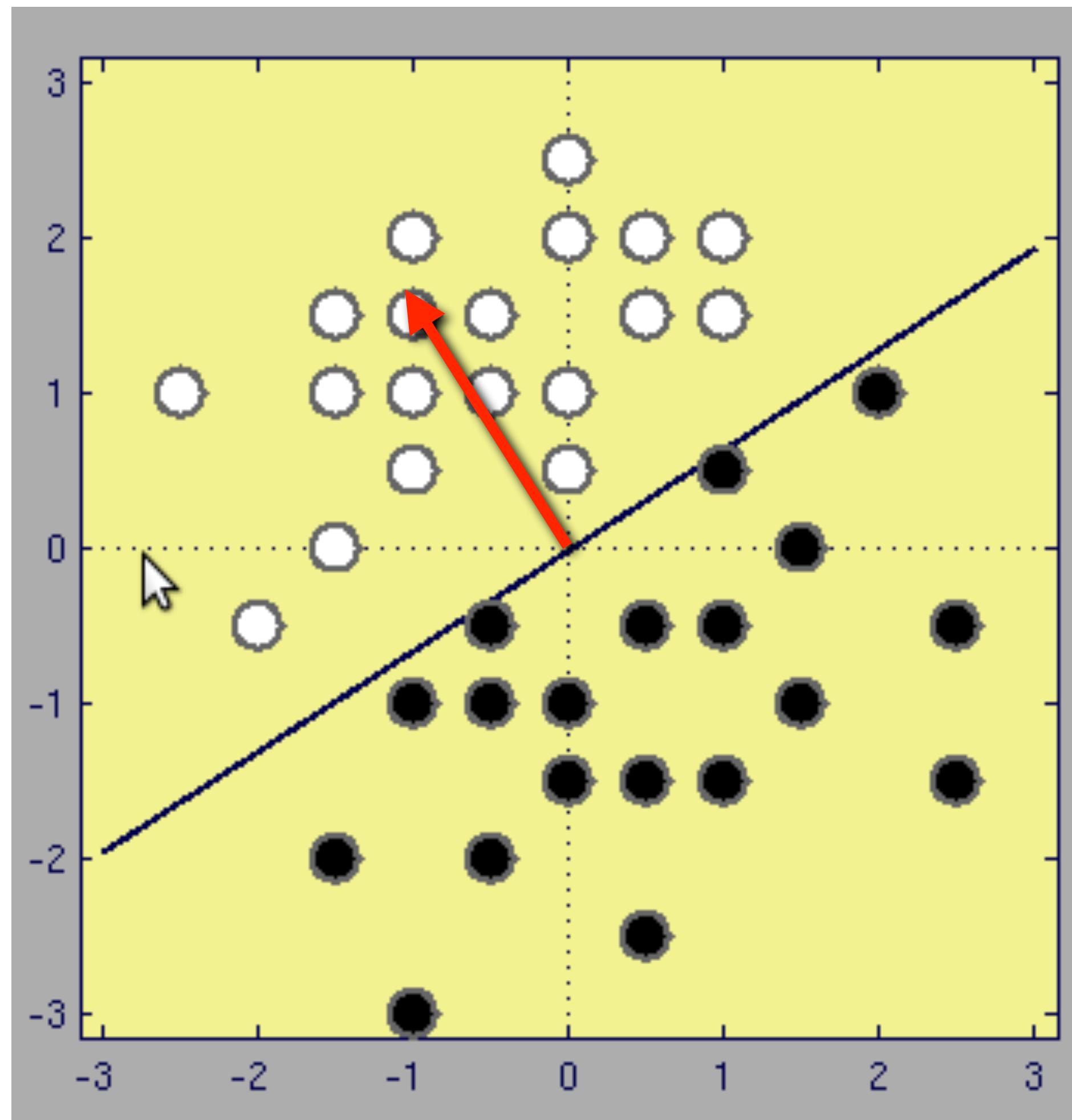
Example of Perceptron Learning



Example of Perceptron Learning

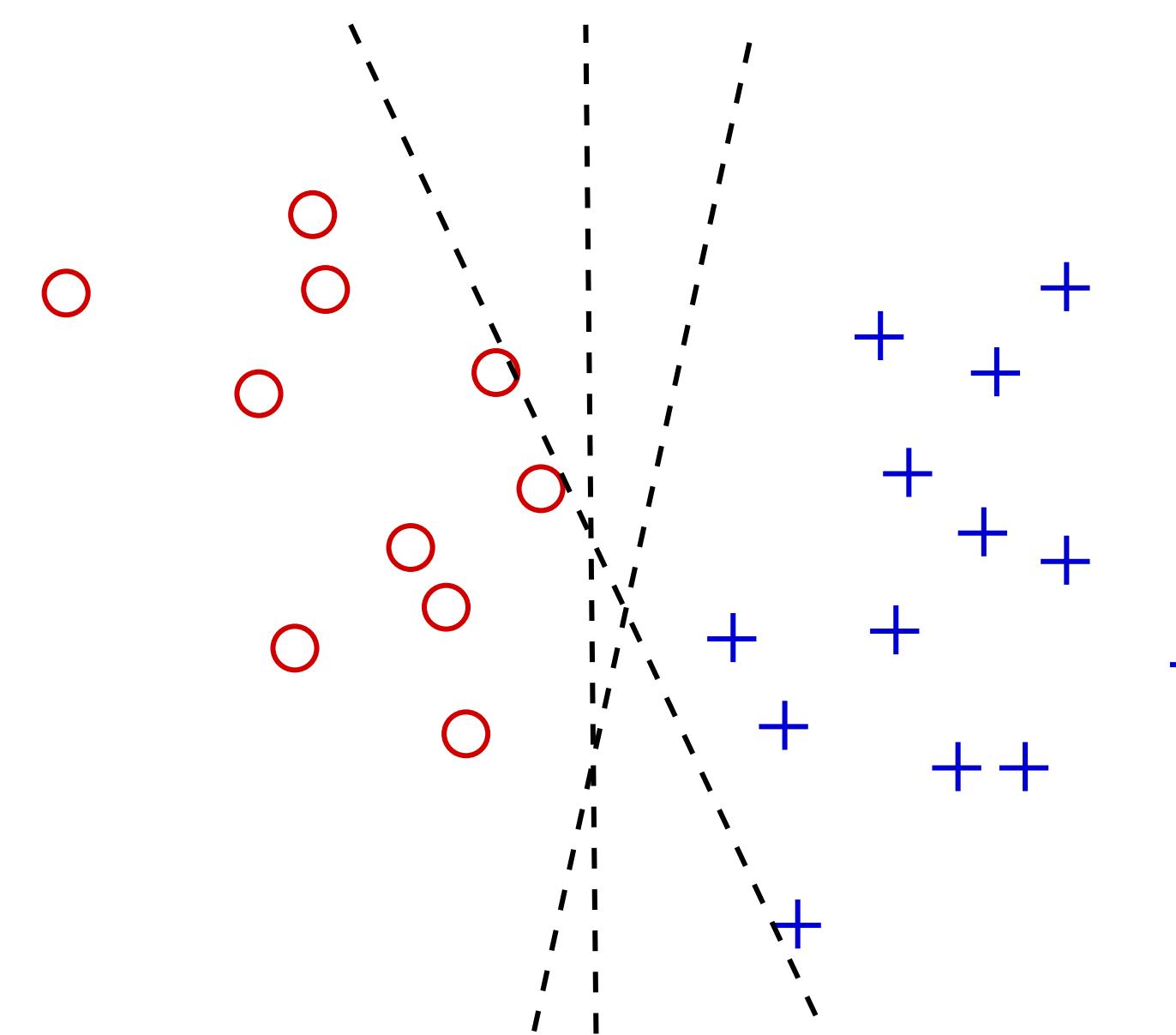


Example of Perceptron Learning

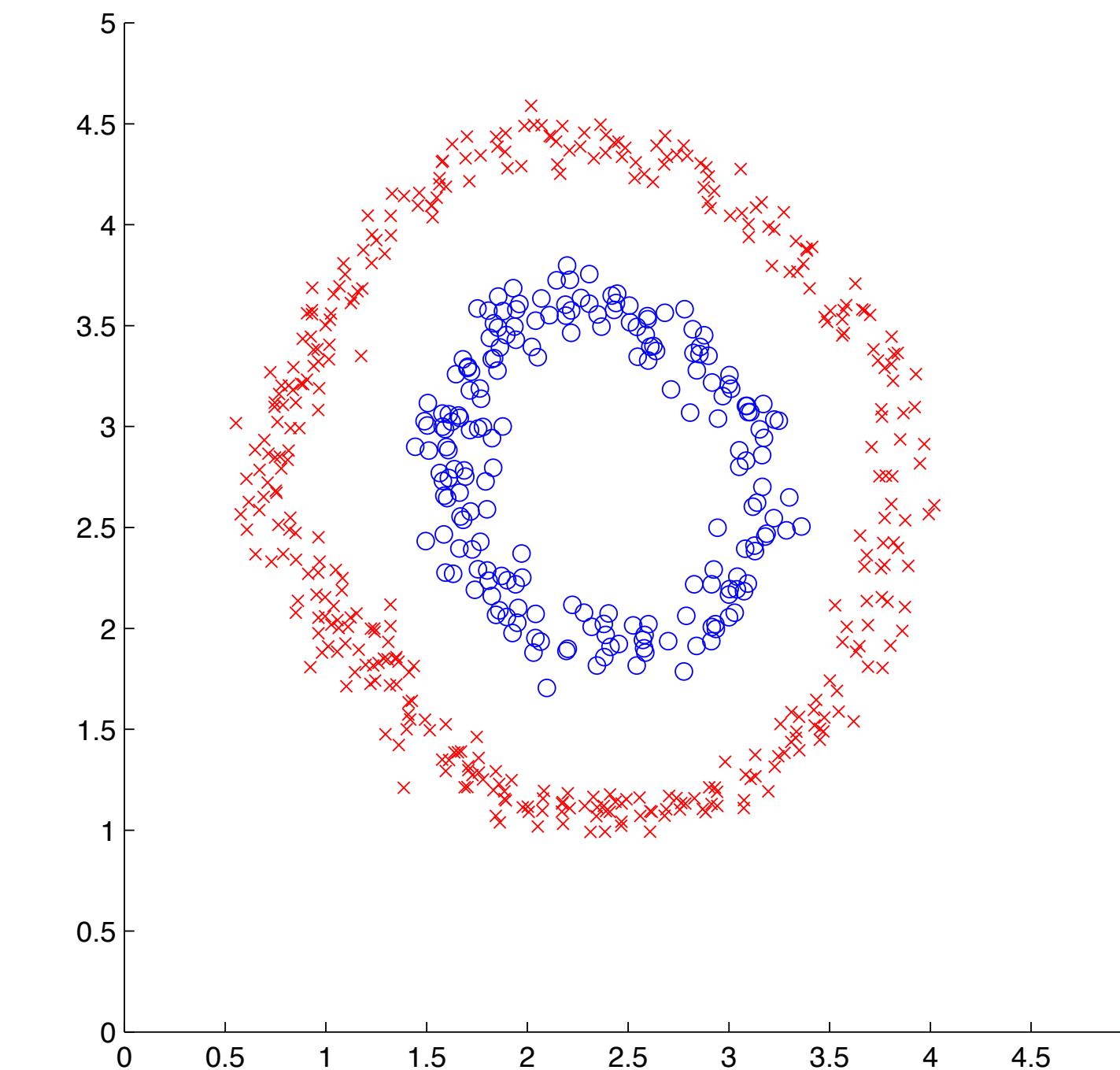


Perceptron Limitations

- Perceptrons + linear + softmax regressors are limited to data that are linearly separable, e.g.,



Linearly separable



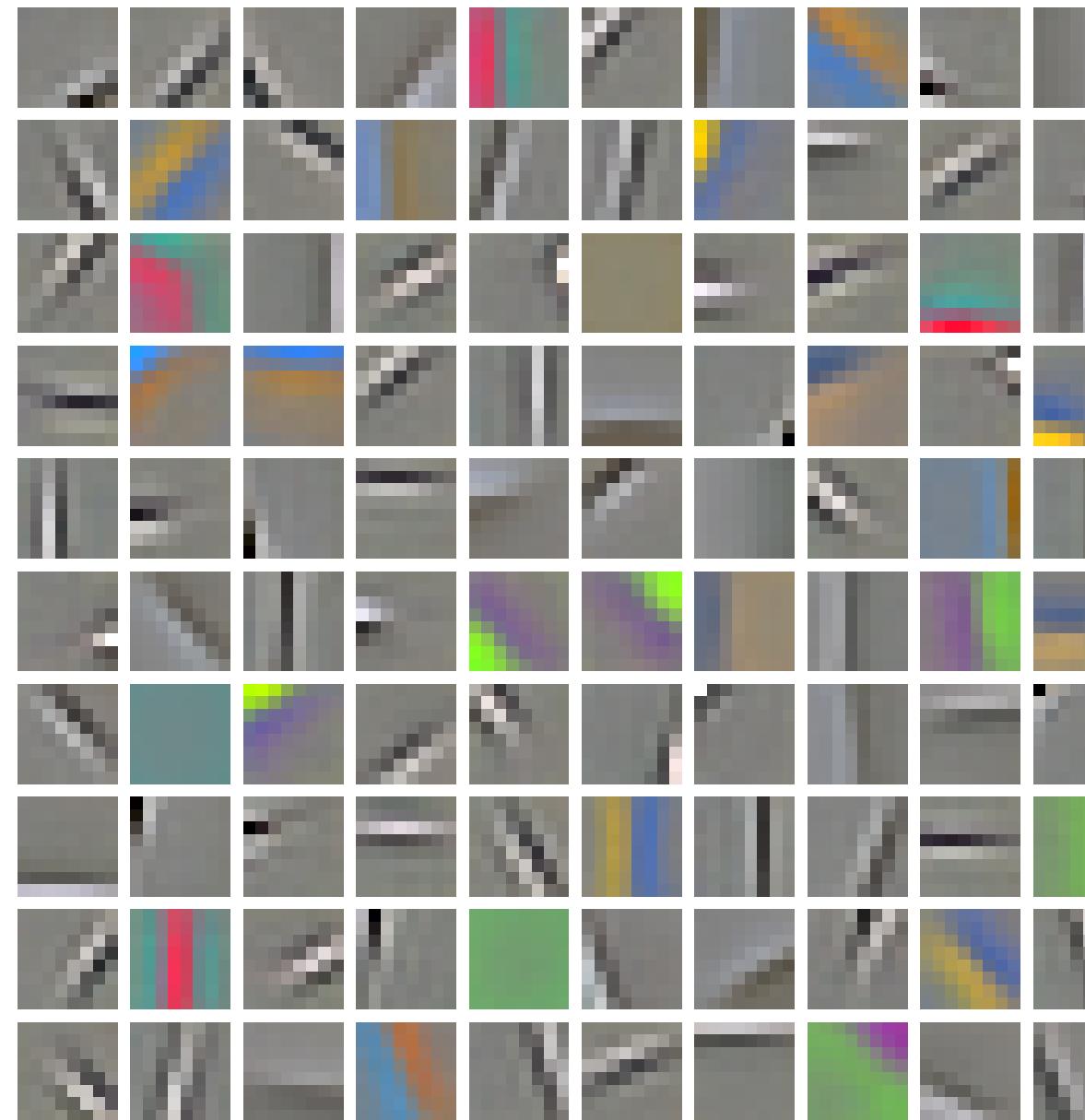
Not linearly separable



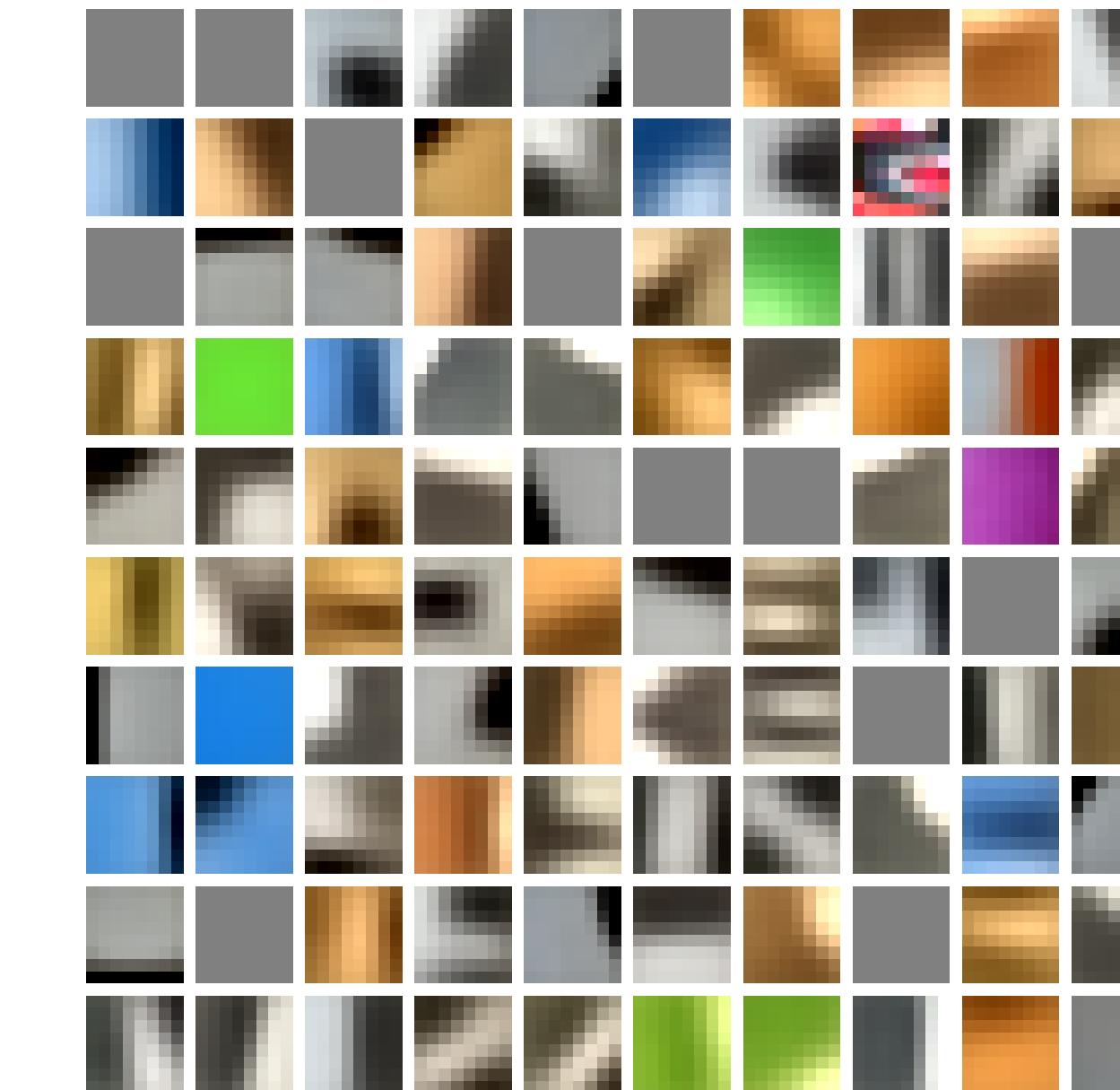
Could we extract features to make the data linearly separable?

CIFAR10 Feature Extraction

- So far, we used RGB pixels as the input to our classifier
- Feature extraction can improve results by a lot
- e.g., Coates et al. achieve 79.6% accuracy on CIFAR10 with features based on k-means of whitened image patches



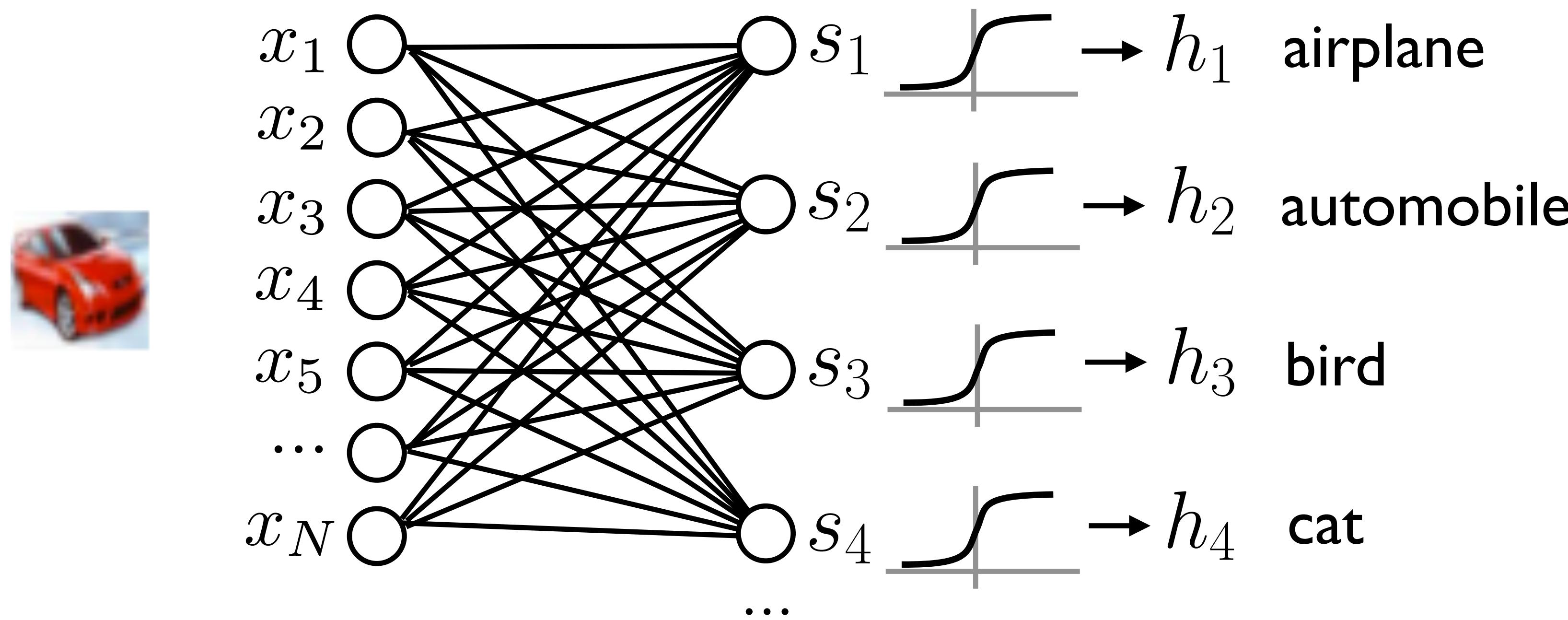
k-means, whitened



k-means, raw RGB

Linear = Fully Connected Layer

- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network

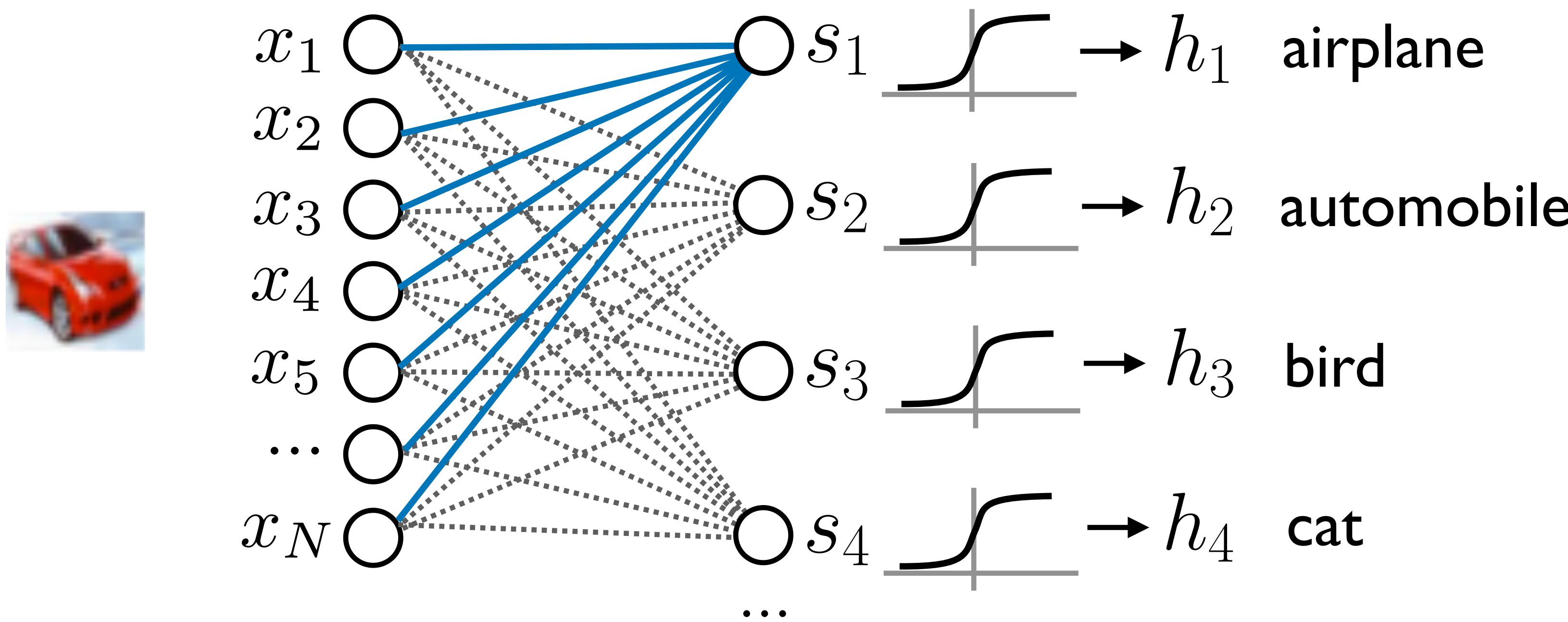


- Typically, we'll also add a bias term b

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

Linear = Fully Connected Layer

- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network

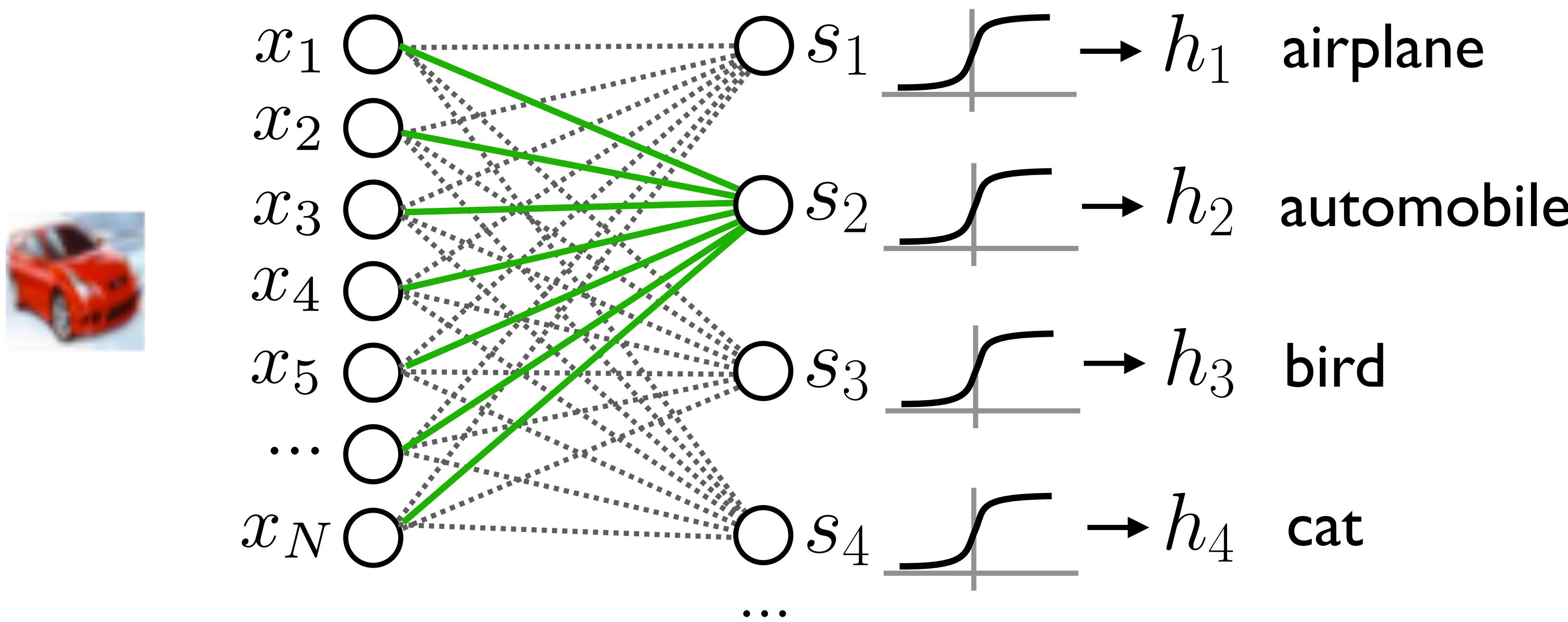


- Typically, we'll also add a bias term b

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

Linear = Fully Connected Layer

- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network

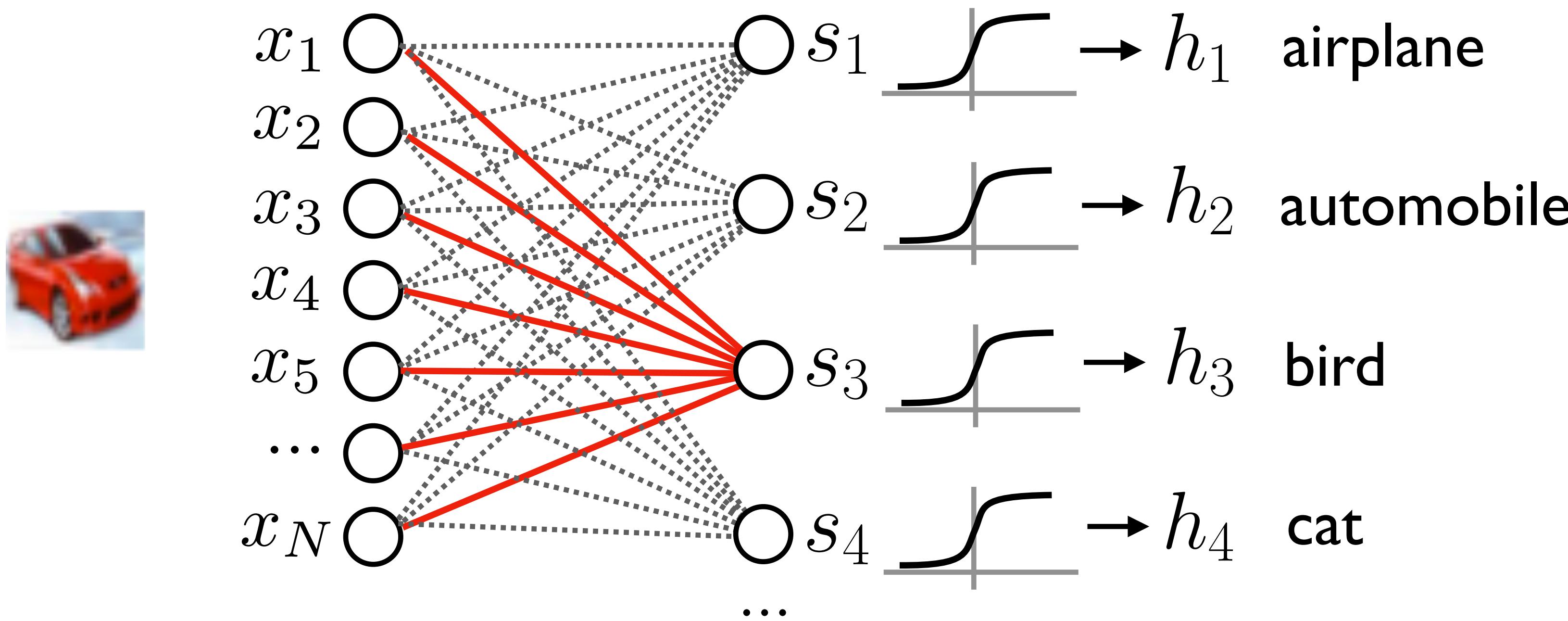


- Typically, we'll also add a bias term b

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

Linear = Fully Connected Layer

- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network

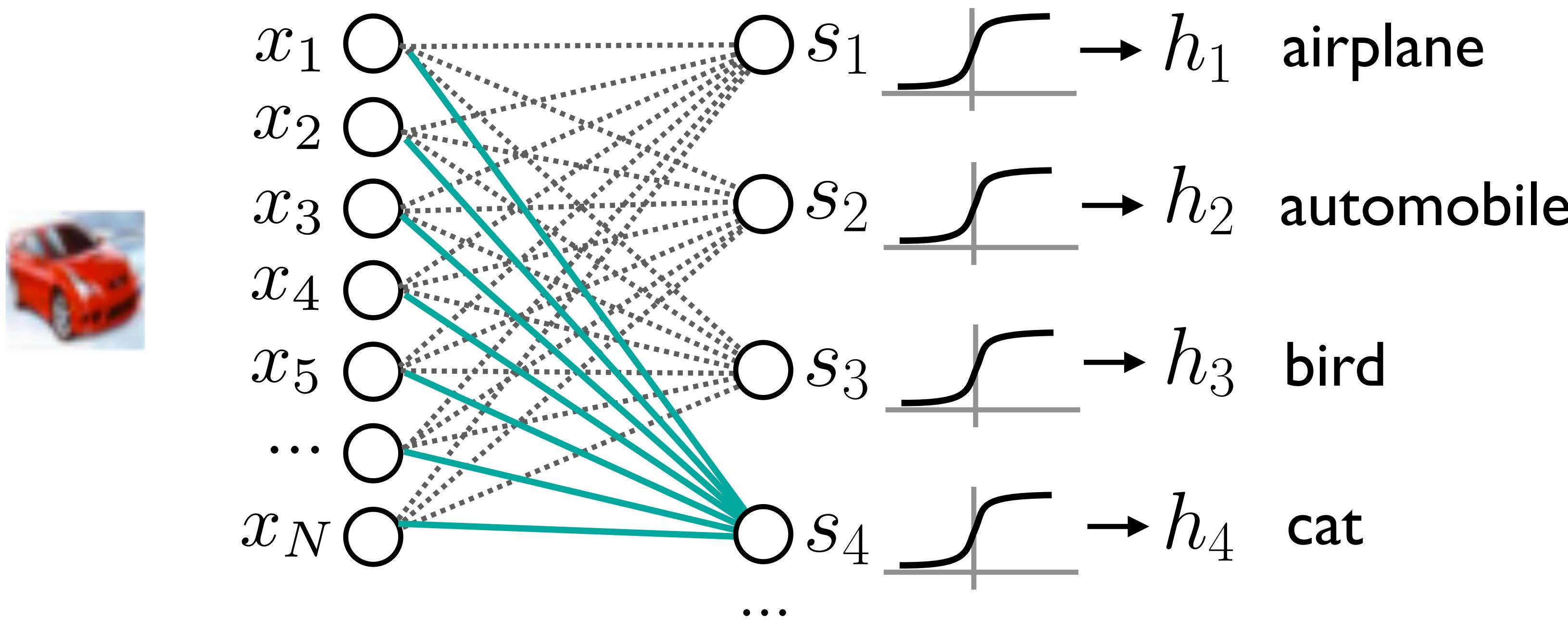


- Typically, we'll also add a bias term b

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

Linear = Fully Connected Layer

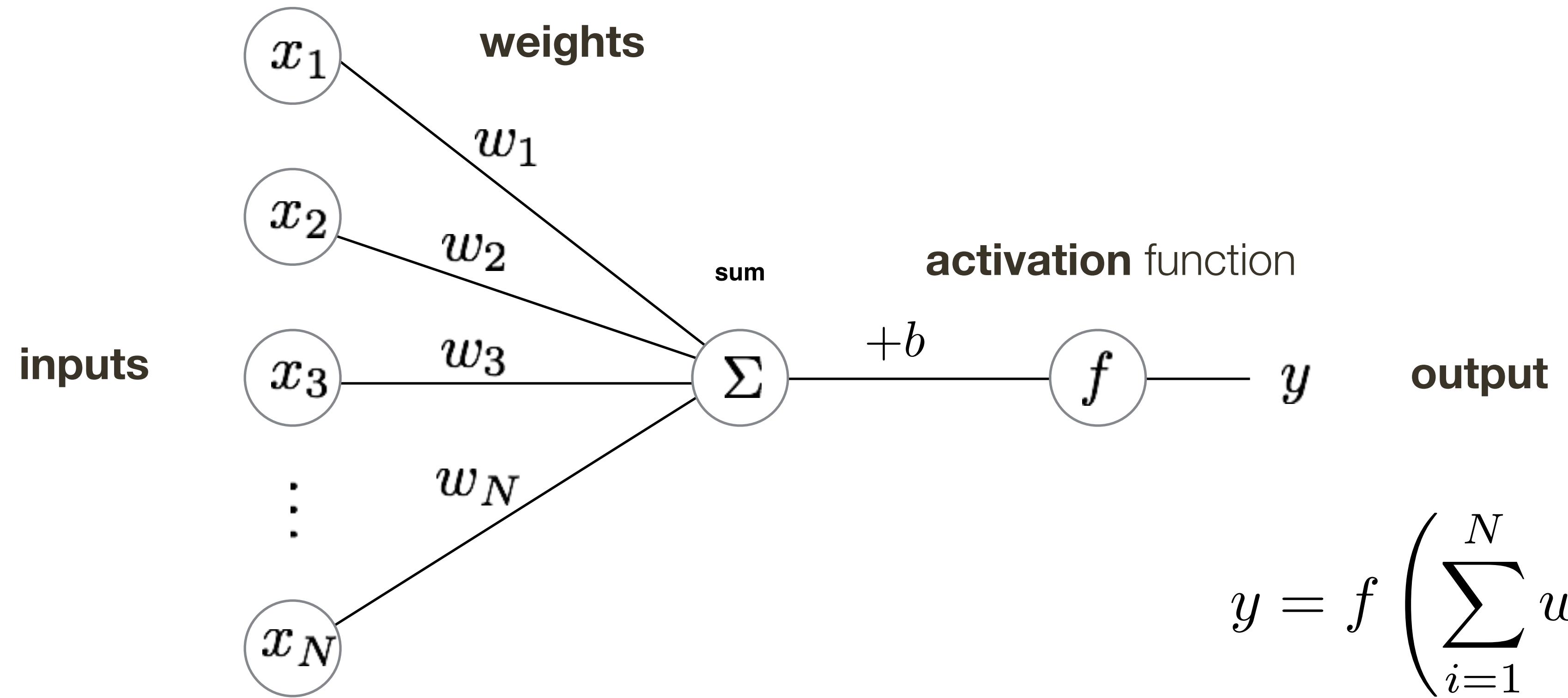
- Note that our linear matrix multiplication classifier is equivalent to a fully connected layer in a neural network



- Typically, we'll also add a bias term b

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

A Neuron



- The basic unit of computation in a neural network is a neuron.
- A neuron accepts some number of input signals, computes their weighted sum, and applies an **activation function** (or **non-linearity**) to the sum.
- Common activation functions include sigmoid and rectified linear unit (ReLU)

Activation Function: Sigmoid

$$f(x) = 1/(1 + e^{-x})$$

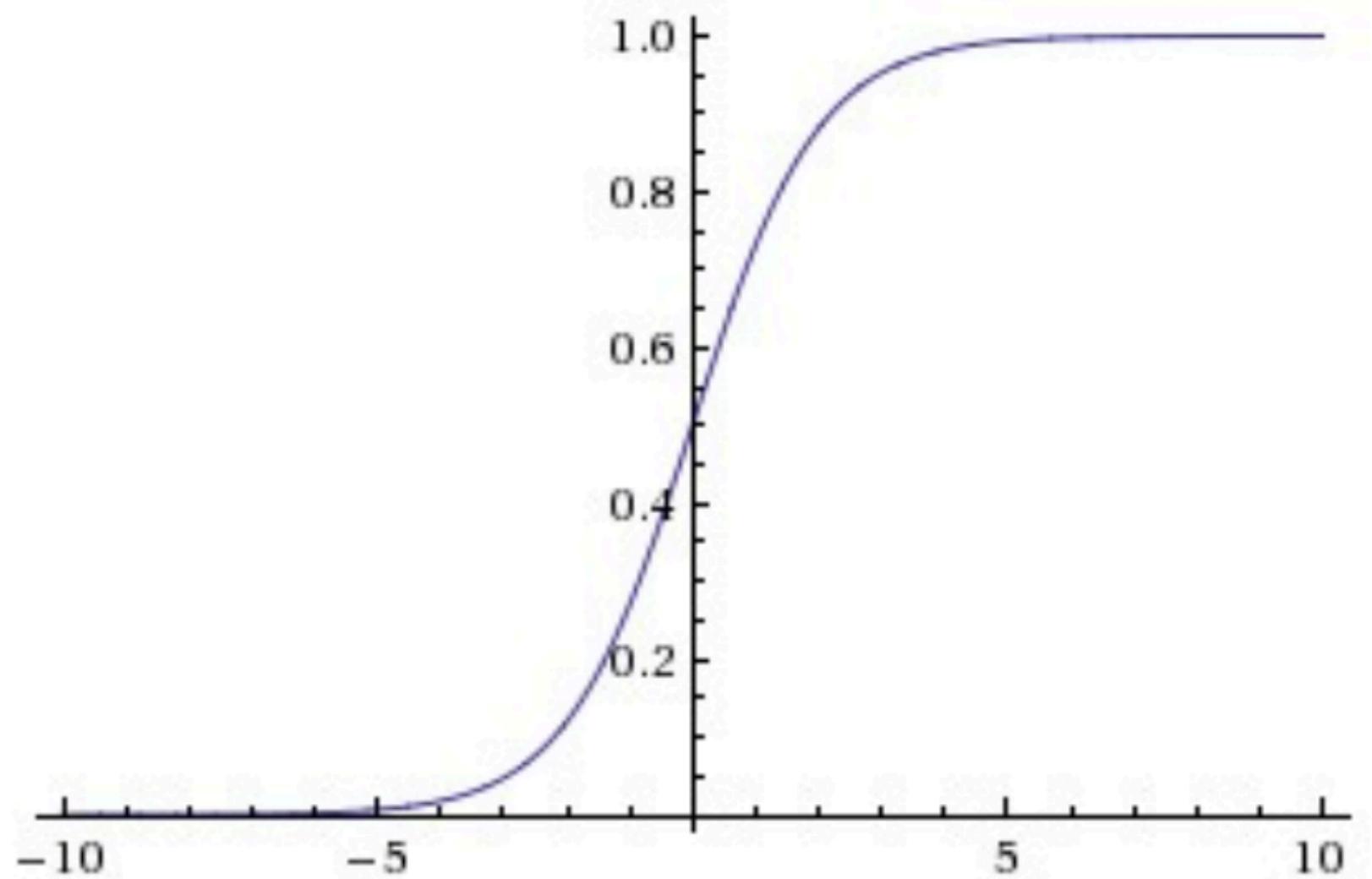


Figure credit: Fei-Fei and Karpathy

Common in many early neural networks

Biological analogy to saturated firing rate of neurons

Maps the input to the range [0,1]

Activation Function: **ReLU** (Rectified Linear Unit)

$$f(x) = \max(0, x)$$

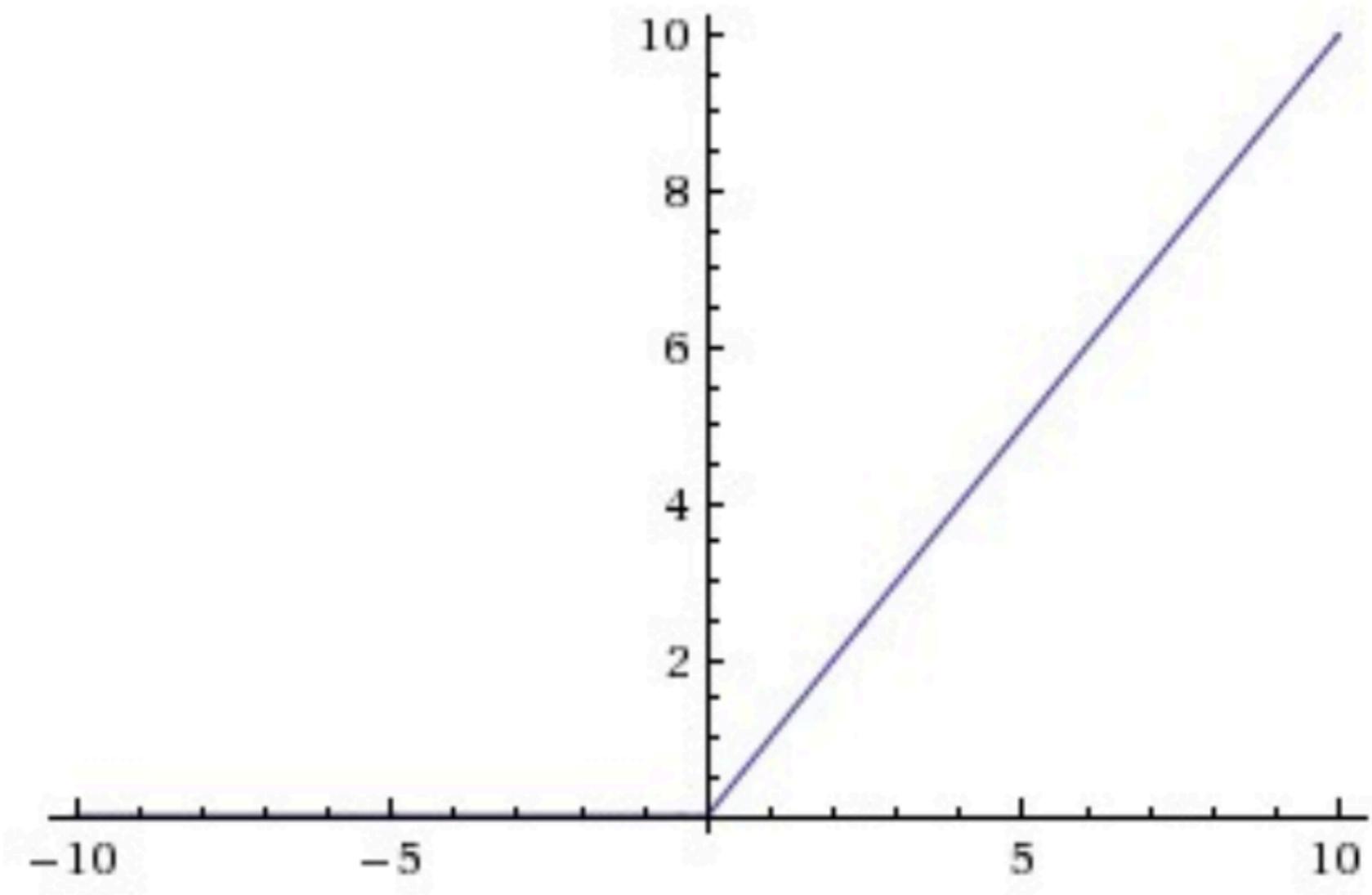
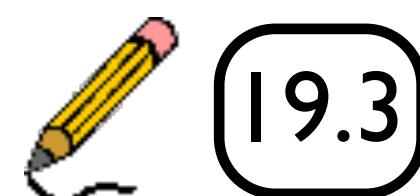


Figure credit: Fei-Fei and Karpathy

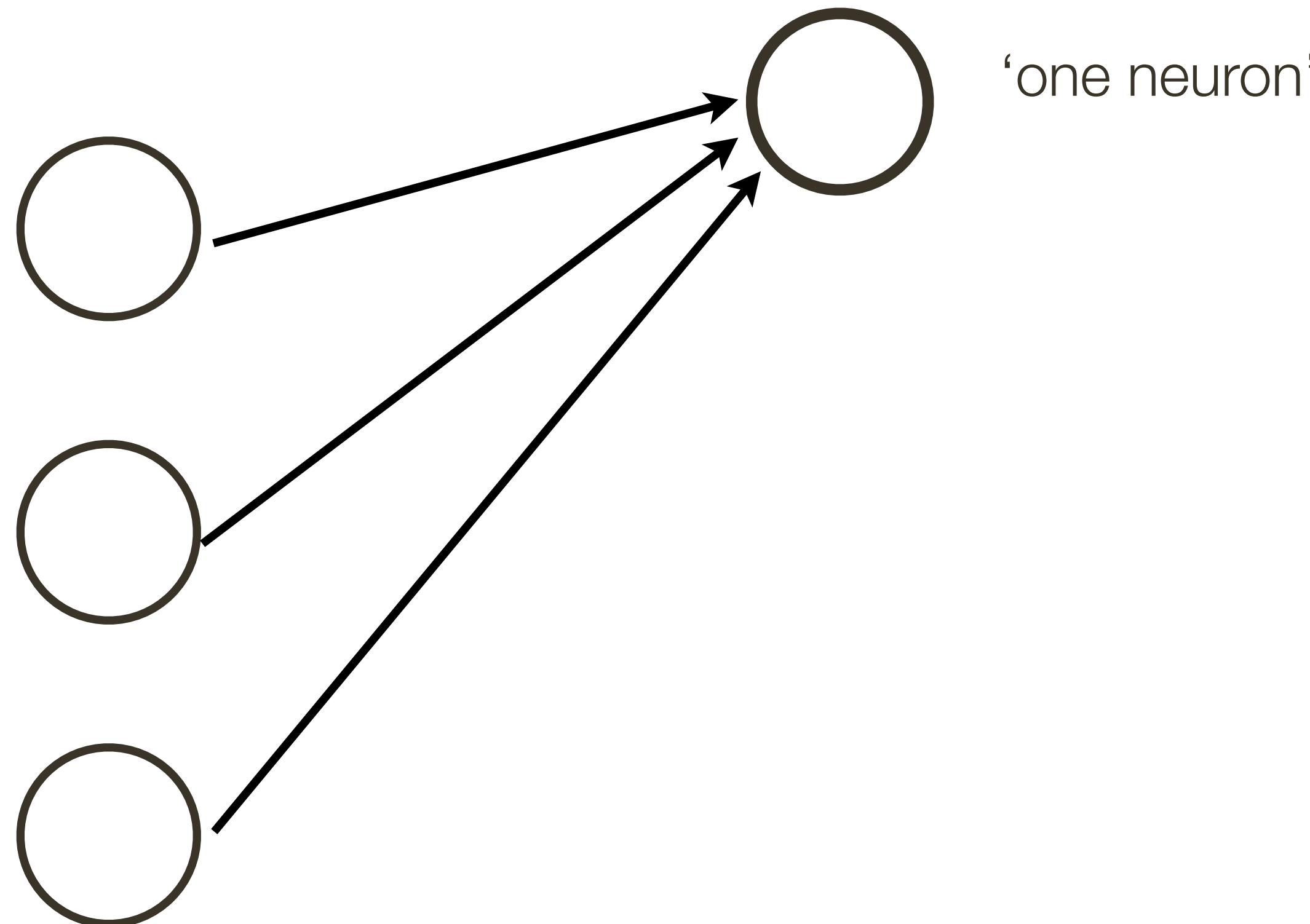
Maintains good gradient flow in networks, prevents vanishing gradient problem
Very commonly used in interior (hidden) layers of neural nets



Why can't we have **linear** activation functions?

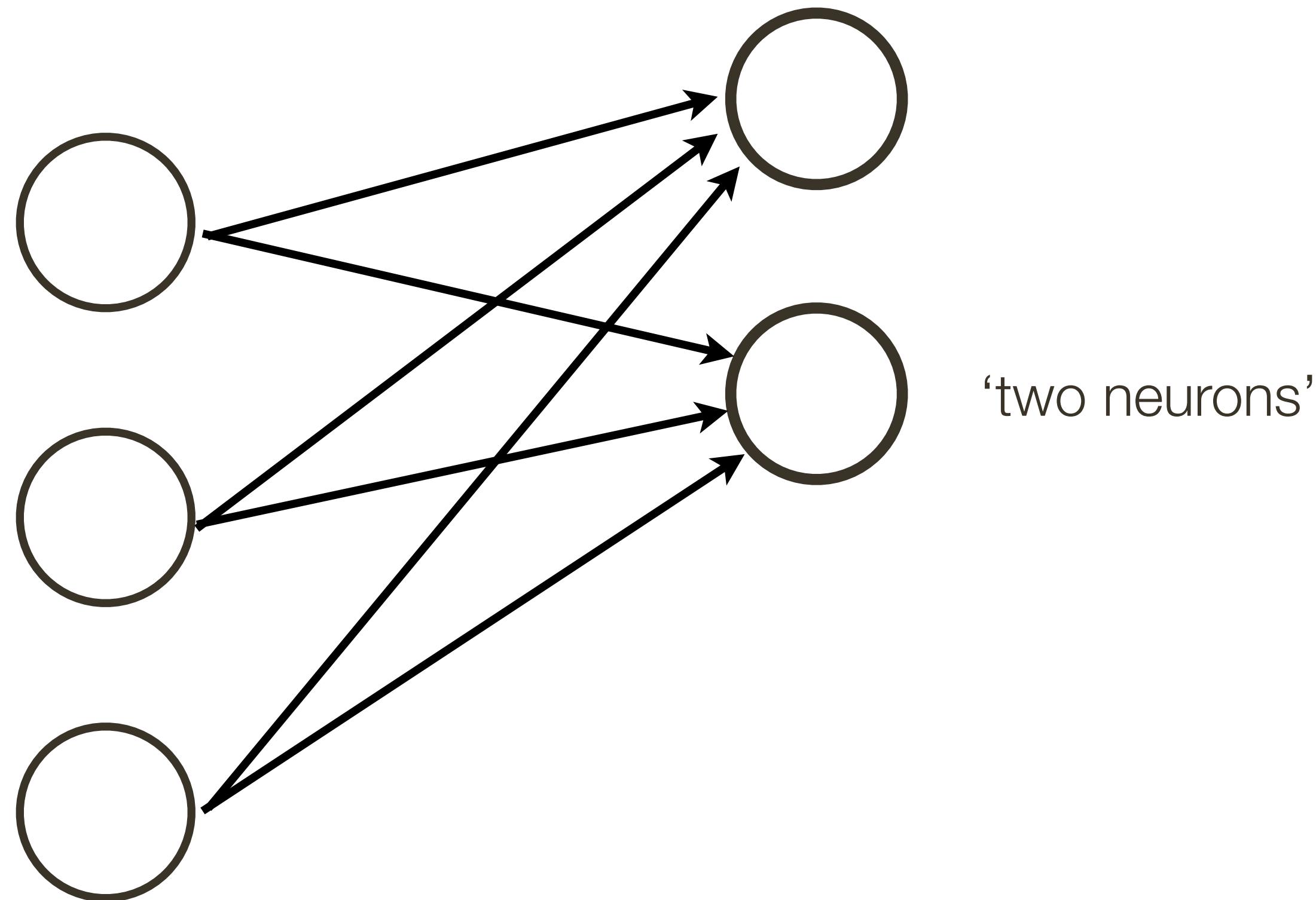
Neural Network

Connect a bunch of neurons together – a collection of connected neurons



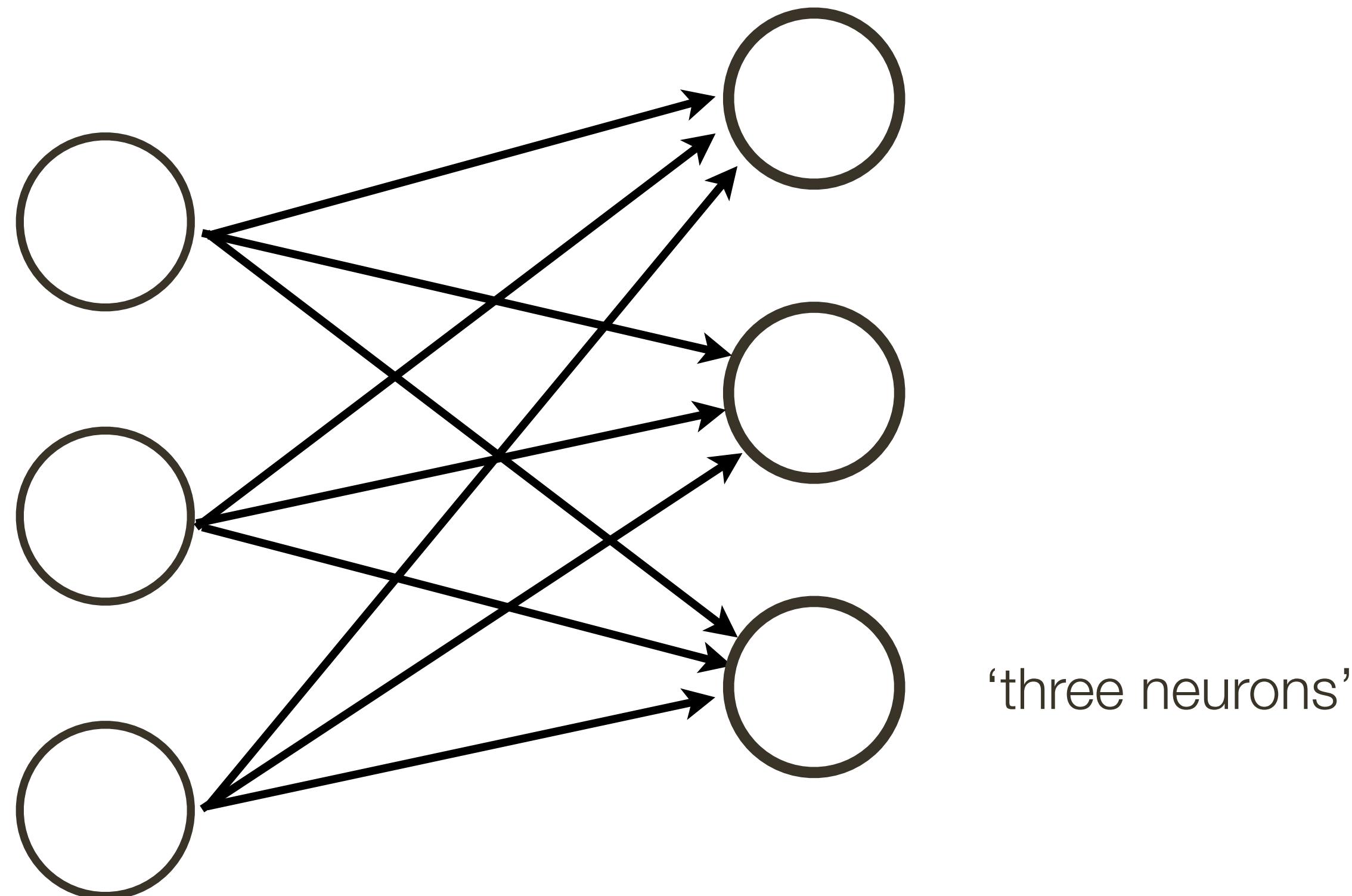
Neural Network

Connect a bunch of neurons together – a collection of connected neurons



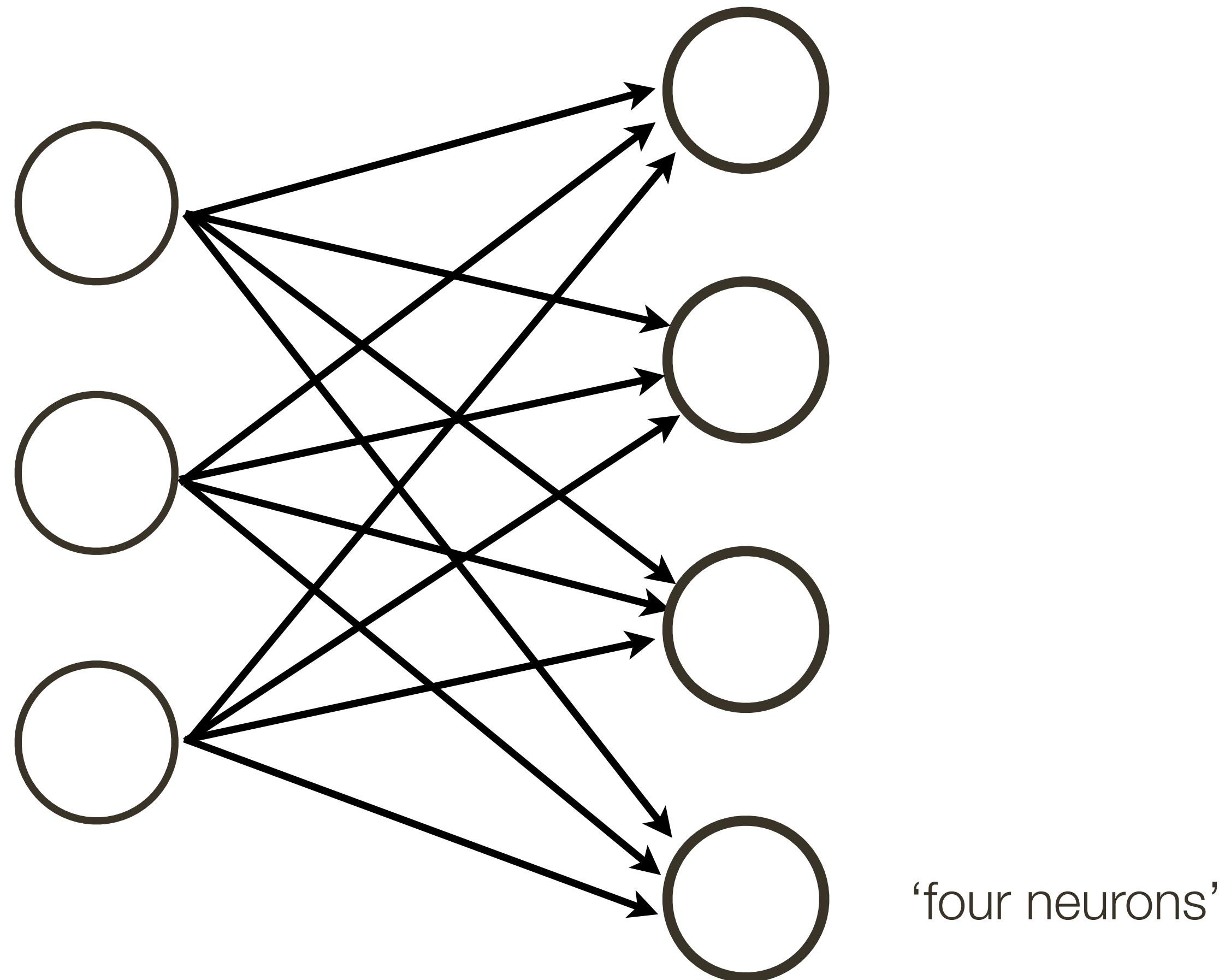
Neural Network

Connect a bunch of neurons together – a collection of connected neurons



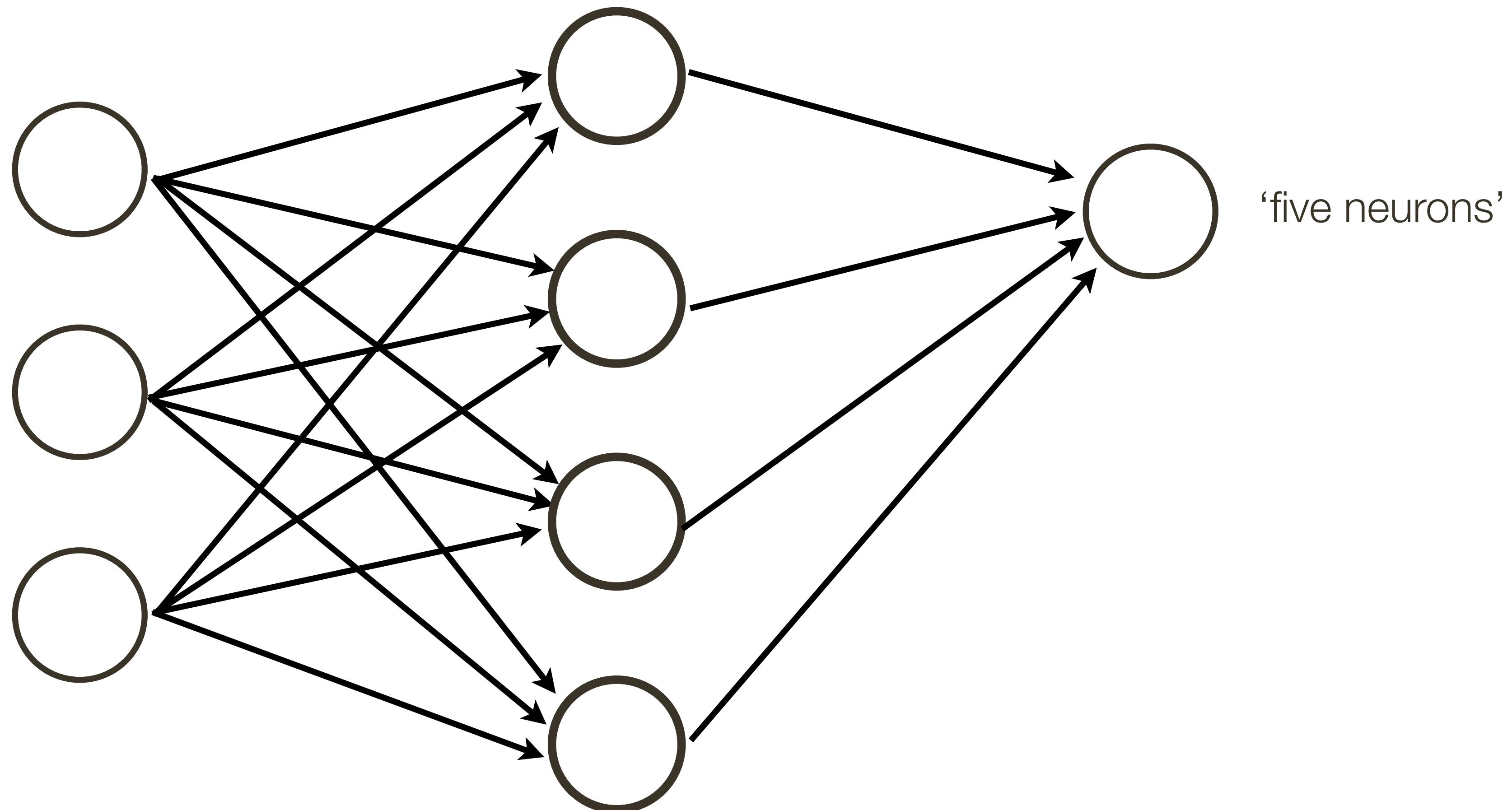
Neural Network

Connect a bunch of neurons together – a collection of connected neurons



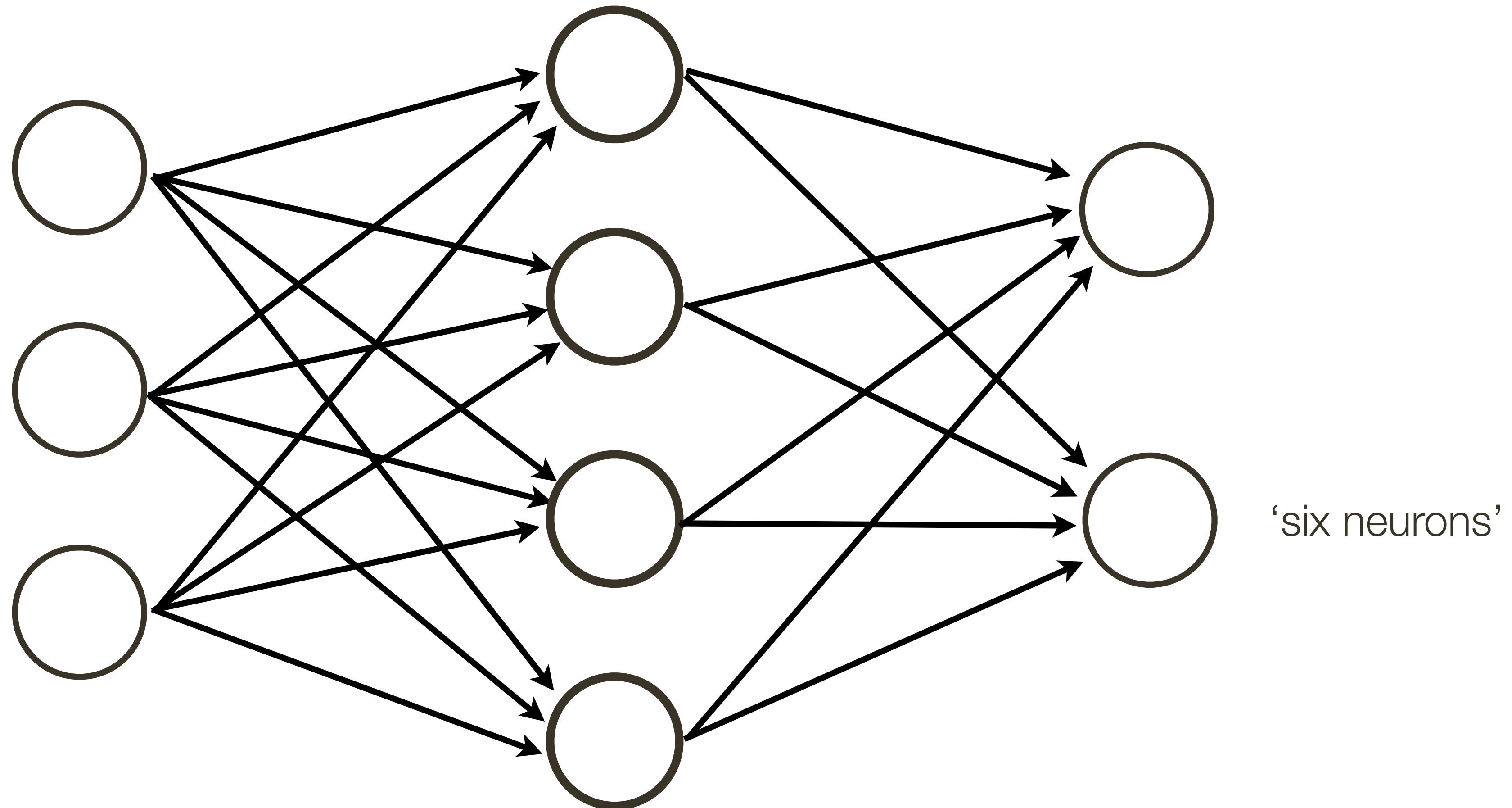
Neural Network

Connect a bunch of neurons together – a collection of connected neurons



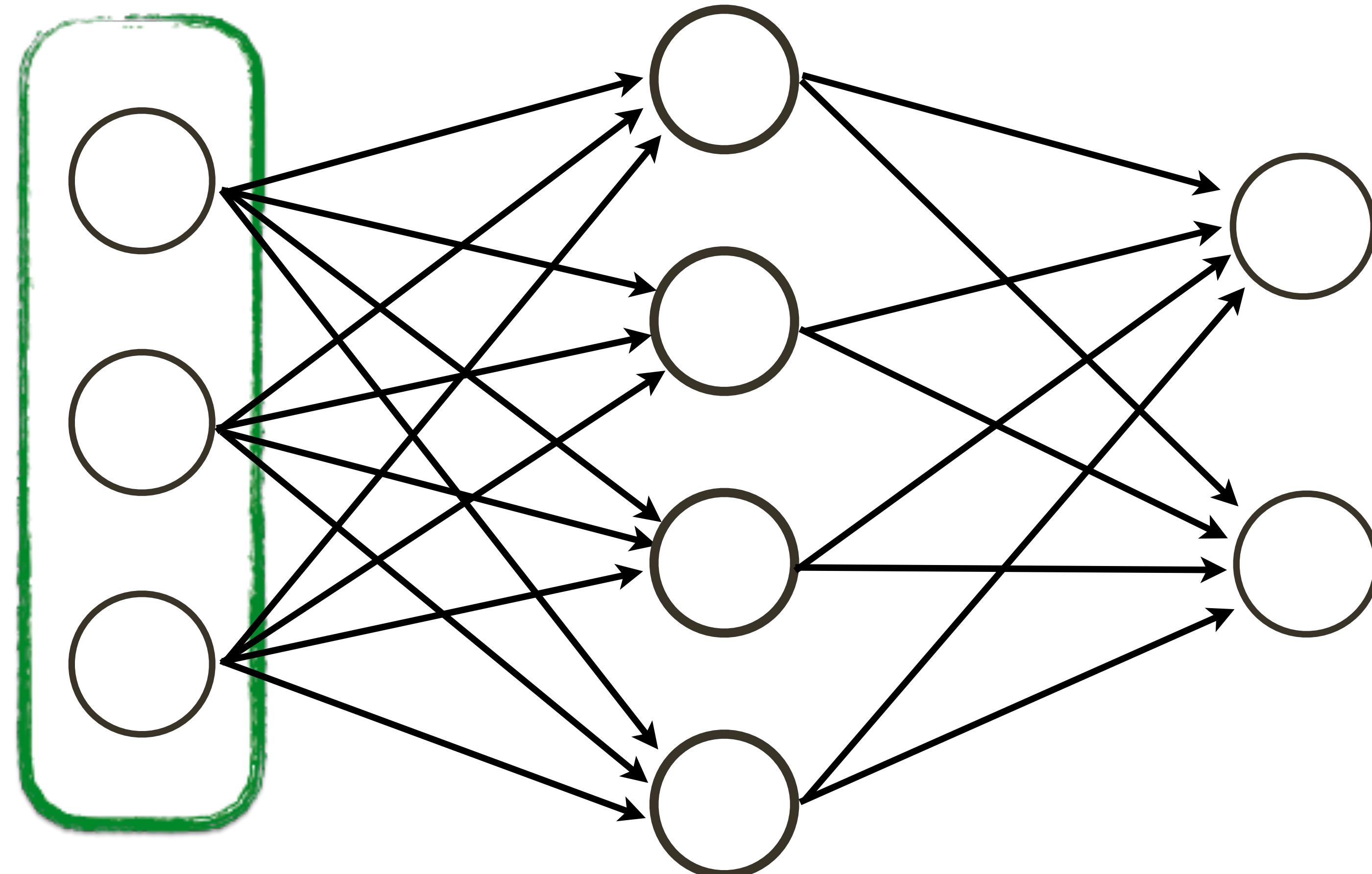
Neural Network

Connect a bunch of neurons together – a collection of connected neurons

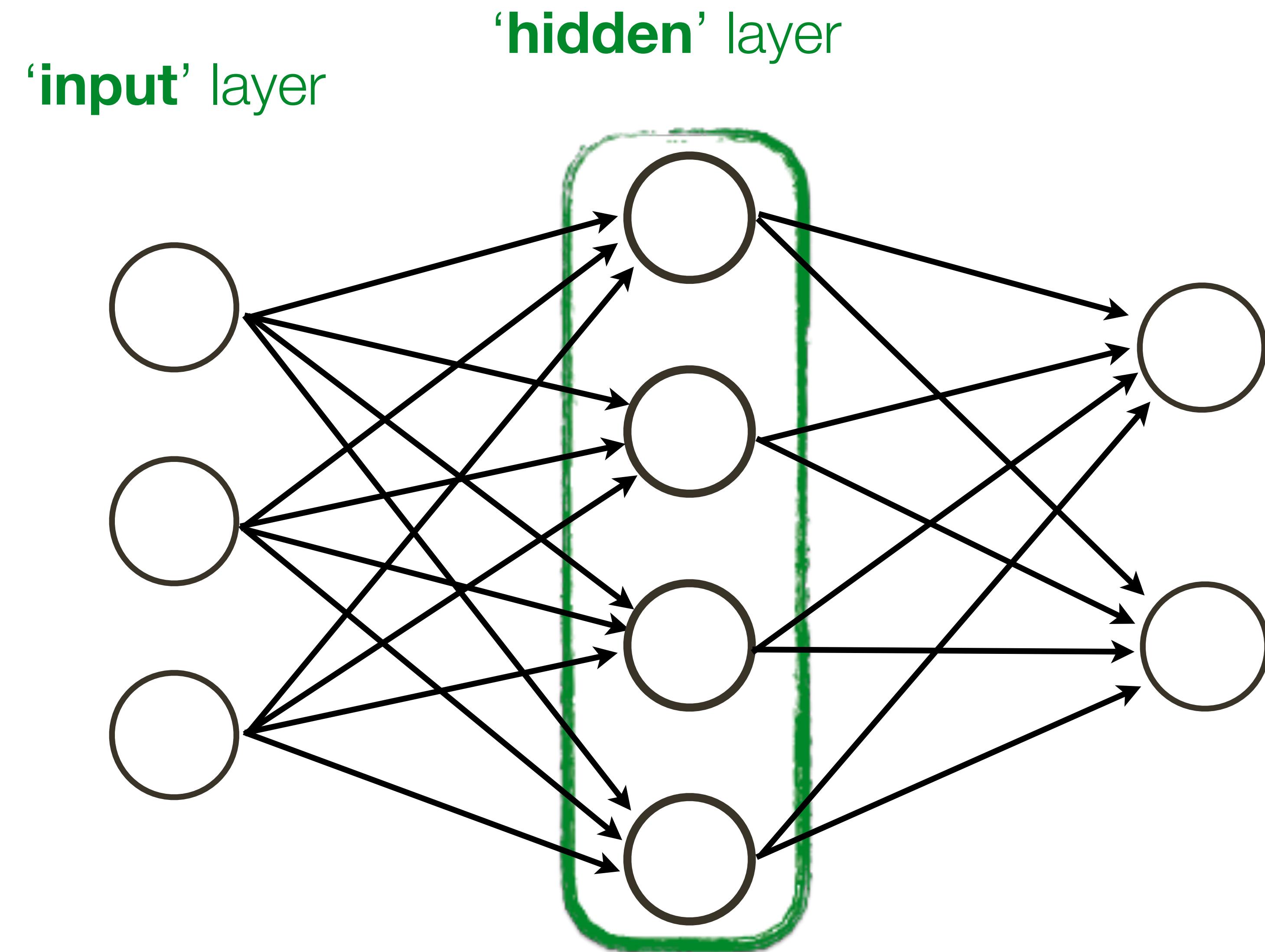


Neural Network: Terminology

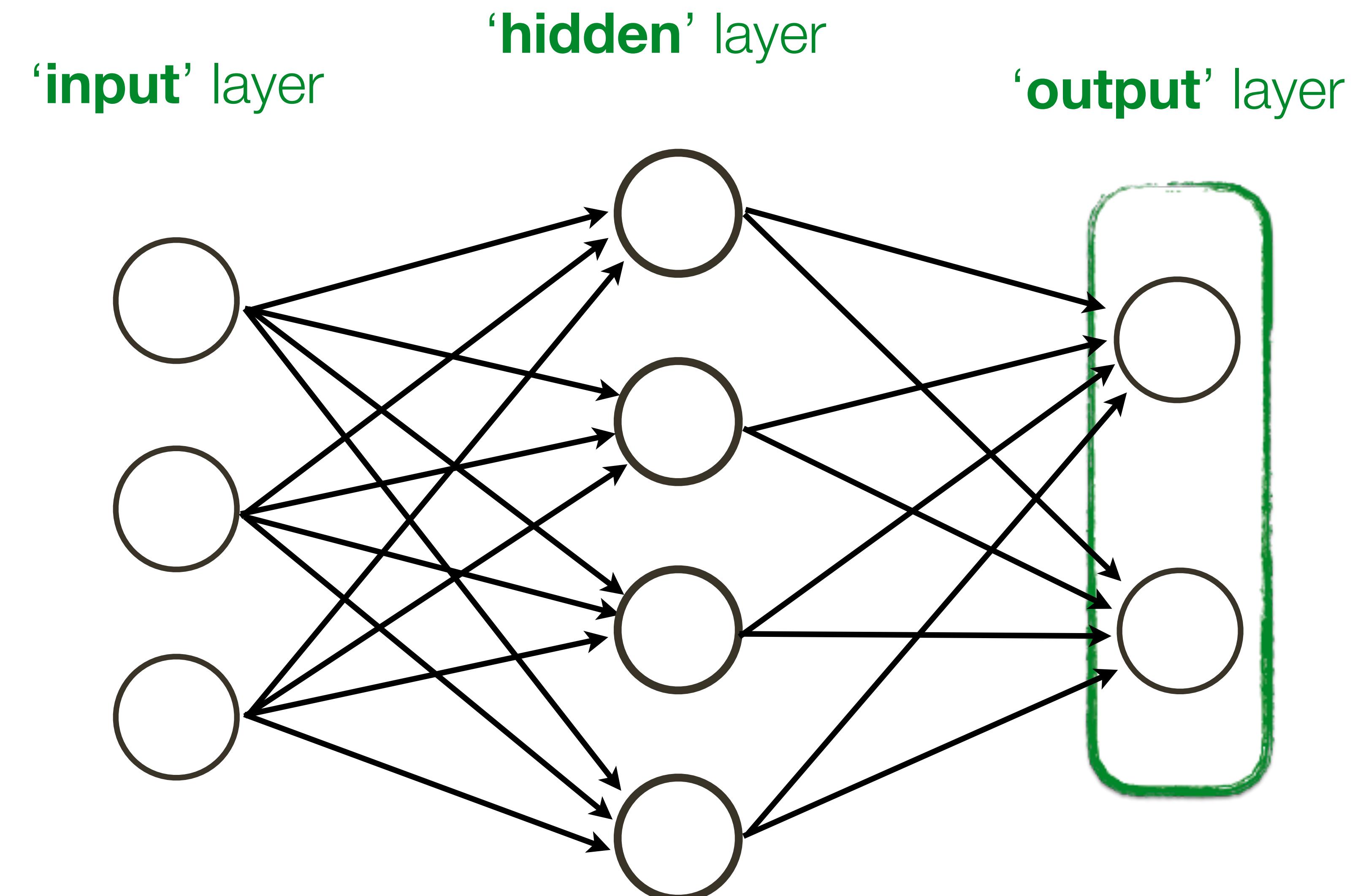
'input' layer



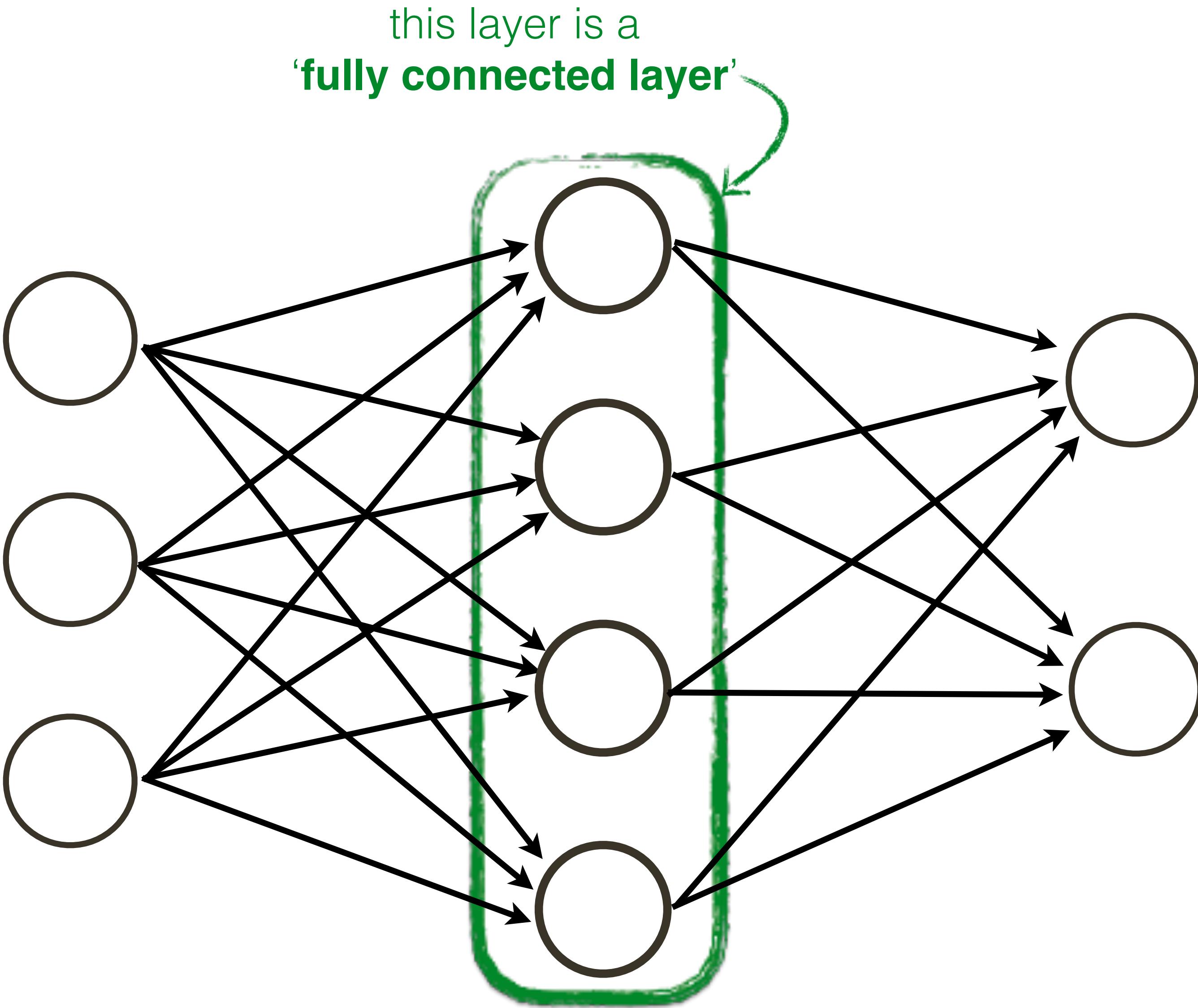
Neural Network: Terminology



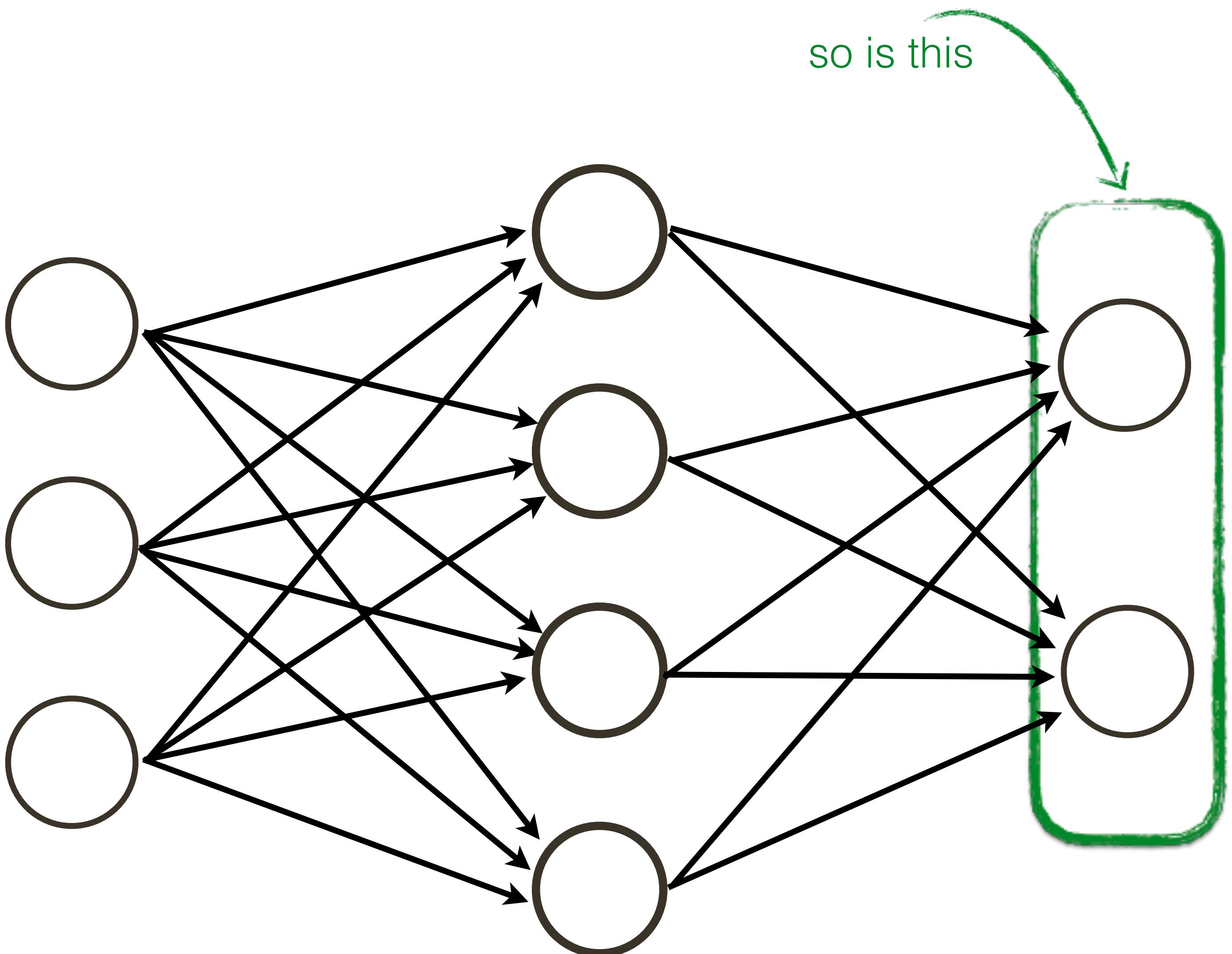
Neural Network: Terminology



Neural Network: Terminology

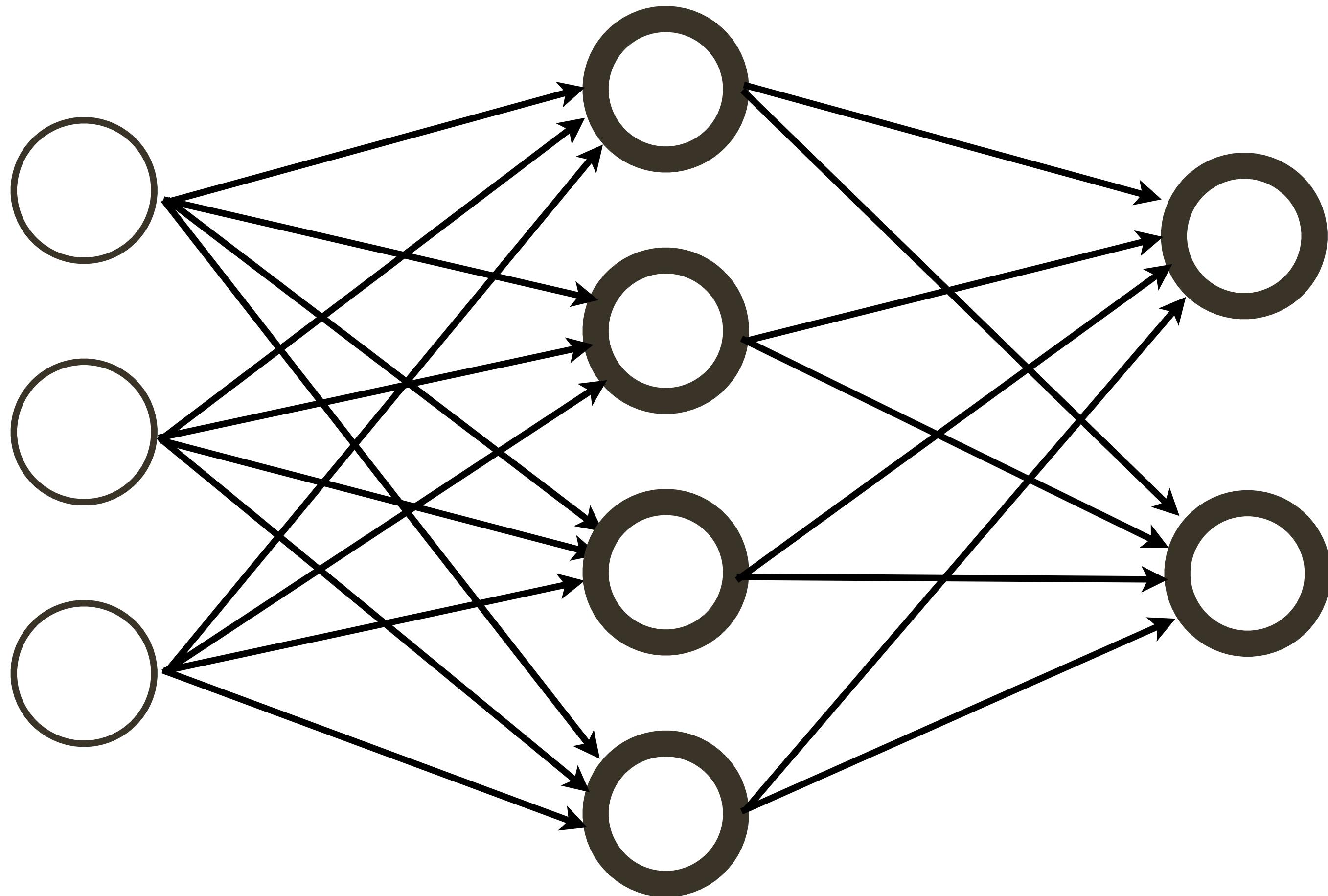


Neural Network: Terminology



Neural Network

How many neurons? $4+2 = 6$

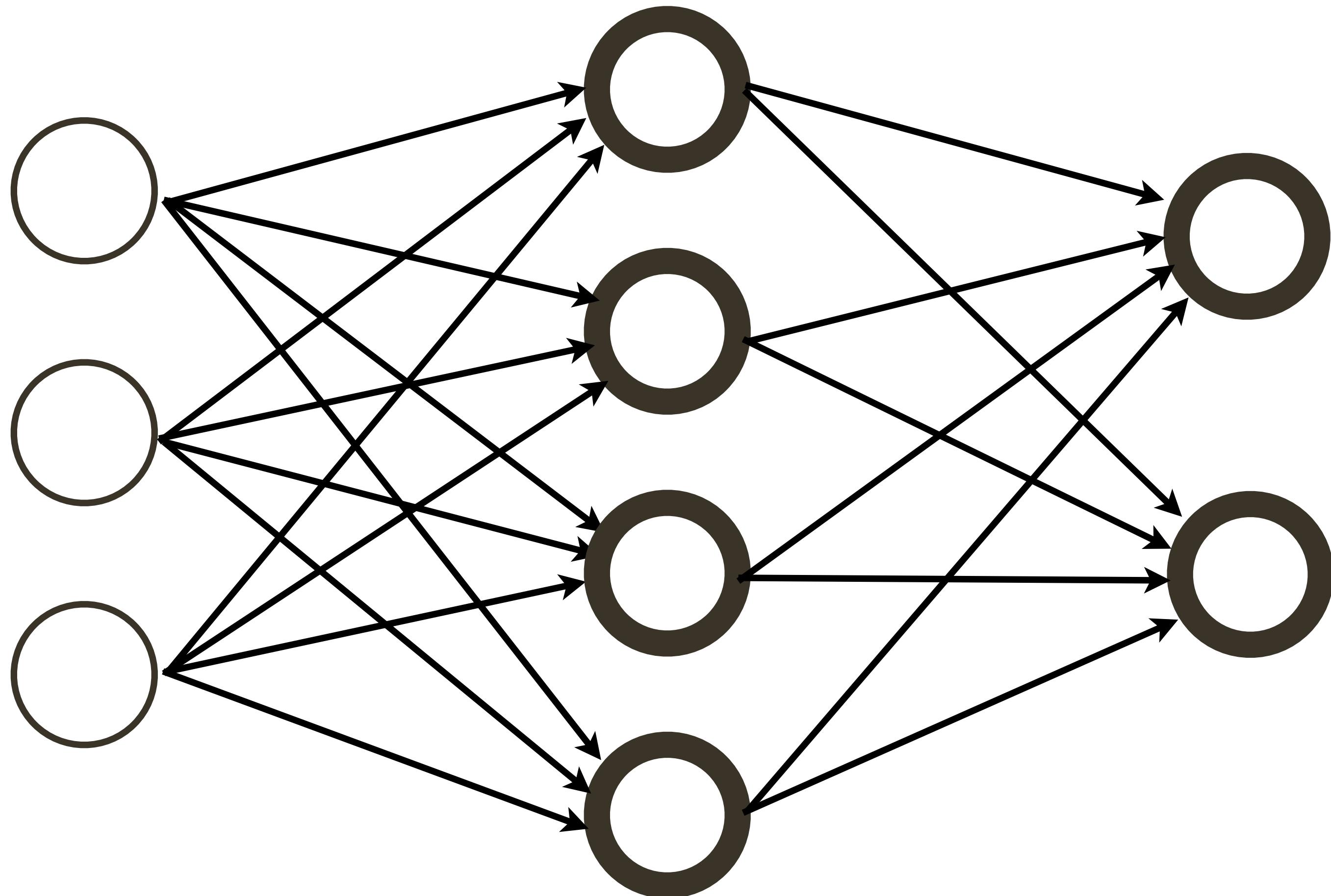


Neural Network

How many neurons?

$$4+2 = 6$$

How many weights?



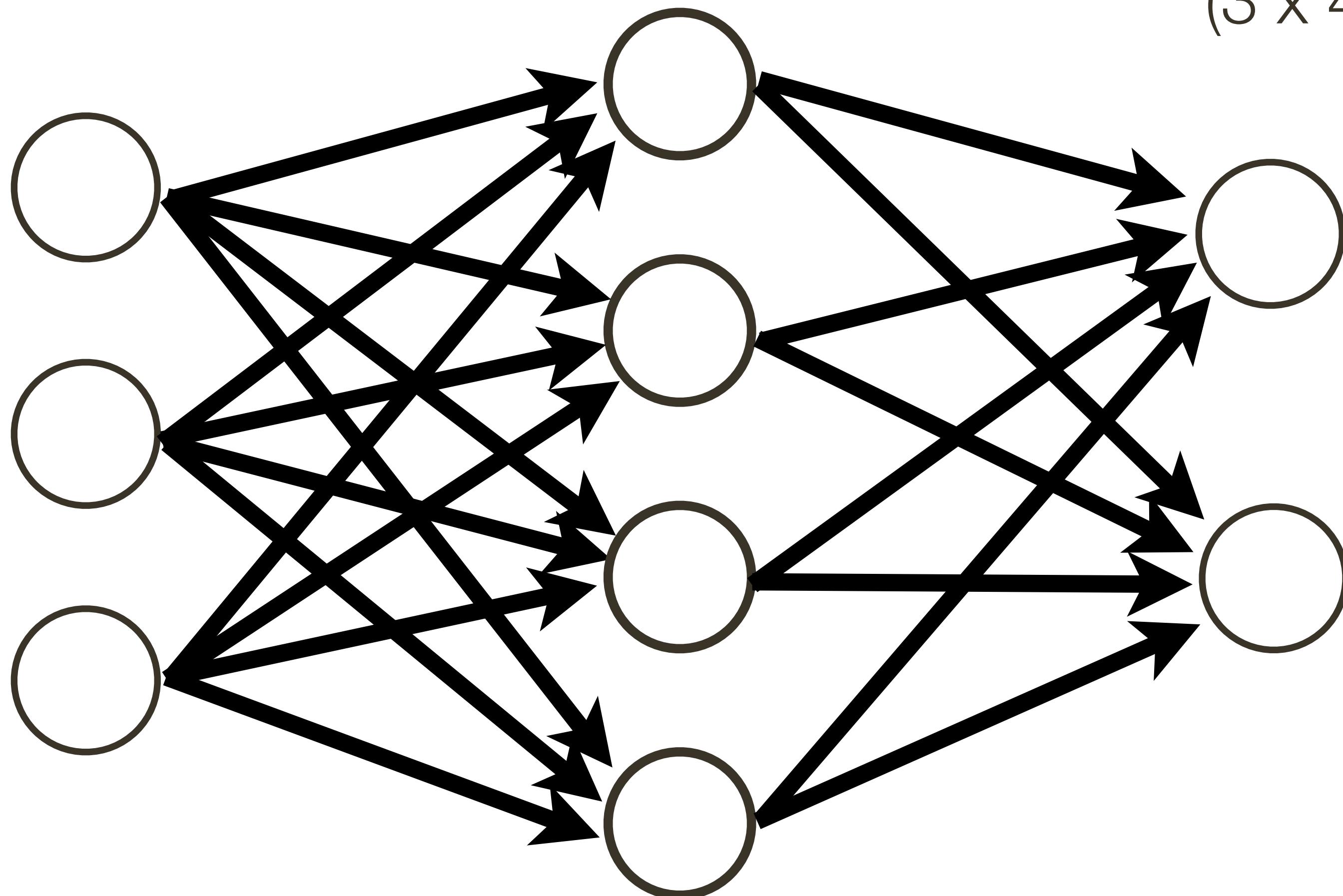
Neural Network

How many neurons?

$$4+2 = 6$$

How many weights?

$$(3 \times 4) + (4 \times 2) = 20$$



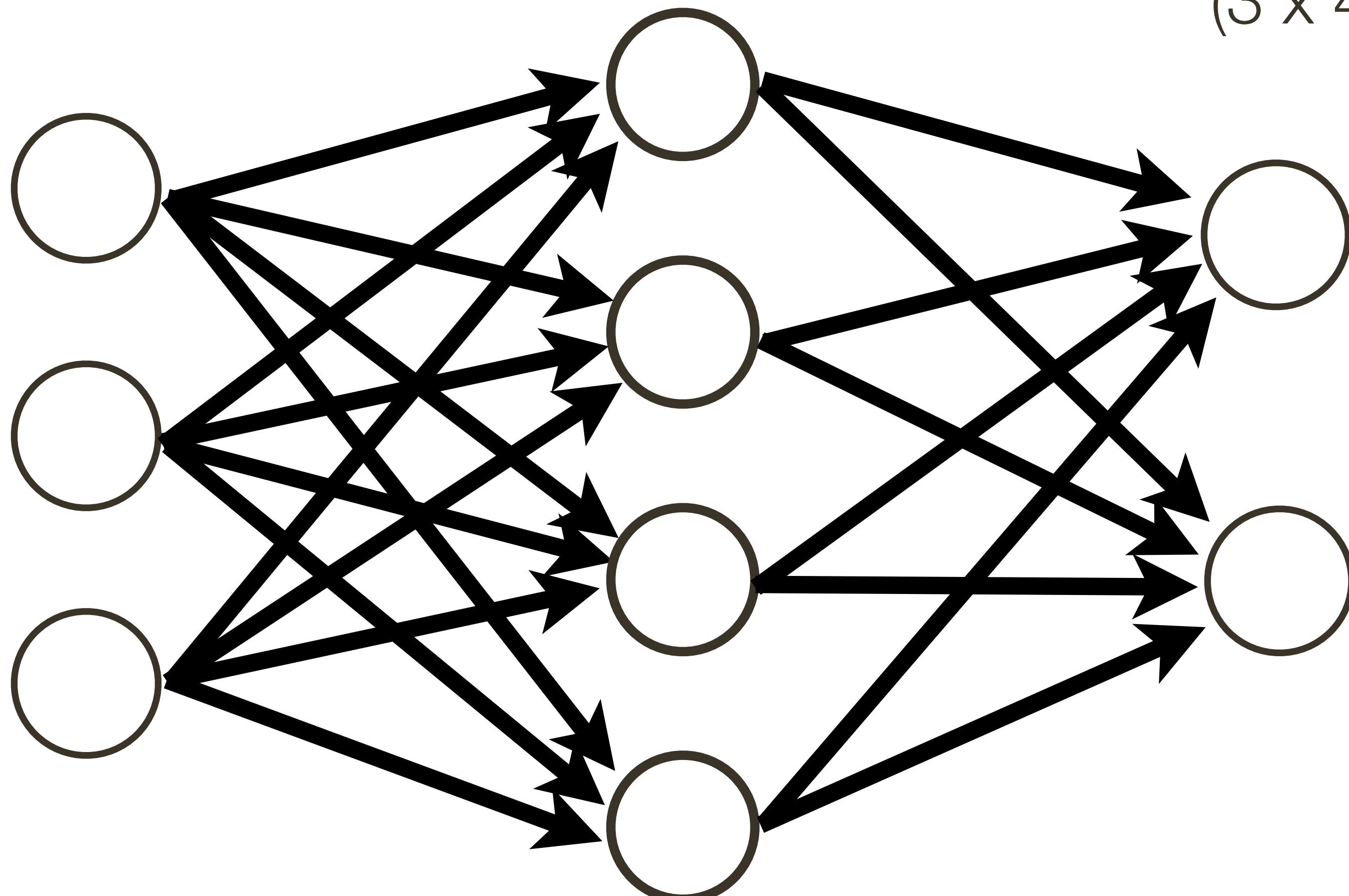
Neural Network

How many neurons?

$$4+2 = 6$$

How many weights?

$$(3 \times 4) + (4 \times 2) = 20$$



How many learnable parameters?

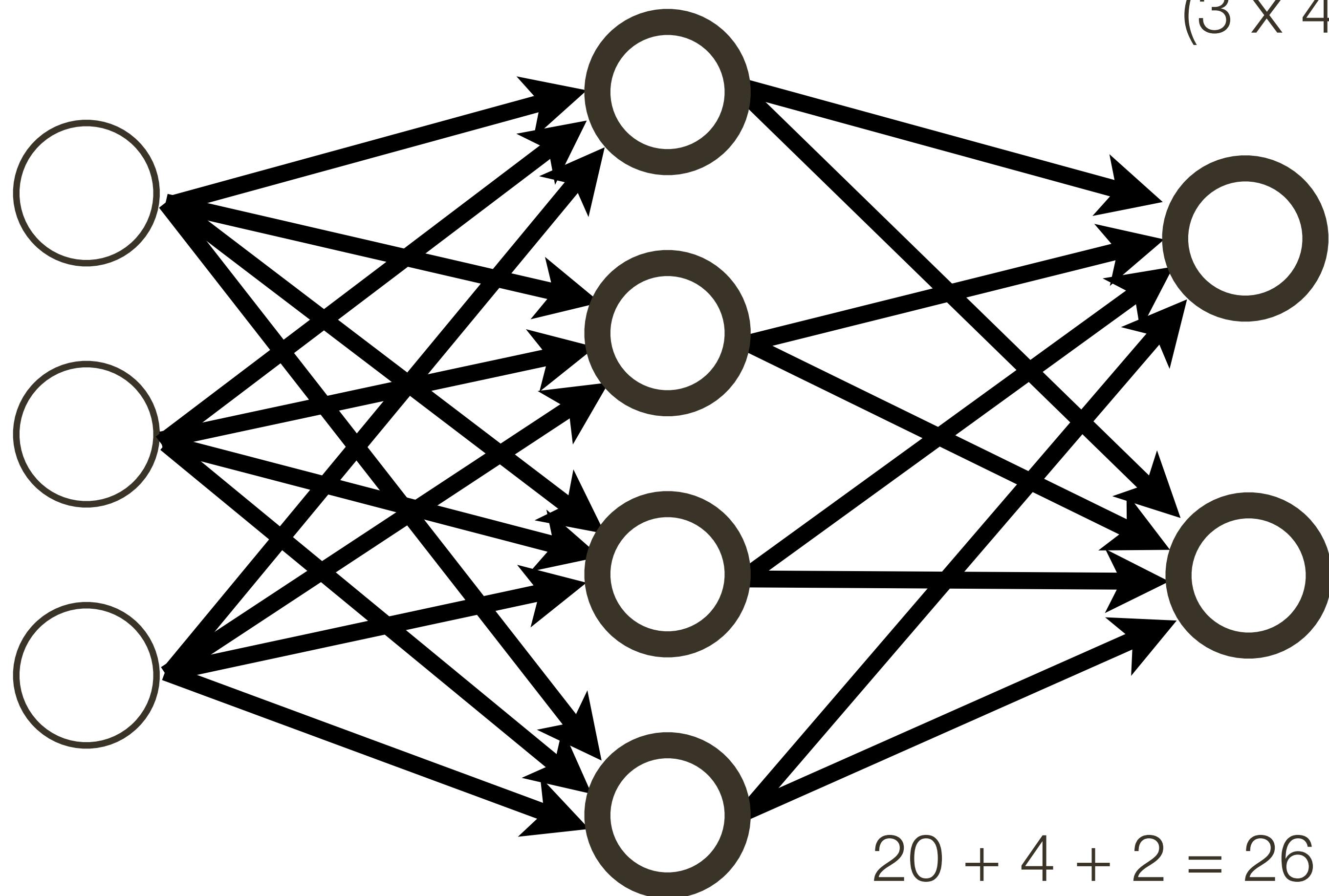
Neural Network

How many neurons?

$$4+2 = 6$$

How many weights?

$$(3 \times 4) + (4 \times 2) = 20$$



How many learnable parameters?

$$20 + 4 + 2 = 26$$

bias terms

Neural Network Intuition

Question: What is a Neural Network?

Answer: Complex mapping from an input (vector) to an output (vector)

Question: What class of functions should be considered for this mapping?

Answer: Compositions of simpler functions (a.k.a. layers)? We will talk more about what specific functions next ...

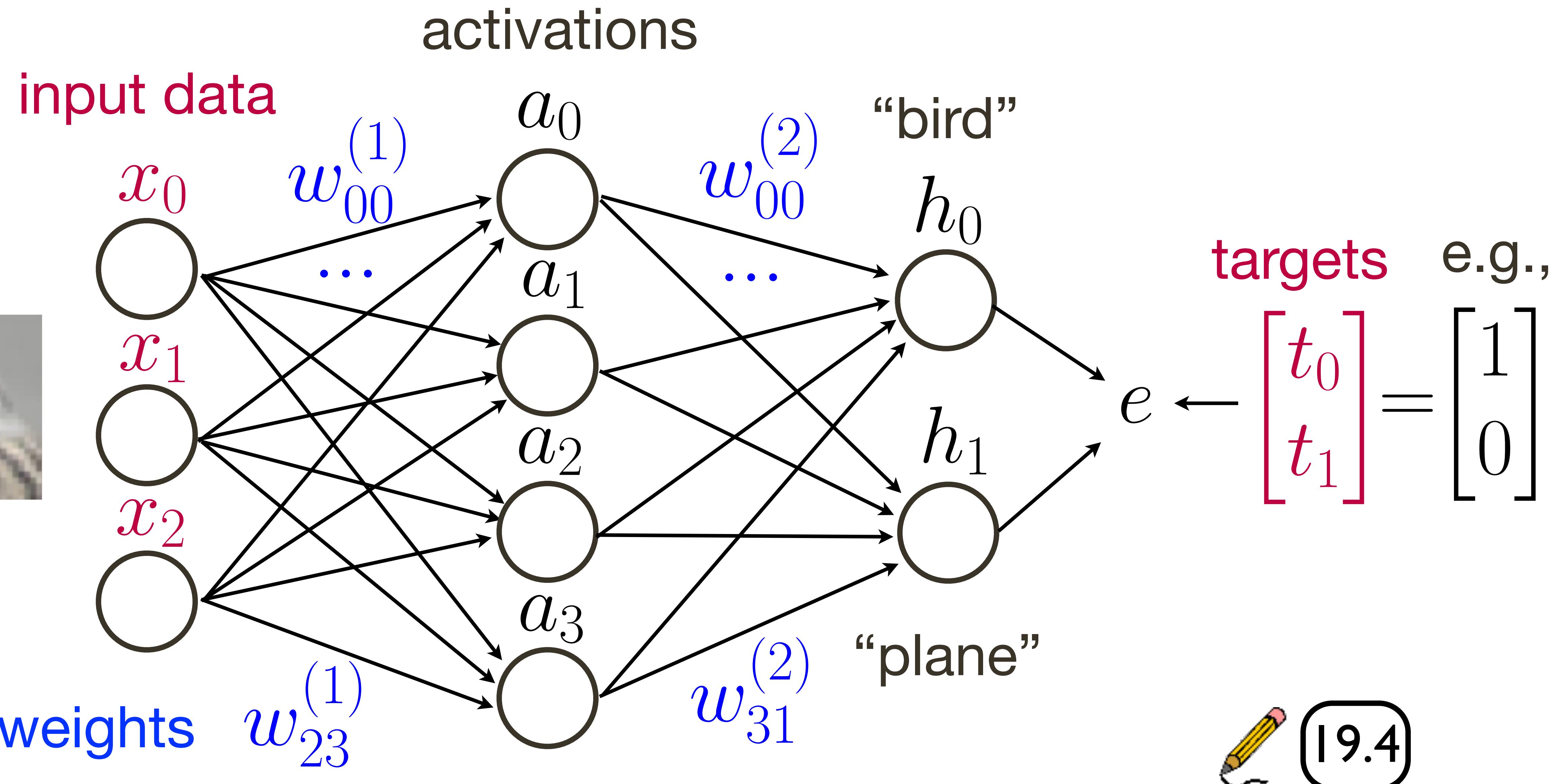
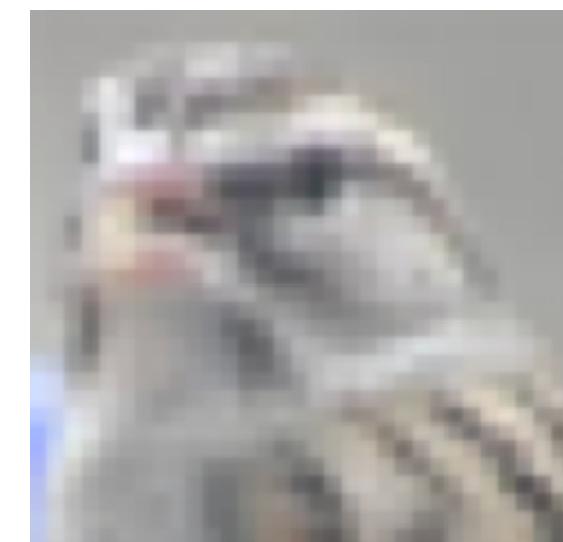
Question: What does a hidden unit do?

Answer: It can be thought of as classifier or a feature.

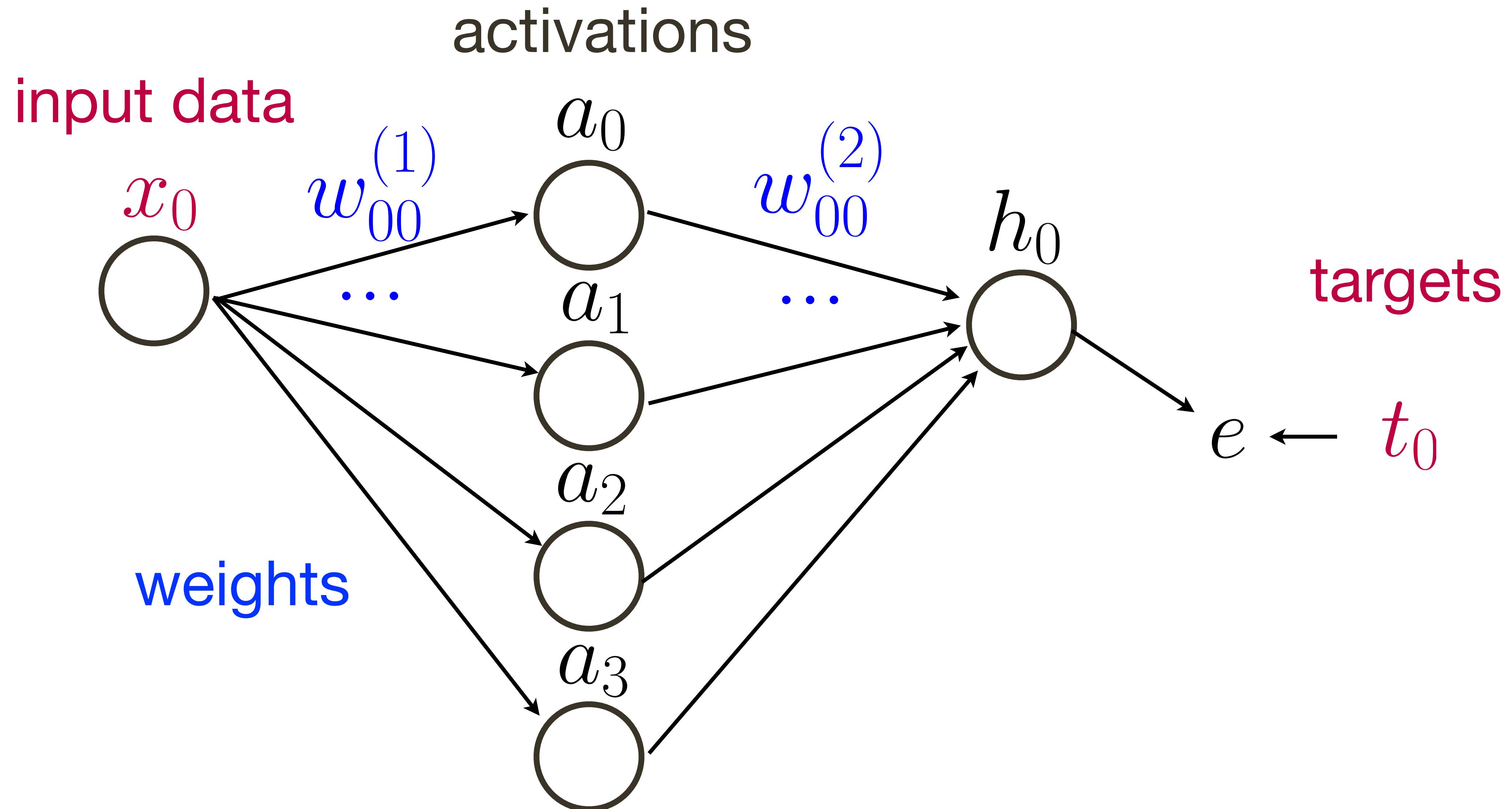
Question: Why have many layers?

Answer: 1) More layers = more complex functional mapping
2) More efficient due to distributed representation

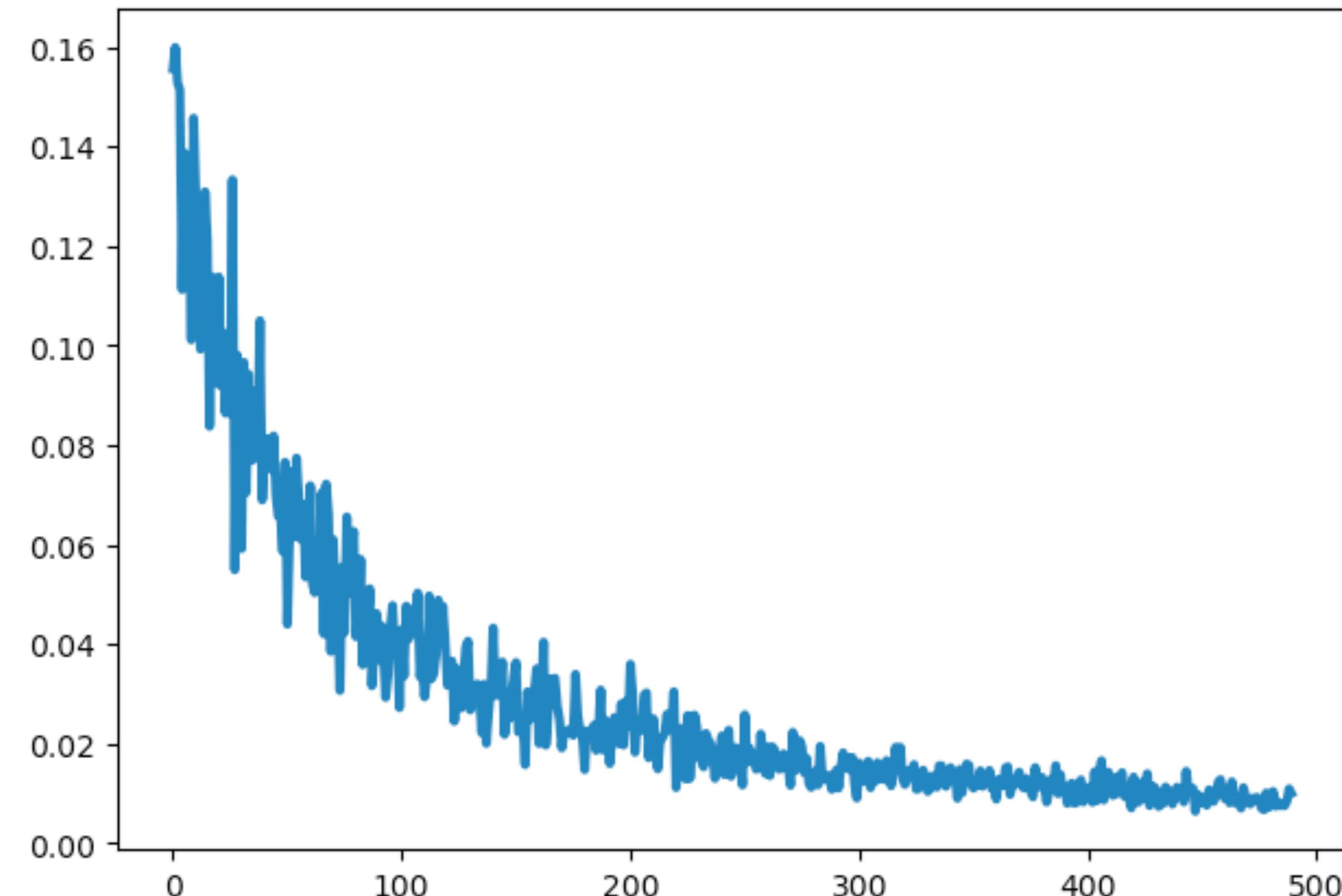
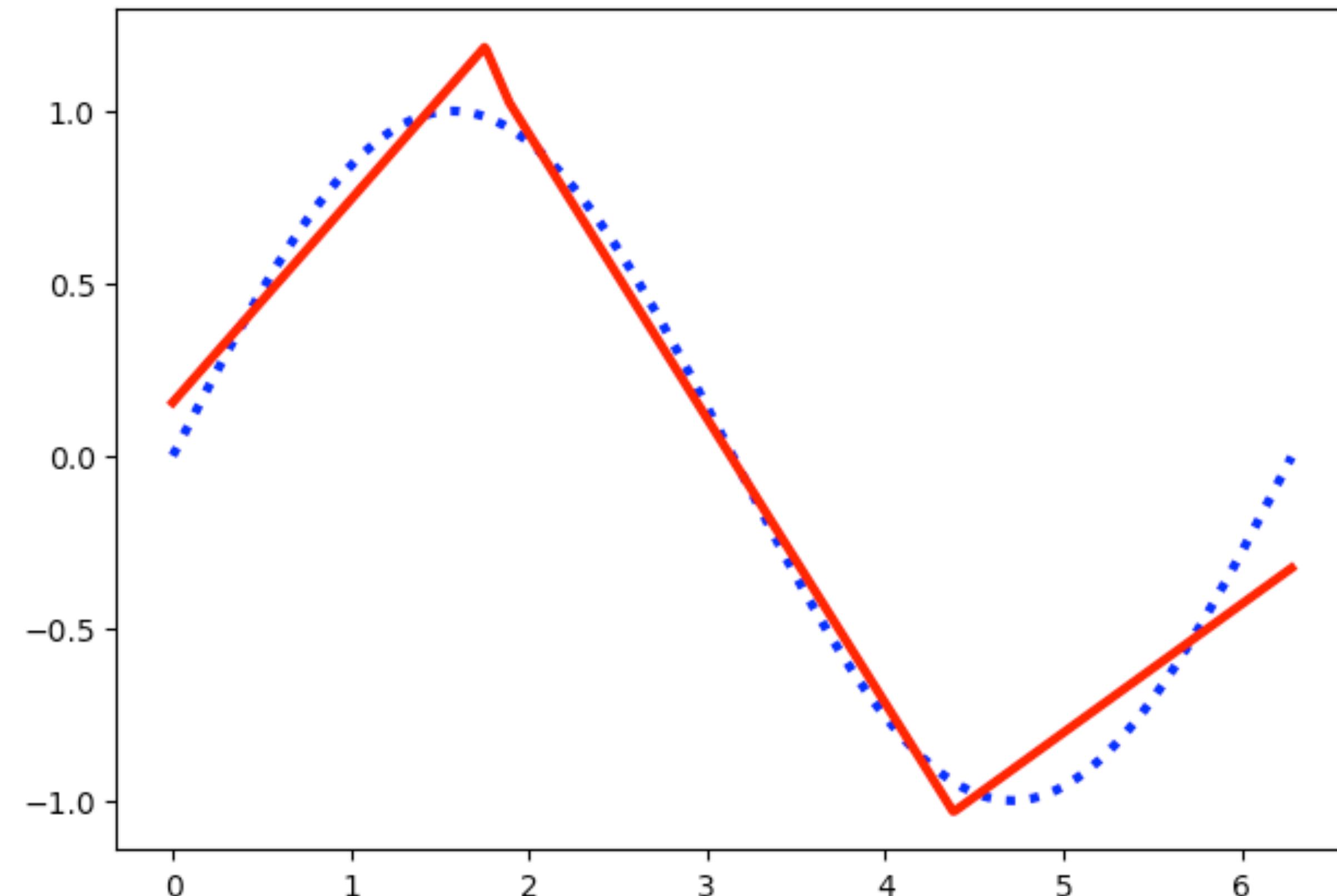
2-Layer Neural Network



2-Layer Neural Network – n hidden, 1 input/output

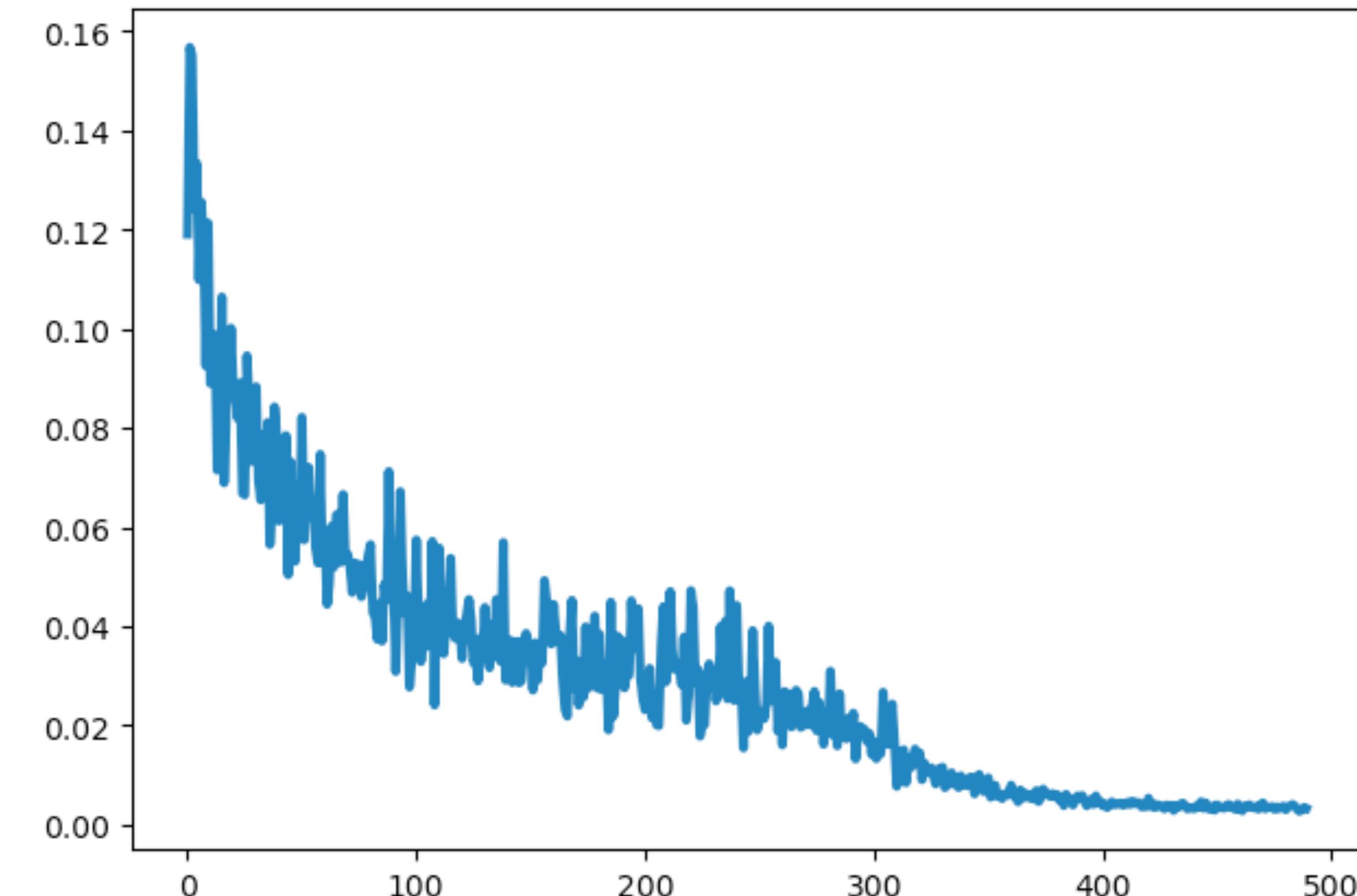
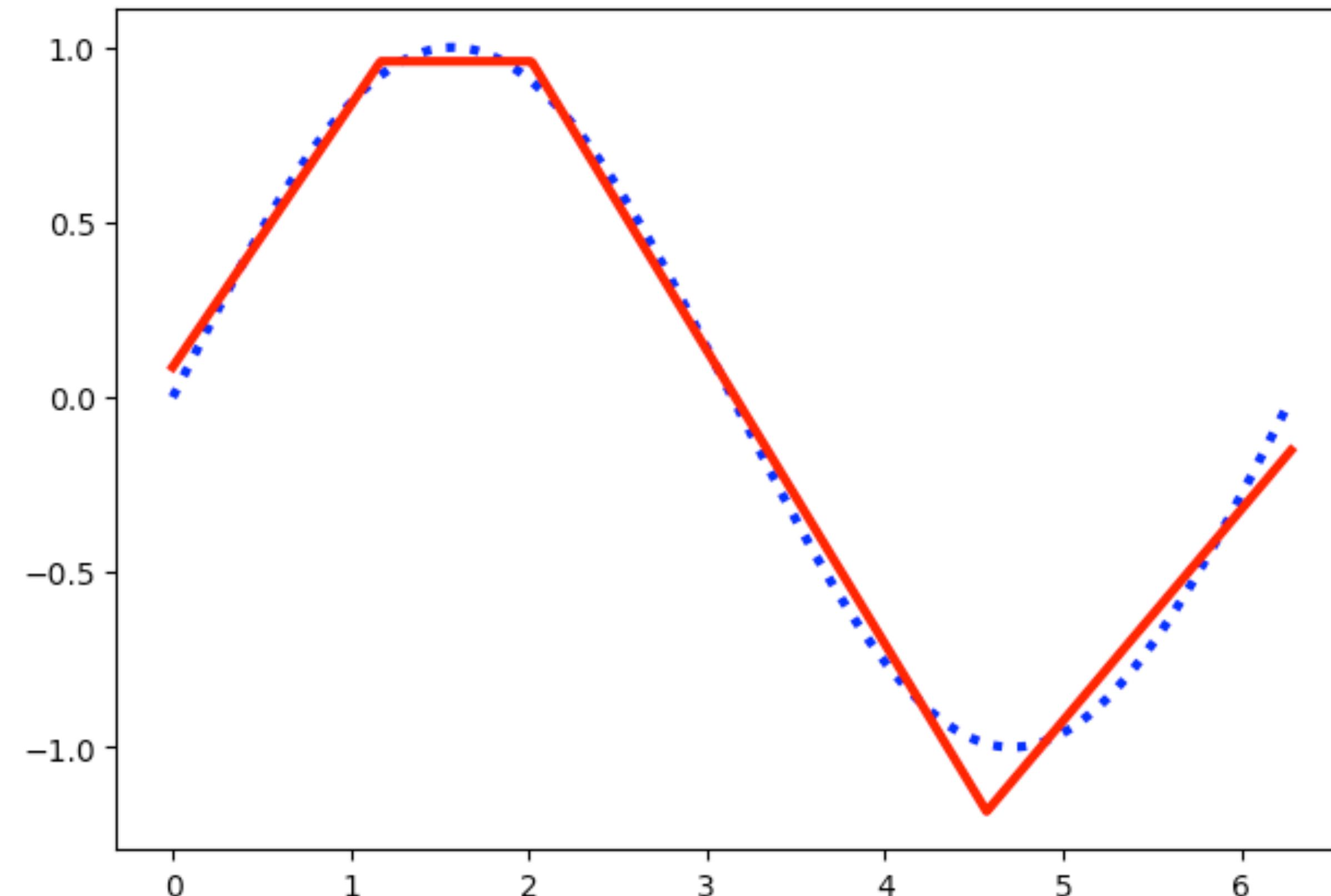


2-Layer Neural Network – n hidden, 1 input/output



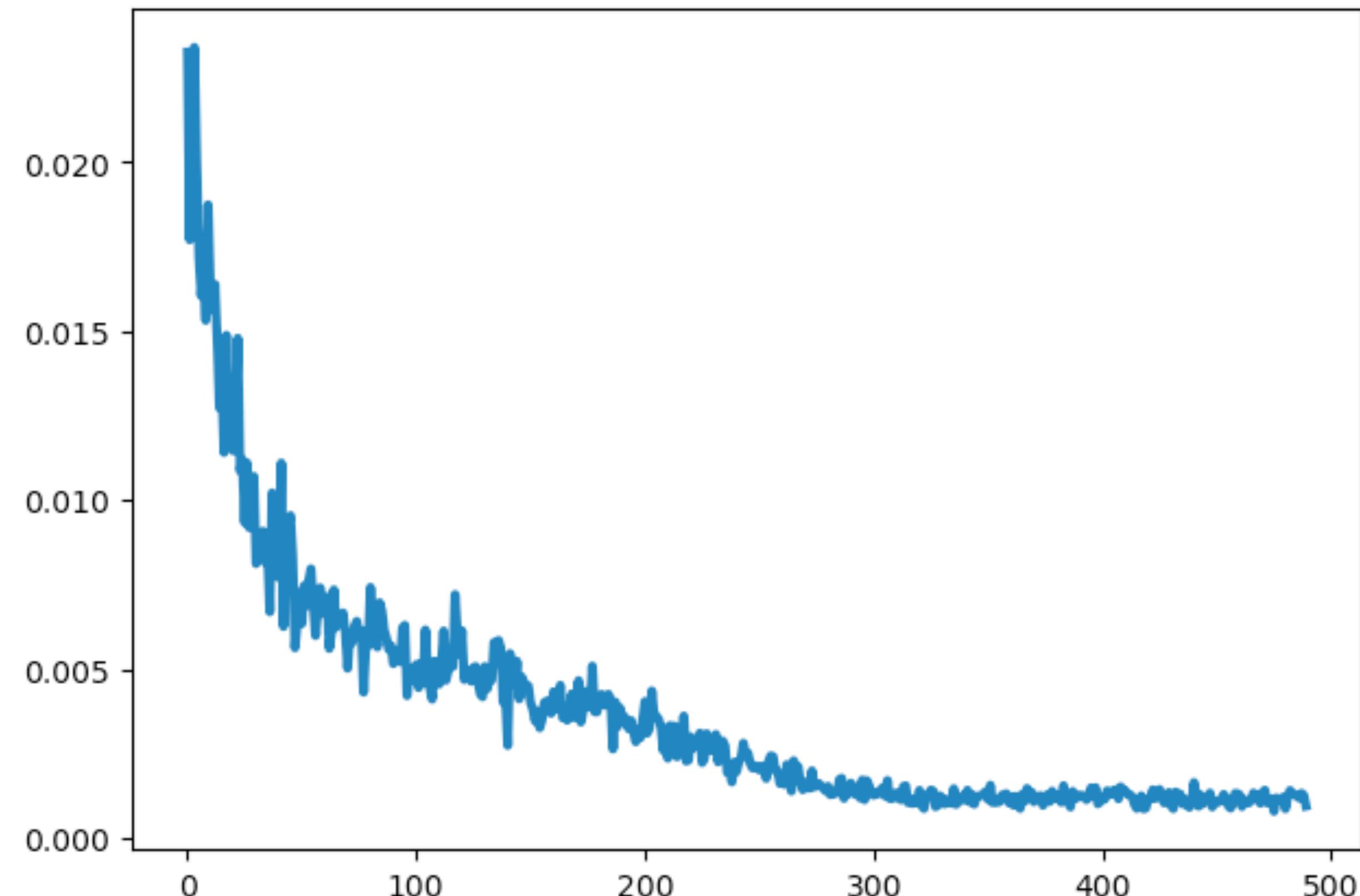
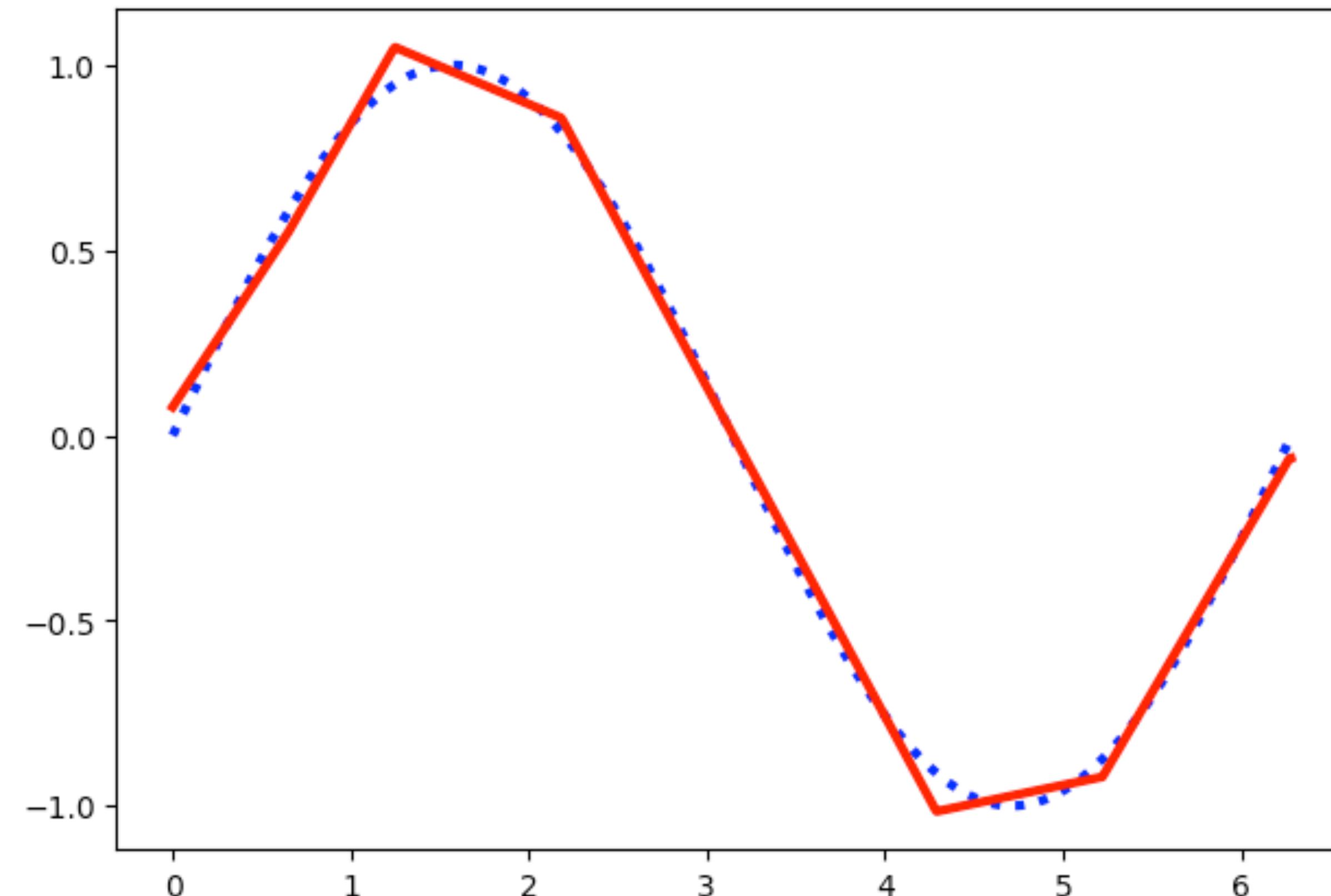
3 hidden units

2-Layer Neural Network – n hidden, 1 input/output



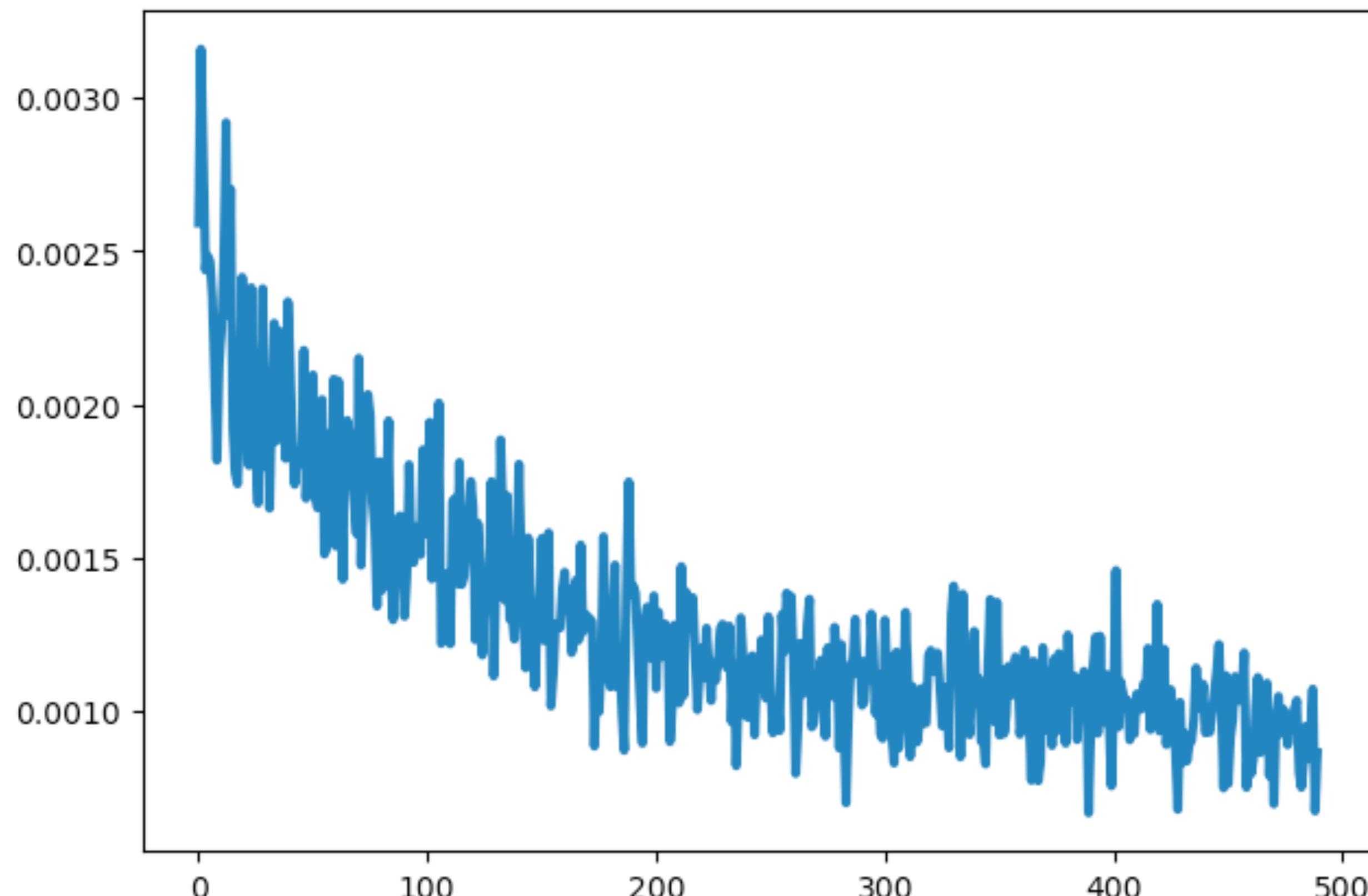
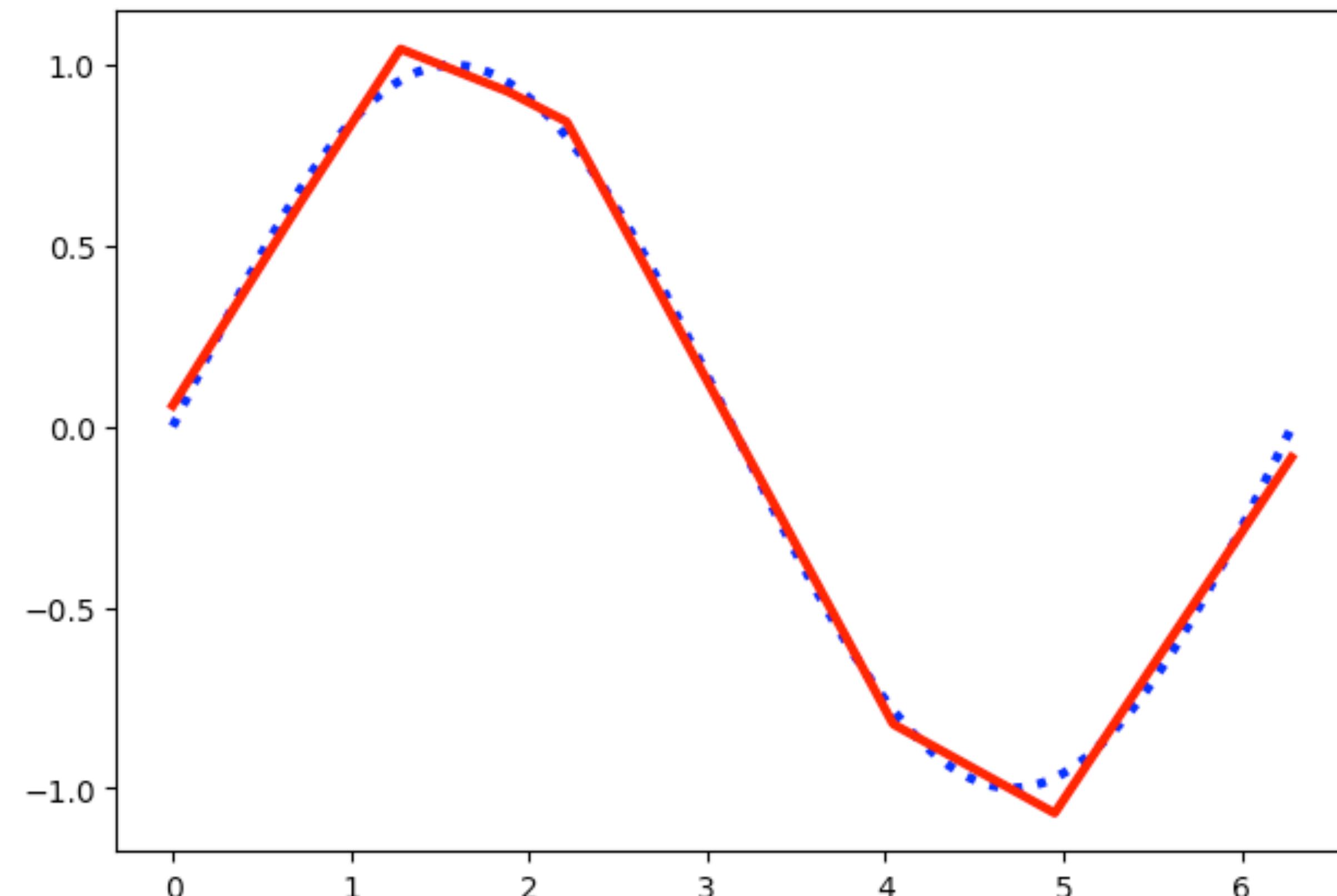
4 hidden units

2-Layer Neural Network – n hidden, 1 input/output



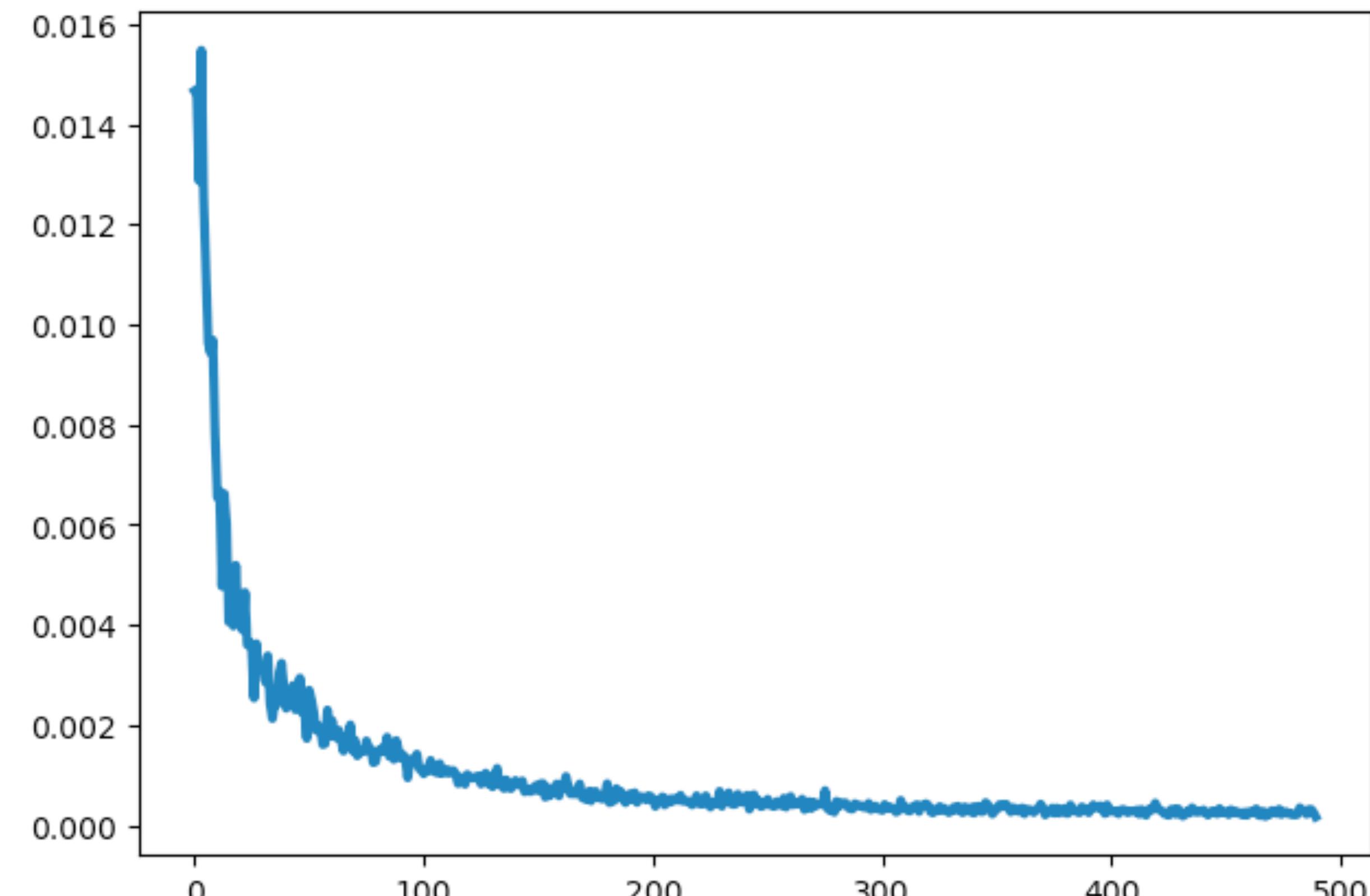
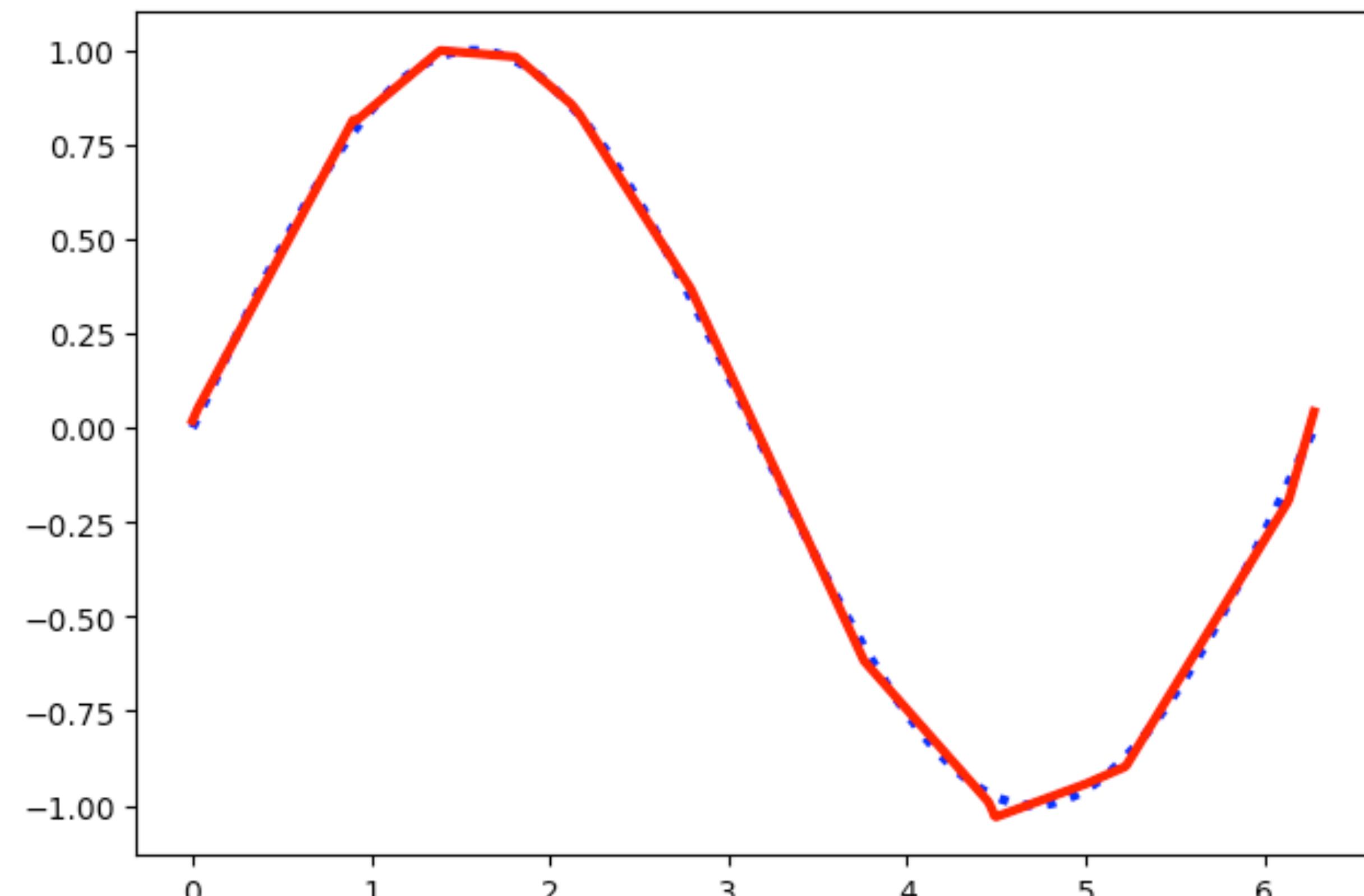
6 hidden units

2-Layer Neural Network – n hidden, 1 input/output



8 hidden units

2-Layer Neural Network – n hidden, 1 input/output



20 hidden units

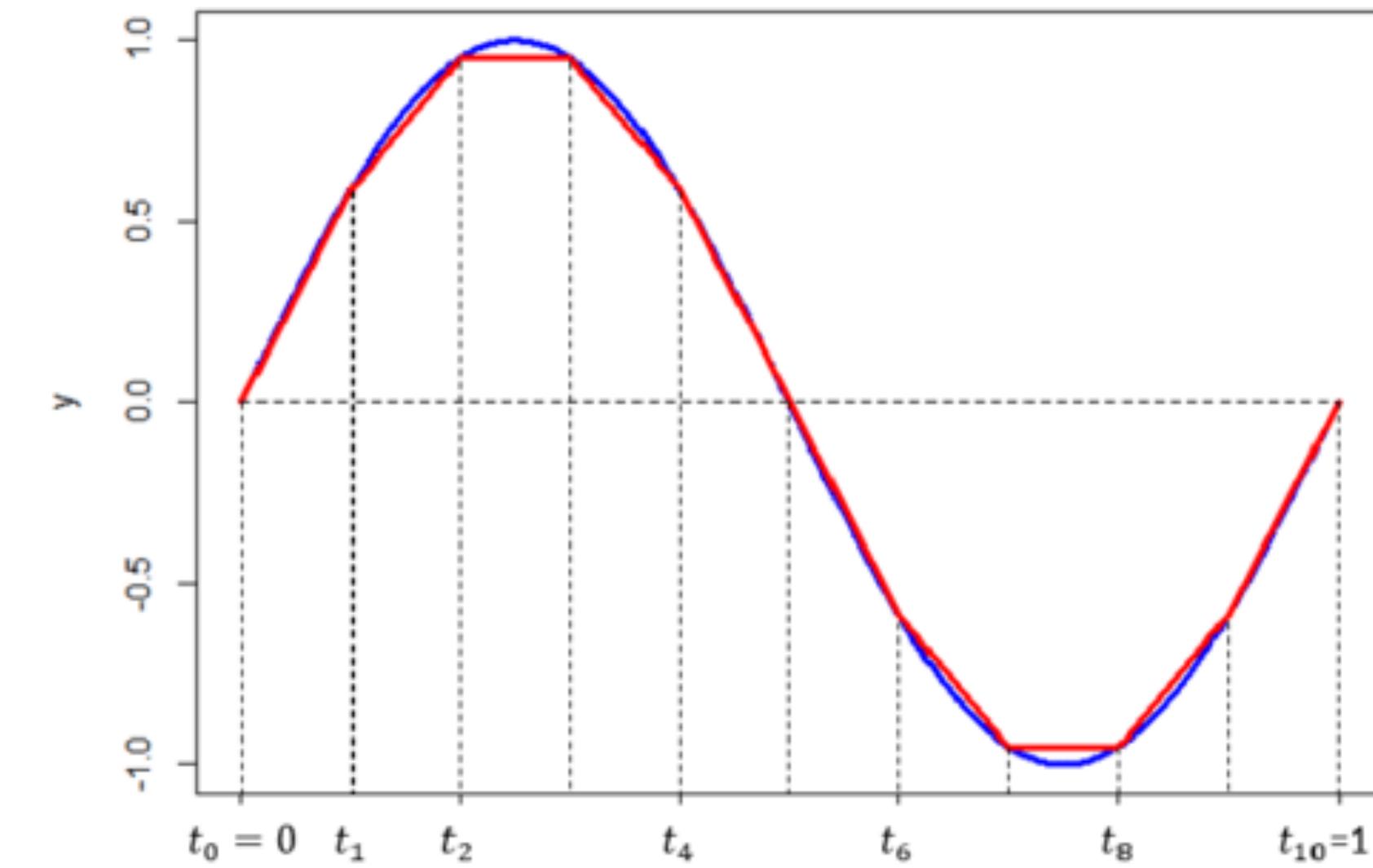
Neural Network as Universal Approximator

Non-linear activation is required to provably make the Neural Net a **universal function approximator**

Intuition: with ReLU activation, we effectively get a linear spline approximation to any function.

Optimization of neural net parameters = finding slopes and transitions of linear pieces

The quality of approximation depends on the number of linear segments



Neural Network as Universal Approximator

Universal Approximation Theorem: Single hidden layer can approximate any continuous function with compact support to arbitrary accuracy, when the width goes to infinity.

[Hornik *et al.*, 1989]

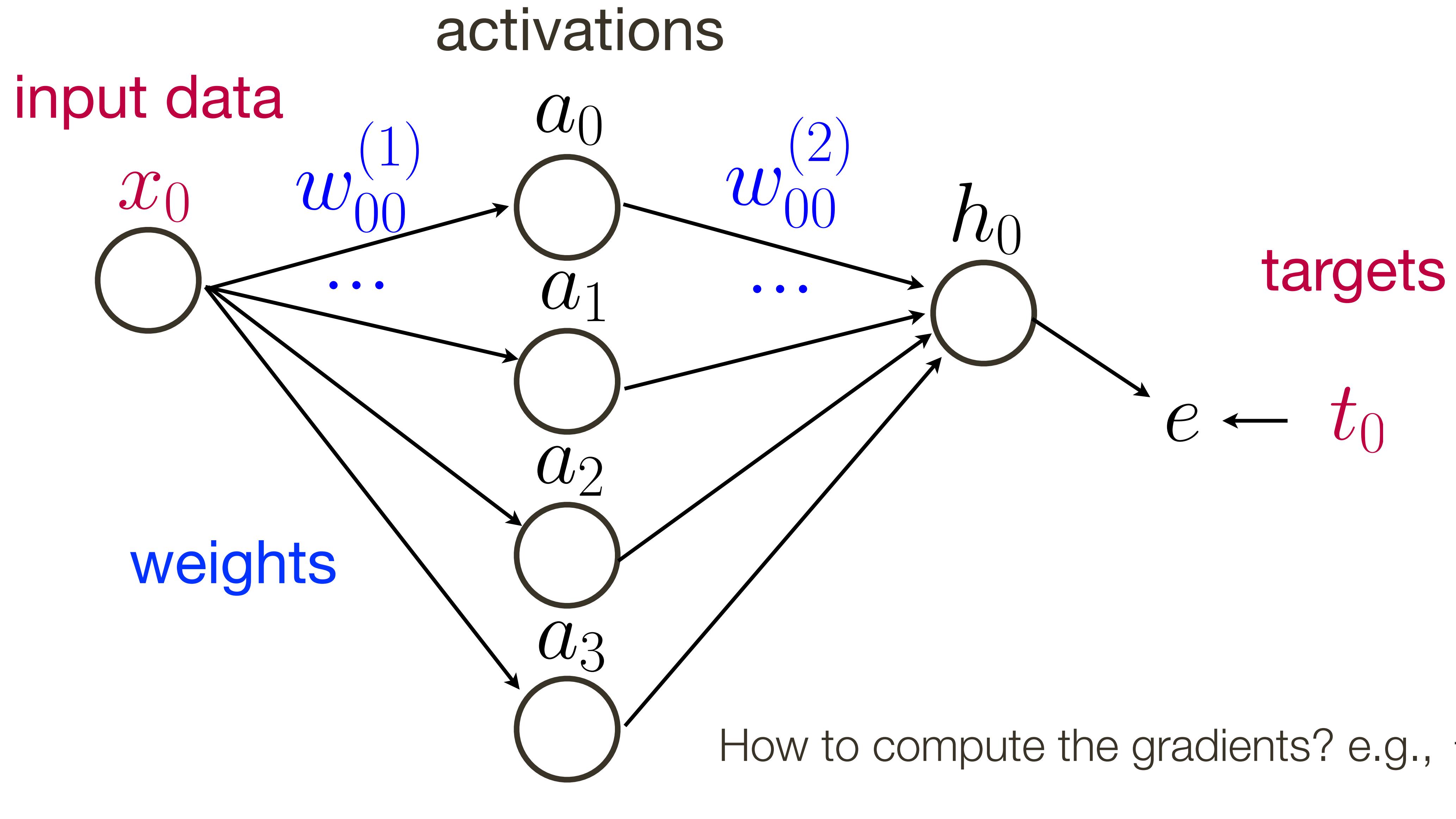
Universal Approximation Theorem (revised): A network of infinite depth with a hidden layer of size $d + 1$ neurons, where d is the dimension of the input space, can approximate any continuous function.

[Lu *et al.*, NIPS 2017]

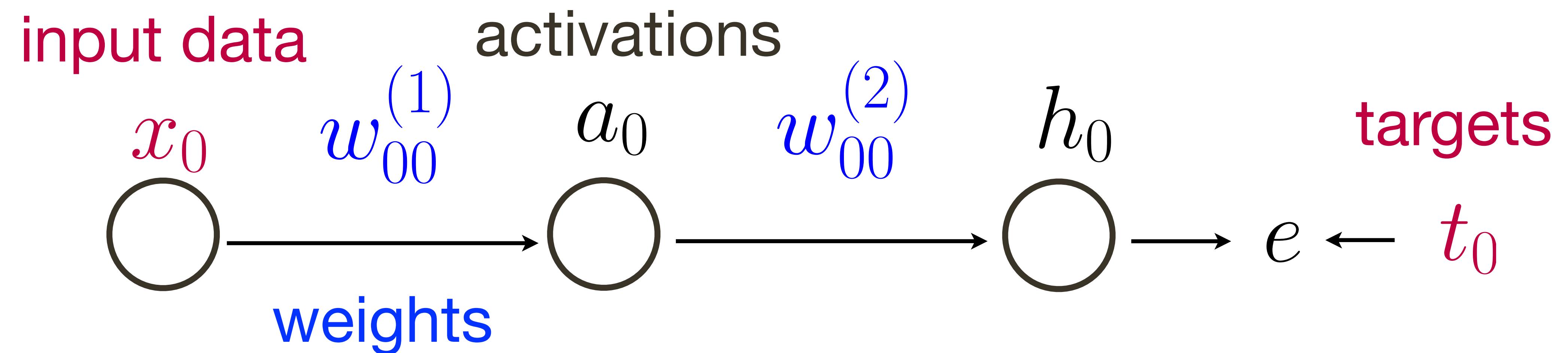
Universal Approximation Theorem (further revised): ResNet with a single hidden unit and infinite depth can approximate any continuous function.

[Lin and Jegelka, NIPS 2018]

2-Layer Neural Network – n hidden, 1 input/output



2-Layer Neural Network – 1 hidden, 1 input/output



2-Layer Neural Network – 1 hidden, 1 input/output

$$y = w_2(\max(0, w_1x + b_1)) + b_2 \quad L = (y - t)^2$$

Optimise by **gradient descent**

$$\begin{bmatrix} w_1 \\ b_1 \\ w_2 \\ b_2 \end{bmatrix} \rightarrow \begin{bmatrix} w_1 \\ b_1 \\ w_2 \\ b_2 \end{bmatrix} - \alpha \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial w_2} \\ \frac{\partial L}{\partial b_2} \end{bmatrix}$$

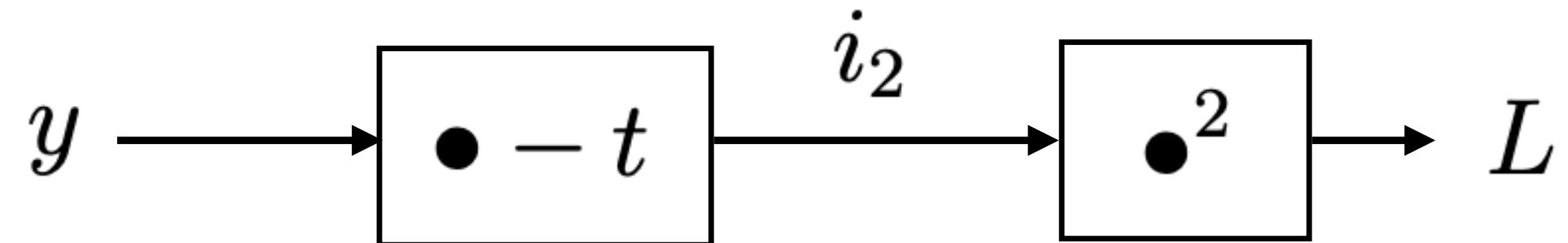
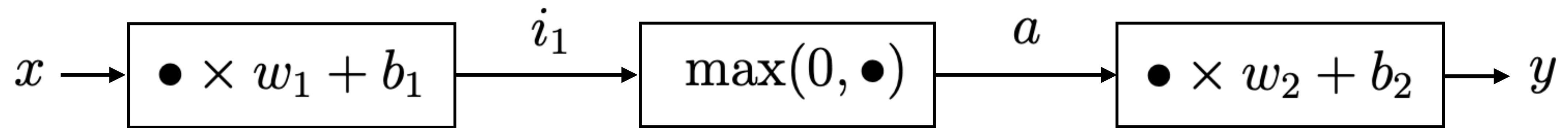


19.5

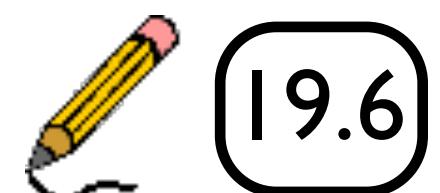
How to compute the gradients? e.g., $\frac{\partial L}{\partial w_1}$

2-Layer Neural Network – 1 hidden, 1 input/output

$$y = w_2(\max(0, w_1x + b_1)) + b_2 \quad L = (y - t)^2$$

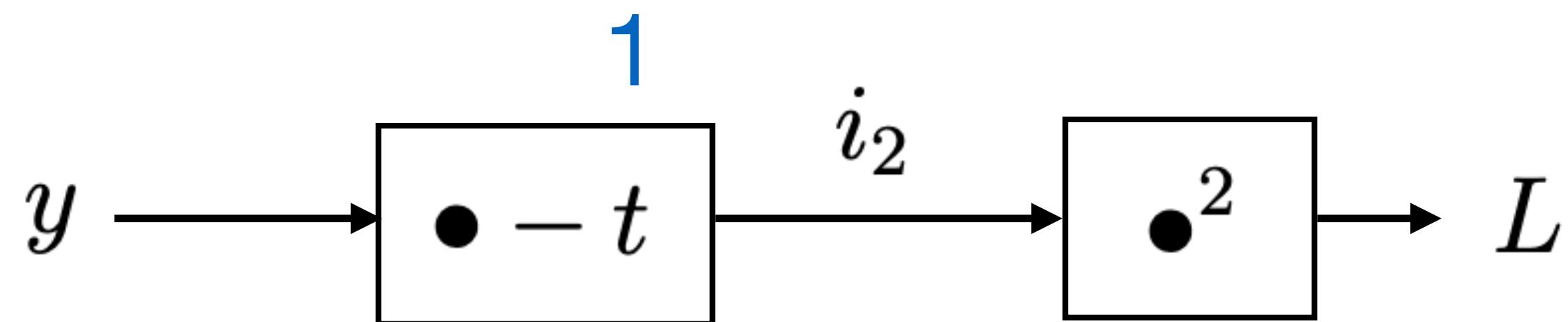
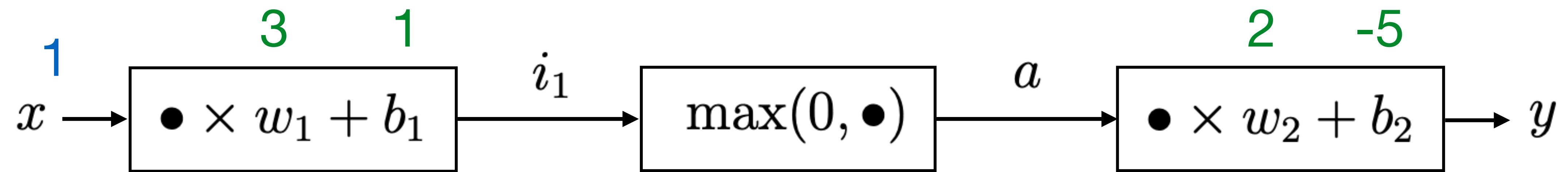


Alternative: build a **computational graph** to apply the **chain rule**



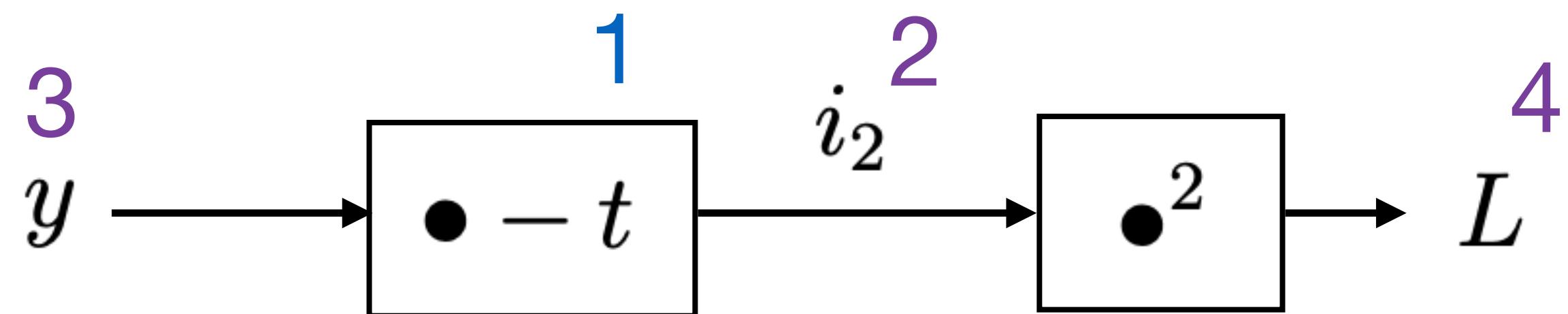
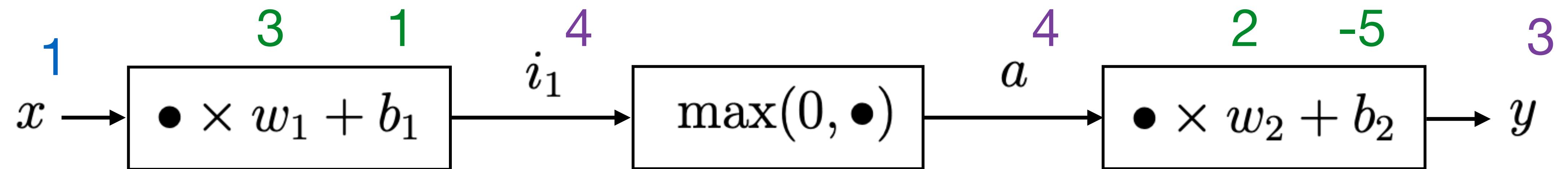
2-Layer Neural Network – 1 hidden, 1 input/output

Input + Initial weights
/target

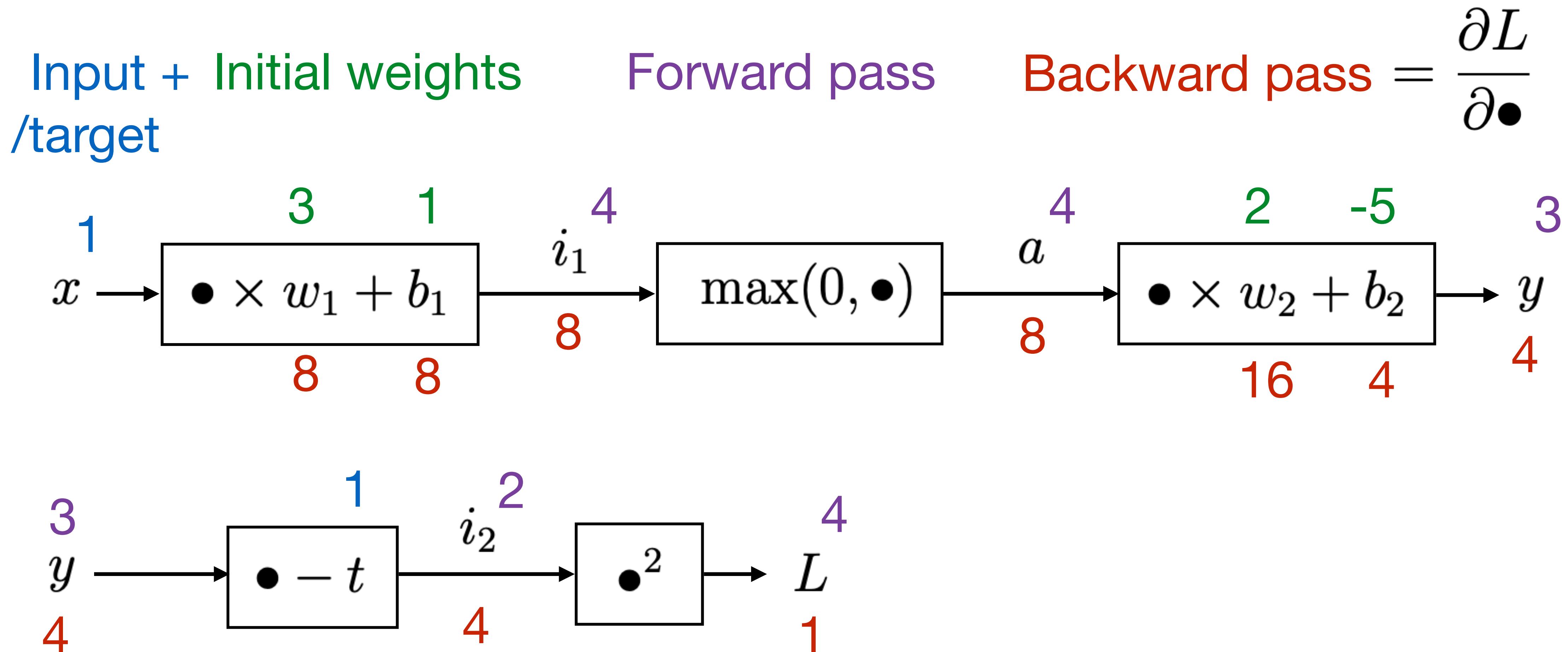


2-Layer Neural Network – 1 hidden, 1 input/output

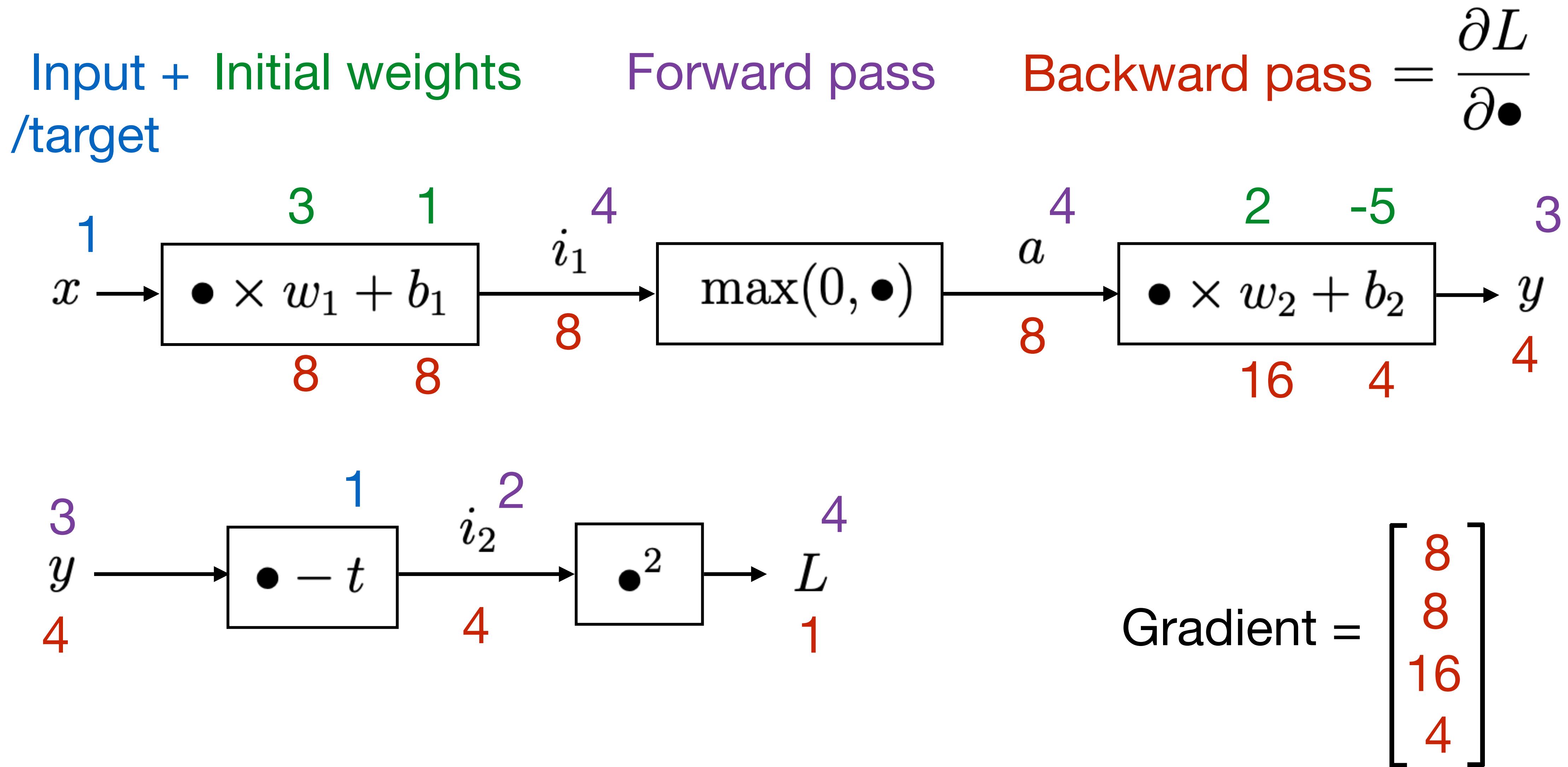
Input + Initial weights
/target



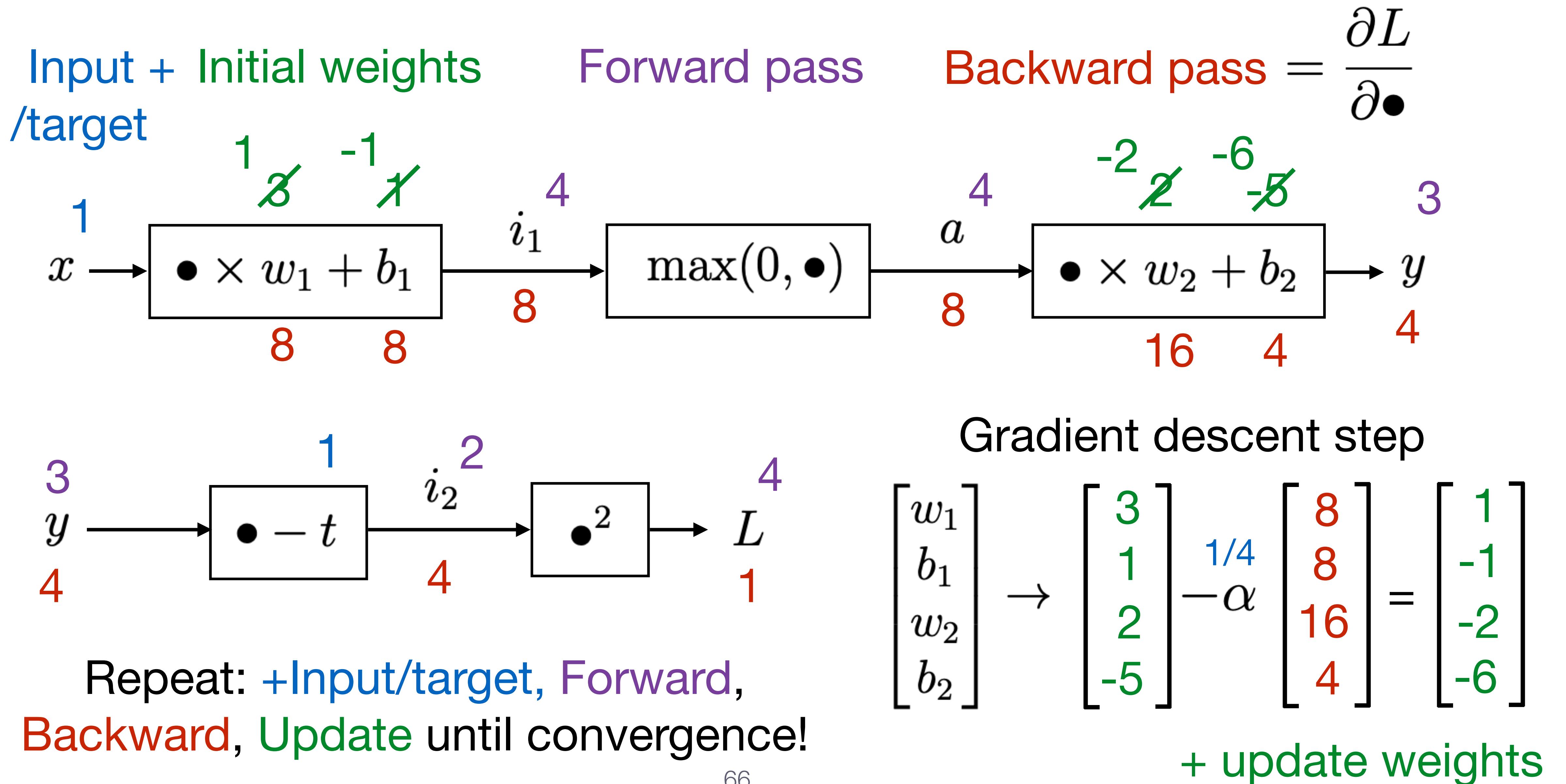
2-Layer Neural Network – 1 hidden, 1 input/output



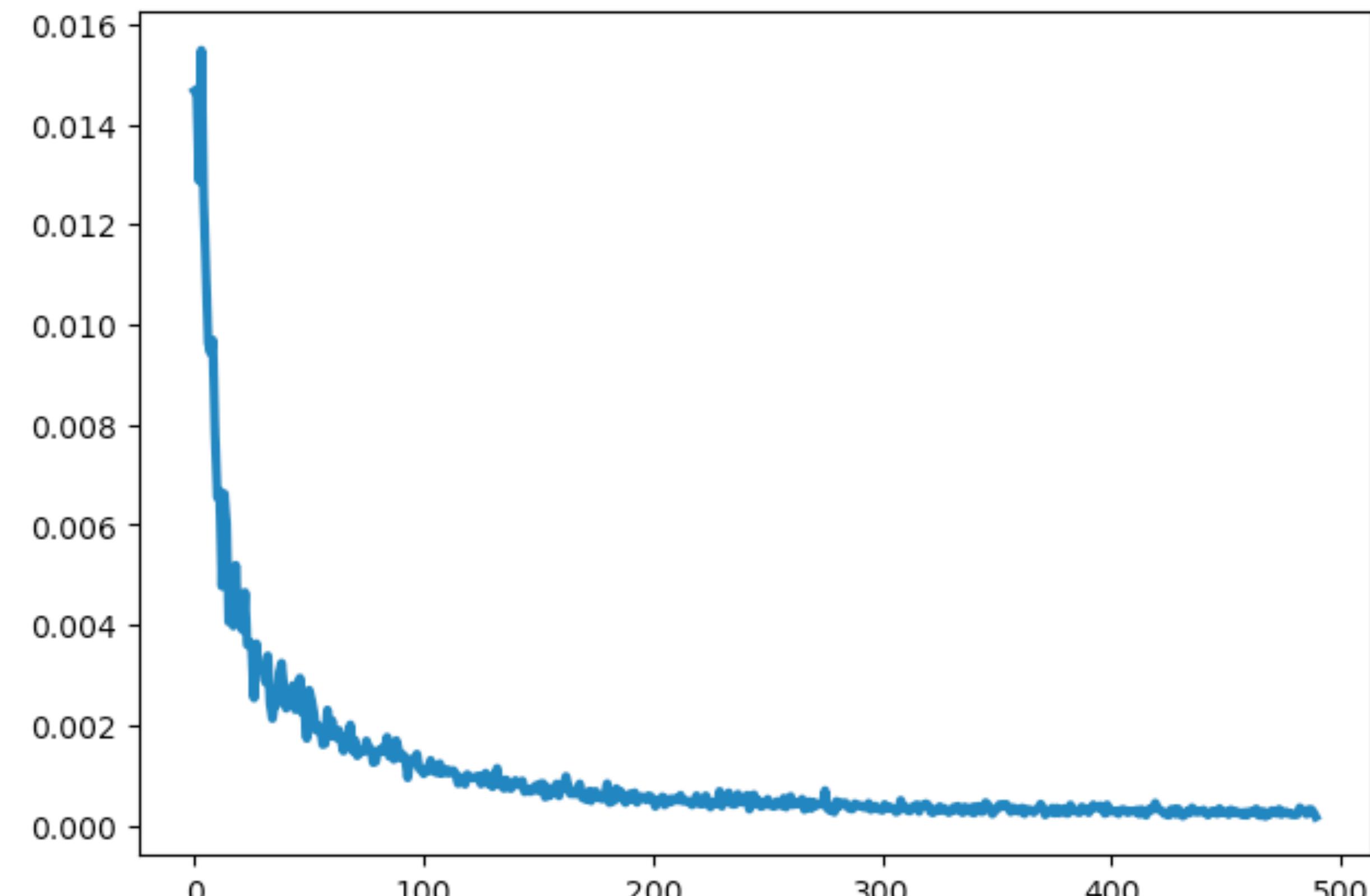
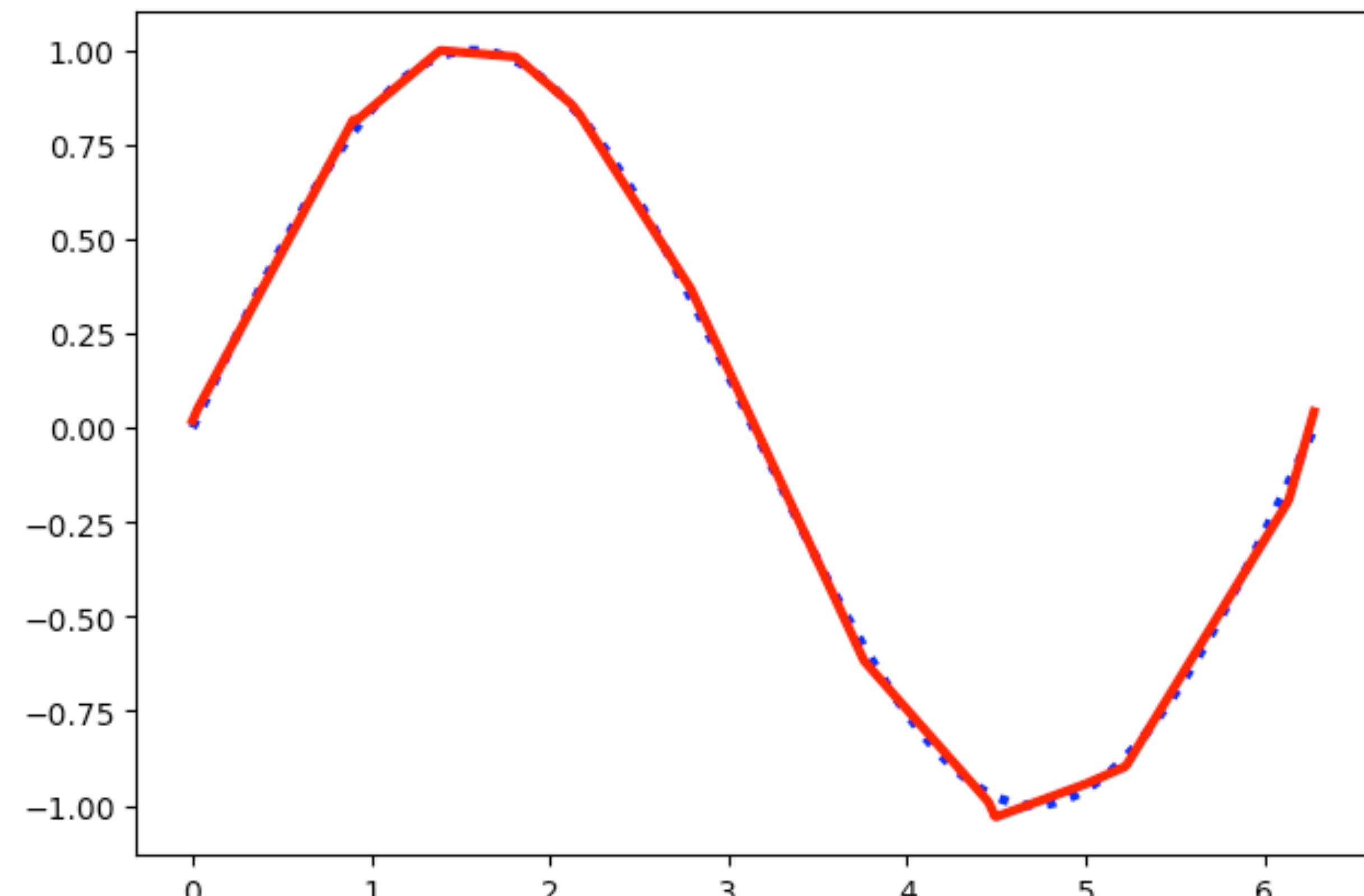
2-Layer Neural Network – 1 hidden, 1 input/output



2-Layer Neural Network – 1 hidden, 1 input/output

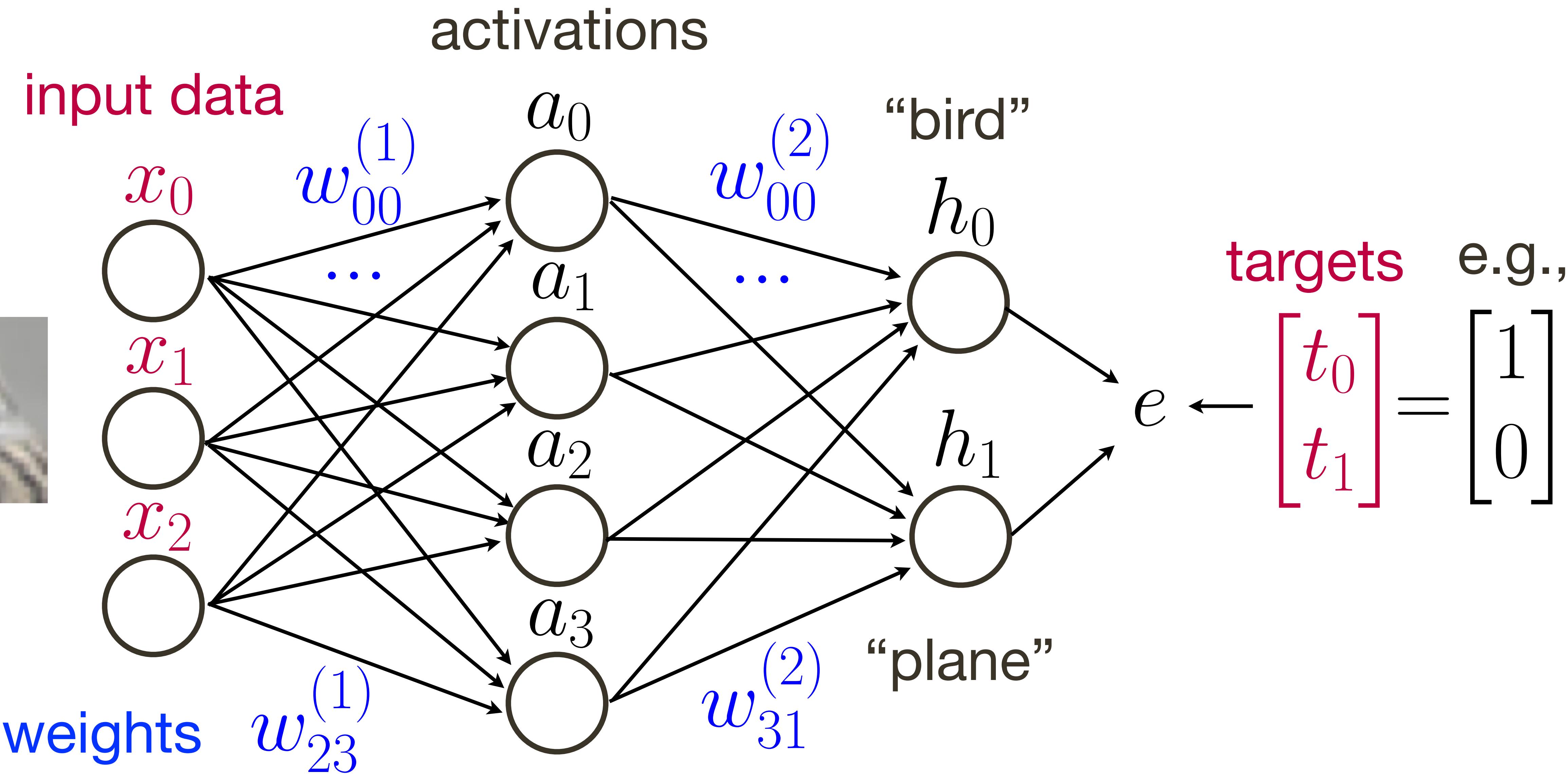
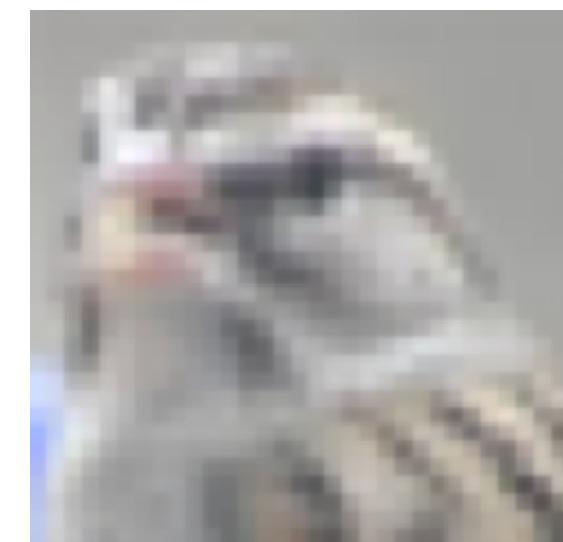


2-Layer Neural Network – n hidden, 1 input/output



20 hidden units

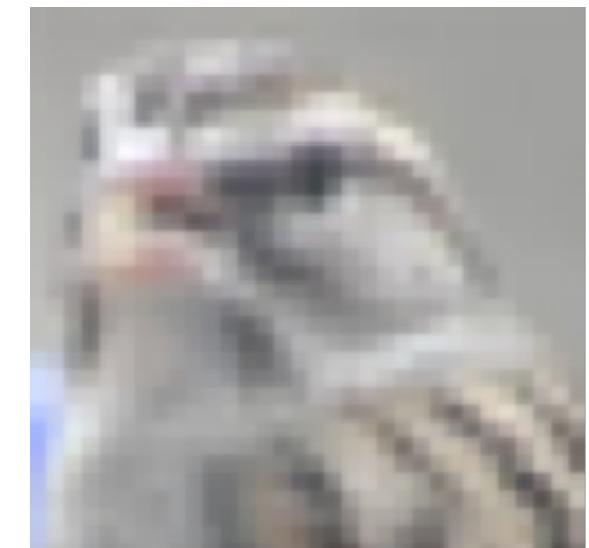
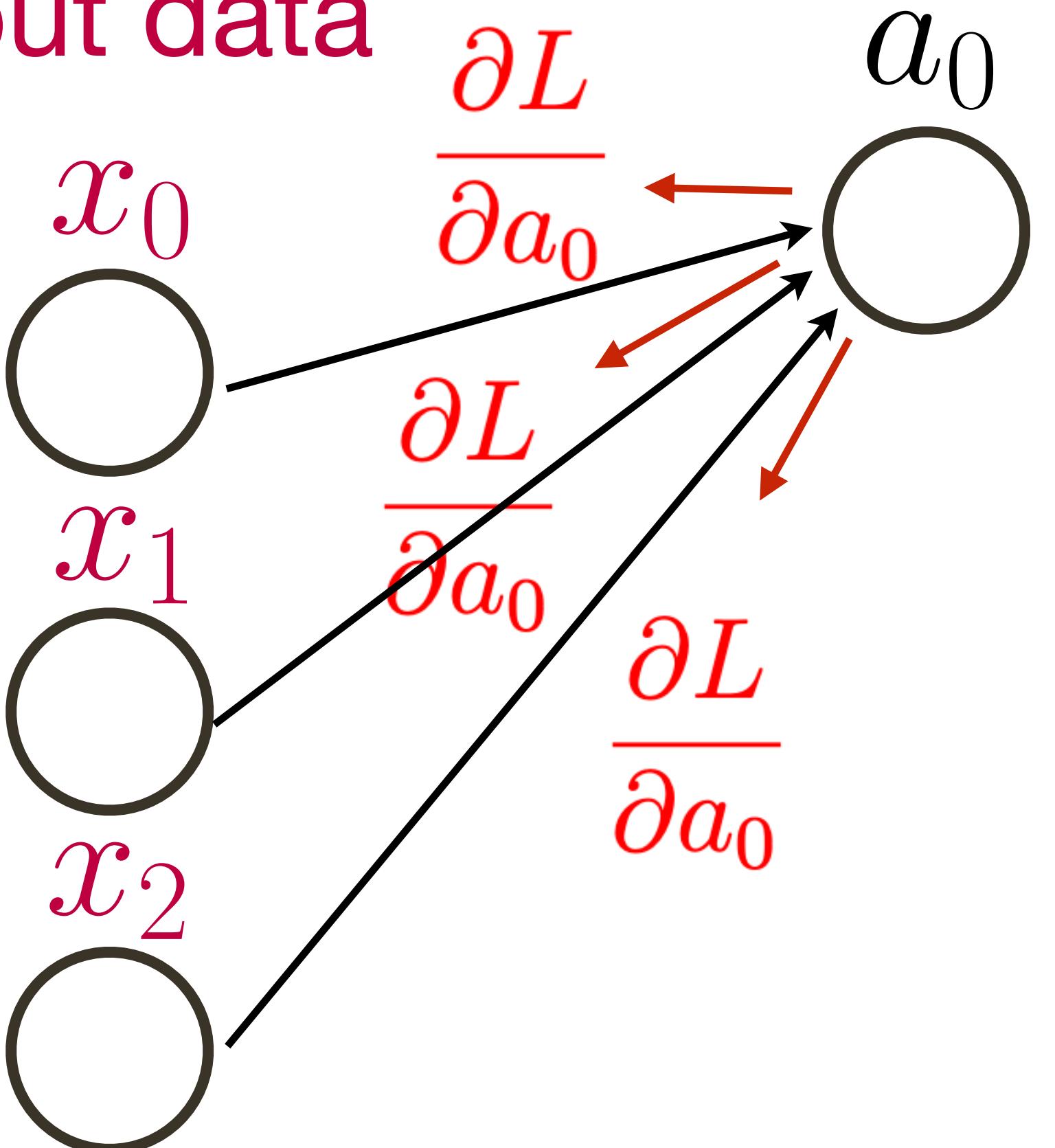
2-Layer Neural Network



2-Layer Neural Network – multiple inputs

activations

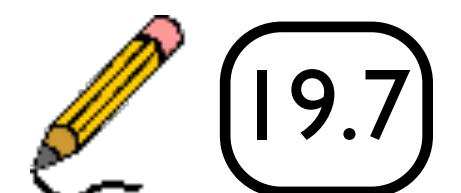
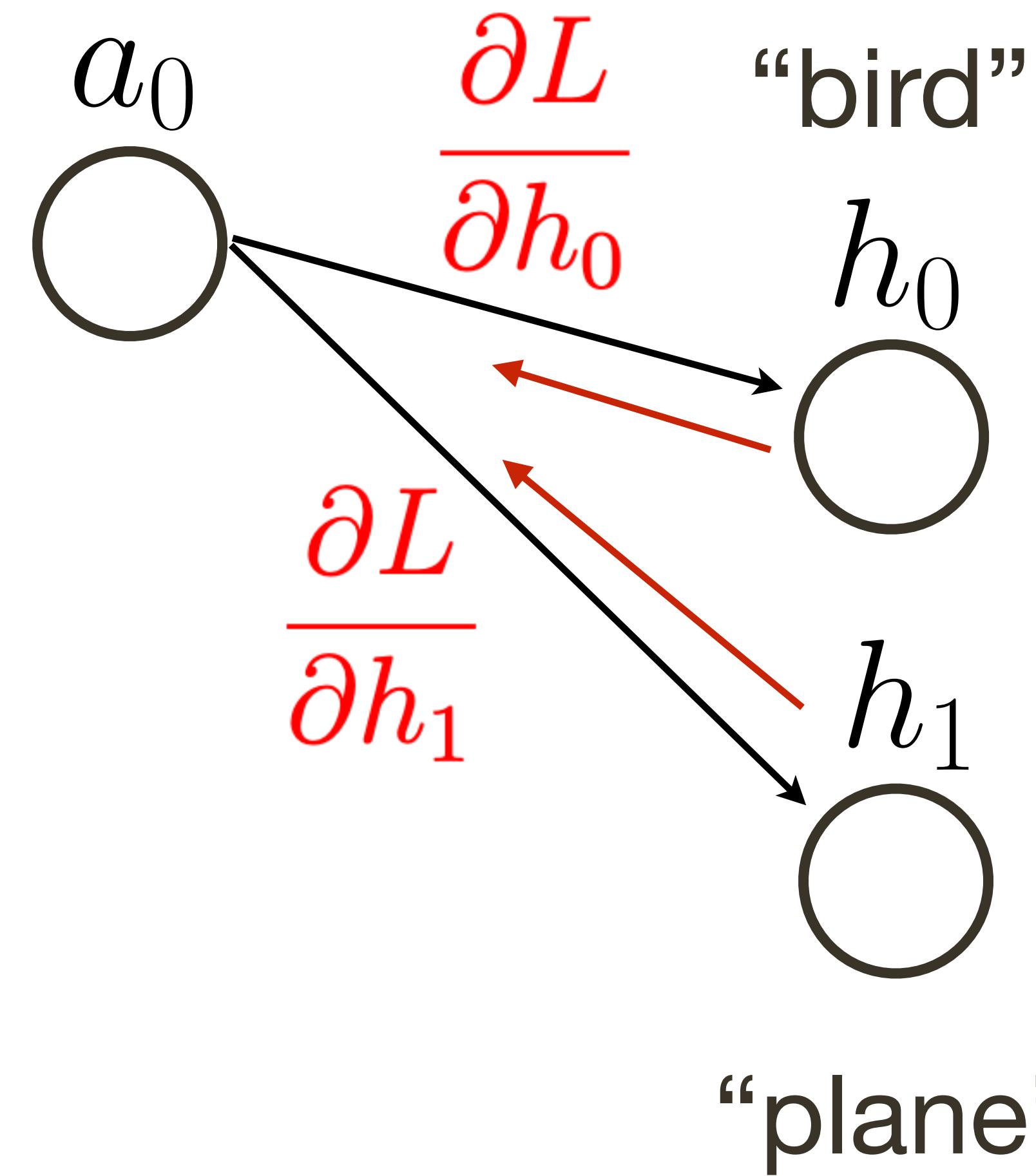
input data



weights

2-Layer Neural Network – multiple outputs

activations



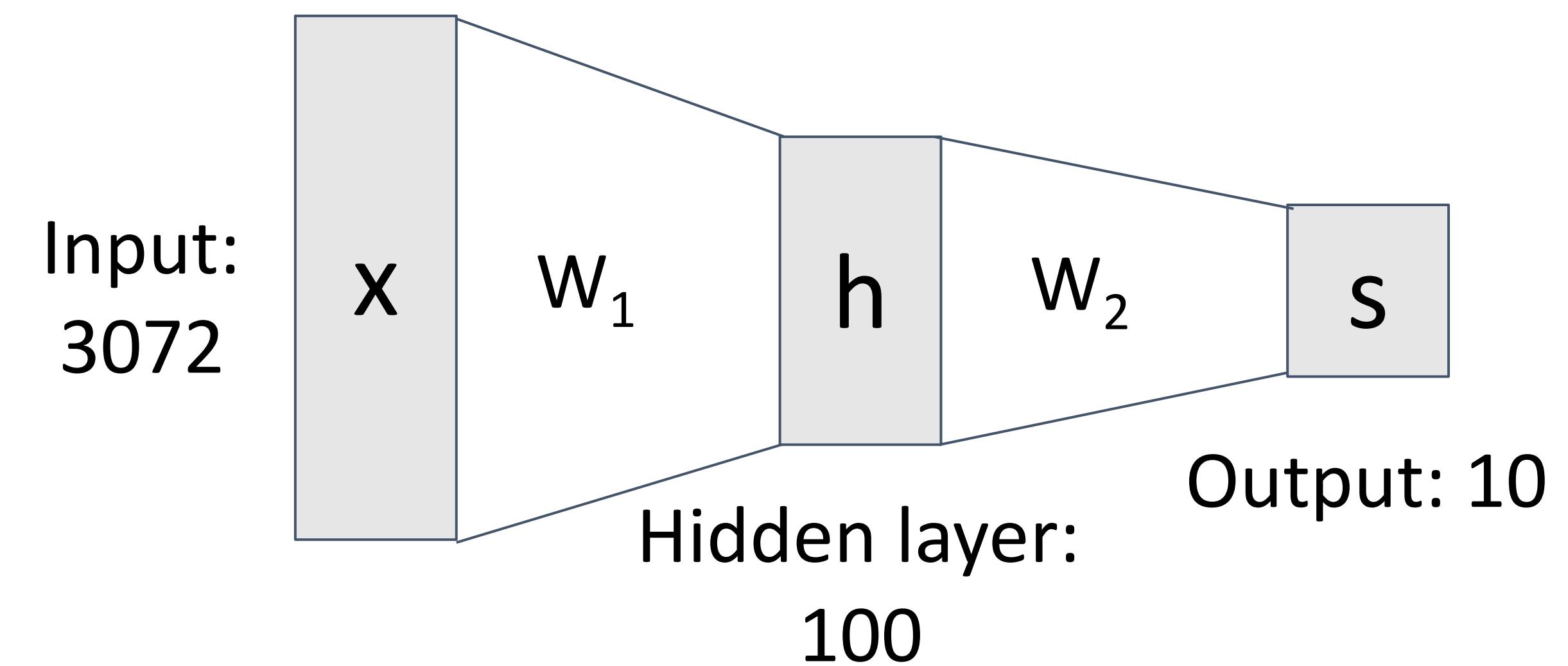
Neural Networks

Linear classifier: One template per class



(Before) Linear score function:

(Now) 2-layer Neural Network



$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

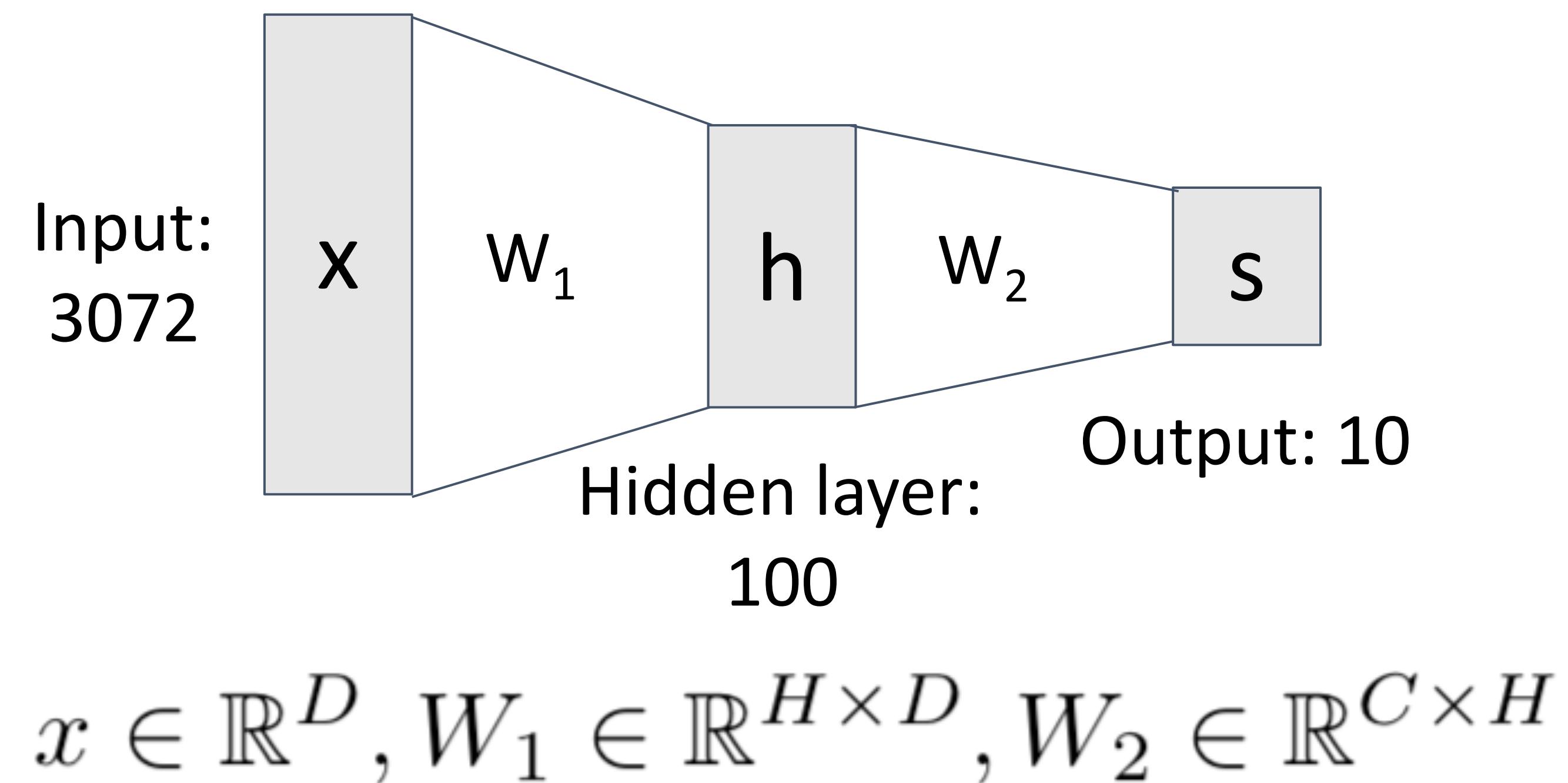
Neural Networks

Neural net: first layer is bank of templates;
Second layer recombines templates



(Before) Linear score function:

(Now) 2-layer Neural Network



Menu for Today

Topics:

- **Neural Networks** introduction
- **Activation functions** softmax, relu
- **2-layer** fully connected net
- **Backprop** intro

Readings:

- **Today's Lecture:** Szeliski 5.1.3, 5.3-5.4, Justin Johnson Michigan EECS 498/598

Reminders:

- **Assignment 5:** due today
- **Assignment 6:** Deep Learning is now available