# Reminders / Announcements

- **Homework 1 is due tonight!**

  o Each student gets 10 late days (total). See the syllabus for details.

  o You **DO NOT** need to email me for permission to use late days!

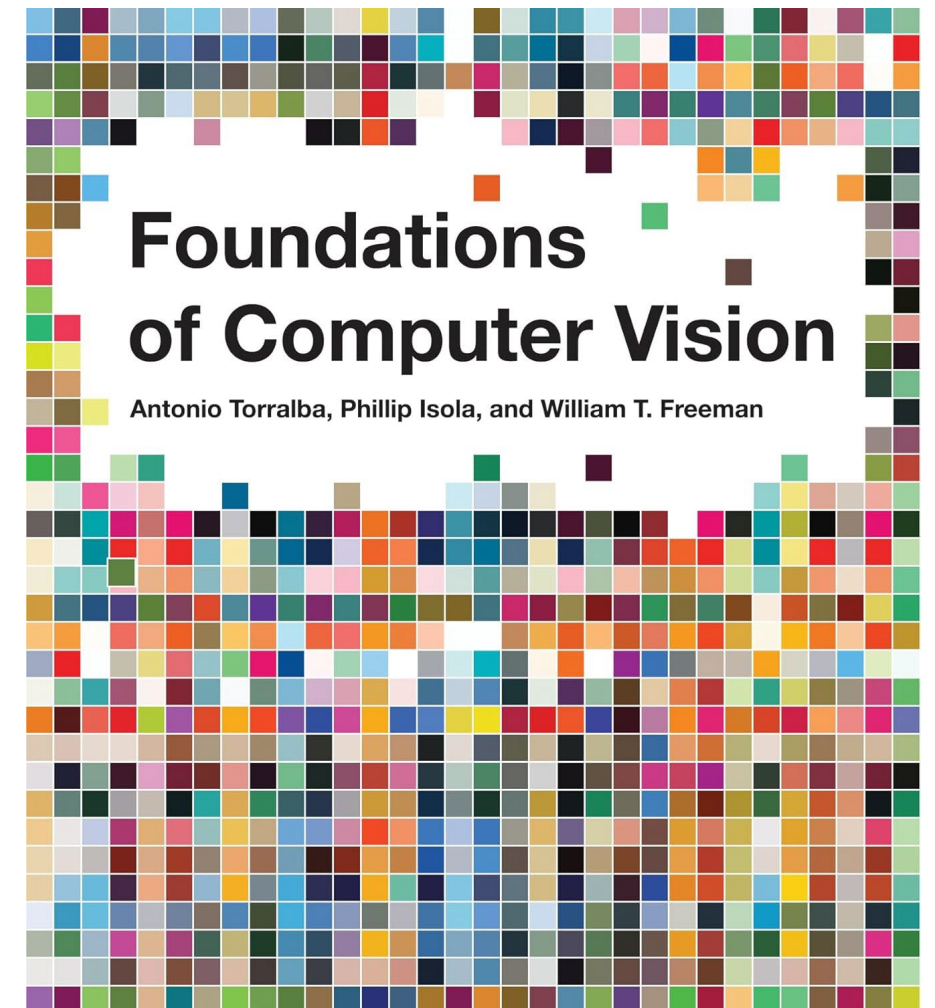- **Project Proposal is due 10/03**

  o Group sizes <3 need my explicit permission!

  o Proposal needs to be turned in on Blackboard by each group member

- **Midterm Exam is on 10/20**

  o In class; closed-book; 1 hour; everything up to and including 10/15 lecture

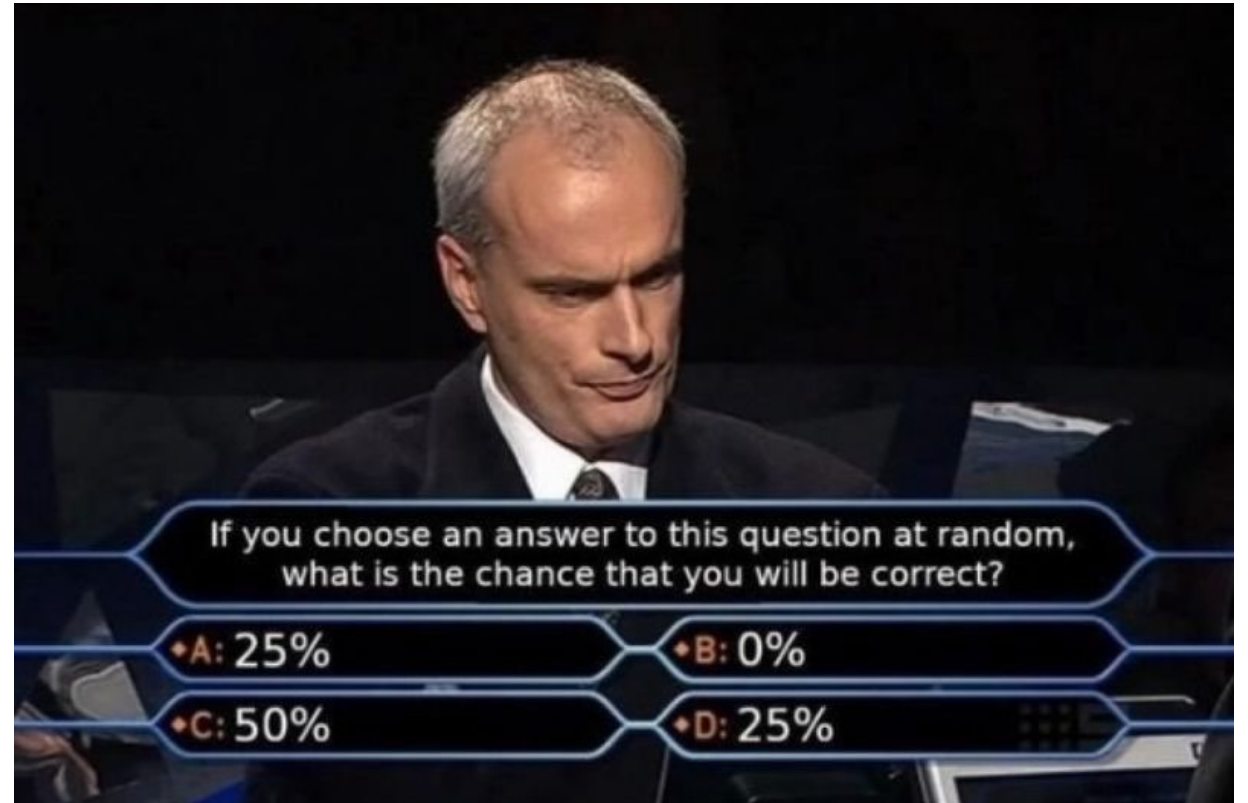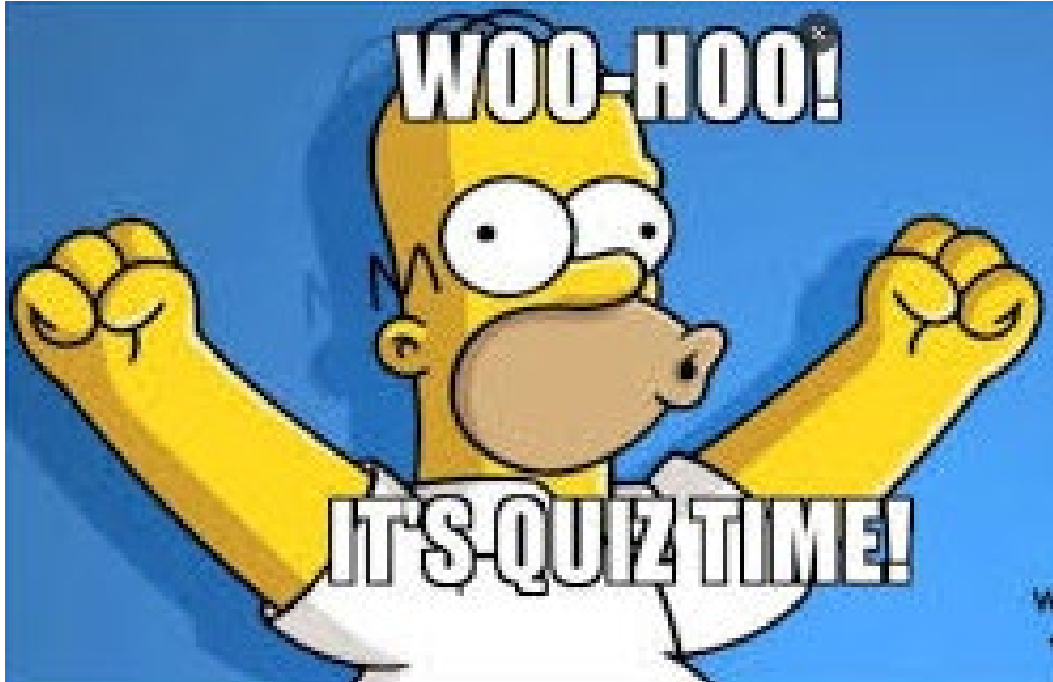  o More details in the next class.

# A New Useful Reference Book

- Available for free on: https://visionbook.mit.edu/

- Published in 2024 – a modern take compared to the other reference books mentioned in the syllabus (these are also available for free)

- As a reminder:  you are encouraged to read relevant chapters of reference books
  - Course website lists book chapters for each lecture
  - This is optional, but encouraged

I've requested UMBC library to buy it, but that might take some time …

# Quiz 3!





If you choose an answer to this question at random, what is the chance that you will be correct?

A: 25%     B: 0%

C: 50%     D: 25%

CMSC 472/672 Computer Vision

# Lecture 8: Neural Networks



Some slides from Suren Jayasuriya (ASU), Phillip Isola (MIT)

UMBC

Artificial Intelligence

$$\hat{y} = w^\top x + b$$

# Limitations to linear classifiers



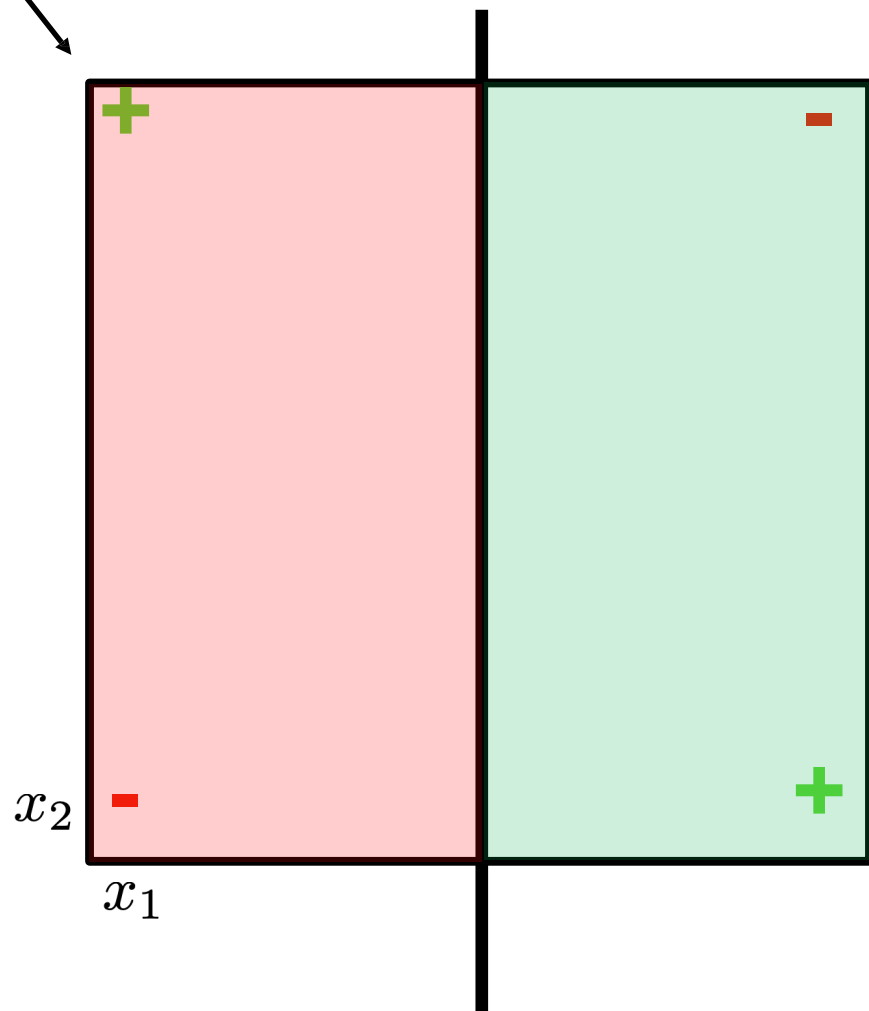|       | $x_2$ | |
|-------|-------|-------|
| $x_1$ |   0   |   1   |
| 0     |   0   |   1   |
| 1     |   1   |   0   |

XOR

# Limitations to linear classifiers

**Wrong!**

**Wrong!**



|       | $x_2$ | |
|-------|-----|-----|
|       | 0   | 1   |
| $x_1$ 0 | 0 | 1 |
| 1     | 1   | 0   |

XOR

# Limitations to linear classifiers

**Wrong!**

**Wrong!**

|       | $x_2$ | |
|-------|-------|---|
|       | 0     | 1 |
| $x_1$ 0 | 0   | 1 |
| 1     | 1     | 0 |

XOR

# Goal: Non-linear decision boundary



|       | $x_2$ |   |
| :---: | :---: | :---: |
|       | 0 | 1 |
| $x_1$ 0 | 0 | 1 |
| 1     | 1 | 0 |

XOR

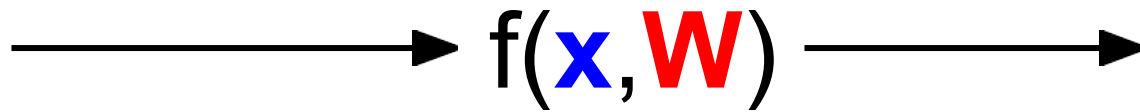# A brief history of Neural Networks

enthusiasm

time

# Parametric Approach

Image



Array of **32x32x3** numbers
(3072 numbers total)

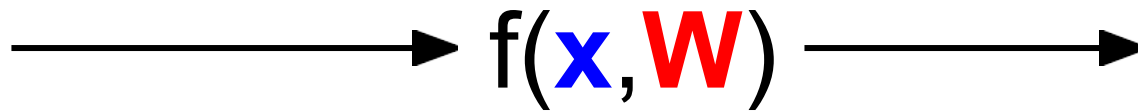f(**x**,**W**) → **10** numbers giving class scores

↑

**W**
parameters
or weights

# Parametric Approach: Linear Classifier

$$f(x,W) = Wx$$
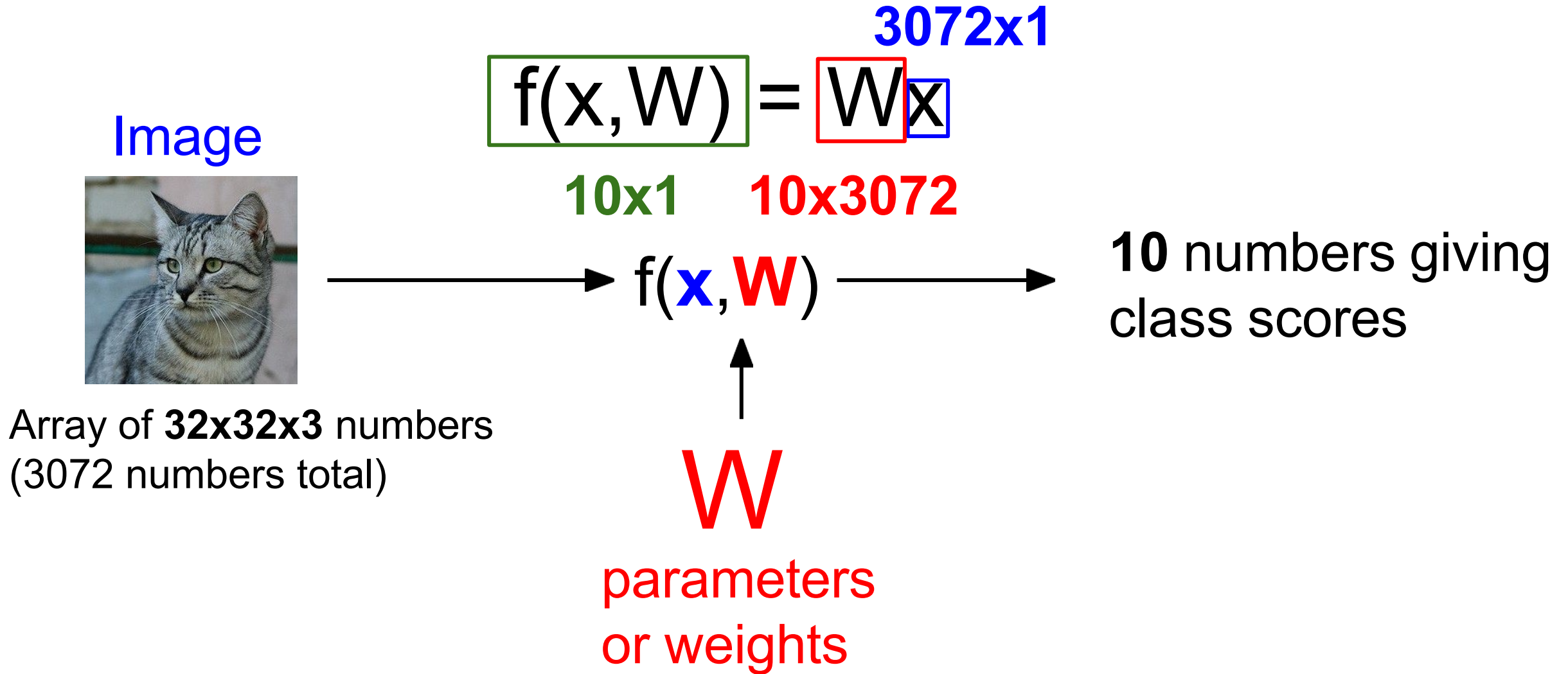
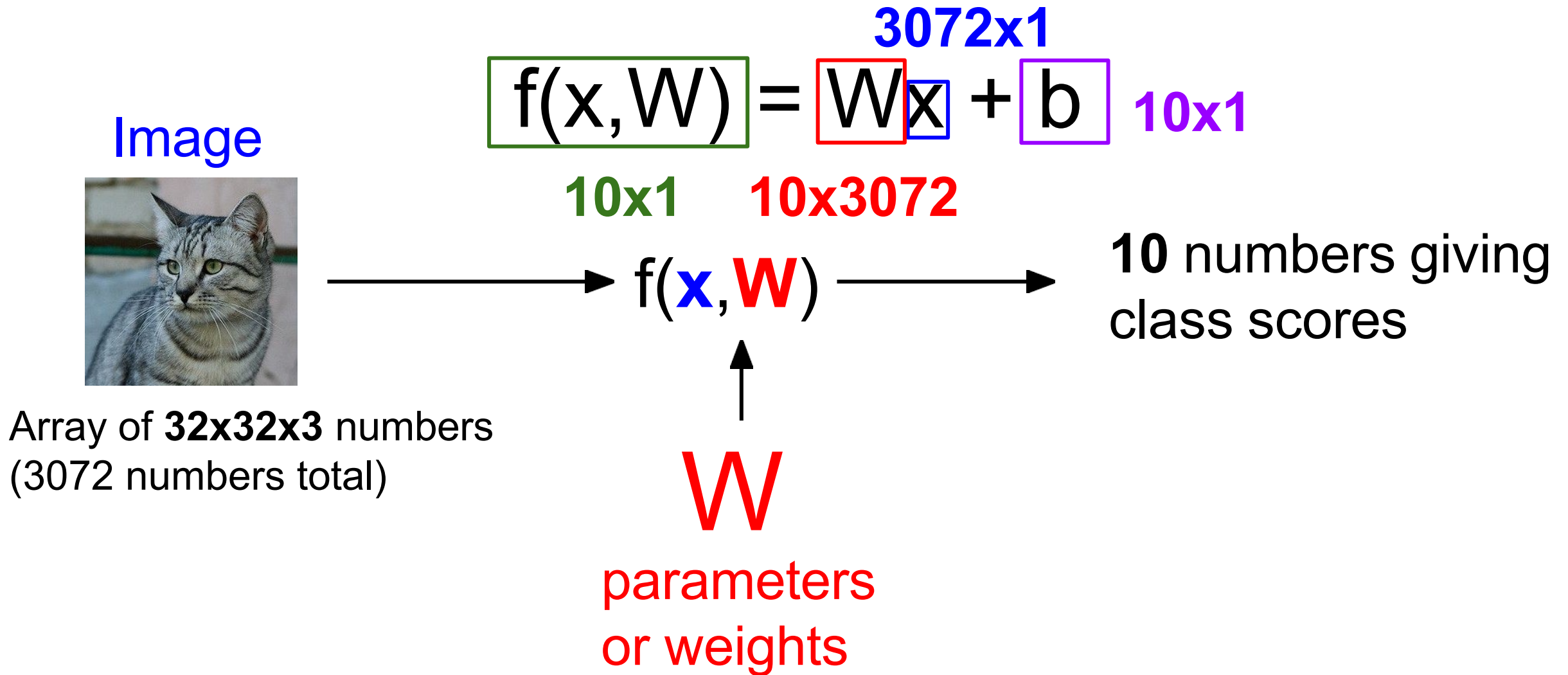Image



Array of **32x32x3** numbers
(3072 numbers total)

$f(\textcolor{blue}{\mathbf{x}},\textcolor{red}{\mathbf{W}})$

**10** numbers giving class scores

$\textcolor{red}{W}$

parameters or weights

# Parametric Approach: Linear Classifier

$$f(x,W) = Wx$$

3072x1

10x1    10x3072

**Image**



Array of **32x32x3** numbers
(3072 numbers total)

f(**x**,**W**)

**10** numbers giving
class scores

**W**

parameters
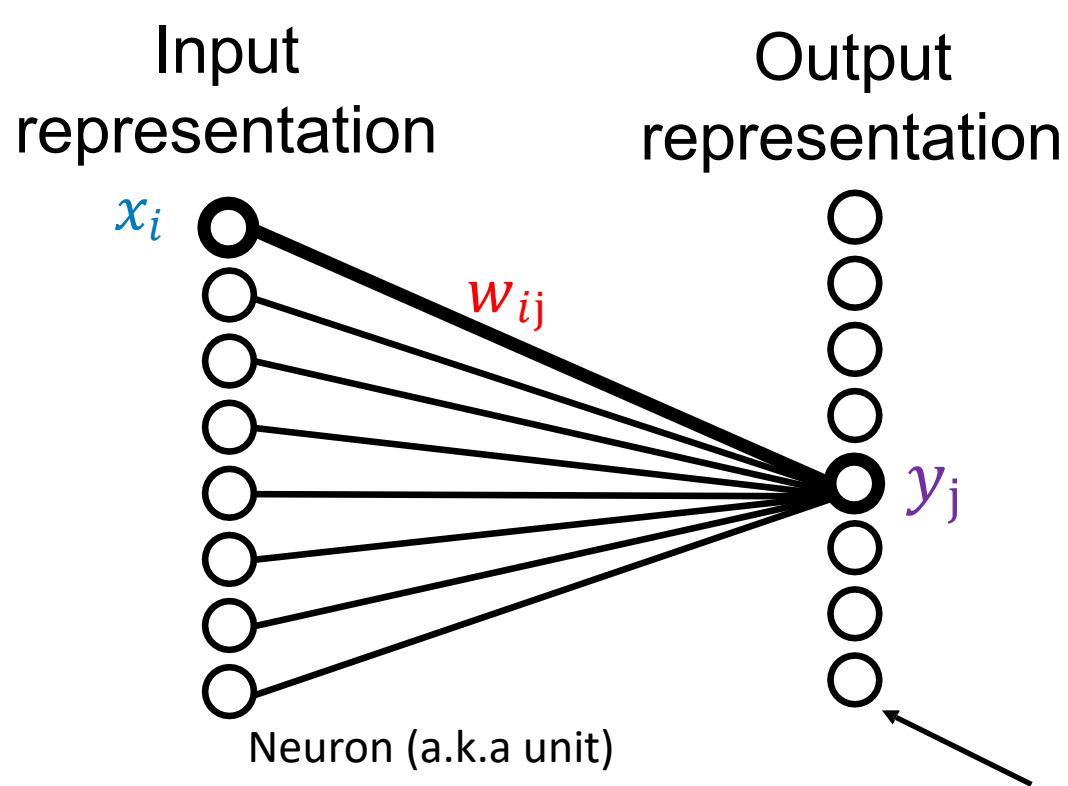or weights

# Parametric Approach: Linear Classifier

$$f(x,W) = Wx + b$$

**3072x1**
**10x1**    **10x3072**
**10x1**

**Image**



Array of **32x32x3** numbers
(3072 numbers total)

$f(\mathbf{x}, \mathbf{W})$

**10** numbers giving class scores

**W**
parameters
or weights

# Computation in a neural net

Let's say we have some 1D input that we want to convert to some new feature space:

**Linear layer**



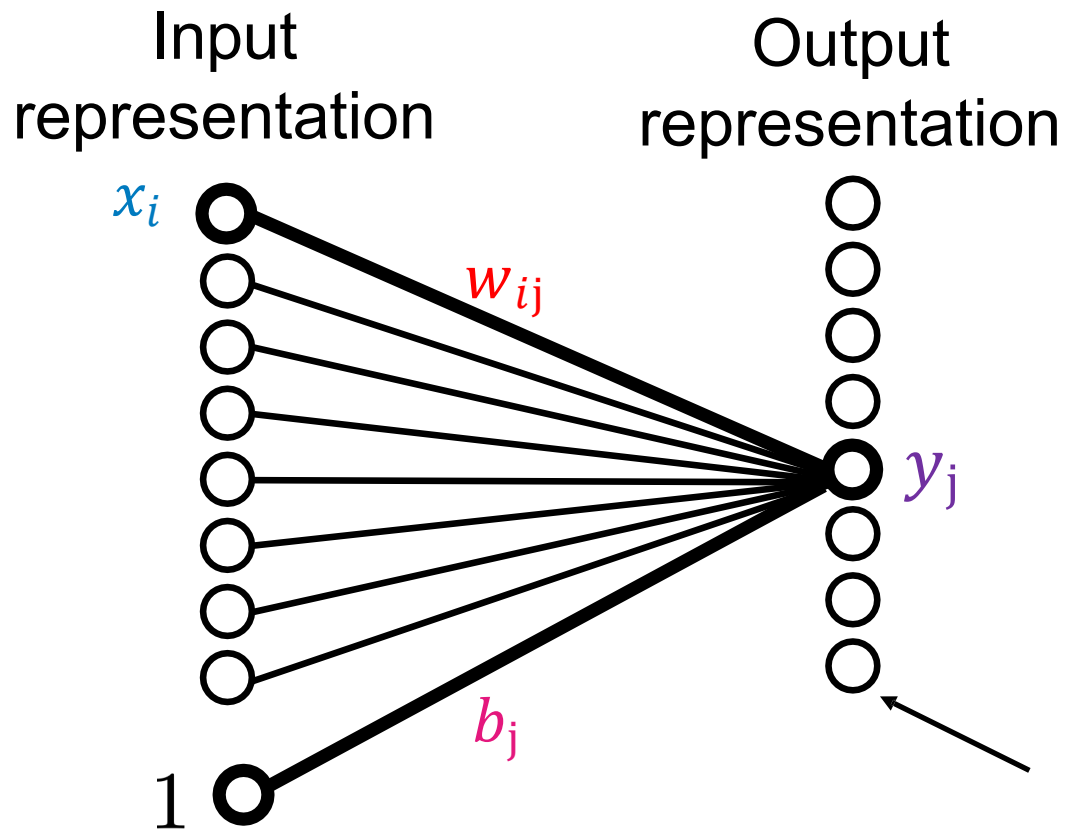$$y_j = \sum_i w_{ij} x_i$$

weights

Input representation

$x_i$

Output representation

$y_j$

$w_{ij}$

Neuron (a.k.a unit)

# Computation in a neural net

Let's say we have some 1D input that we want to convert to some new feature space
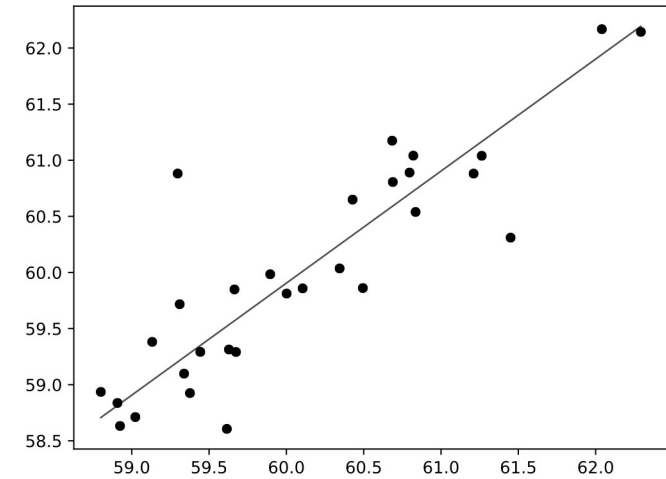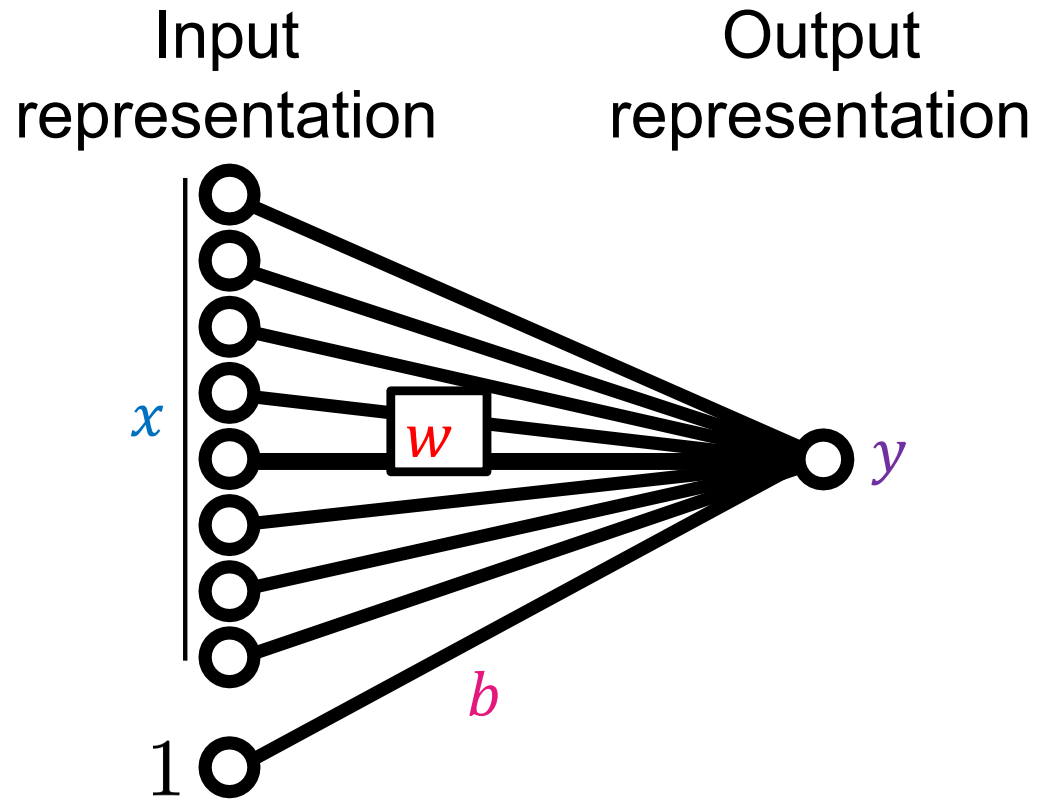
## <u>Linear layer</u>



$$y_{\text{j}} = \sum_i w_{i\text{j}} x_i + b_{\text{j}}$$

weights

bias

# Example: Linear Regression

**Linear layer**

Input
representation

Output
representation

$x$

$w$

$b$

$y$

$1$
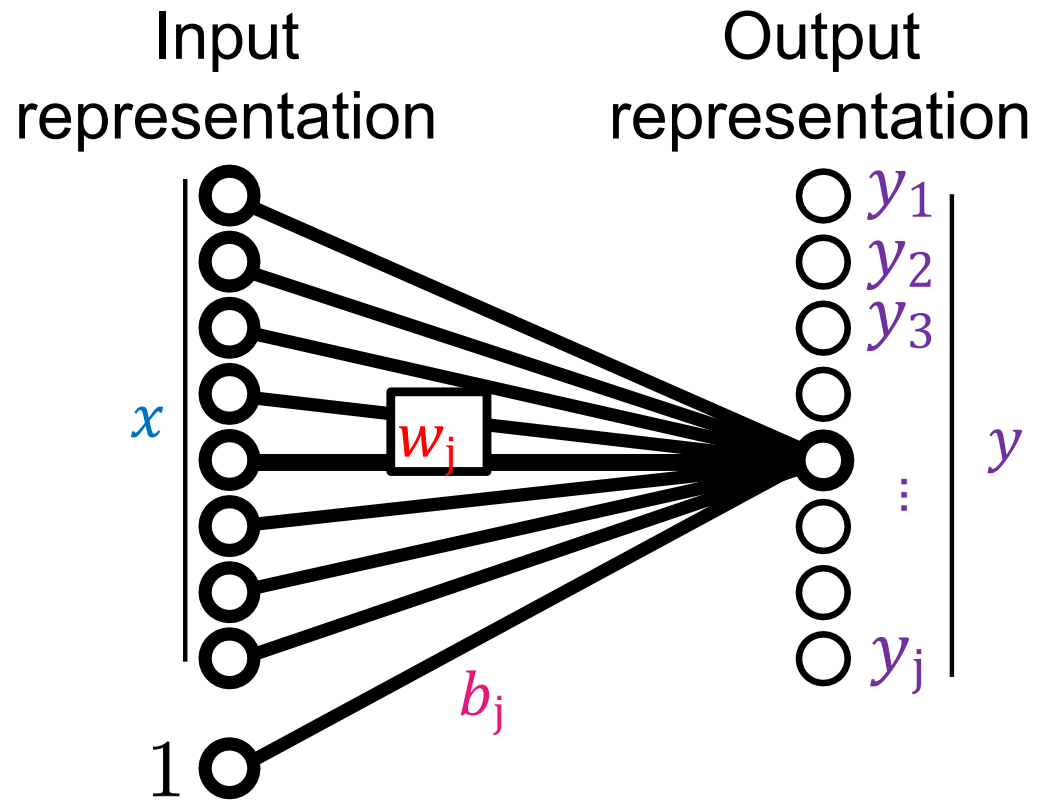
$$f_{\mathbf{w},b}(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b$$

# Computation in a neural net – Full Layer

**Linear layer**



Input representation

Output representation

$x$

$w_j$

$b_j$

1

$y_1$
$y_2$
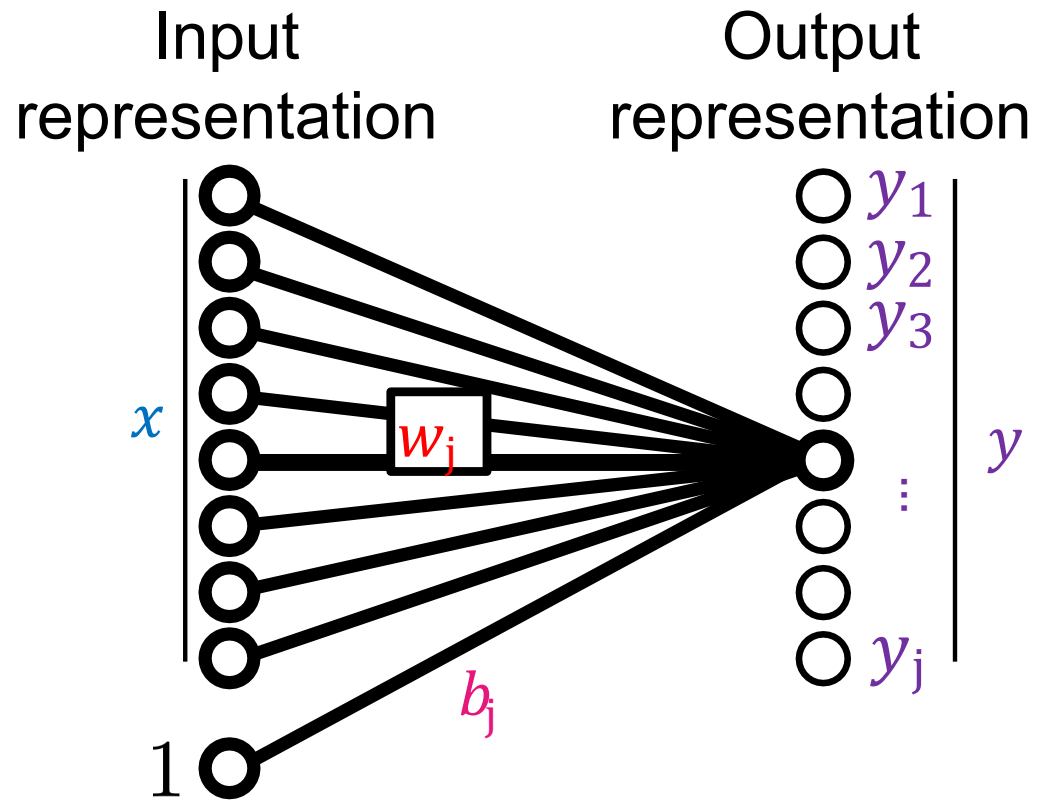$y_3$
$\vdots$
$y_j$

$y$

$$y = Wx + b$$

$$\begin{bmatrix} w_{11} & \cdots & w_{1n} \\ \vdots & \ddots & \vdots \\ w_{j1} & \cdots & w_{jn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_n \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \cdots \\ b_j \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_j \end{bmatrix}$$

**parameters of the model: $\boldsymbol{\theta} = \{\boldsymbol{W}, \boldsymbol{b}\}$**

# Computation in a neural net – Full Layer

**Linear layer**

Input
representation

Output
representation



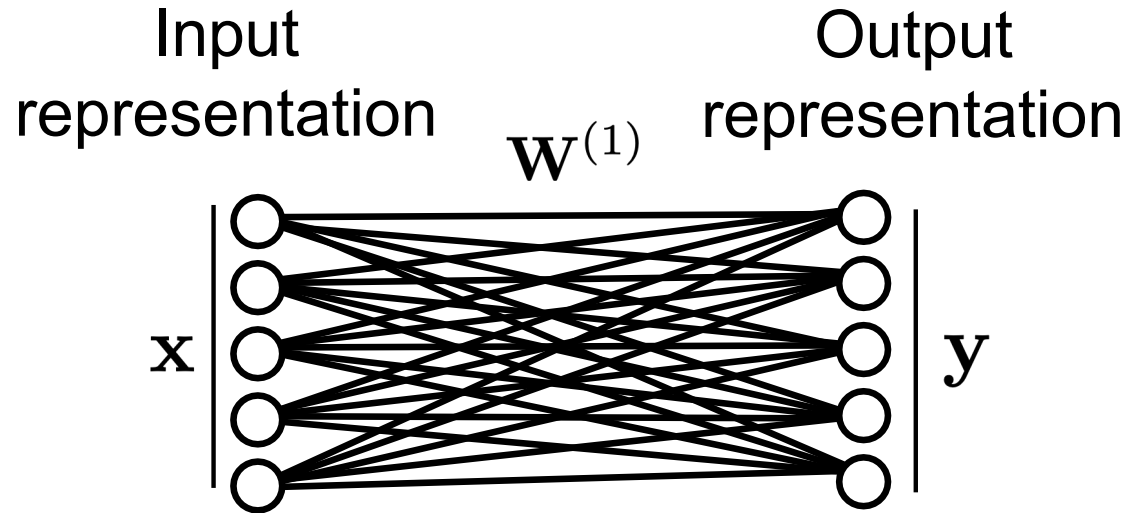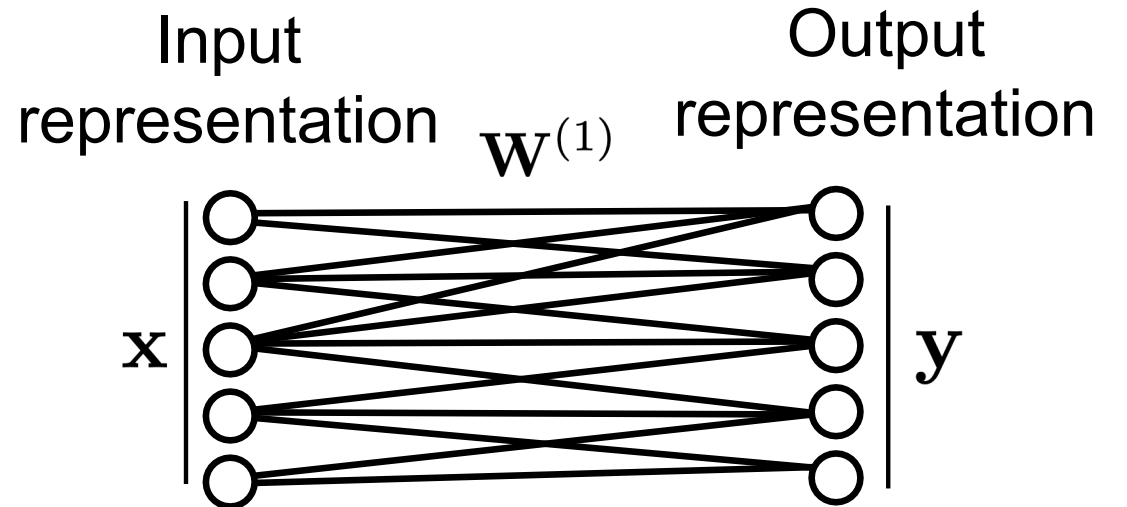**Full layer**

$$y = Wx + b$$

$$\begin{bmatrix} w_{11} & \cdots & w_{jn} & b_1 \\ \vdots & \ddots & \vdots & \vdots \\ w_{j1} & \cdots & w_{jn} & b_j \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \ldots \\ x_n \\ 1 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \ldots \\ y_j \end{bmatrix}$$

Can again simplify notation by
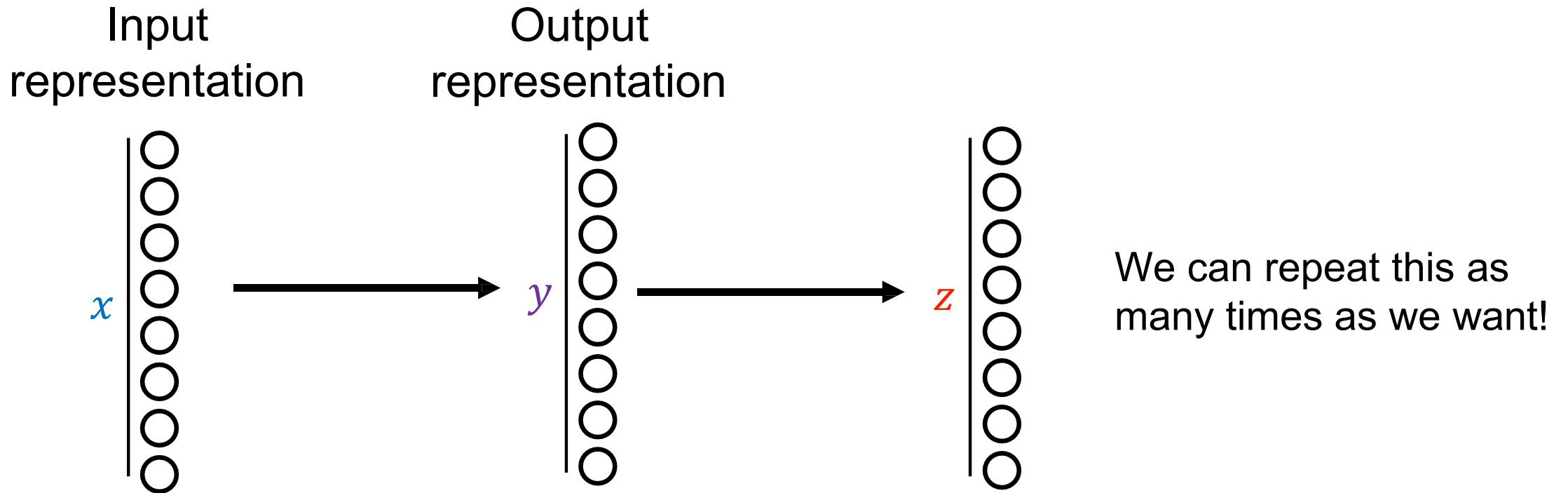appending a 1 to **x**

# Connectivity patterns



Input representation $\mathbf{W}^{(1)}$ Output representation

$\mathbf{x}$ $\mathbf{y}$

*Fully connected layer*

Input representation $\mathbf{W}^{(1)}$ Output representation

$\mathbf{x}$ $\mathbf{y}$
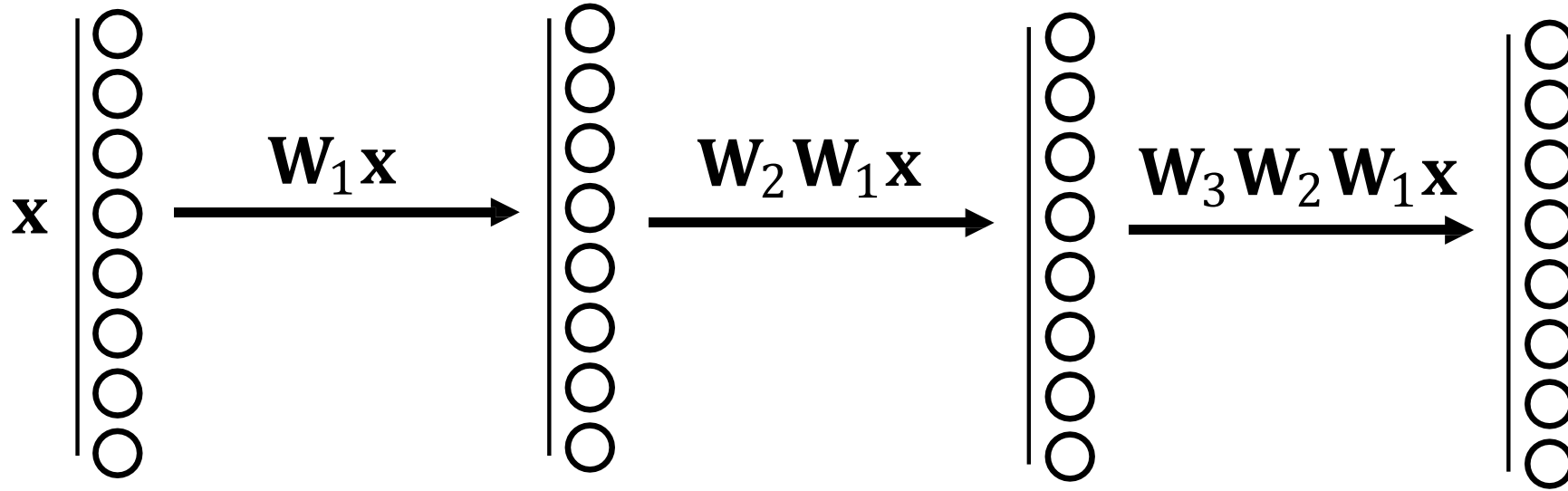
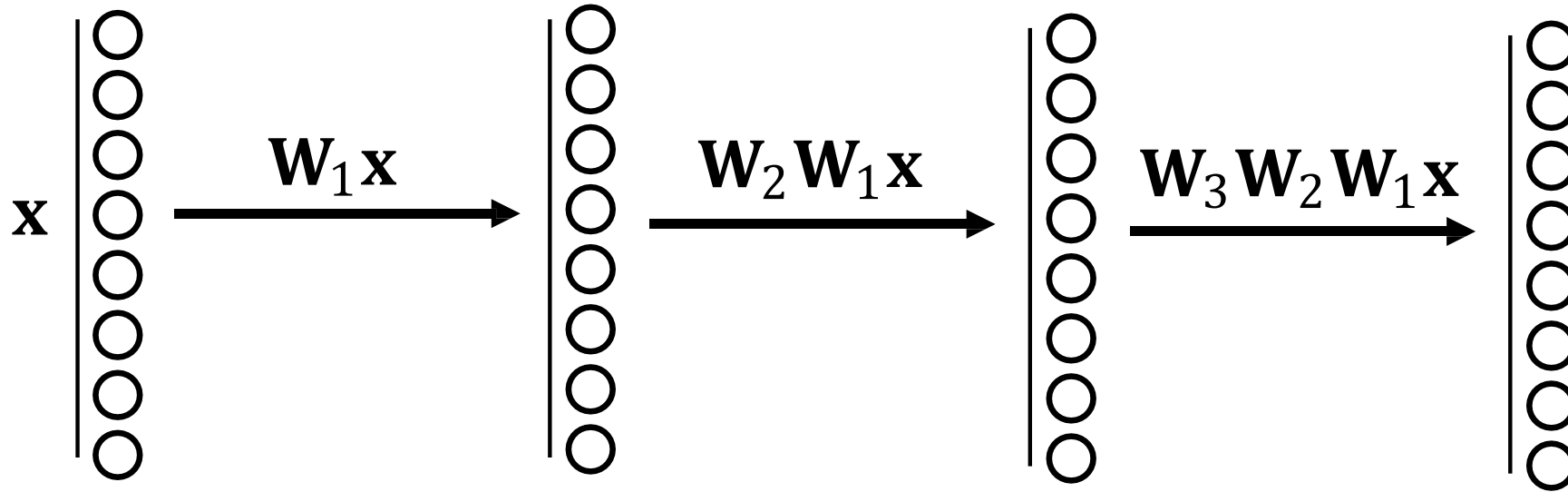*Locally connected layer (Sparse W)*

# Computation in a neural network

We can now transform our input representation vector into some output representation vector using a bunch of linear combinations of the input:

Input representation

Output representation

$x$

$y$

$z$

We can repeat this as many times as we want!

# What is the problem with this idea?

# What is the problem with this idea?

$$\mathbf{x} \quad \xrightarrow{\mathbf{W_1 x}} \quad \xrightarrow{\mathbf{W_2 W_1 x}} \quad \xrightarrow{\mathbf{W_3 W_2 W_1 x}}$$

Can be expressed as single linear layer!

$$\widehat{W} x$$

Limited power: can't solve XOR ☹

# Recall
# Goal: Non-linear decision boundary



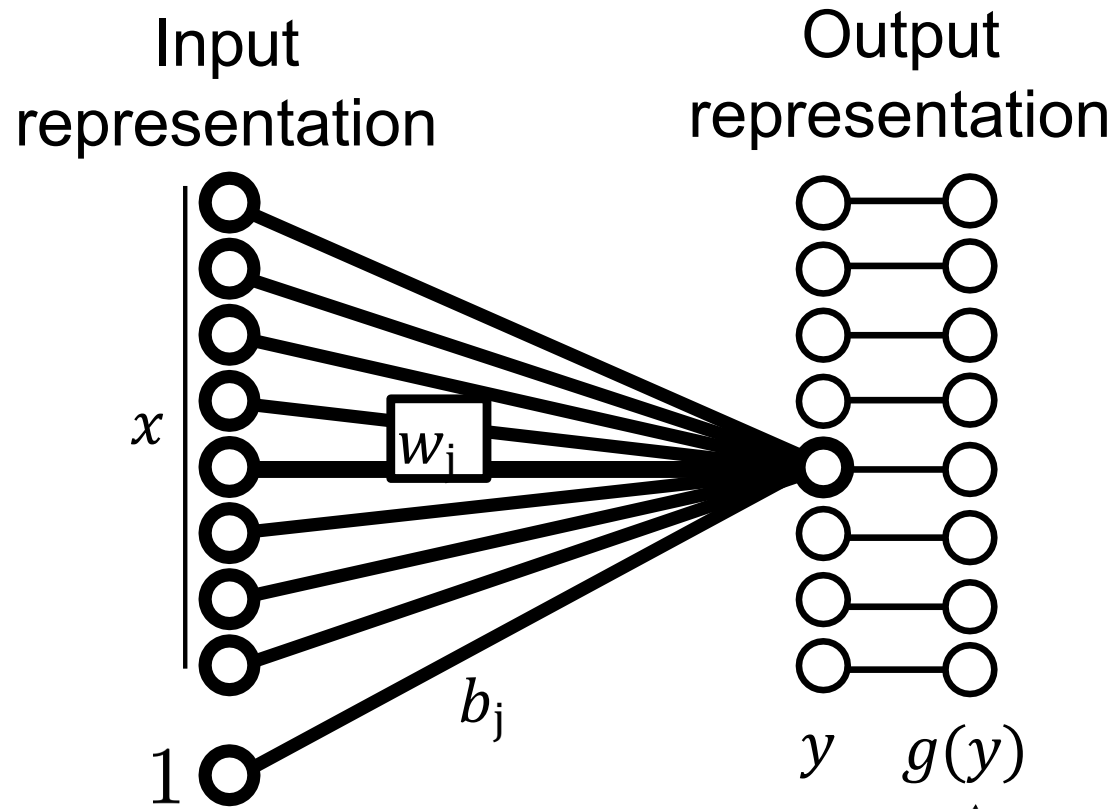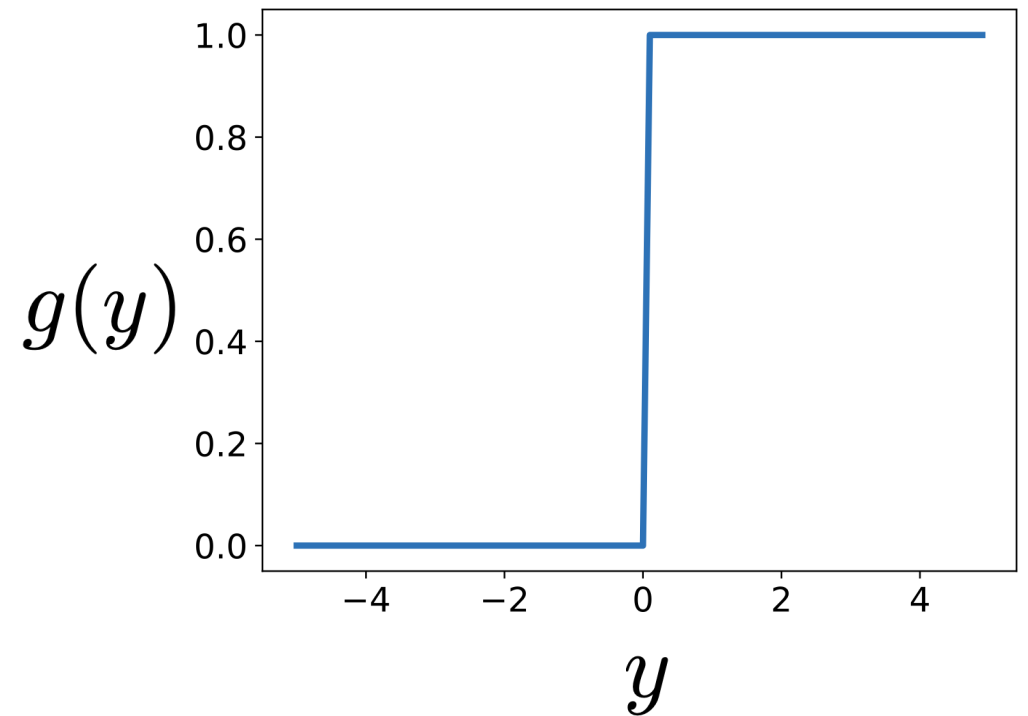|       |   | $x_2$ | |
|-------|---|---|---|
|       |   | 0 | 1 |
| $x_1$ | 0 | 0 | 1 |
|       | 1 | 1 | 0 |

XOR

# Solution: simple nonlinearity

**Linear layer**

Input representation

Output representation

$x$

$w_i$

$b_j$

1

$y$ $g(y)$

Pointwise Non-linearity

$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$

$g(y)$

# The Perceptron

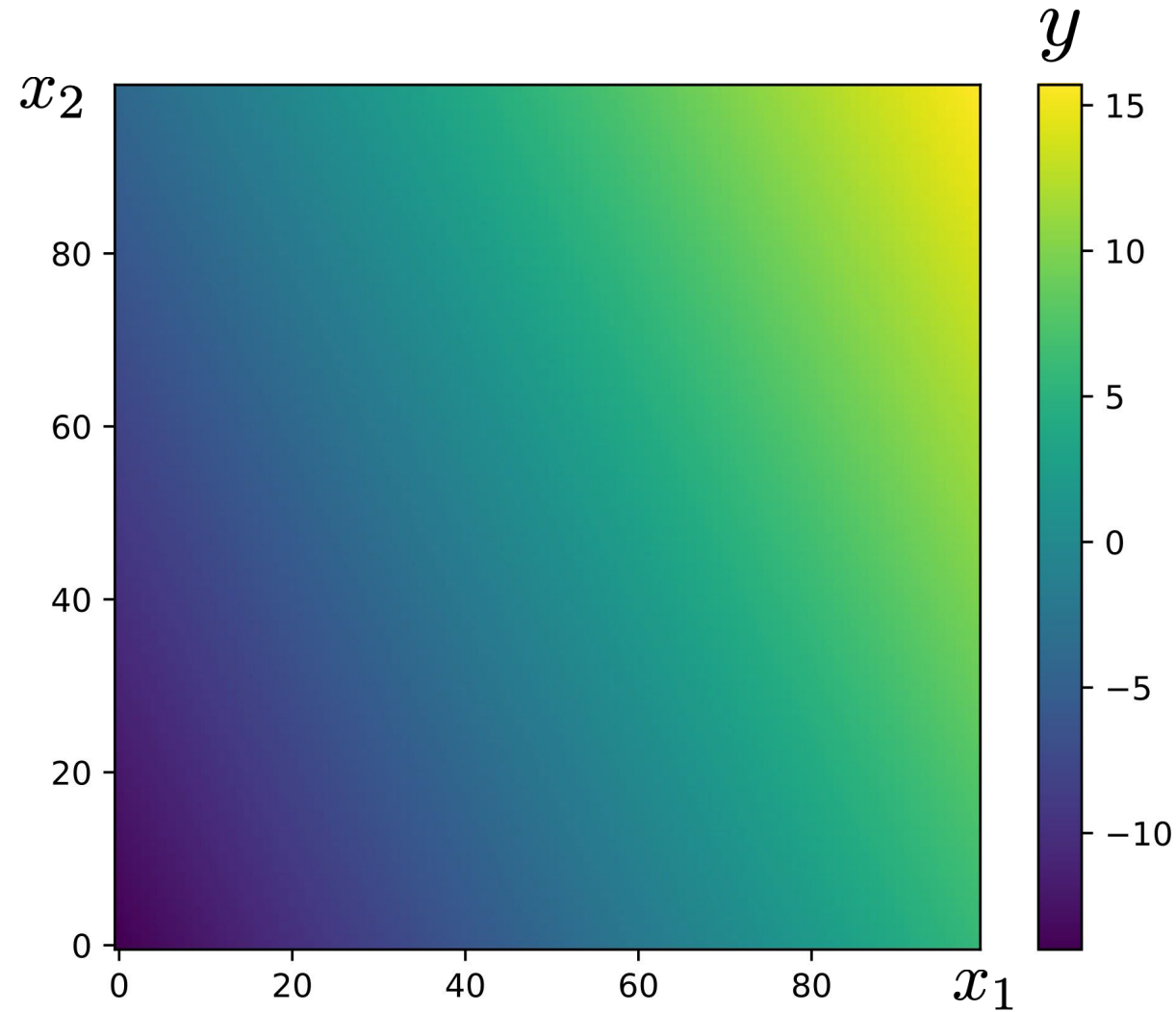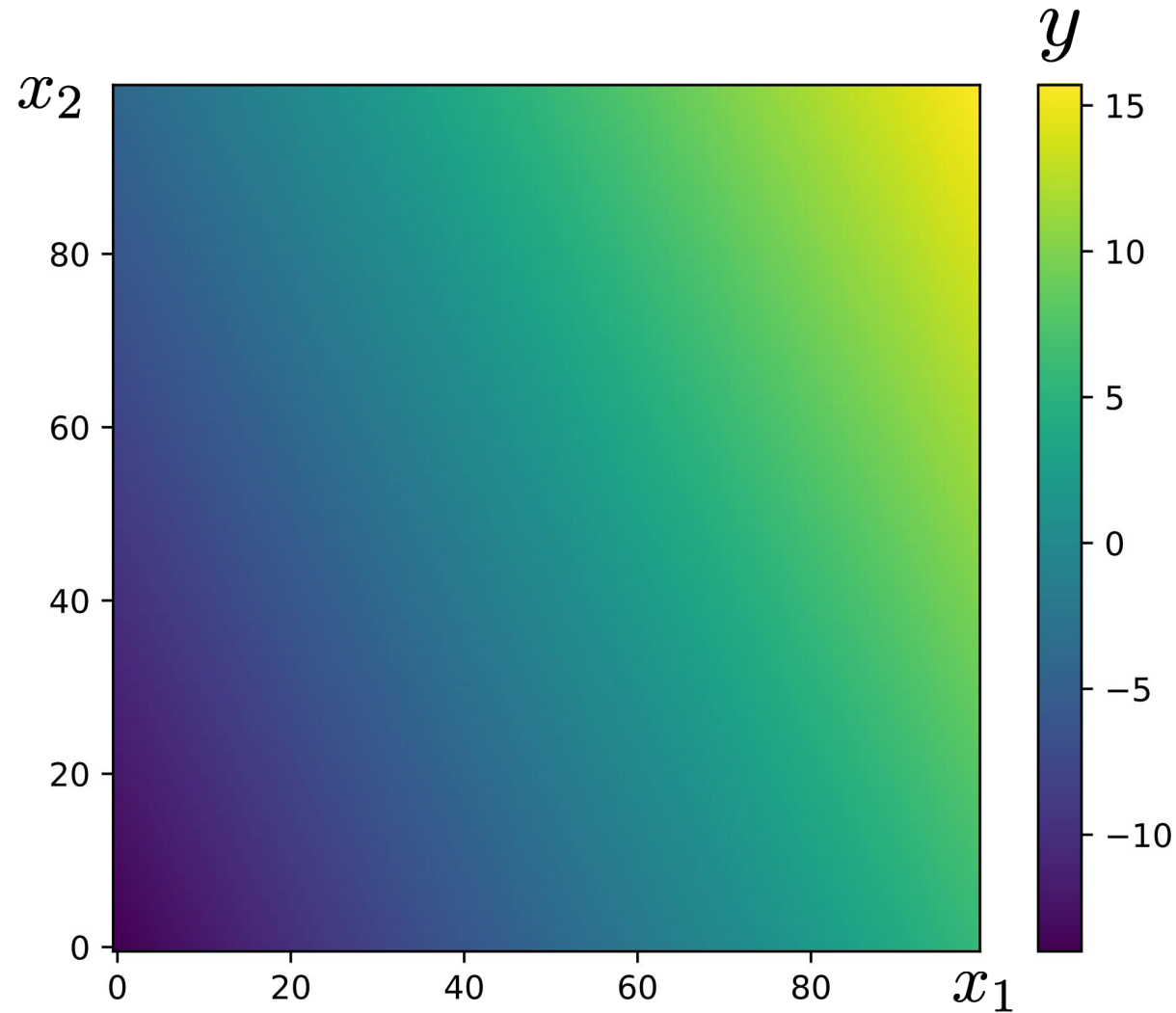# Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

# Example: linear classification with a perceptron



$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$
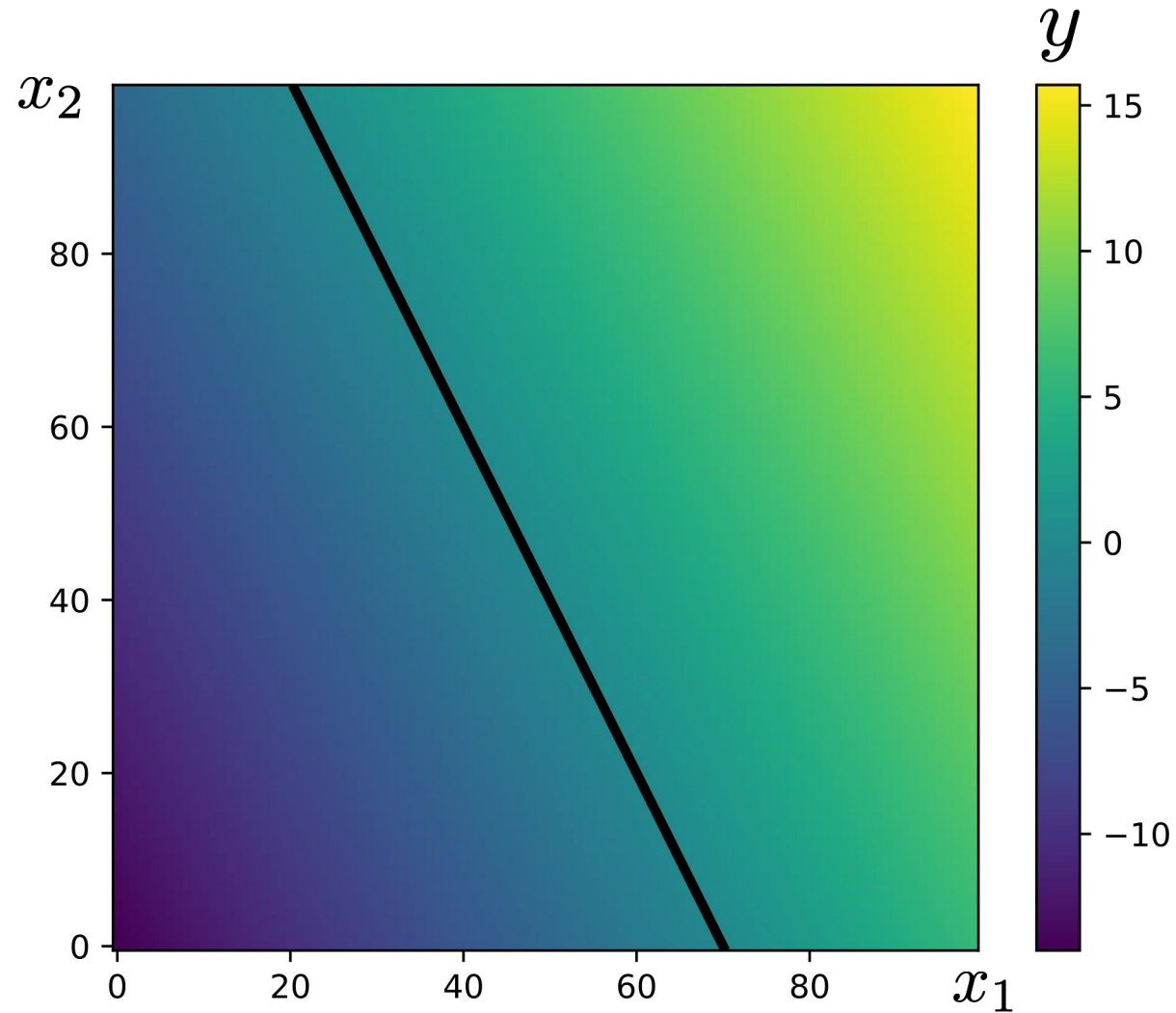
# Example: linear classification with a perceptron
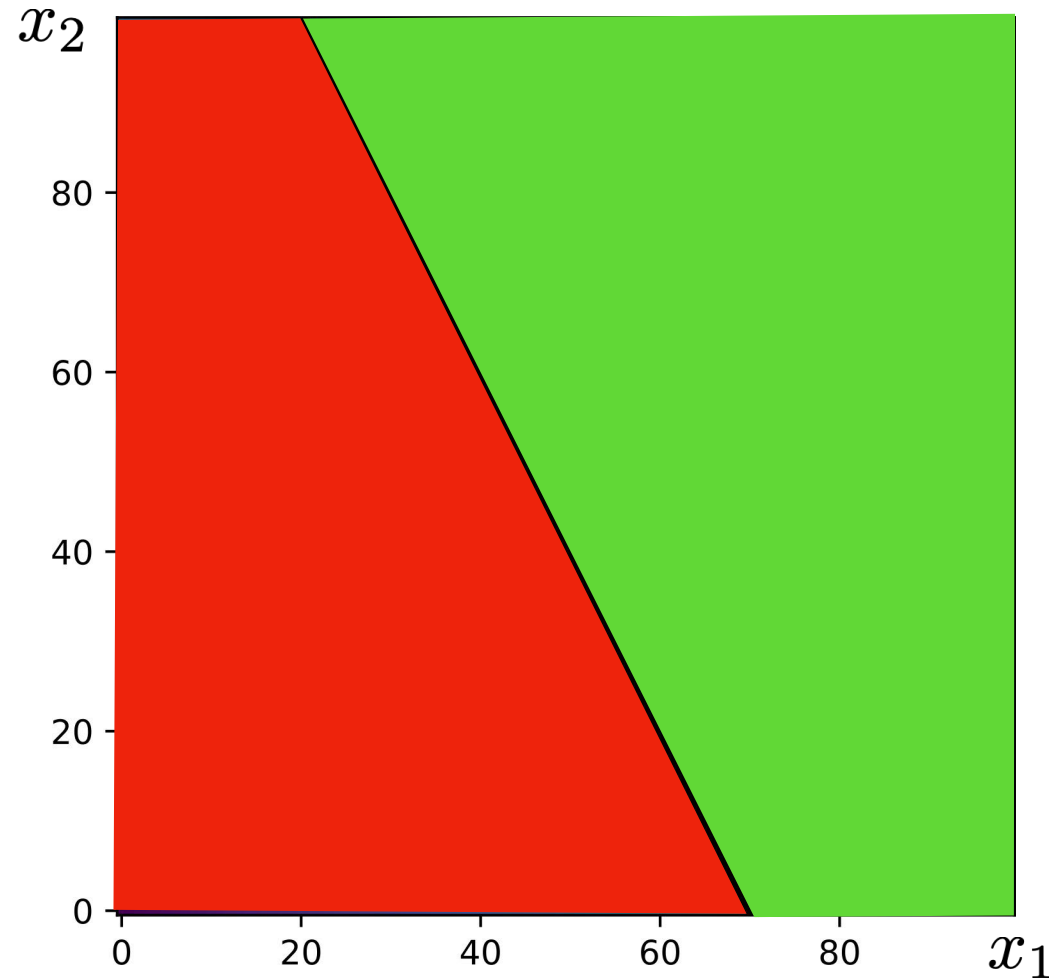


$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$

"when y is greater than 0, set all pixel values to 1 (green), otherwise, set all pixel values to 0 (red)"

# Example: linear classification with a perceptron
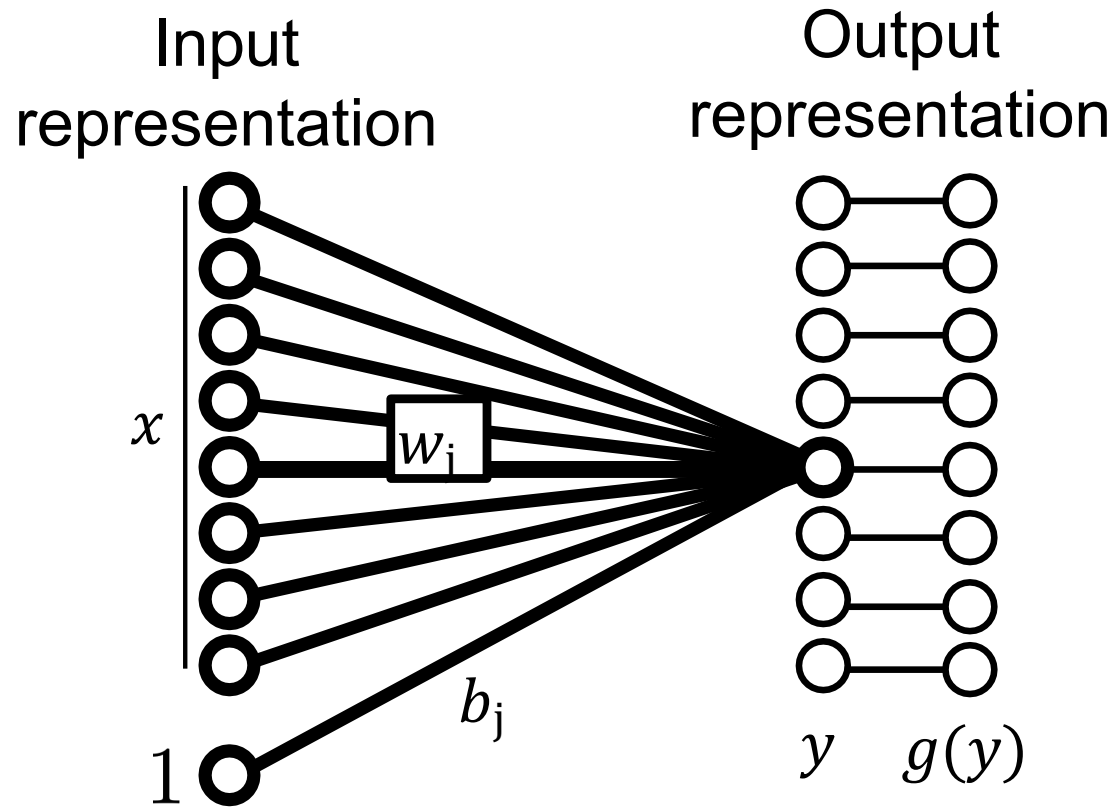
$g(y)$



$$y = \mathbf{x}^T \mathbf{w} + b$$

$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$

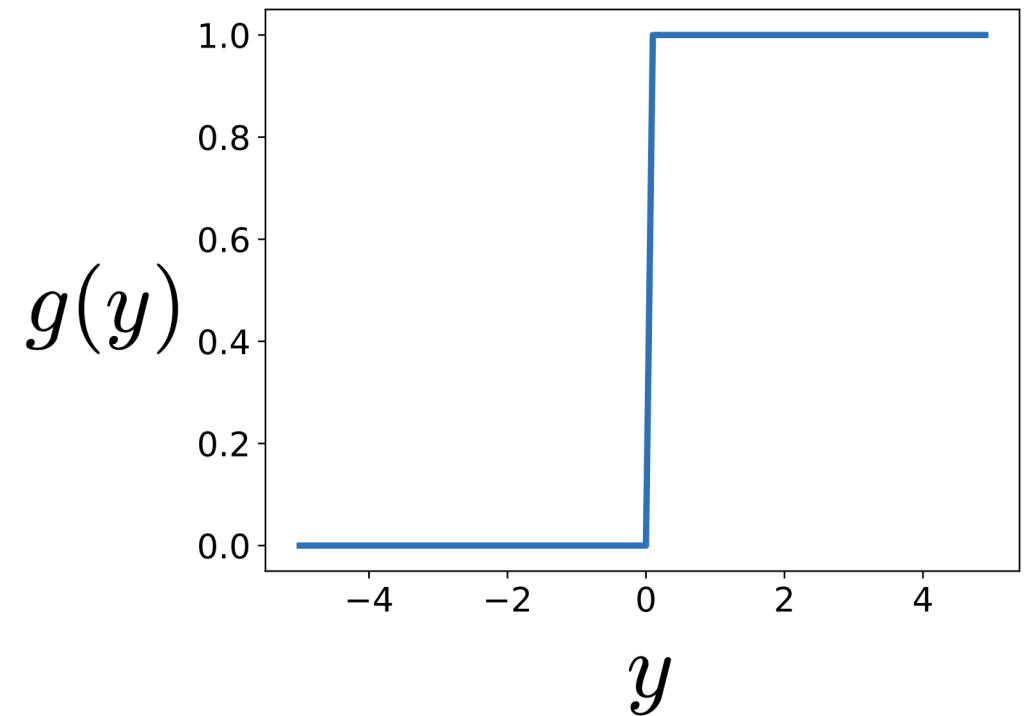"when y is greater than 0, set all pixel values to 1 (green), otherwise, set all pixel values to 0 (red)"

# Computation in a neural net - nonlinearity

**Linear layer**

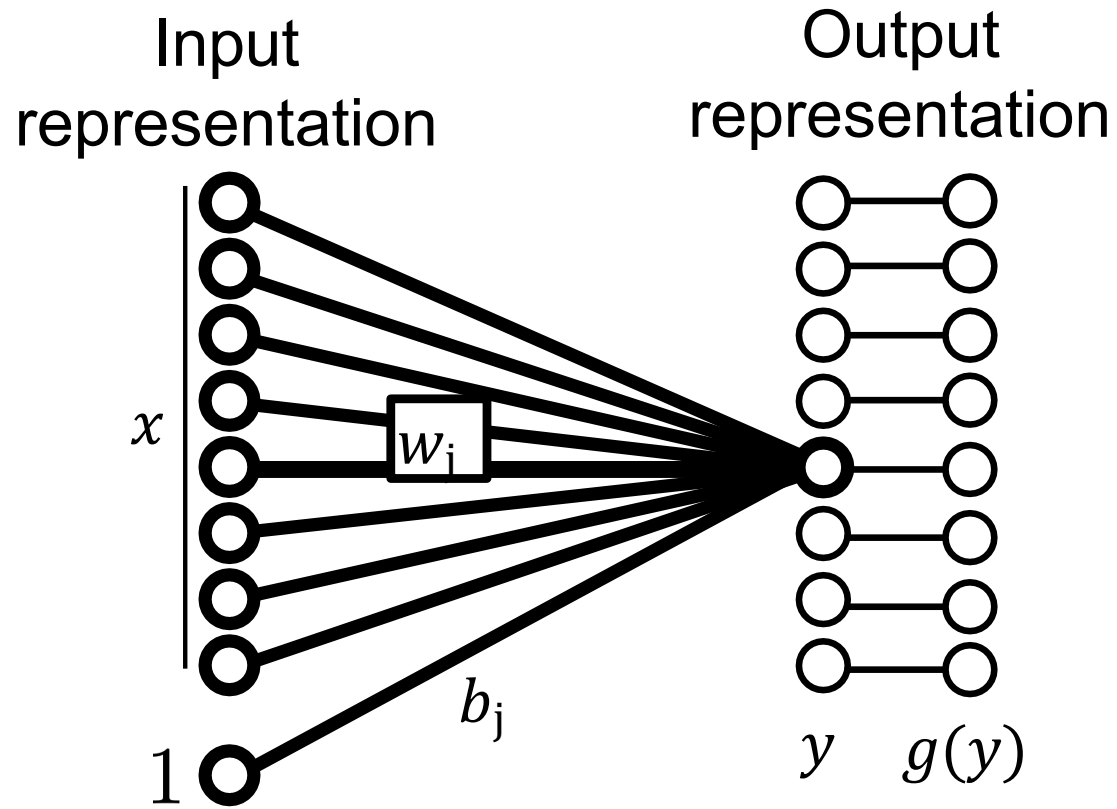$$g(y) = \begin{cases} 1, & \text{if} \quad y > 0 \\ 0, & \text{otherwise} \end{cases}$$



Input representation

Output representation

$x$

$w_i$

$b_j$

$1$

$y \quad g(y)$

$g(y)$

$y$

Can't use gradient-based optimization, $\dfrac{\partial}{\partial y} g = 0$

# Computation in a neural net - nonlinearity

**Linear layer**

Input representation

Output representation

$x$

$w_\mathrm{i}$

$b_\mathrm{j}$

$1$

$y$  $g(y)$

**Sigmoid**
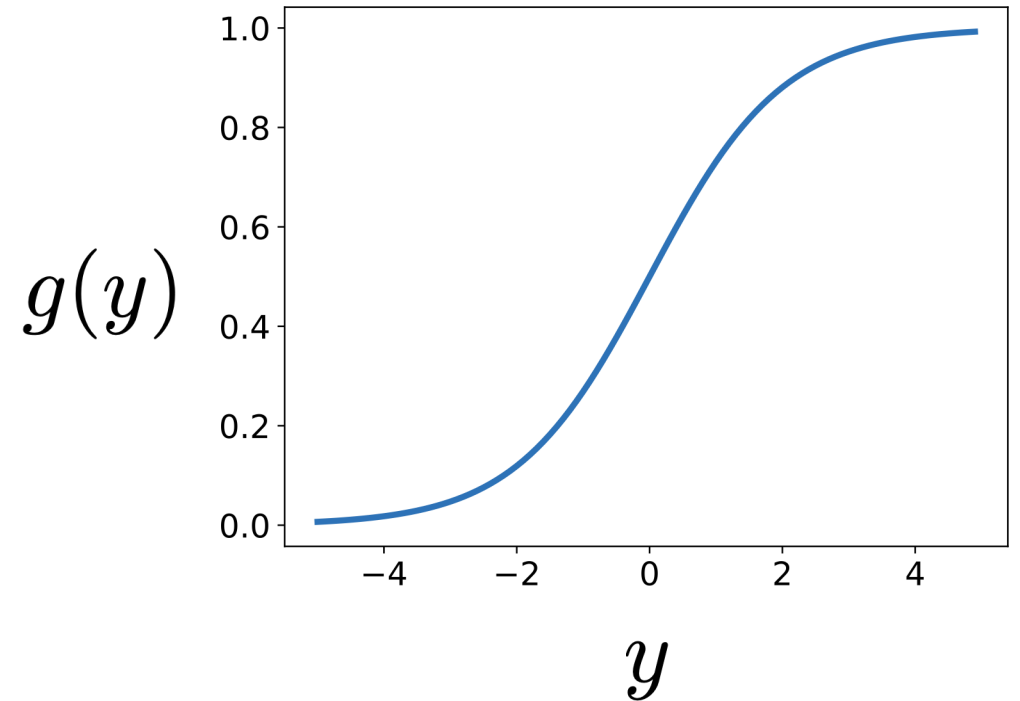
$$g(y) = \sigma(y) = \frac{1}{1 + e^{-y}}$$

# Computation in a neural net - nonlinearity

- Bounded between [0,1]

- Saturation for large +/- inputs

- Gradients go to zero

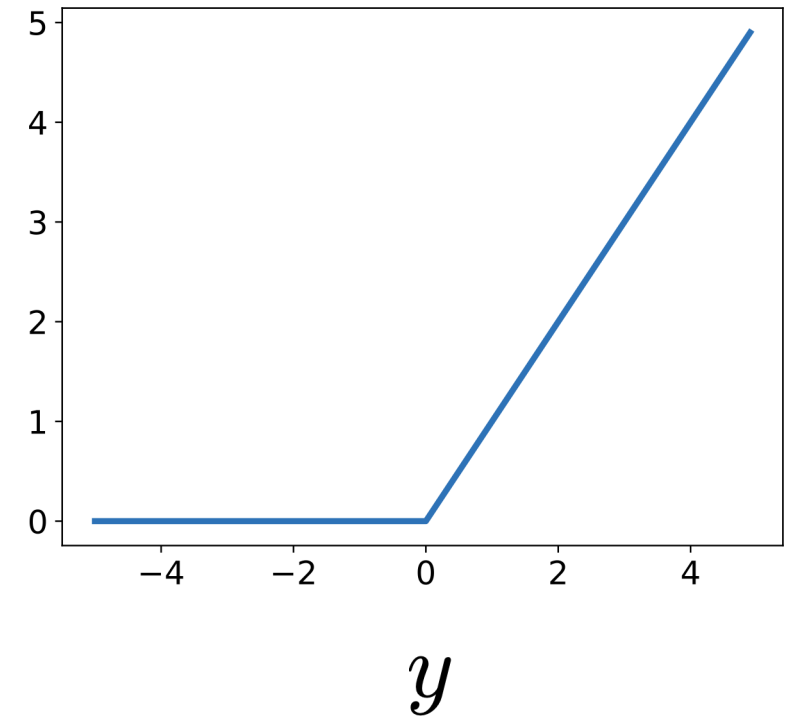**Sigmoid**

$$g(y) = \sigma(y) = \frac{1}{1 + e^{-y}}$$

$g(y)$

$y$

# Computation in a neural net — nonlinearity

- Unbounded output (on positive side)

- Efficient to implement: $\dfrac{\partial g}{\partial y} = \begin{cases} 0, & \text{if} \quad y < 0 \\ 1, & \text{if} \quad y \geq 0 \end{cases}$

- Also seems to help convergence (6x speedup vs. tanh in [Krizhevsky et al. 2012])

- Drawback: if strongly in negative region, unit is dead forever (no gradient).

- Default choice: widely used in current models!

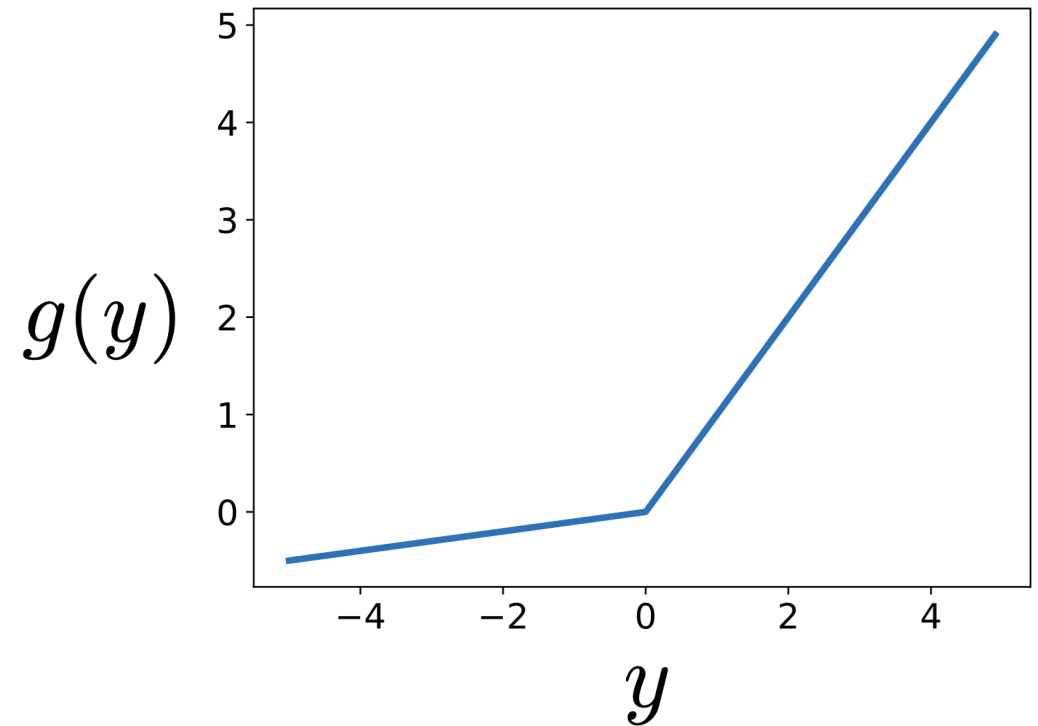**Rectified linear unit (ReLU)**

$$g(y) = \max(0, y)$$



$g(y)$

$y$

# Computation in a neural net — nonlinearity

- where a is small (e.g., 0.02)

- Efficient to implement:

- Has non-zero gradients everywhere (unlike ReLU)

$$\frac{\partial g}{\partial y} = \begin{cases} -a, & \text{if} \quad y < 0 \\ 1, & \text{if} \quad y \geq 0 \end{cases}$$

**Leaky ReLU**

$$g(y) = \begin{cases} \max(0, y), & \text{if} \quad y \geq 0 \\ a\min(0, y), & \text{if} \quad y < 0 \end{cases}$$

# Perceptron: Old Idea!

Late 1950s video on Rosenblatt's perceptron research

"While promising, this approach to machine intelligence virtually died out …"
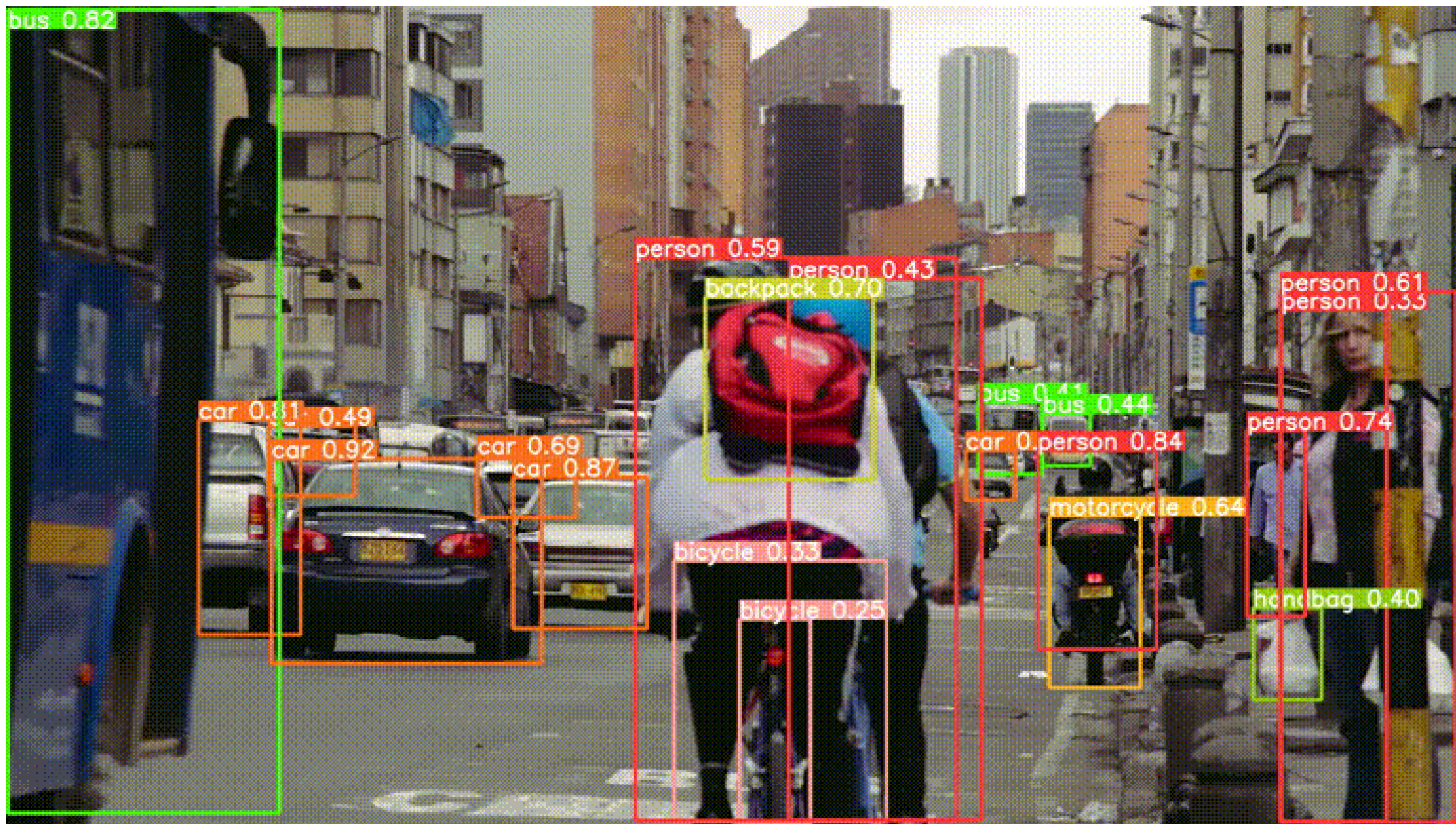
# Perceptron: Old Idea!

Late 1950s video on Rosenblatt's perceptron research

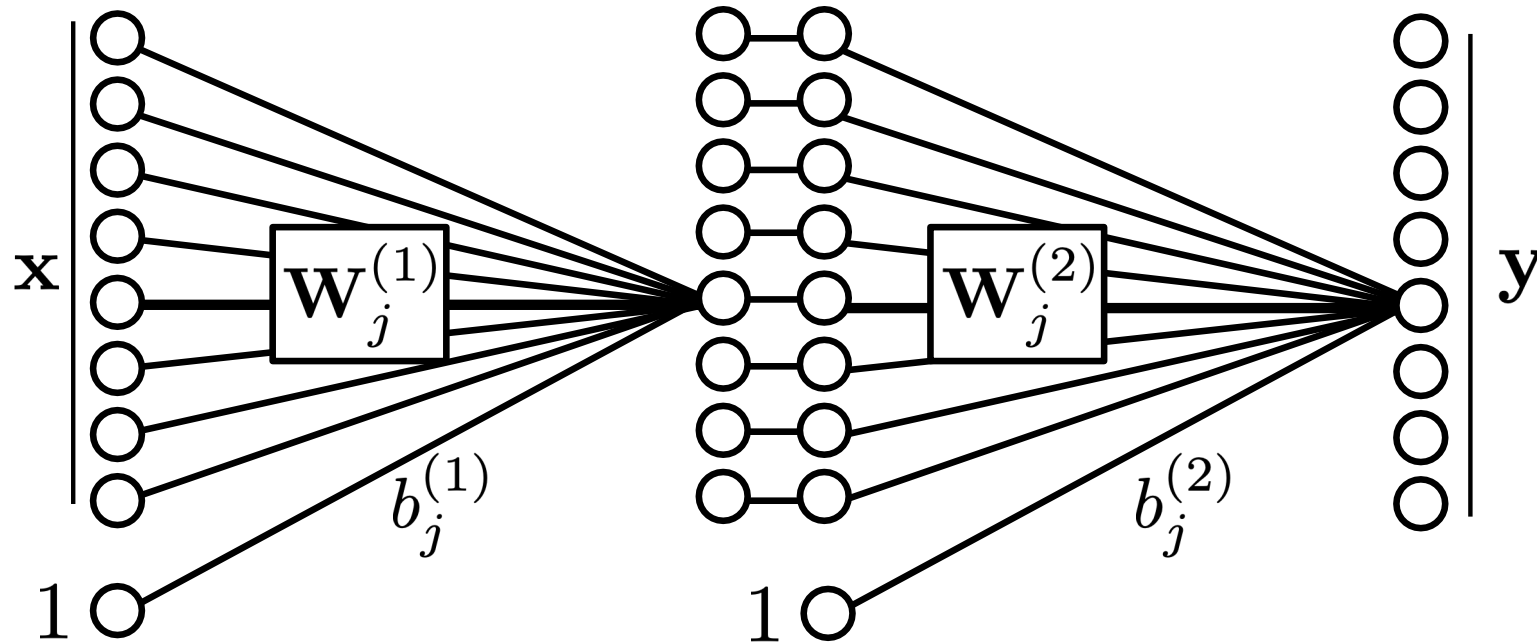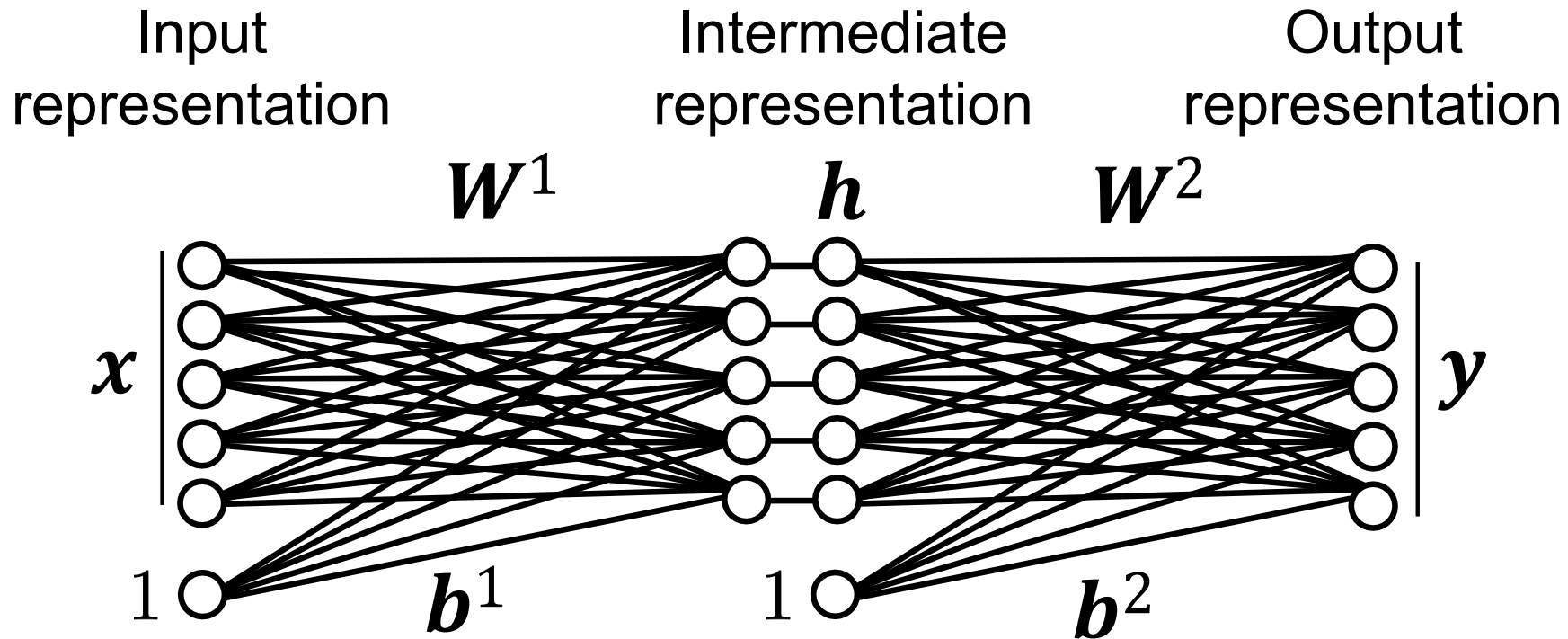*"While promising, this approach to machine intelligence virtually died out …"*



I guess you guys aren't ready for that yet. But your kids are gonna love it.

# Stacking layers

Input representation

Intermediate representation

Output representation



$\mathbf{h}$ = "hidden units"

# Stacking layers
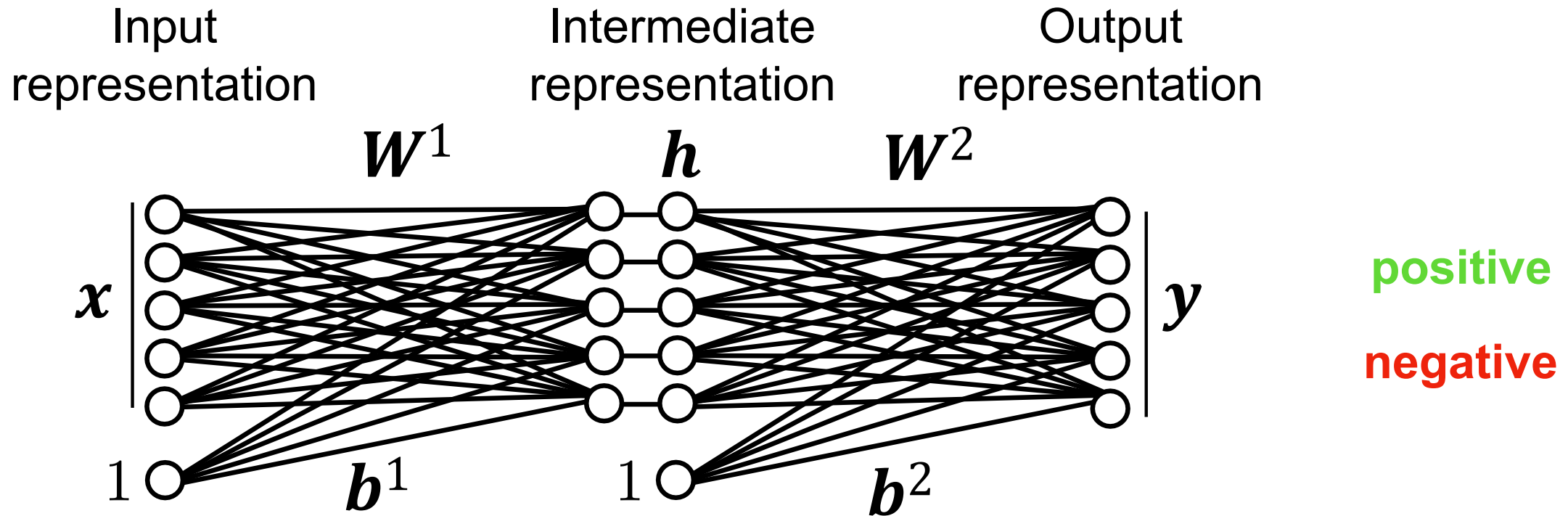
Input
representation

Intermediate
representation

Output
representation

$$\boldsymbol{W}^1 \qquad \boldsymbol{h} \qquad \boldsymbol{W}^2$$



$$x \qquad\qquad 1 \qquad\qquad y$$

$$1 \qquad \boldsymbol{b}^1 \qquad 1 \qquad \boldsymbol{b}^2$$

$$\boldsymbol{h} = g(\boldsymbol{W}^1 \boldsymbol{x} + \boldsymbol{b}^1) \qquad \boldsymbol{y} = g(\boldsymbol{W}^2 \boldsymbol{h} + \boldsymbol{b}^2)$$

**ReLU**

$$\theta = \{\boldsymbol{W}^1, ..., \boldsymbol{W}^L, \boldsymbol{b}^1, ..., \boldsymbol{b}^L\}$$

# Stacking layers



Input representation     Intermediate representation     Output representation

$W^1$    $h$    $W^2$

$x$    $y$

positive

negative

$1$   $b^1$    $1$   $b^2$

$$h = g(W^1 x + b^1) \qquad y = g(W^2 h + b^2)$$

**ReLU**

$$\theta = \{W^1, ..., W^L, b^1, ..., b^L\}$$

# Stacking layers



Input representation

Intermediate representation

Output representation

$W^1$    $h$    $W^2$

$x$

$1$   $b^1$    $1$    $b^2$

$y$

**positive**

**negative**

$$h = g(W^1 x + b^1) \quad y = g(W^2 h + b^2)$$

**ReLU**

$$\theta = \{W^1, ..., W^L, b^1, ..., b^L\}$$

# Stacking layers



Input representation

Intermediate representation

Output representation

$W^1$

$h$

$W^2$

$x$

$y$

positive

negative

$1$

$b^1$

$1$

$b^2$

$$h = g(W^1 x + b^1) \quad y = g(W^2 h + b^2)$$

ReLU

$$\theta = \{W^1, ..., W^L, b^1, ..., b^L\}$$

# Stacking layers

Input representation

Intermediate representation

Output representation

$$\boldsymbol{W}^1 \qquad \boldsymbol{h} \qquad \boldsymbol{W}^2$$

$\boldsymbol{x}$

positive

negative

$\boldsymbol{y}$

$1 \qquad \boldsymbol{b}^1 \qquad 1 \qquad \boldsymbol{b}^2$

$$\boldsymbol{h} = g(\boldsymbol{W}^1 \boldsymbol{x} + \boldsymbol{b}^1) \qquad \boldsymbol{y} = g(\boldsymbol{W}^2 \boldsymbol{h} + \boldsymbol{b}^2)$$

**ReLU**

$$\theta = \{\boldsymbol{W}^1, ..., \boldsymbol{W}^L, \boldsymbol{b}^1, ..., \boldsymbol{b}^L\}$$

# Stacking layers



Input representation

$$\boldsymbol{W}^1 \qquad \boldsymbol{h} \qquad \boldsymbol{W}^2$$

Intermediate representation

Output representation

$\boldsymbol{x}$

$1 \qquad \boldsymbol{b}^1 \qquad 1 \qquad \boldsymbol{b}^2$

$\boldsymbol{y}$

**positive**

**negative**

$$\boldsymbol{h} = g(\boldsymbol{W}^1 \boldsymbol{x} + \boldsymbol{b}^1) \qquad \boldsymbol{y} = g(\boldsymbol{W}^2 \boldsymbol{h} + \boldsymbol{b}^2)$$

**ReLU**

$$\theta = \{\boldsymbol{W}^1, ..., \boldsymbol{W}^L, \boldsymbol{b}^1, ..., \boldsymbol{b}^L\}$$

WE NEED TO GO

DEEPER

DEEP Neural Nets?

# Stacking layers - What's actually happening?



Low level features: e.g., edge, texture, …

higher level features: e.g., shape

even higher level features: e.g., "paw", "fur"
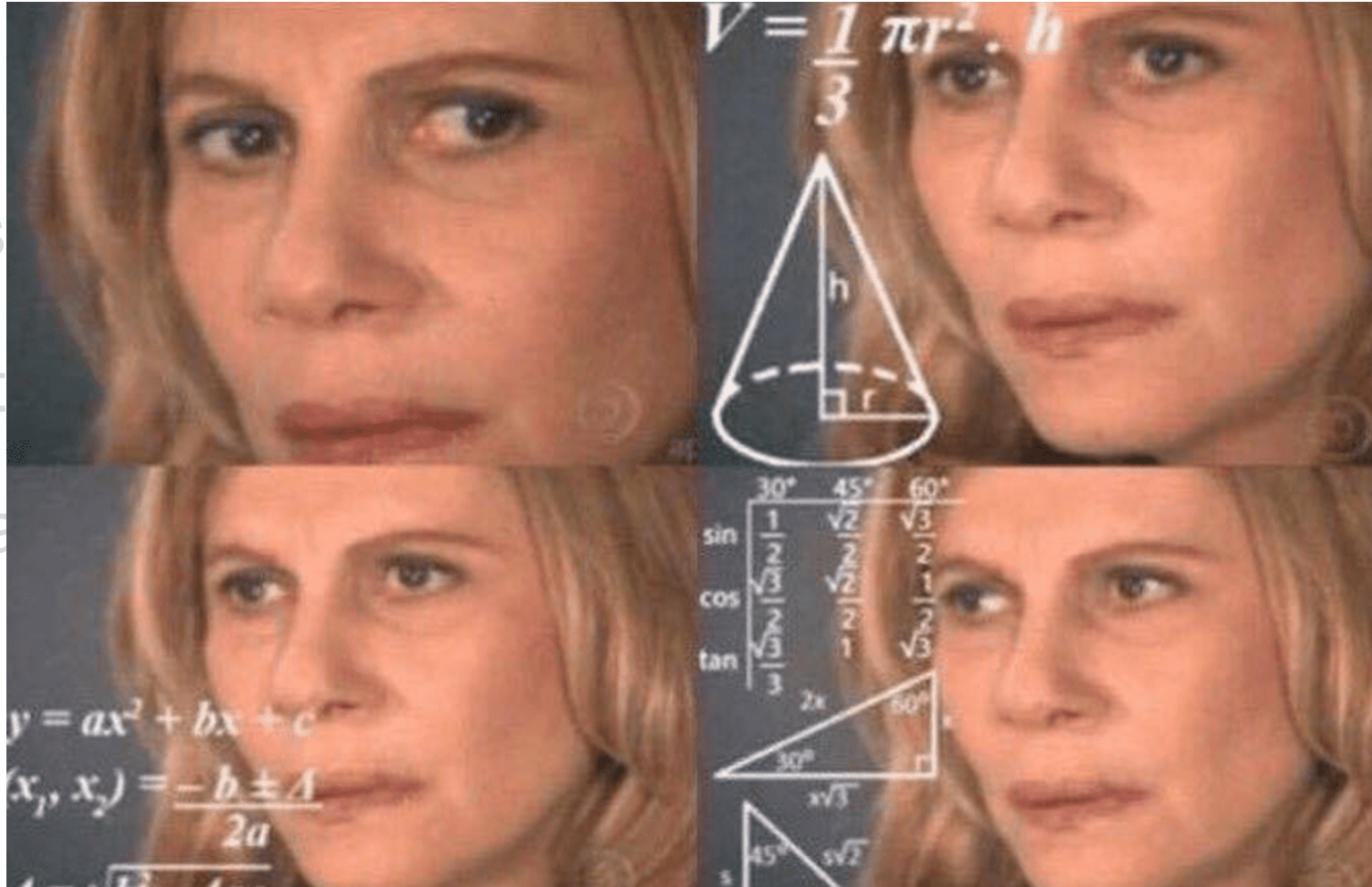
# Deep nets



$$f(x) = f_L( \ldots f_3(f_2(f_1(x))))$$

# Computation has a simple form

- Composition of linear functions with nonlinearities in between

- E.g. matrix multiplications with ReLU, $max(0, \mathbf{x})$ afterwards

- Do a matrix multiplication, set all negative values to 0, repeat

But where do we get the weights from?

# Computation has a simple form



- Compos ... in between
- E.g. mat... afterwards
- Do a ma... o 0, repeat

But where do we get the weights from?

# Where do we get the weights from ?

# How would we learn the parameters?

$\mathbf{y}_1$
"dog"

$\mathbf{x}_1$



$\mathcal{L}(f_\theta(\mathbf{x}_1), \mathbf{y}_1)$

predicted          ground truth

Learned

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

$$\theta^* = \arg\min_\theta \sum_{i=1}^{N} \mathcal{L}(f_\theta(\mathbf{x}_i), \mathbf{y}_i)$$

# Training neural networks

Let's start easy

# world's smallest neural network!
# (a "perceptron")



$$y = wx$$

(a.k.a. line equation, linear regression)

# Training a Neural Network

Given a set of samples and a Perceptron

$$\{x_i, y_i\}$$
$$y = f_{\mathrm{PER}}(x; w)$$

Estimate the parameter of the Perceptron

$$w$$

Given training data:

| $x$ | $y$ |
| --- | --- |
| 10 | 10.1 |
| 2 | 1.9 |
| 3.5 | 3.4 |
| 1 | 1.1 |

*What do you think the weight parameter is?*

$$y = wx$$

Given training data:

| $x$ | $y$ |
| --- | --- |
| 10 | 10.1 |
| 2 | 1.9 |
| 3.5 | 3.4 |
| 1 | 1.1 |

*What do you think the weight parameter is?*

$$y = wx$$

not so obvious as the network gets more complicated so we use …

# An Incremental Learning Strategy
## (gradient descent)

Given several examples

$$\{(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)\}$$
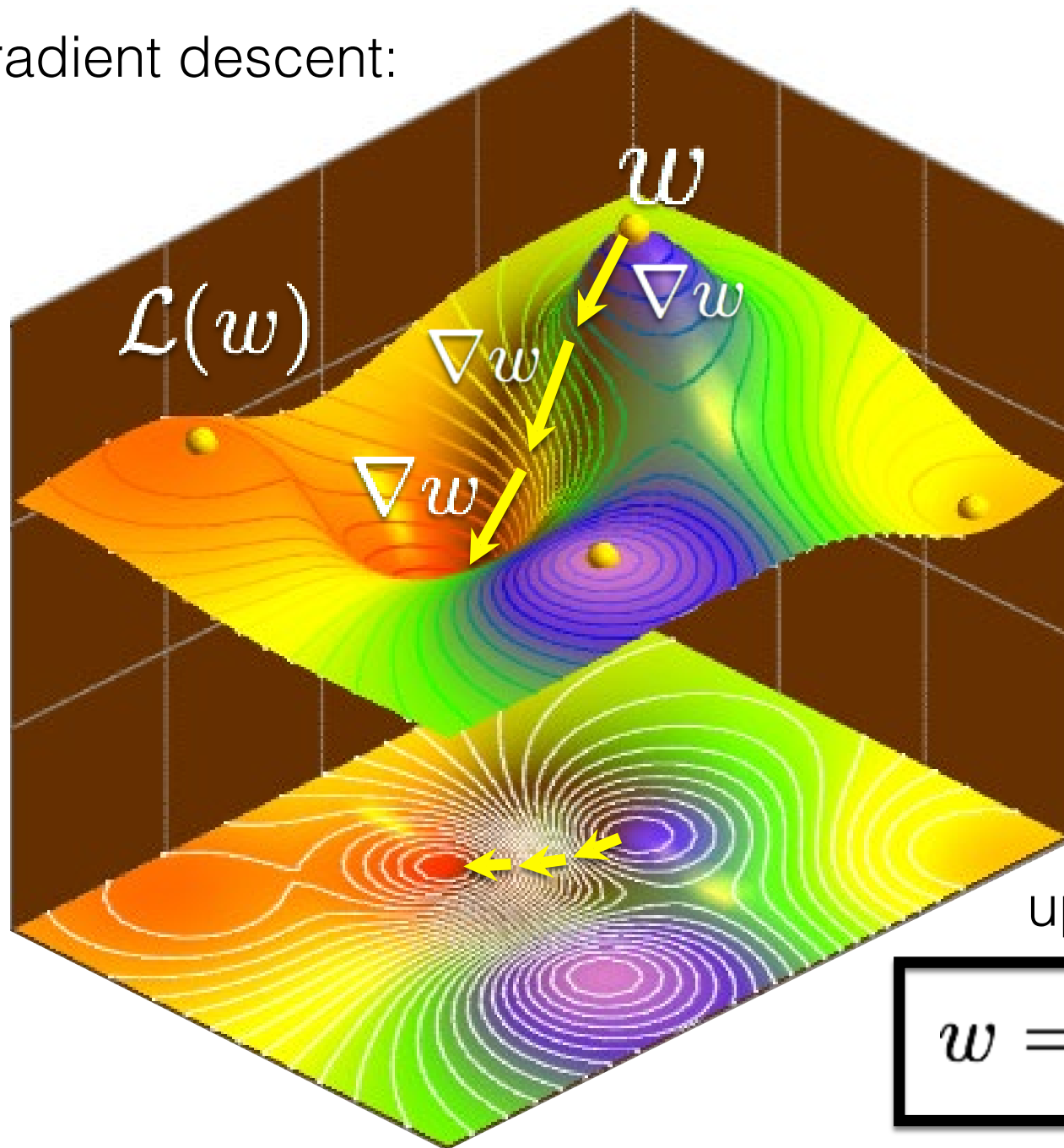
and a perceptron

$$\hat{y} = wx$$

# An Incremental Learning Strategy
### (gradient descent)

Given several examples

$$\{(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)\}$$

and a perceptron

$$\hat{y} = wx$$

Modify weight $w$ such that $\hat{y}$ gets **'closer'** to $y$

# An Incremental Learning Strategy
(gradient descent)

Given several examples

$$\{(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)\}$$

and a perceptron

$$\hat{y} = wx$$

Modify weight $w$ such that $\hat{y}$ gets **'closer'** to $y$

perceptron parameter

perceptron output

true label

Gradient descent:



update rule:

$$w = w - \nabla w$$

$$\frac{d\mathcal{L}}{dw}$$

…is the rate at which **this** will change…

$$\mathcal{L} = \frac{1}{2}(y - \hat{y})^2$$

the loss function

… per unit change of **this**

$$y = \textcircled{w}x$$

the weight parameter

Let's compute the derivative…

# L1 Loss

$$\ell(\hat{y}, y) = |\hat{y} - y|$$
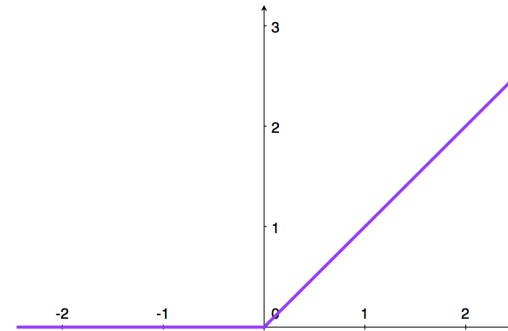
# L2 Loss

$$\ell(\hat{y}, y) = (\hat{y} - y)^2$$

# Zero-One Loss

$$\ell(\hat{y}, y) = \mathbf{1}[\hat{y} = y]$$

# Hinge Loss

$$\ell(\hat{y}, y) = \max(0, 1 - y \cdot \hat{y})$$

$$\frac{d\mathcal{L}}{dw}$$ …is the rate at which **this** will change…

$$\mathcal{L} = \frac{1}{2}(y - \hat{y})^2$$

the loss function

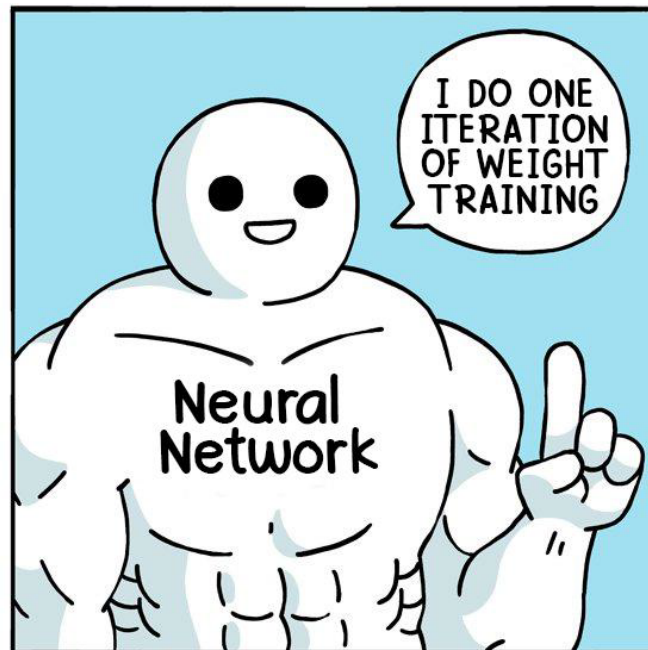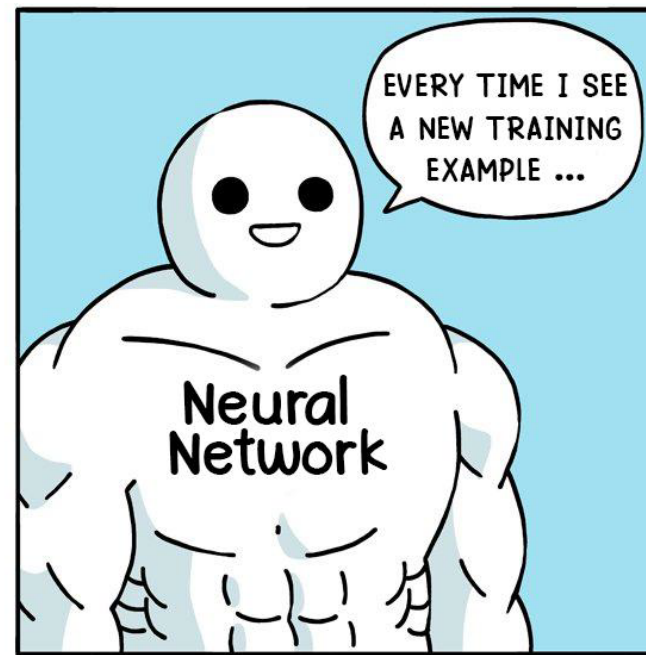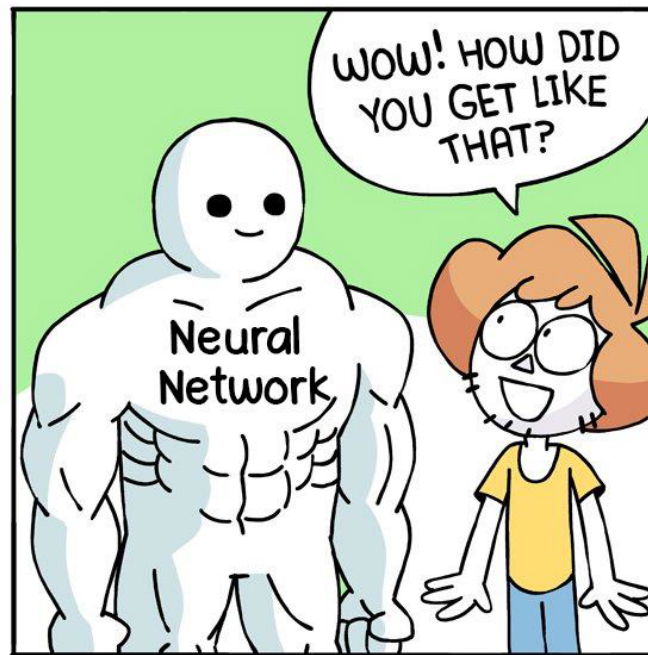… per unit change of **this**

$$y = wx$$

the weight parameter

Let's compute the derivative…

Compute the derivative

$$\frac{d\mathcal{L}}{dw} = \frac{d}{dw}\left\{\frac{1}{2}(y - \hat{y})^2\right\}$$

$$= -(y - \hat{y})\frac{dwx}{dw}$$

$$= -(y - \hat{y})x = \nabla w$$

That means the weight update for **gradient descent** is:

$$w = w - \nabla w \quad \text{move in direction of negative gradient}$$

$$= w + (y - \hat{y})x$$

SHEN COMIX

**Gradient Descent (world's smallest perceptron)**

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

1. Predict

   a. Forward pass $\hat{y} = wx_i$

   b. Compute Loss $\mathcal{L}_i = \dfrac{1}{2}(y_i - \hat{y})^2$
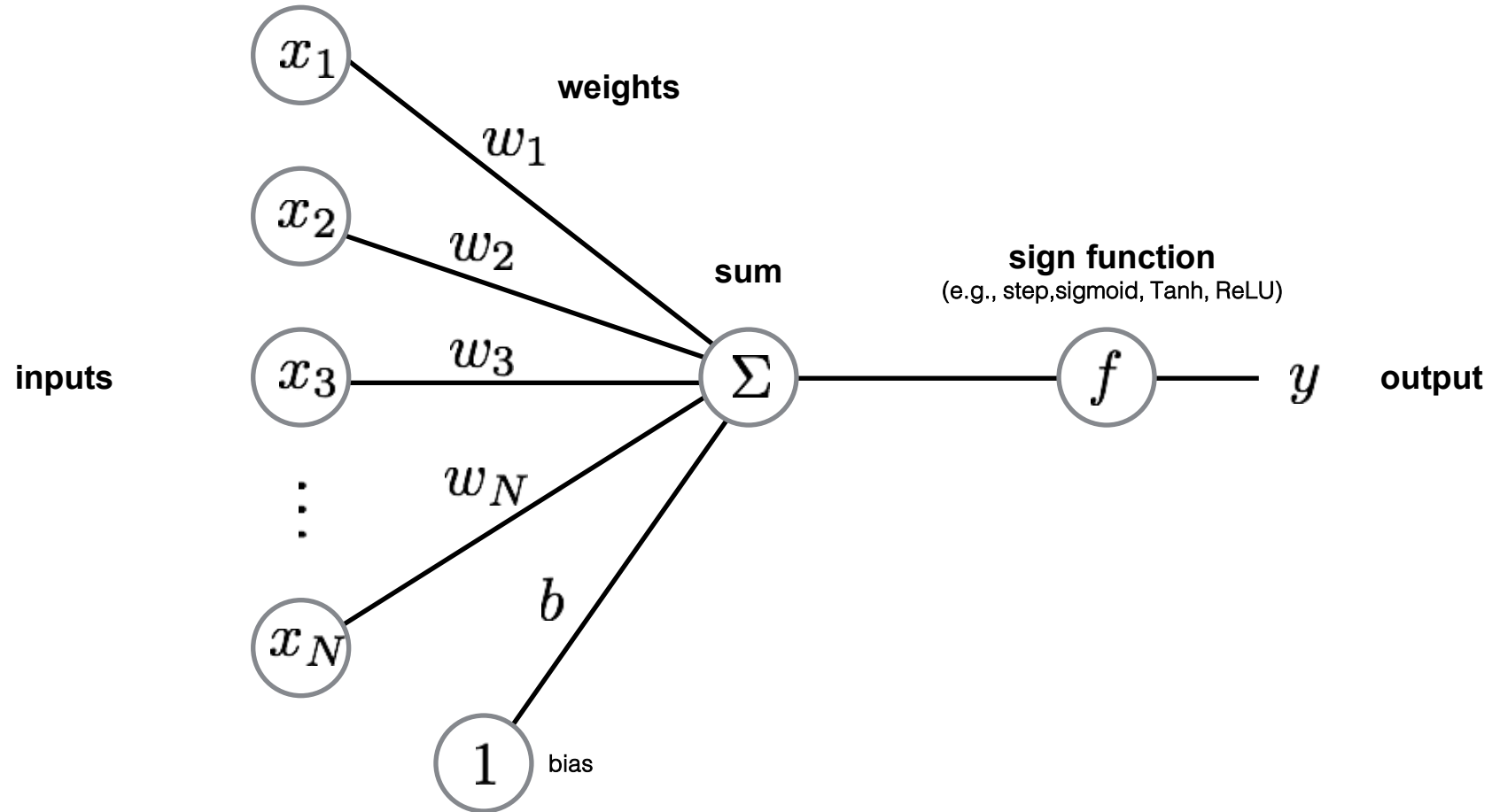
2. Update

   a. Back Propagation $\dfrac{d\mathcal{L}_i}{dw} = -(y_i - \hat{y})x_i = \nabla w$
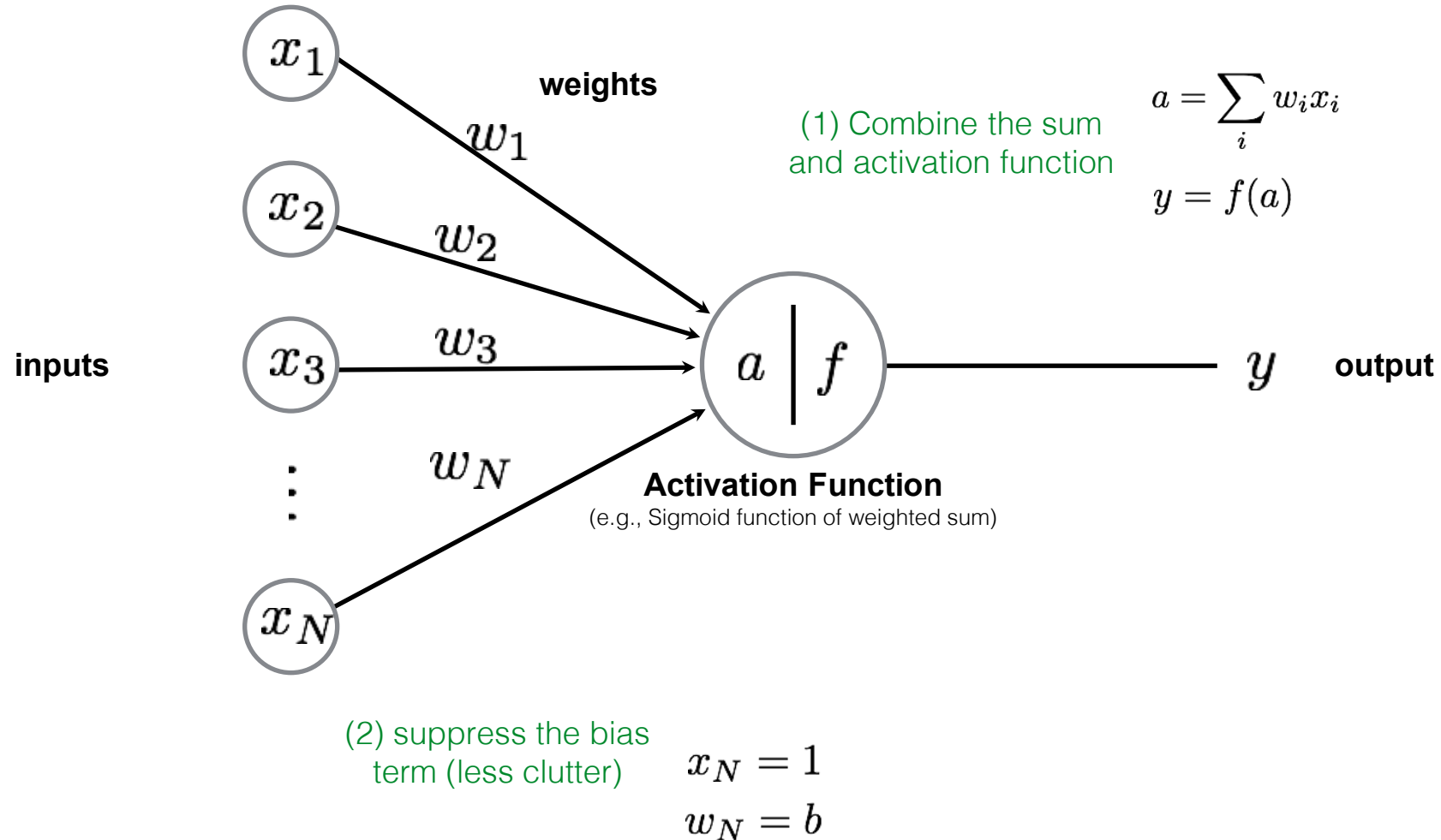
   b. Gradient update $w = w - \nabla w$

# The Perceptron
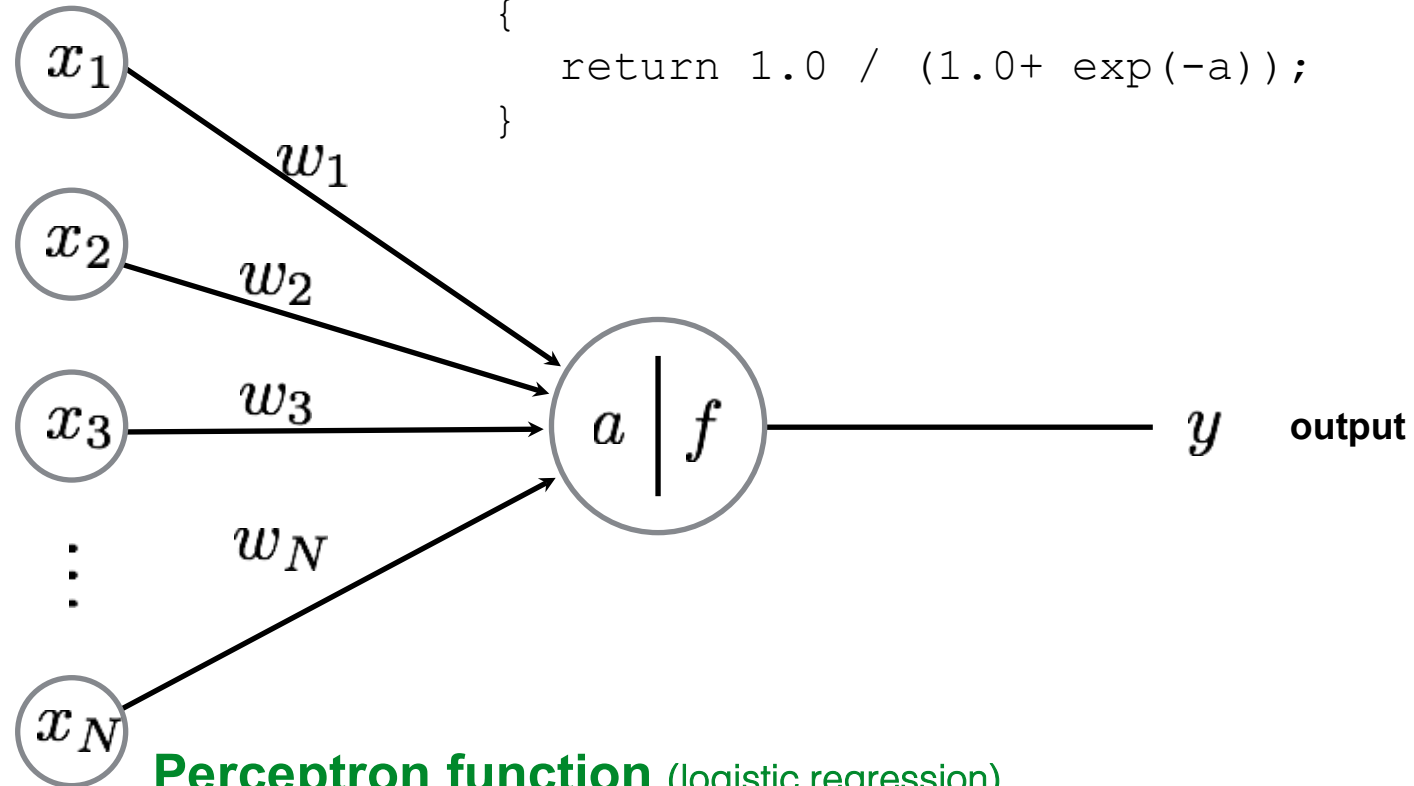


inputs

$x_1$

$x_2$

$x_3$

$x_N$

**weights**

$w_1$

$w_2$

$w_3$

$w_N$

$b$

$1$ bias

**sum**

$\Sigma$

**sign function**
(e.g., step,sigmoid, Tanh, ReLU)

$f$

$y$ **output**

# Another way to draw it...



inputs

$x_1$

$x_2$

$x_3$

$x_N$

weights

$w_1$

$w_2$

$w_3$

$w_N$

$a \mid f$

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

$y$    output

(1) Combine the sum
and activation function

$$a = \sum_i w_i x_i$$

$$y = f(a)$$

(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));
}
```



$x_1$

$x_2$

$x_3$

$x_N$

$w_1$

$w_2$

$w_3$

$w_N$

$a \mid f$

$y$ **output**

**Perceptron function** (logistic regression)

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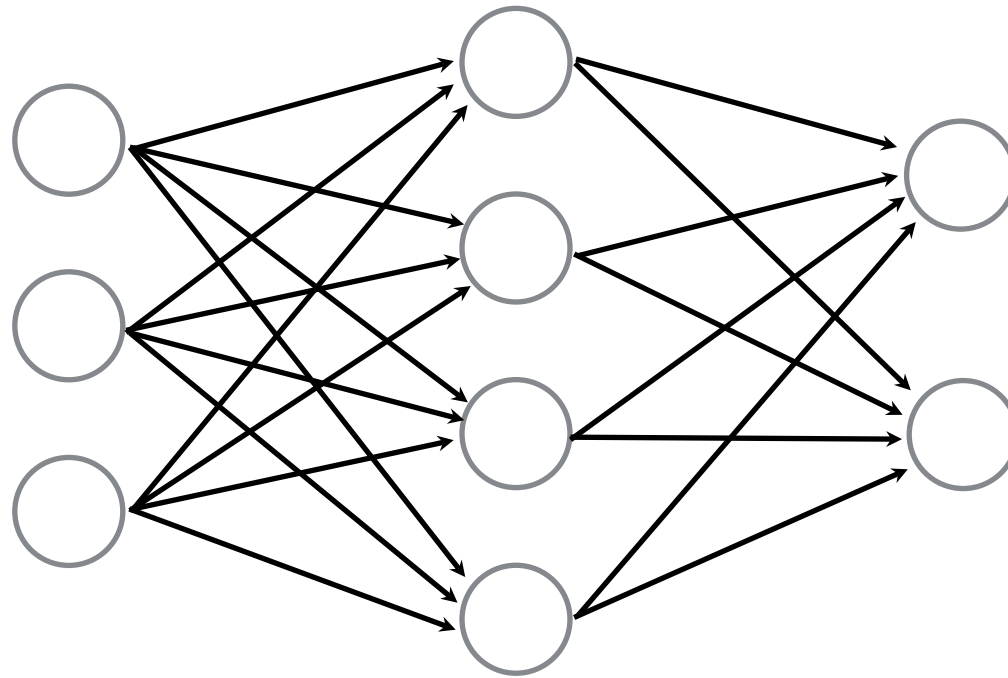
# Neural networks

Connect a bunch of perceptrons together …

# Neural Network

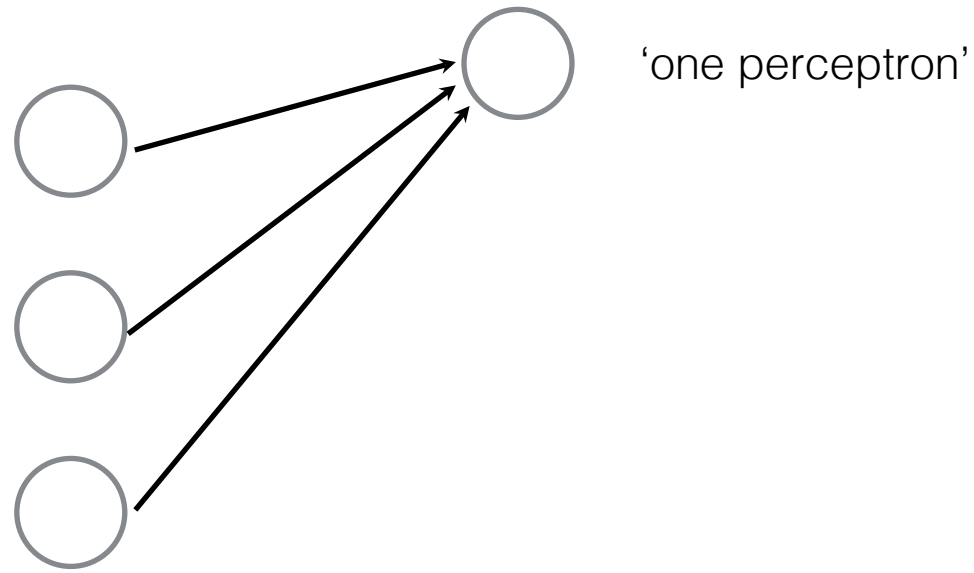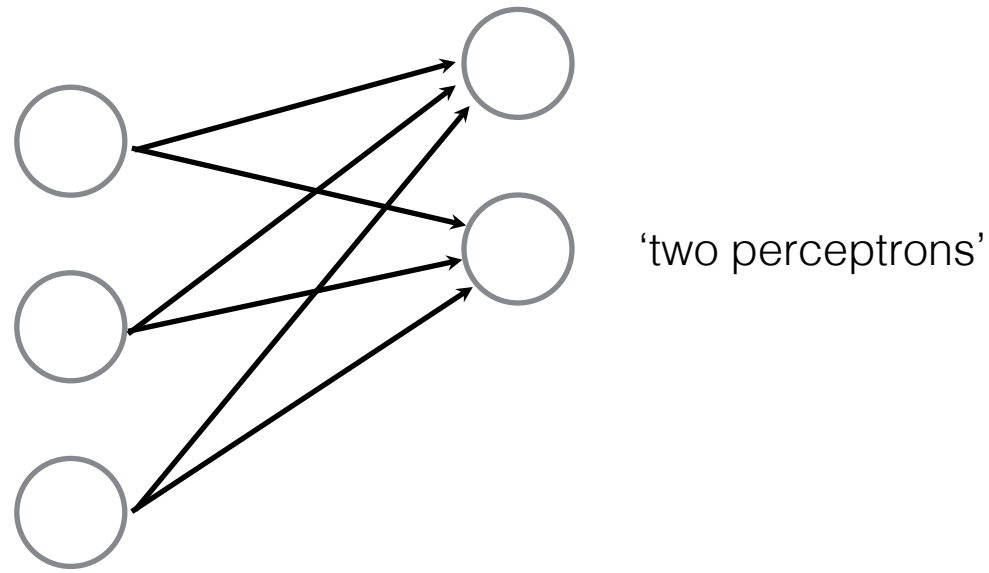Connect a bunch of perceptrons together …

a collection of connected perceptrons

# Connect a bunch of perceptrons together …

# Neural Network

a collection of connected perceptrons

Connect a bunch of perceptrons together …
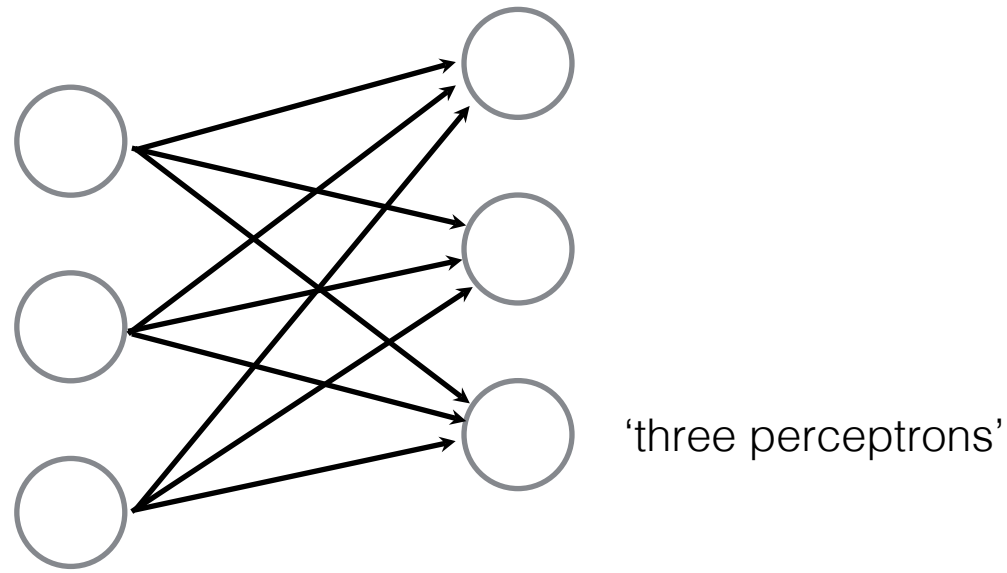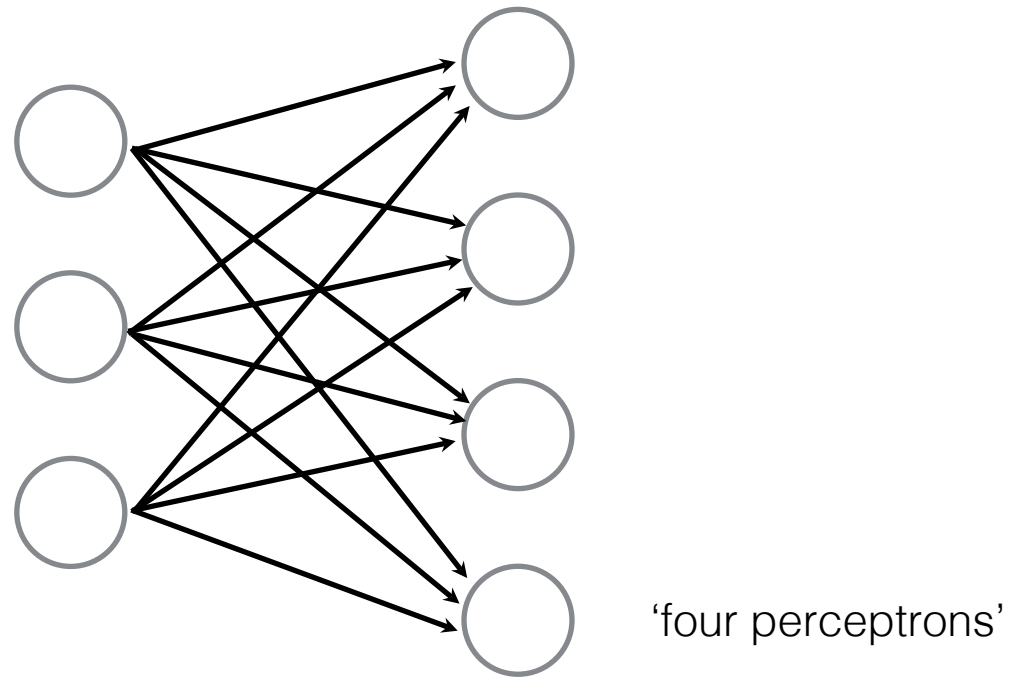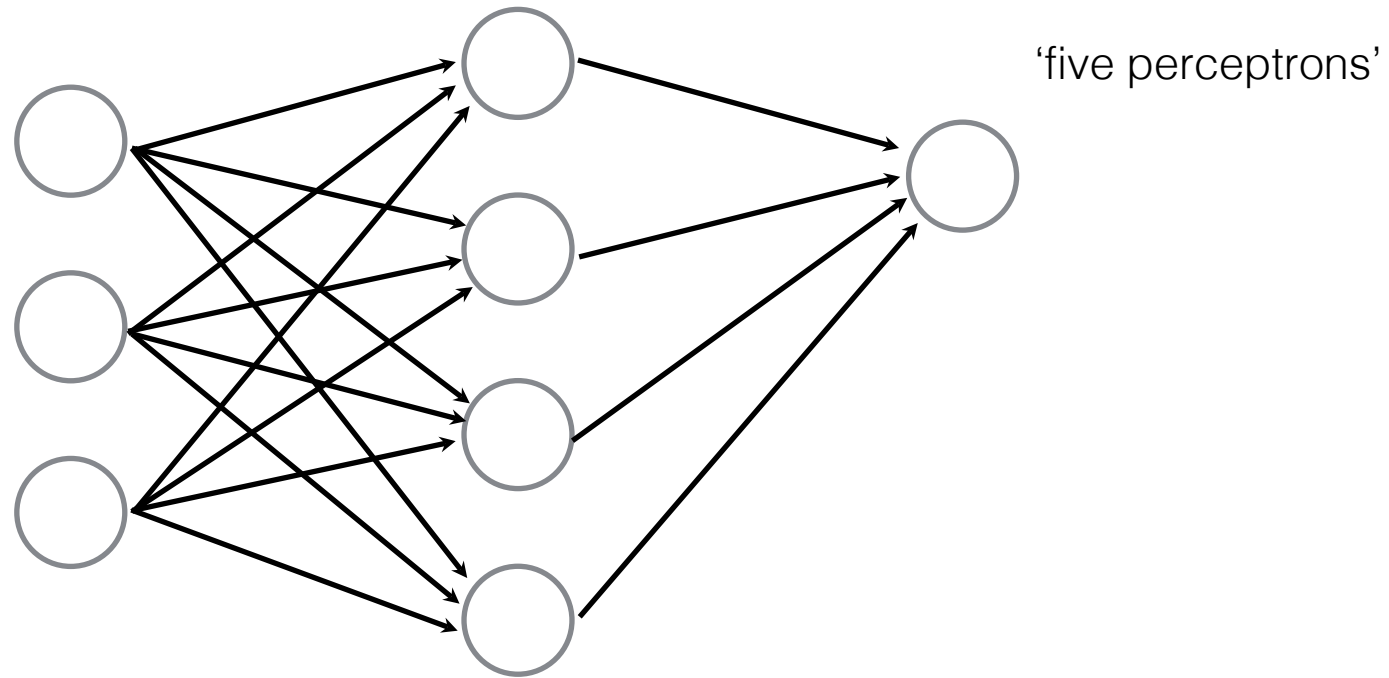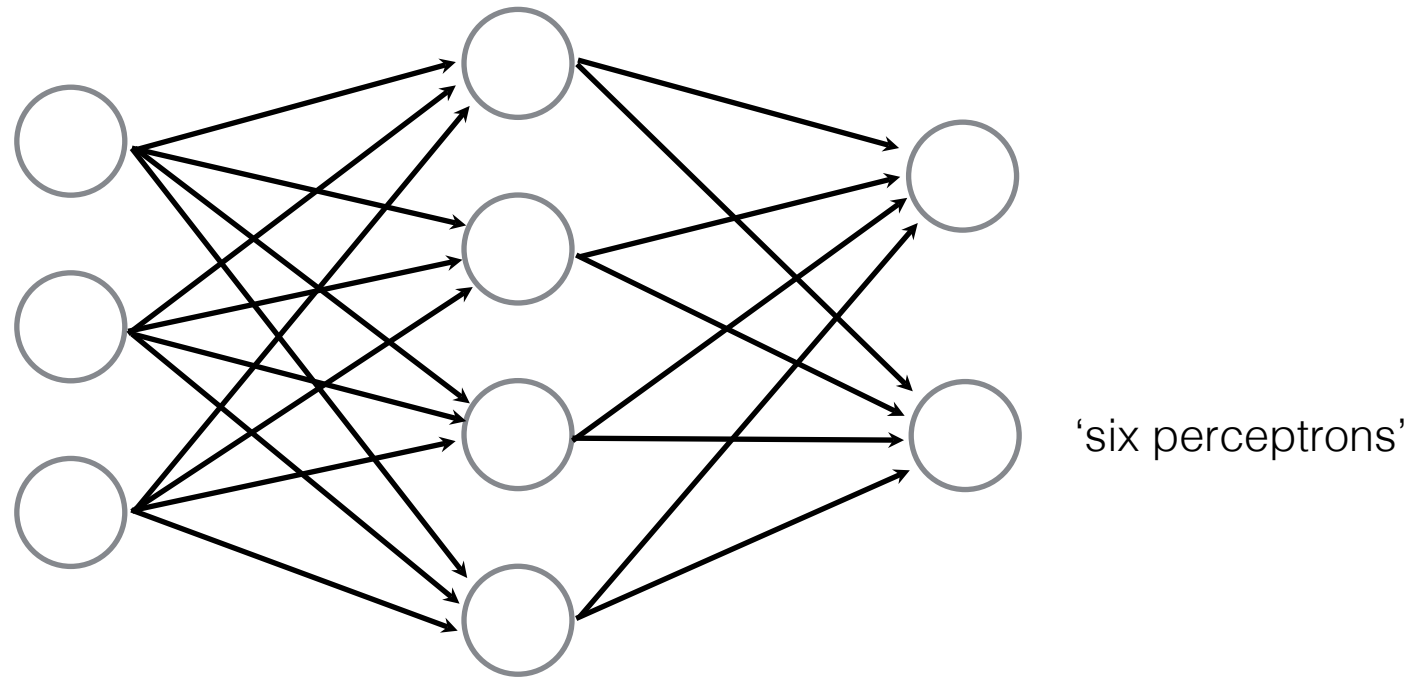
# Neural Network

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 …

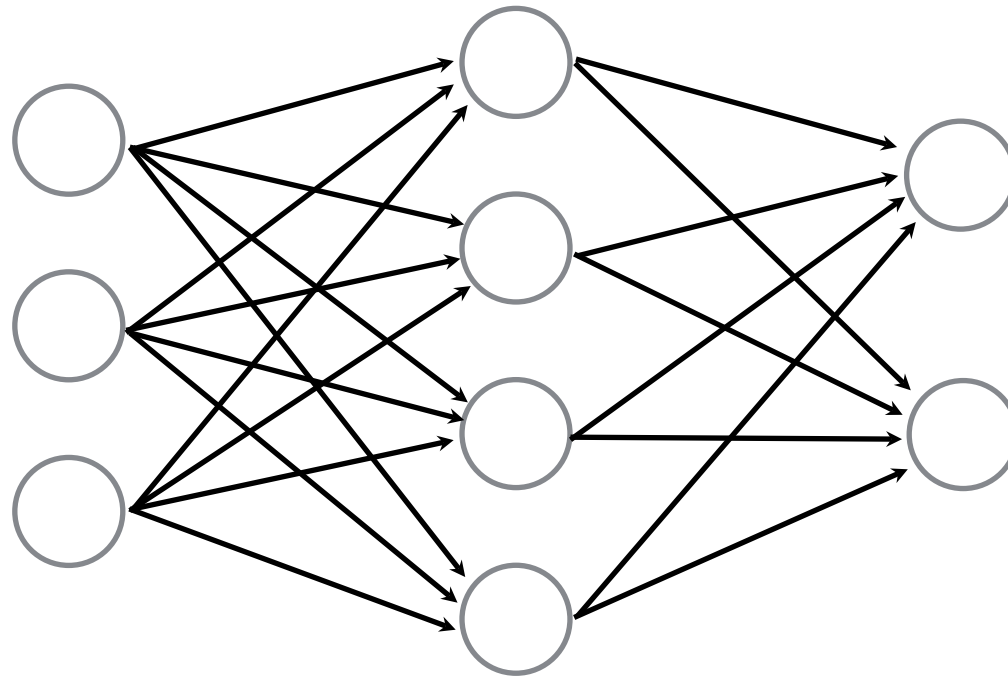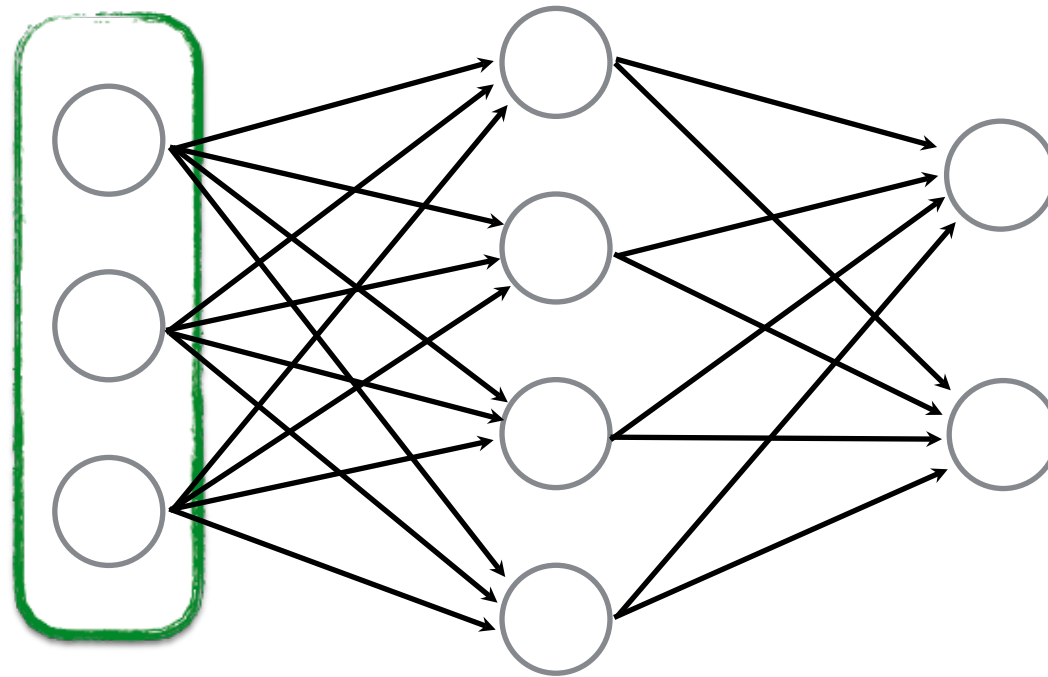# Neural Network

a collection of connected perceptrons



'six perceptrons'

Some terminology…



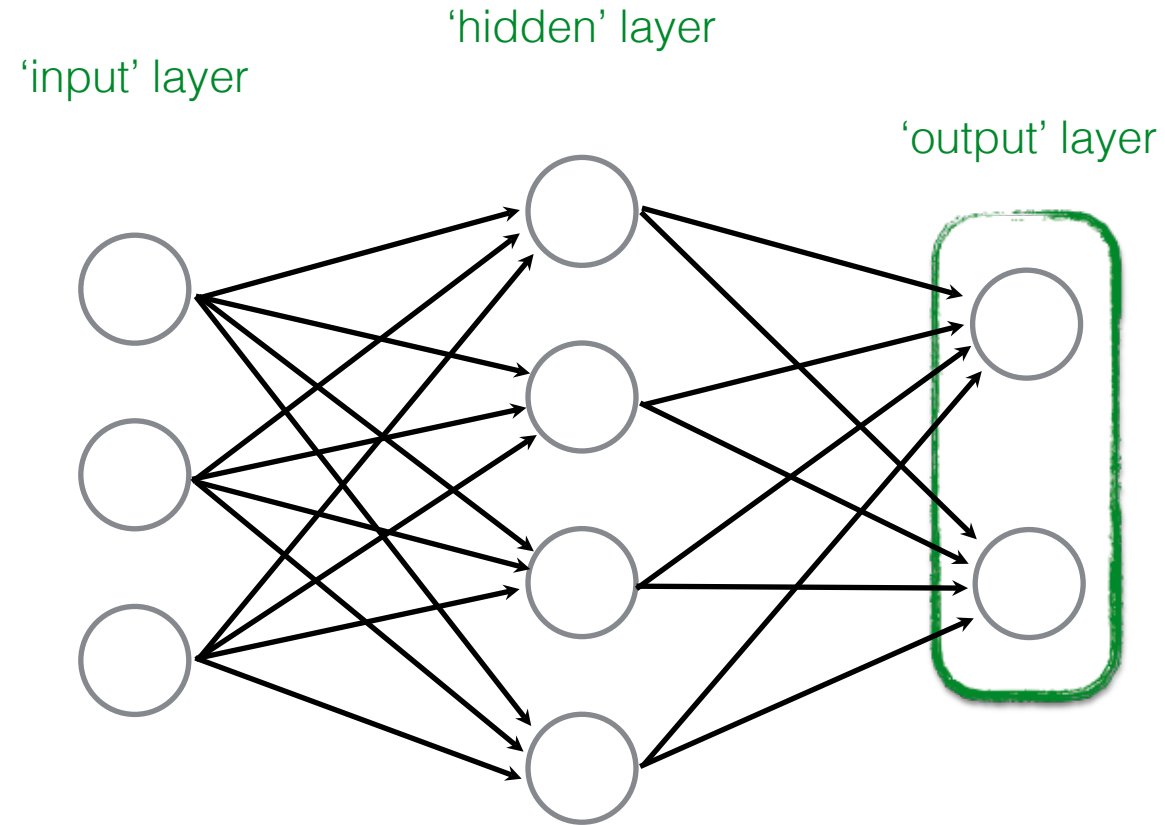…also called a **Multi-layer Perceptron** (MLP)

# Some terminology…

‘input’ layer



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

# Some terminology…



'hidden' layer

'input' layer

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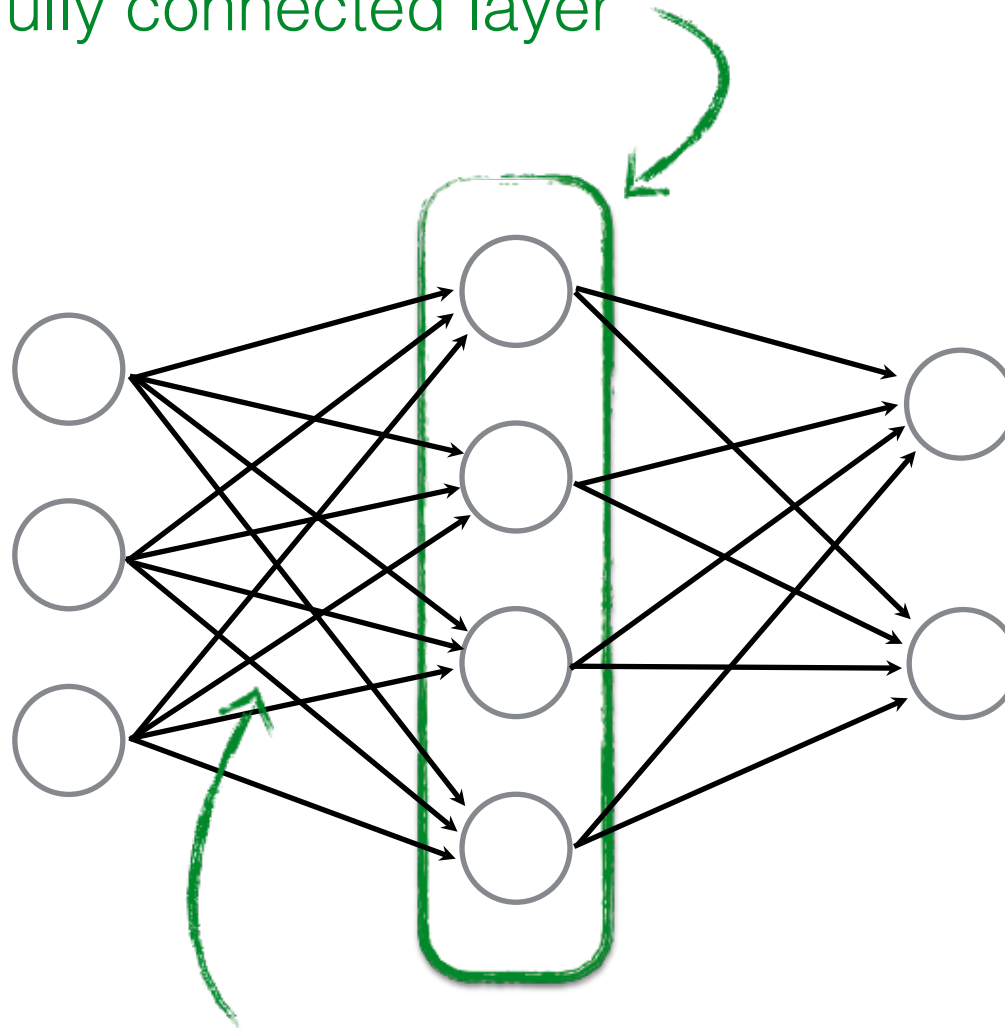
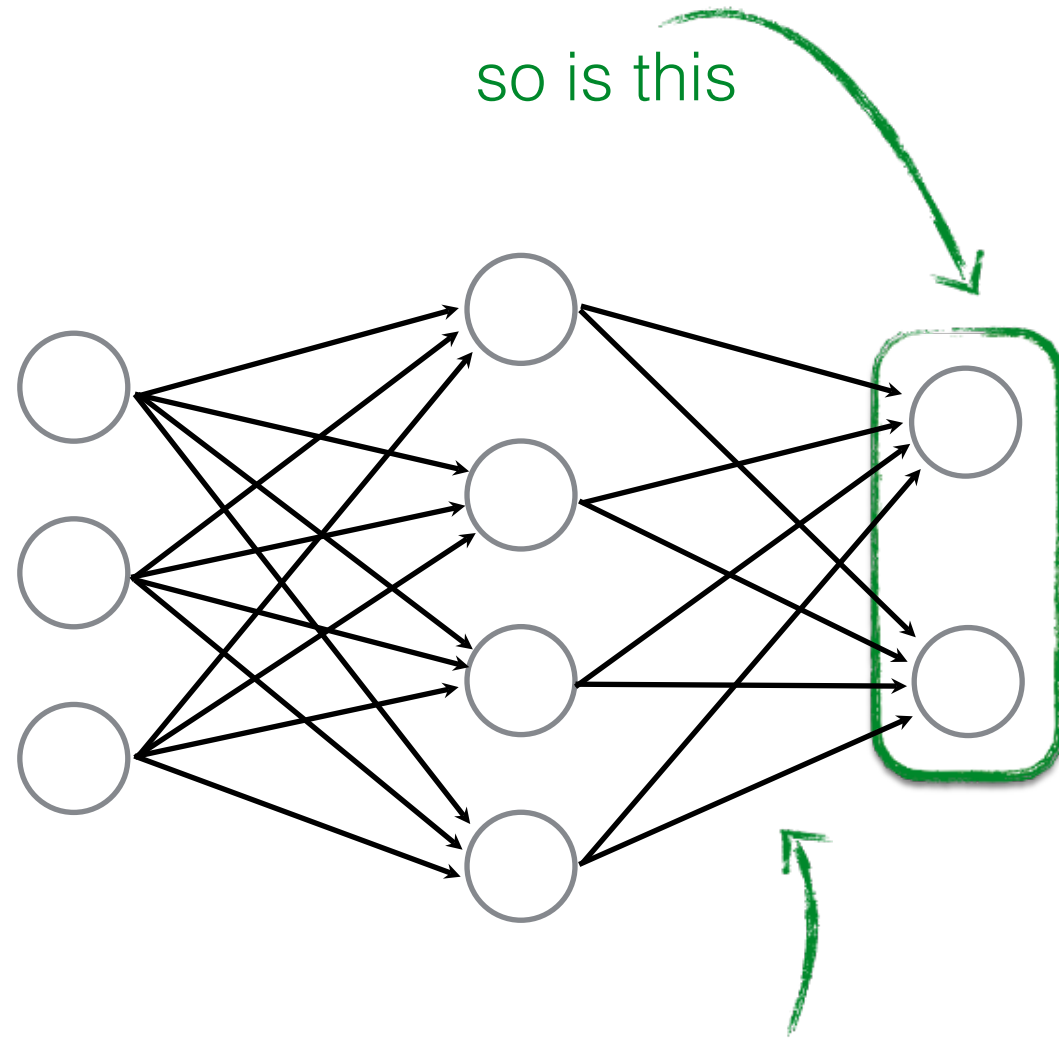# Some terminology…

'hidden' layer

'input' layer

'output' layer

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

this layer is a
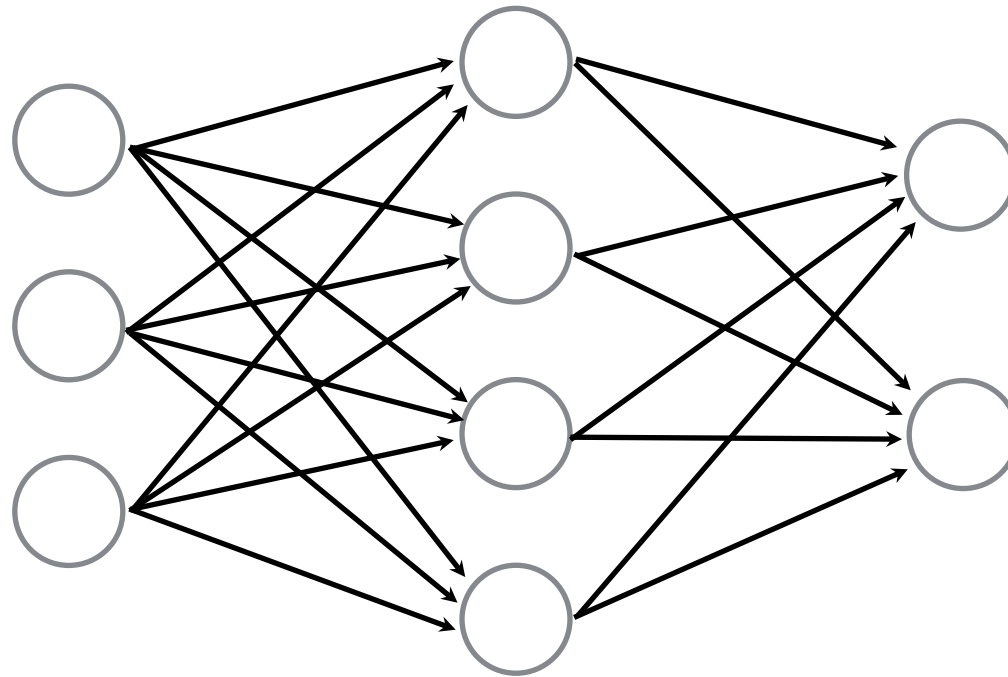'fully connected layer'

all pairwise neurons between layers are connected

so is this

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

*How many neurons (perceptrons)?*

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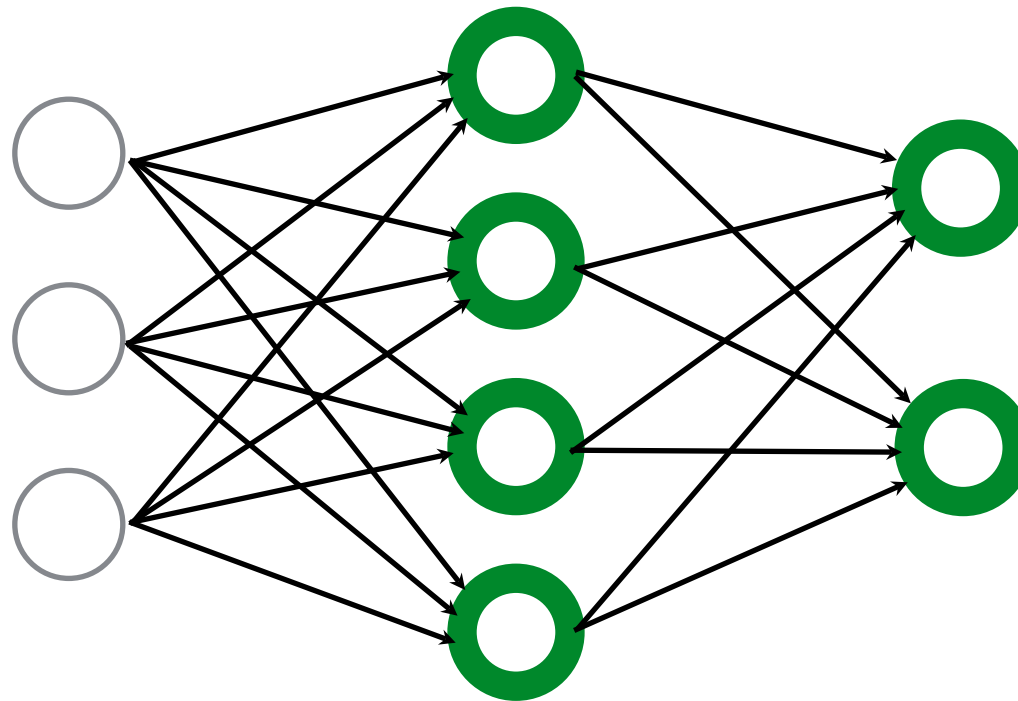
*How many learnable parameters total?*

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

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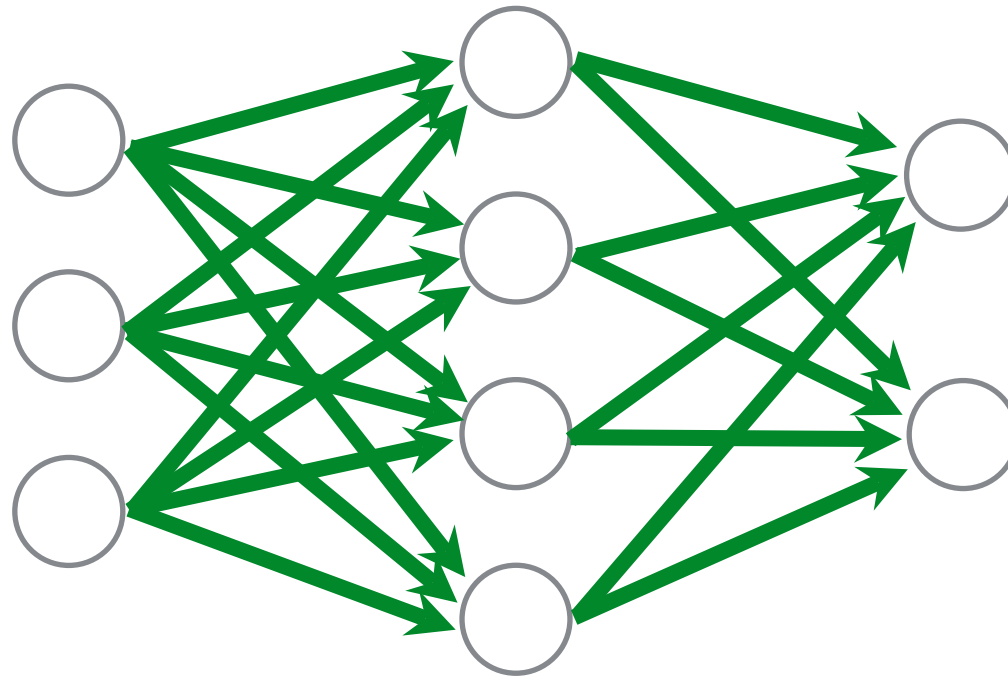
*How many learnable parameters total?*

How many neurons (perceptrons)?

$4 + 2 = 6$

How many weights (edges)?

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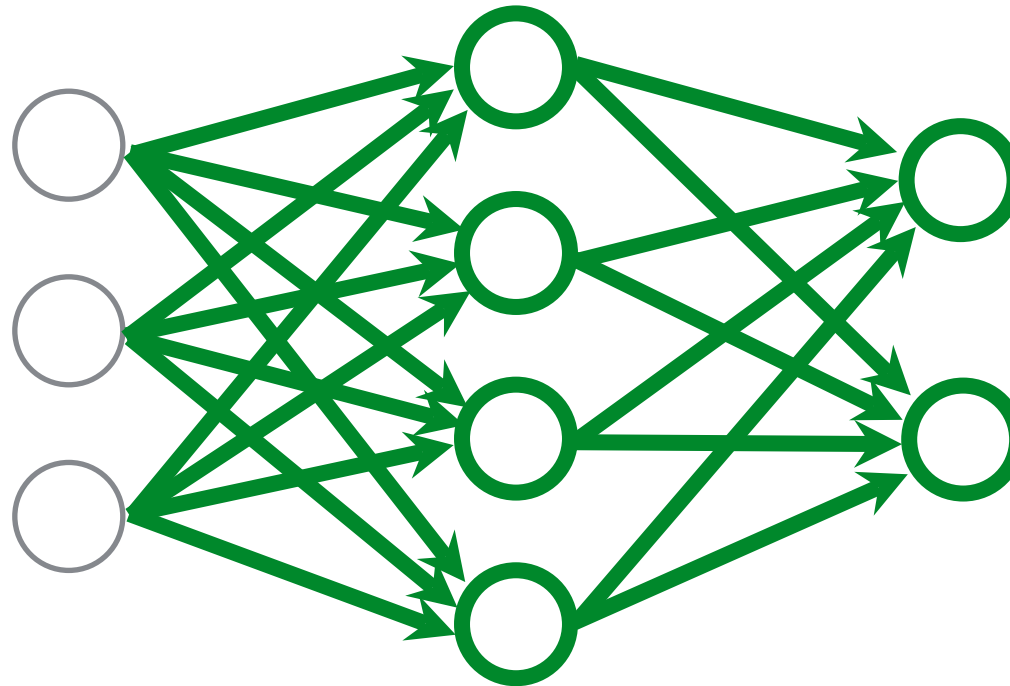


How many learnable parameters total?

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

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

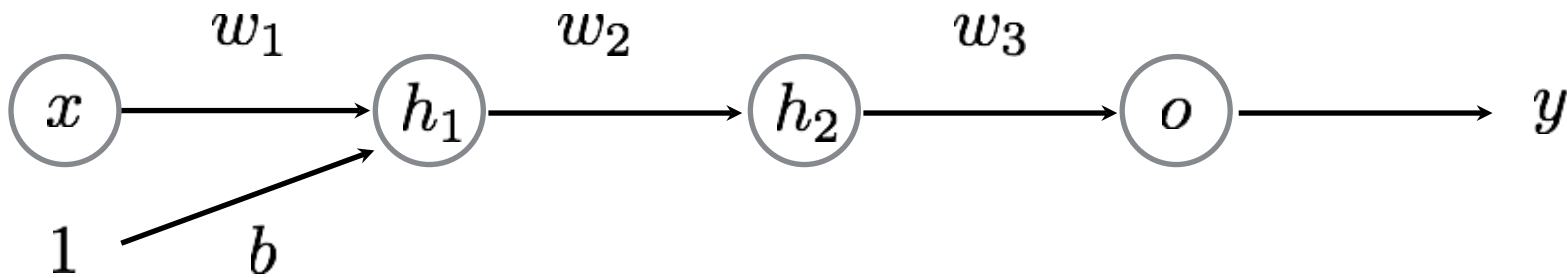*How many learnable parameters total?*          20 + 4 + 2 = 26
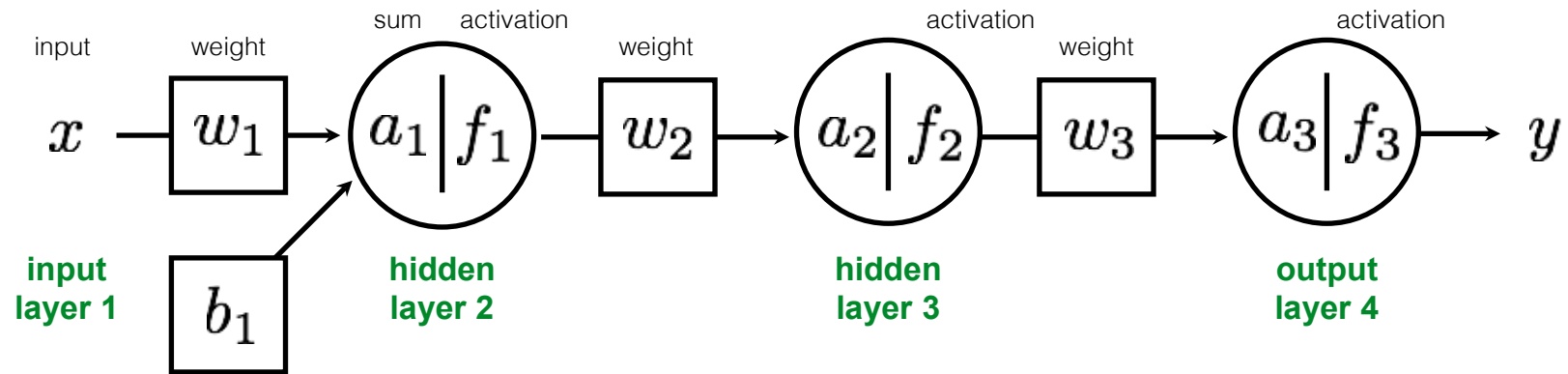
bias terms

# Example of a multi-layer perceptron



function of **FOUR** parameters and **FOUR** layers!

input
weight
sum    activation
weight
activation
weight
activation

$x$ — $w_1$ → $a_1 \vert f_1$ — $w_2$ → $a_2 \vert f_2$ — $w_3$ → $a_3 \vert f_3$ → $y$

$b_1$

**input
layer 1**

**hidden
layer 2**

**hidden
layer 3**

**output
layer 4**

Diagram showing a four-layer neural network. The green-boxed region labeled "input layer 1" and "hidden layer 2" contains: input $x$, weight $w_1$, bias $b_1$, and the node $a_1 \mid f_1$ (sum and activation). This connects through weight $w_2$ to node $a_2 \mid f_2$ (hidden layer 3, activation), then through weight $w_3$ to node $a_3 \mid f_3$ (output layer 4, activation), producing output $y$.

$$a_1 = w_1 \cdot x + b_1$$

$$a_1 = w_1 \cdot x + b_1$$

input layer 1    hidden layer 2    hidden layer 3    output layer 4

$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

$$a_1 = w_1 \cdot x + b_1$$
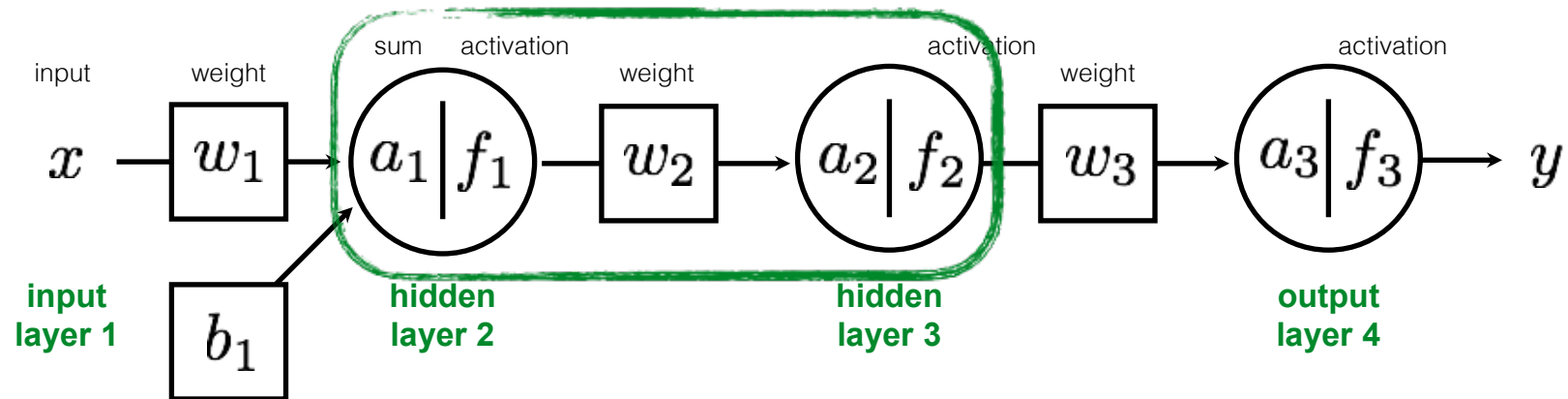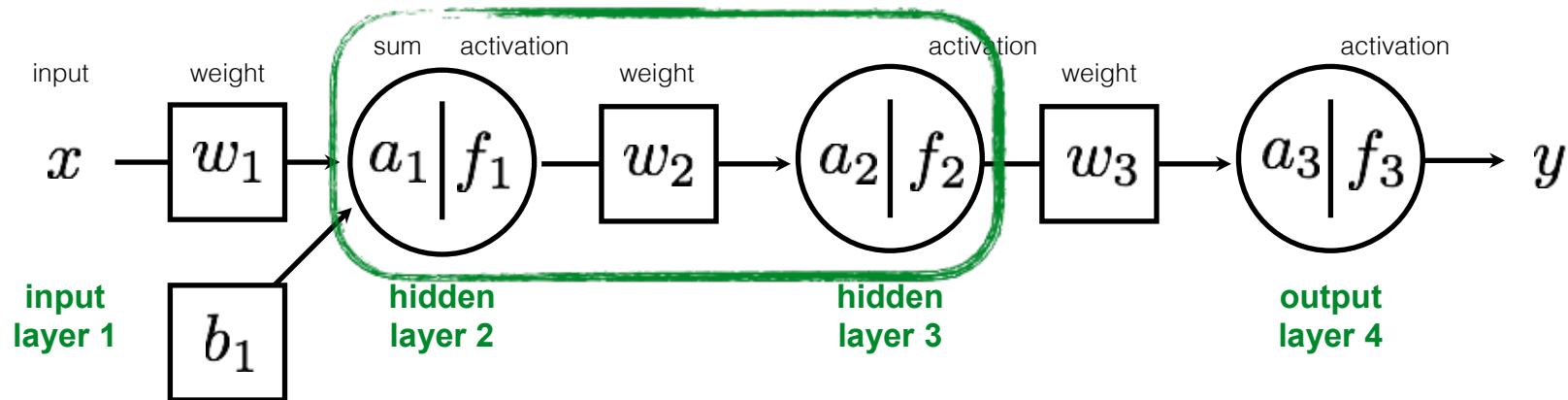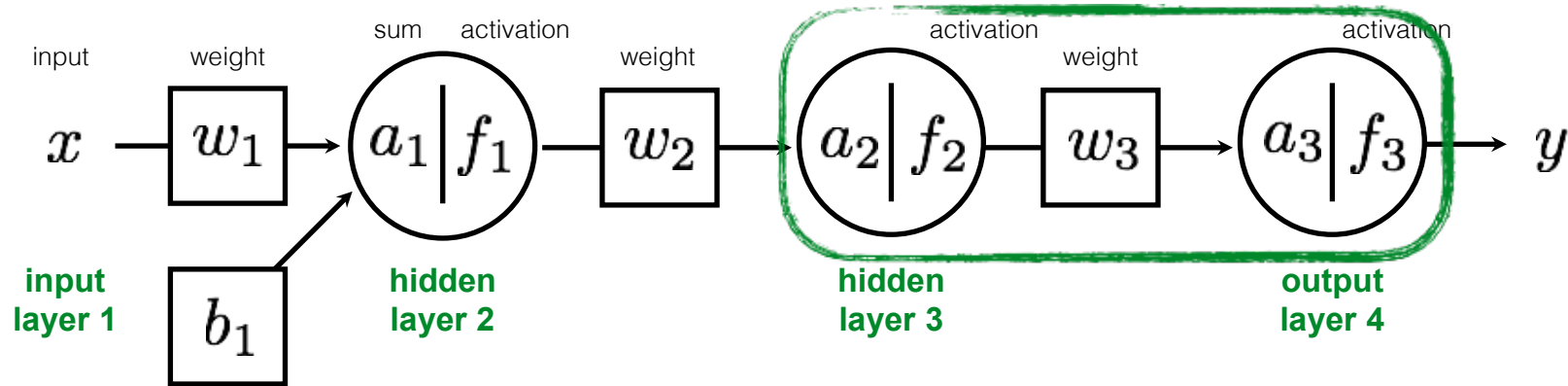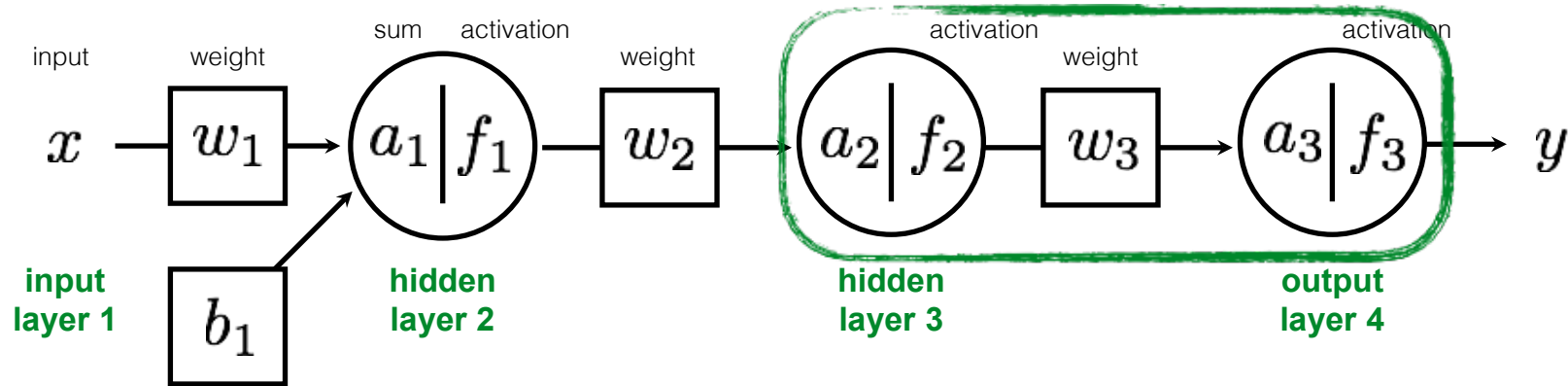
$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$
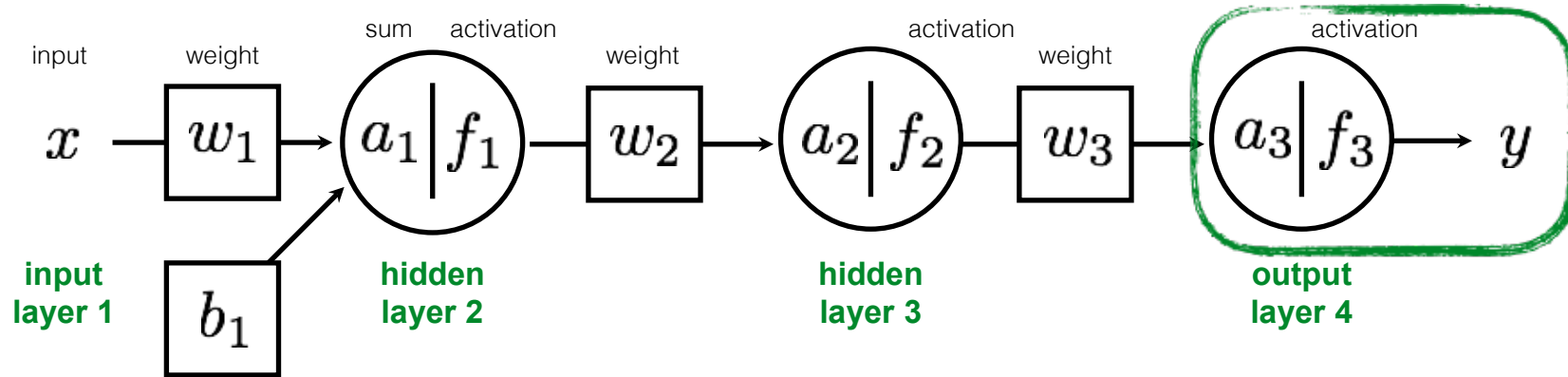
$$a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$$

$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

$$a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$$

$$a_1 = w_1 \cdot x + b_1$$

$$a_2 = w_2 \cdot f_1(w_1 \cdot x + b_1)$$

$$a_3 = w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1))$$

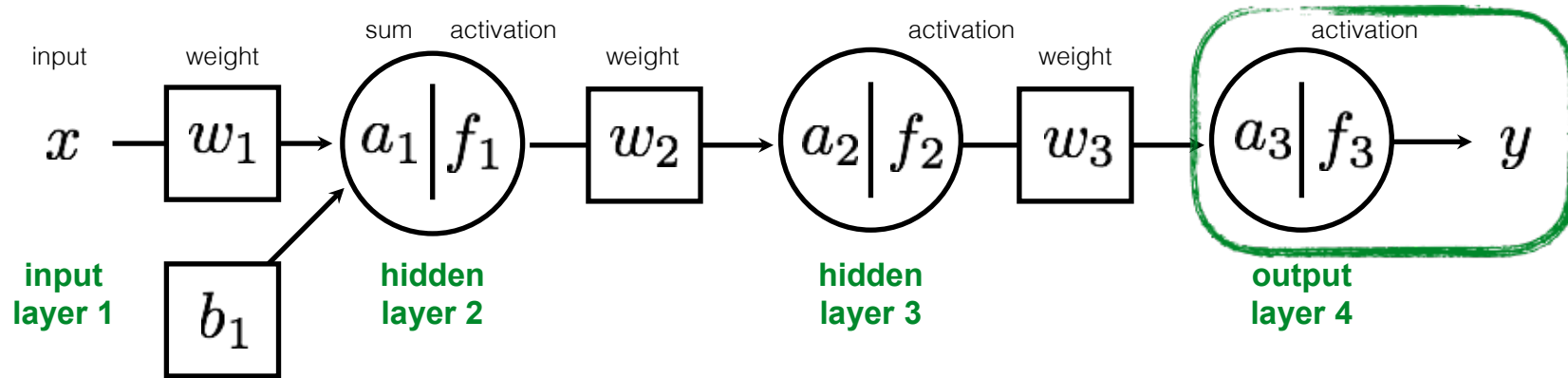$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

Entire network can be written out as one long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

We need to train the network:

*What is known? What is unknown?*

Entire network can be written out as one long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

**known**

We need to train the network:

*What is known? What is unknown?*

Entire network can be written out as one long equation

$$y = f_3(w_3 \cdot f_2(w_2 \cdot f_1(w_1 \cdot x + b_1)))$$

activation function
sometimes has
parameters

**unknown**

We need to train the network:

*What is known? What is unknown?*

# Learning an MLP

Given a set of samples and a MLP

$$\{x_i, y_i\}$$

$$y = f_{\text{MLP}}(x; \theta)$$

Estimate the parameters of the MLP

$$\theta = \{f, w, b\}$$

## Gradient Descent

For each **random** sample $\{x_i, y_i\}$

1. Predict

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

   b. Compute Loss

2. Update

   a. Back Propagation $\quad \dfrac{\partial \mathcal{L}}{\partial \theta}$    vector of parameter partial derivatives

   b. Gradient update $\quad\quad \theta \leftarrow \theta - \eta \nabla \theta$
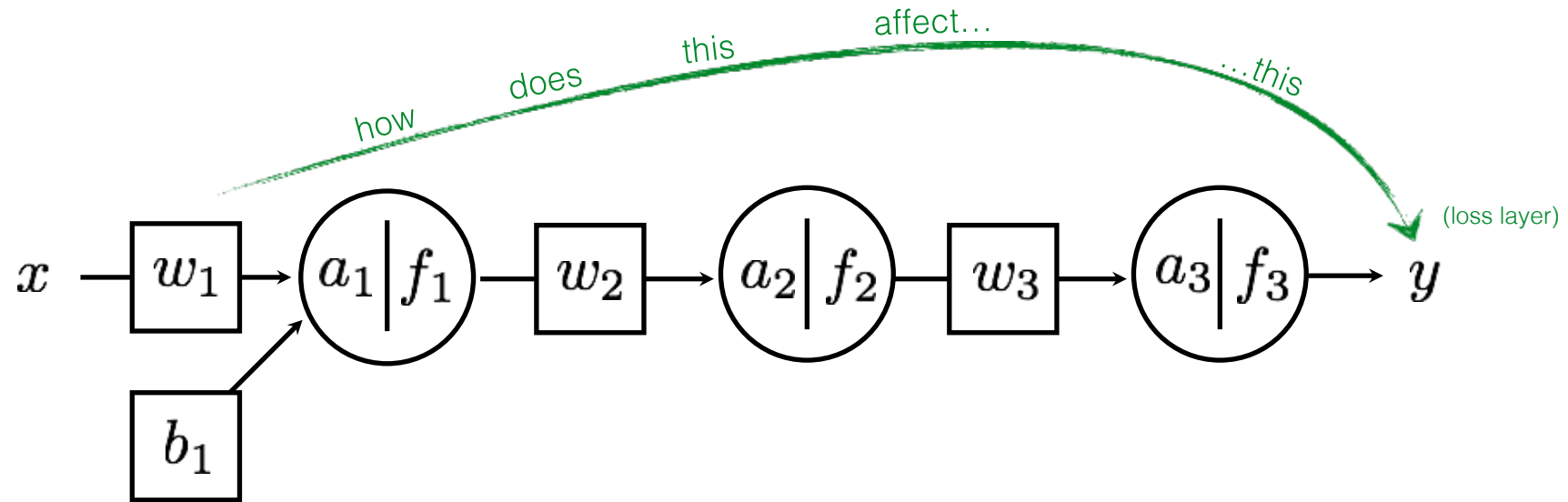
                     vector of parameter update equations

So we need to compute the partial derivatives

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \left[ \frac{\partial \mathcal{L}}{\partial w_3} \; \frac{\partial \mathcal{L}}{\partial w_2} \; \frac{\partial \mathcal{L}}{\partial w_1} \; \frac{\partial \mathcal{L}}{\partial b} \right]$$
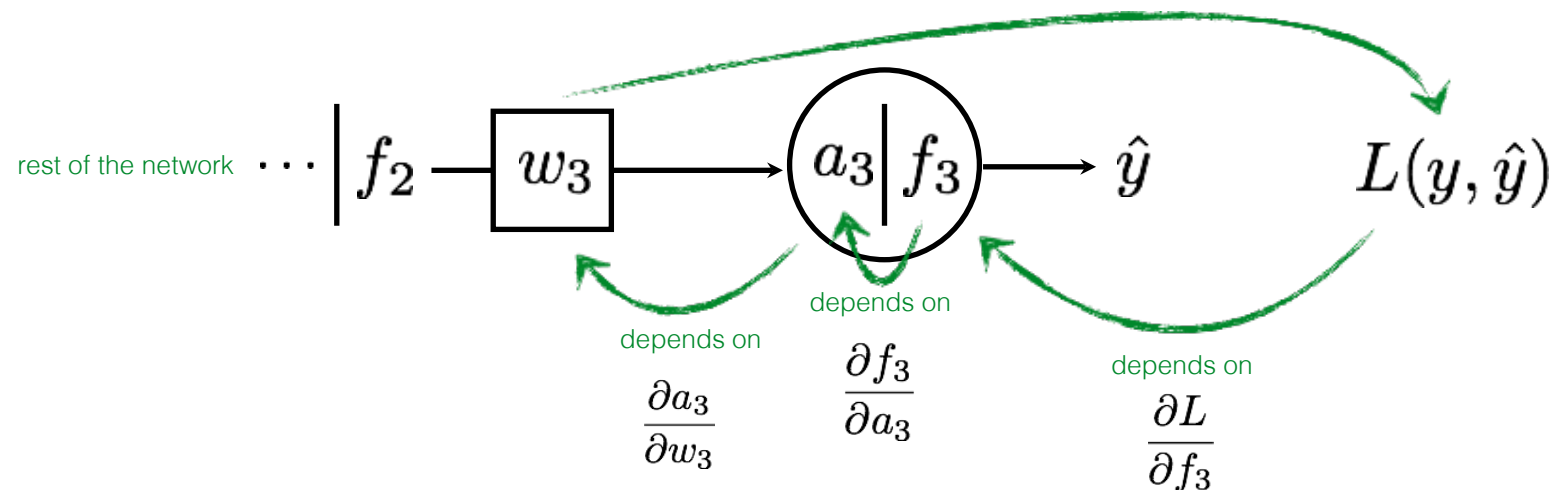
Remember,

Partial derivative $\dfrac{\partial L}{\partial w_1}$ describes…



how does this affect… …this

$x \longrightarrow \boxed{w_1} \longrightarrow \left(a_1 \middle| f_1\right) \longrightarrow \boxed{w_2} \longrightarrow \left(a_2 \middle| f_2\right) \longrightarrow \boxed{w_3} \longrightarrow \left(a_3 \middle| f_3\right) \longrightarrow y$
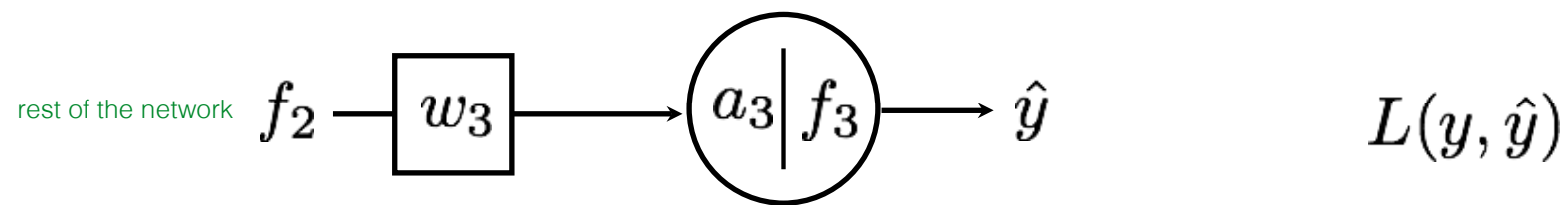
$\boxed{b_1}$

(loss layer)

According to the chain rule…

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

Intuitively, the effect of weight on loss function : $\dfrac{\partial L}{\partial w_3}$

rest of the network $f_2$ — $w_3$ → $a_3 \big| f_3$ → $\hat{y}$

$$L(y, \hat{y})$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial a_3} \frac{\partial a_3}{\partial w_3}$$

Chain Rule!

# Next Class:

# How to use chain rule and "backpropagation" to train any neural network