# CSC4007 Advanced Machine Learning

## Lesson 08: Neural Networks

by Vien Ngo
EEECS / ECIT / DSSC

# Outline

- Neural network basics and representation

- Perceptron learning, multi-layer perceptron

- Nueral network training: Backpropagation

- Modern neural network architecture (a.k.a Deep learning):

  - Convolutional neural network (CNN)

  - Recurrent neural network (RNN), long-short term memory network (LSTM)
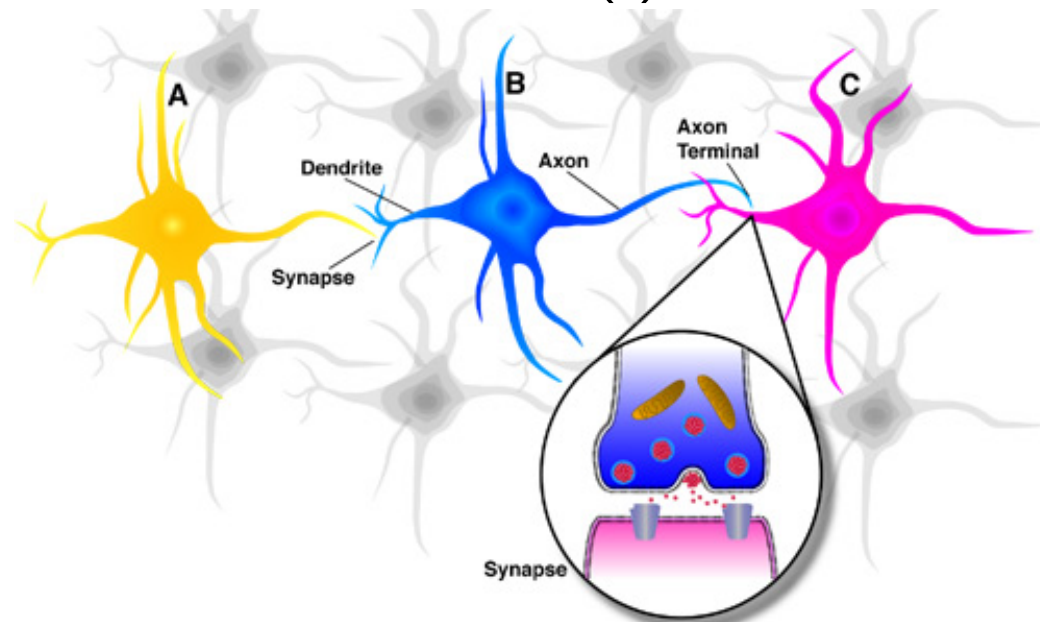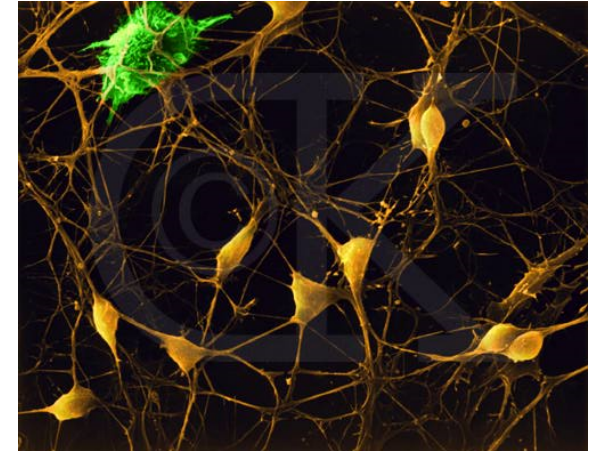
# Neural Network History

- History traces back to the 40's but became popular in the 80's with work by Hopfield, Rumelhart, Hinton, and Mclelland

- Peaked in the 90's. Today:
  - Hundreds of variants
  - Less a model of the actual brain than a useful tool, but still some debate

- Numerous applications
  - Handwriting, face, speech recognition
  - Vehicles that drive themselves
  - Models of reading, sentence production, dreaming

- Recent major resurgence
  - NNs are computationally expensive, so only recently large scale neural networks became computationally feasible

# Reasons to study neural computation

- To understand how the brain actually works.
  - It's very big and very complicated and made of stuff that dies when you poke it around. So we need to use computer simulations.
- To understand a style of parallel computation inspired by neurons and their adaptive connections.
  - Very different style from sequential computation.
    - should be good for things that brains are good at (e.g. vision)
    - should be bad for things that brains are bad at (e.g. 23 x 71)
- To solve practical problems by using novel learning algorithms inspired by the brain (this lesson)
  - Learning algorithms can be very useful even if they are not how the brain actually works.
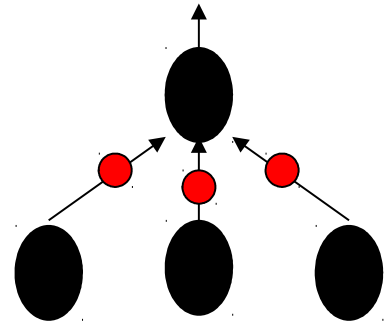
# Neurons in the Brain

- Although heterogeneous, at a low level the brain is composed of neurons

  - A neuron receives input from other neurons (generally thousands) from its synapses (can be adapted ~ learning)

  - Inputs are approximately summed

  - When the input exceeds a threshold the neuron sends an electrical spike that travels from the body, down the axon, to the next neuron(s)

# How the brain works on one slide!

- Each neuron receives inputs from other neurons
  - A few neurons also connect to receptors.
  - Cortical neurons use spikes to communicate.
- The effect of each input line on the neuron is controlled by a synaptic weight
  - The weights can be positive or negative.
- The synaptic weights adapt so that the whole network learns to perform useful computations
  - Recognizing objects, understanding language, making plans, controlling the body.
- You have about $10^{11}$ neurons each with about $10^{4}$ weights.
  - A huge number of weights can affect the computation in a very short time. Much better bandwidth than a workstation.

# Idealized Neurons

- To model things we have to idealize them (e.g. atoms)
  - Idealization removes complicated details that are not essential for understanding the main principles.
  - It allows us to apply mathematics and to make analogies to other, familiar systems.
  - Once we understand the basic principles, its easy to add complexity to make the model more faithful.
- It is often worth understanding models that are known to be wrong (but we must not forget that they are wrong!)
  - E.g. neurons that communicate real values rather than discrete spikes of activity.
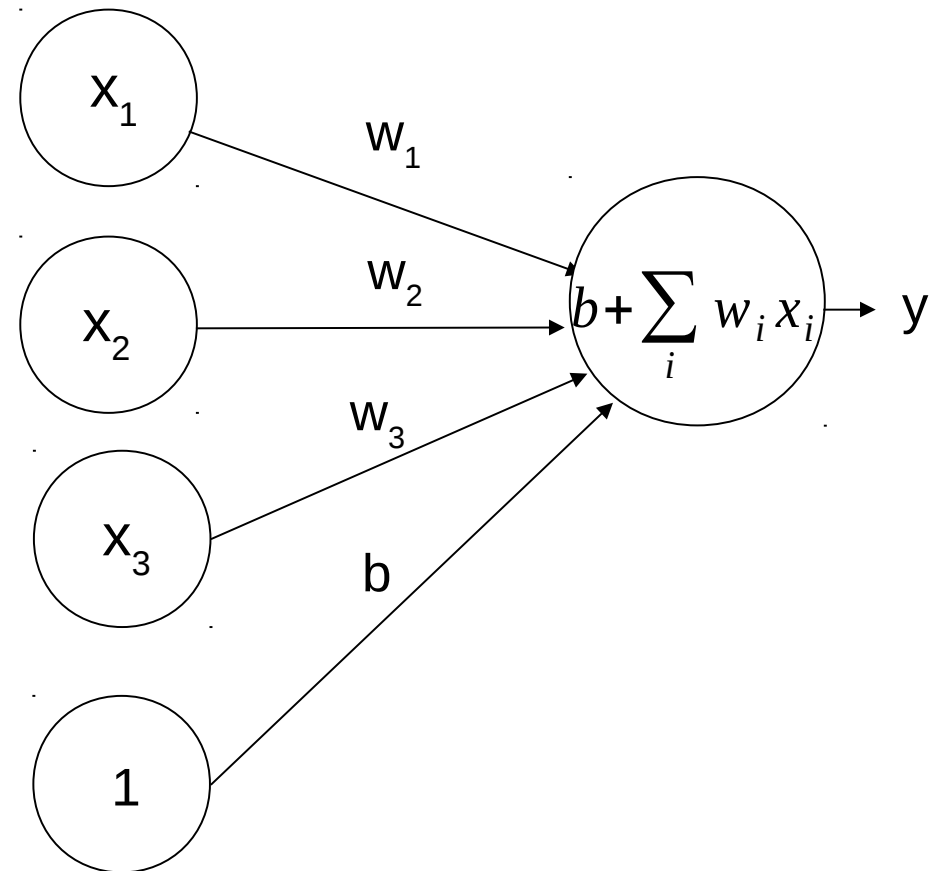
# Linear Neurons

- These are simple but computationally limited
  - If we can make them learn we may get insight into more complicated neurons.

$$y = b + \sum_i x_i w_i$$

bias — i$^{th}$ input dimension

output

index over input connections

weight on i$^{th}$ dimension

# Binary threshold neurons

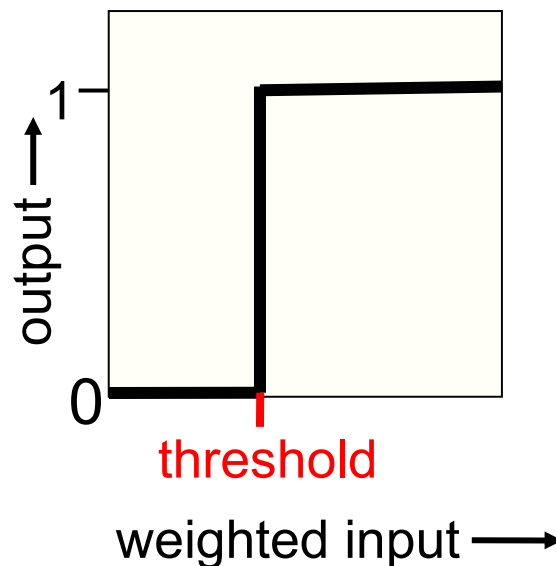- There are two equivalent ways to write the equations for a binary threshold neuron:

$$z = \sum_i x_i w_i$$

$$z = b + \sum_i x_i w_i$$

$$\boxed{\theta = -b}$$

$$y = \begin{cases} 1 \text{ if } z > \theta \\ 0 \text{ otherwise} \end{cases}$$

$$y = \begin{cases} 1 \text{ if } z > 0 \\ 0 \text{ otherwise} \end{cases}$$
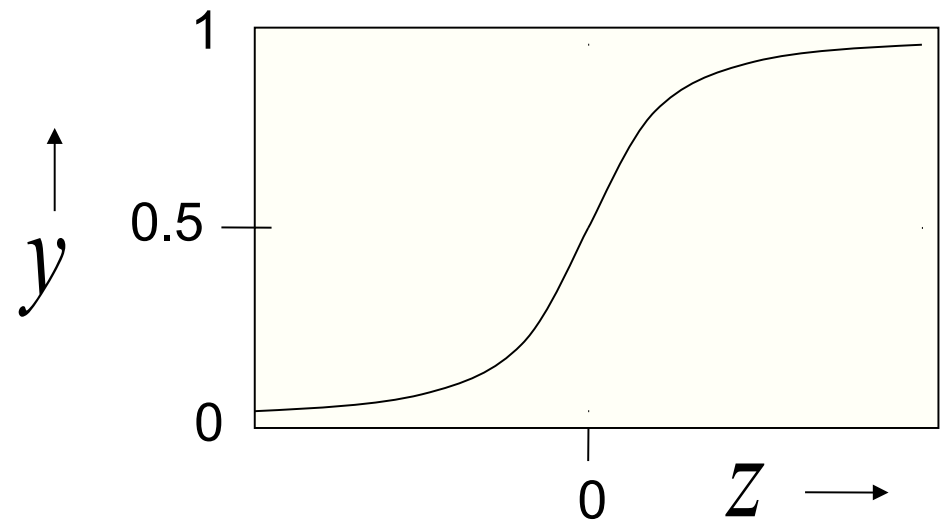
# Sigmoid neurons

- These give a real-valued output that is a smooth and bounded function of their total input.
  - Typically they use the logistic function
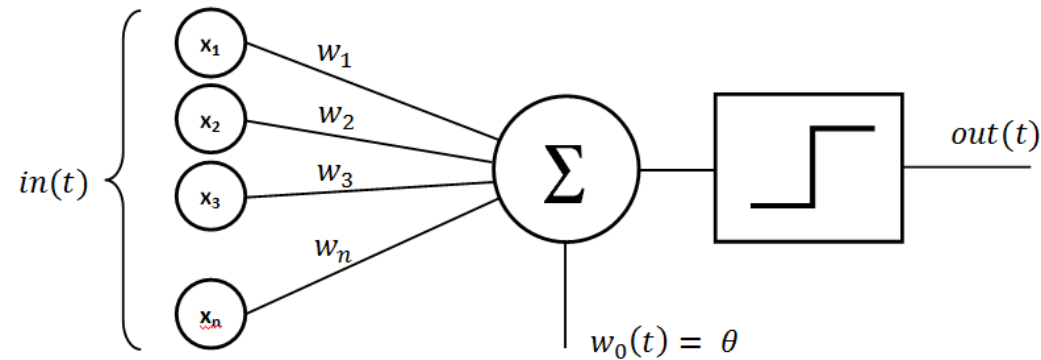  - They have nice derivatives which make learning easy (we will learn backpropagation).

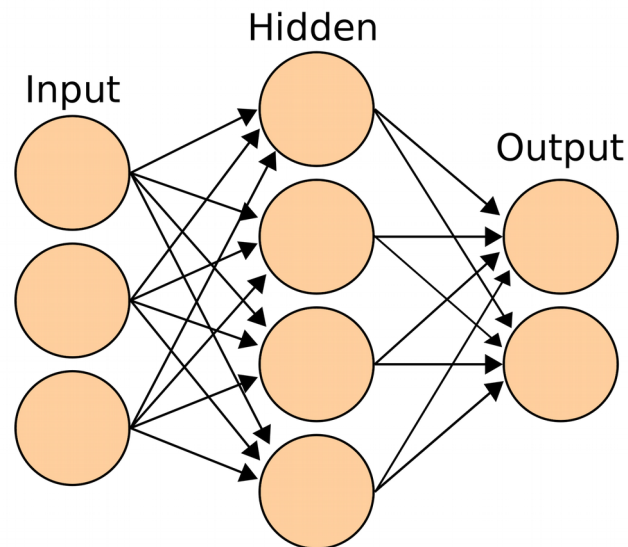$$z = b + \sum_i x_i w_i \qquad y = \frac{1}{1 + e^{-z}}$$

# Simple Neural Networks

- Perceptron



- Multi-layer Perceptron

# Perceptron
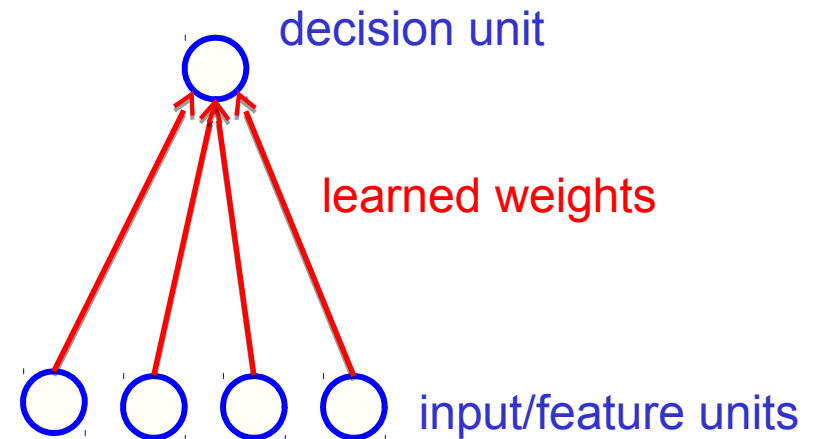
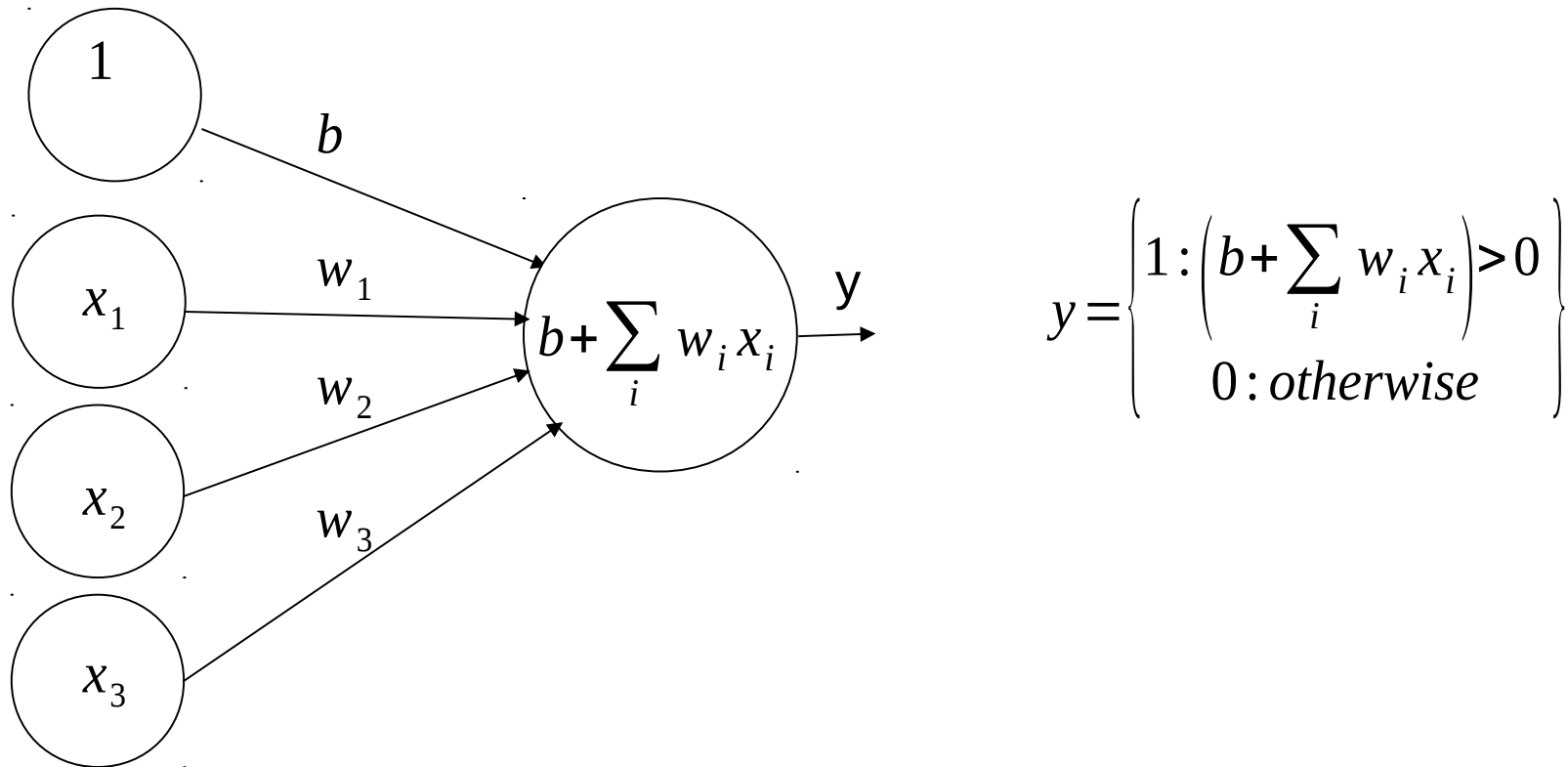# The standard Perceptron architecture

- Learn how to weight each of the input/feature variables to get a single scalar quantity.

- If this quantity is above some threshold, decide that the input vector is a positive example of the target class.

- The simplest neural network algorithm for supervised learning of **binary classification**

decision unit

learned weights

input/feature units

# Perceptrons

- Essentially a linear discriminant composed of nodes, weights
- It is based on a linear neuron and a binary threshold neuron.



$$y = \begin{cases} 1 : \left( b + \sum_i w_i x_i \right) > 0 \\ 0 : otherwise \end{cases}$$

E.g.: Input x=[$x_1$, $x_2$, $x_3$] is a 3-dimensional vector

The weights are w=[b, $w_1$, $w_2$, $w_3$] is a 4-dimensional vector

The bias is b

# Perceptrons: Binary threshold neurons

- McCulloch-Pitts (1943)
  - First compute a weighted sum of the inputs from other neurons (plus a bias).
  - Then output a 1 if the weighted sum exceeds zero.

$$z = b + \sum_i x_i w_i$$

$$y = \begin{cases} 1 \text{ if } z > 0 \\ 0 \text{ otherwise} \end{cases}$$

# Perceptron: Example

Input x=[2,-1]   weights w=[0.3,0.25,0.3]

weights w=[0.3,-0.25,0.3]

bias



z = 1(0.3) + 2(0.25) - 1(0.3)= 0.5

z>0 ⟶ y=1

Z = 1(0.3) - 2(0.25) - 1(0.3) = -0.5

z<0 ⟶ y=0

# The Perceptron: Learning

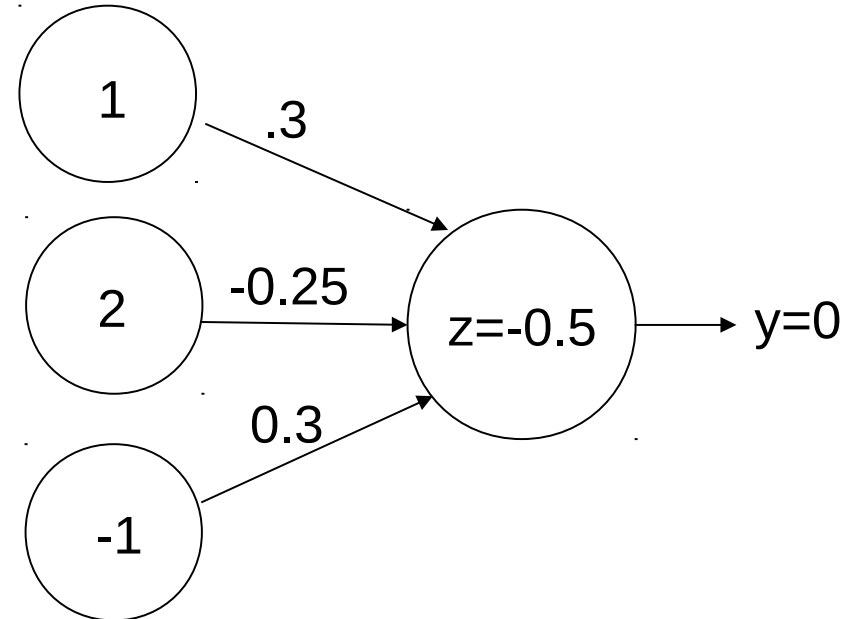- Given a dataset with N labeled data $(x_1,y_1)$, $(x_2,y_2)$, $(x_3,y_3)$,....$(x_N,y_N)$
- Learning consists of modifying the weights.

- We initialise weights to be random.
  - e.g. randomize $w=[b, w_1, w_2, w_3]$
- Training
  - Loop:
    - For each training item ***i***:
      1. Present $x_i$ as input to the perceptron
      2. Compute the output $\bar{y}_i$ (predicted)
      3. Compute the **error** (mismatch between predicted and correct output):
      $$e_i = y_i - \bar{y}_i$$
      4. Adjust the weights according to the error:

$$w = w + \alpha e_i \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$

$\alpha$ is a learning rate (step-size)

$1$

$b$

$x_{i1}$

$w_1$

output

$\bar{y}_i$

$x_{i2}$

$w_2$

$w_3$

$x_{i3}$

$$\bar{y}_i = \begin{cases} 1 : (b + w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3}) > 0 \\ 0 : otherwise \end{cases}$$

e.g.

$$\begin{bmatrix} 1 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix} = \begin{bmatrix} 1 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix} + \alpha e_i \begin{bmatrix} 1 \\ x_{i1} \\ x_{i2} \\ x_{i3} \end{bmatrix}$$

# The Perceptron: Learning

- Given a dataset with N labeled data $(x_1,y_1)$, $(x_2,y_2)$, $(x_3,y_3)$,….$(x_N,y_N)$
- Learning consists of modifying the weights.

- We initialise weights to be random.
  $$w=[b, w_1, w_2, w_3]$$
- Training
  - Loop:
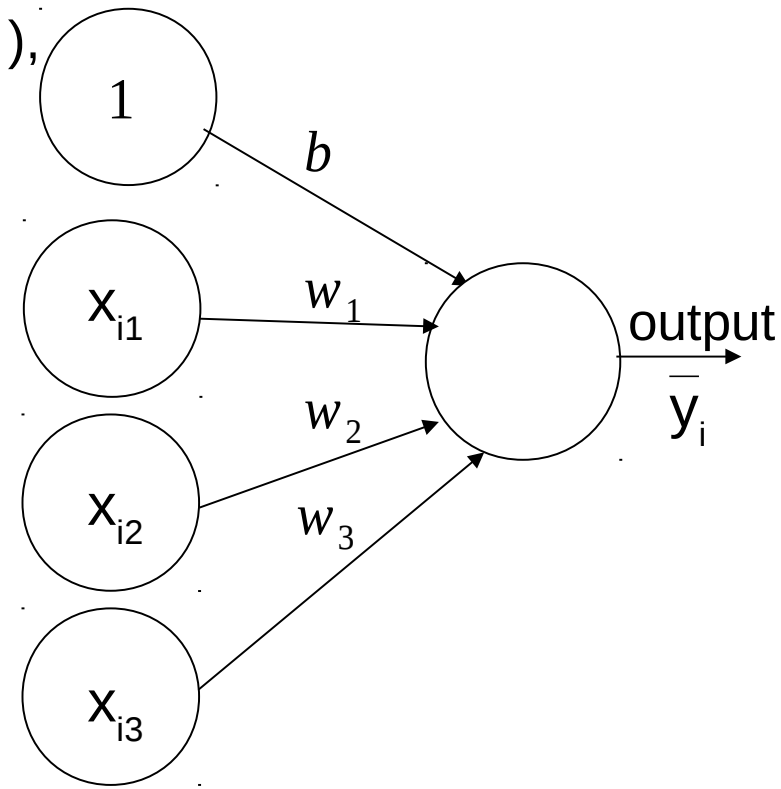    - For each training item *i*:
      1. Present $x_i$ as input to the perceptron
      2. Compute the output $\bar{y}_i$ (predicted)
      3. Compute the **error** (mismatch between predicted and correct output):
         $$e_i = y_i - \bar{y}_i$$
      4. Adjust the weights according to the error:
         $$w = w + \alpha\, e_i \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$
  - Terminate when certain accuracy reached.

# The Perceptron: Learning

- Given a dataset with N labeled data $(x_1, y_1)$, $(x_2, y_2)$, $(x_3, y_3)$,….$(x_N, y_N)$
- Learning consists of modifying the weights.

- We initialise weights to be random.
    $$w = [b, w_1, w_2, w_3]$$
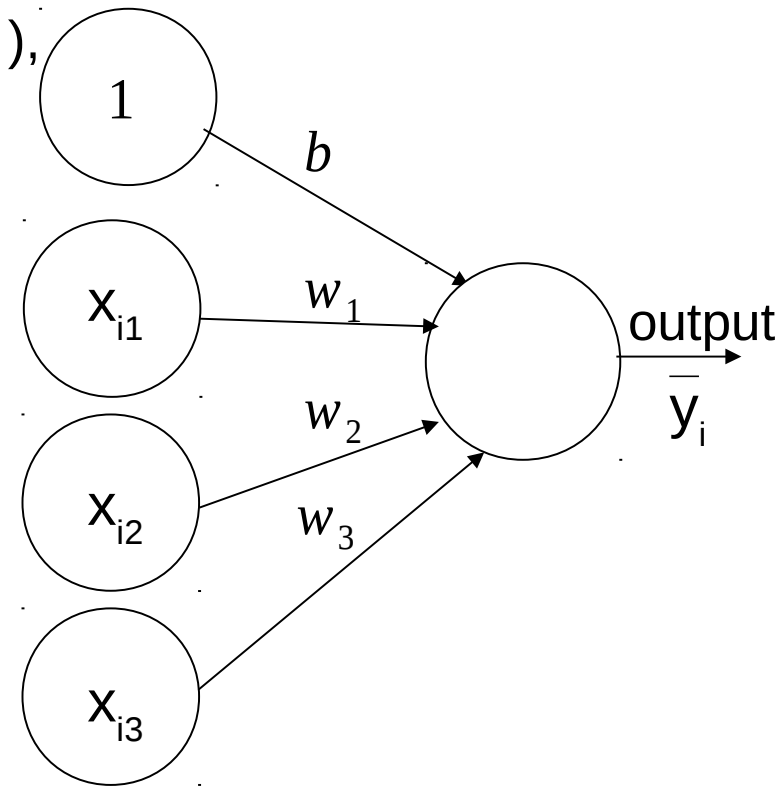
- Training
    - Loop:
        - For each training item $i$:
        1. Present $x_i$ as input to the perceptron
        2. Compute the output $\bar{y}_i$ (predicted)
        3. Compute the **error** (mismatch between predicted and correct output):
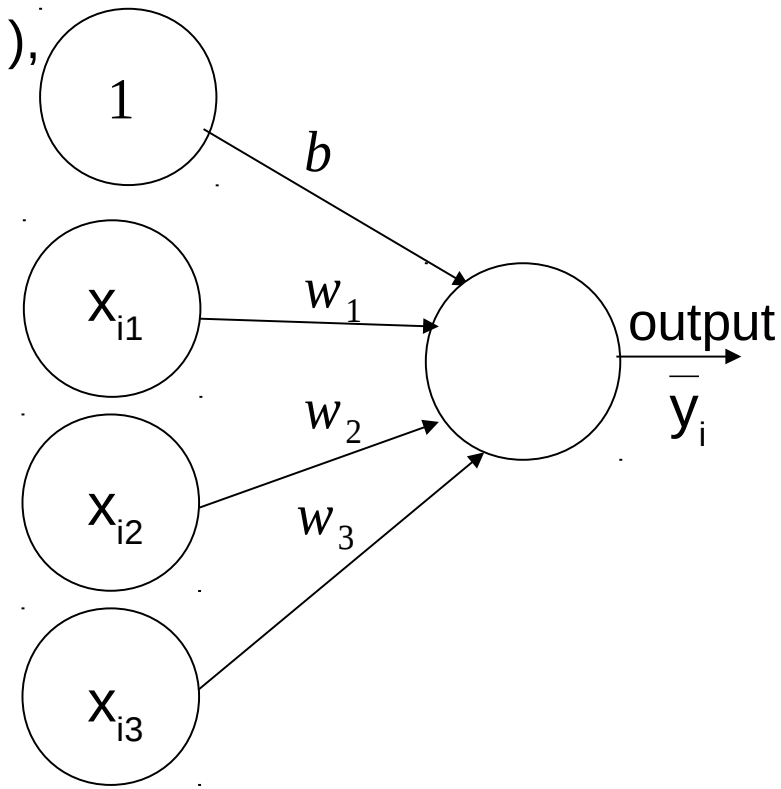            $$e_i = y_i - \bar{y}_i$$
        4. Adjust the weights according to the error:
            $$w = w + \alpha\, e_i \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$

1 — $b$

$x_{i1}$ — $w_1$

output $\bar{y}_i$

$x_{i2}$ — $w_2$

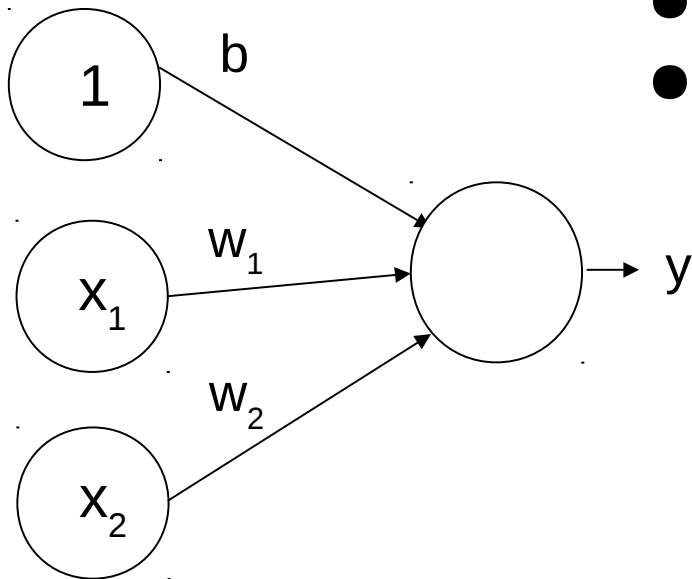$x_{i3}$ — $w_3$

**One iteration: called an epoch**
(when the training process has one complete pass through the training dataset

# Perceptron Example: AND function

- Two input nodes (with bias term): x = [1,$x_1$, $x_2$]
- No hidden layer
- One output node: y
- The weights of this network: w=[b, $w_1$, $w_2$]

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



- Initialize: w=[0.1, 0.2, -0.2], set $\alpha = 0.5$
- Loop:
  - For each training instance:
    - x=[1,0,0],
      - Compute 1+0*0.2-0*0.2=1>0, predict $\bar{y}$ = 1.
      - The error is e = y- $\bar{y}$ = 0-1 = -1
      - Update w = w – 0.5*1 * [1,0,0] = [-0.4,0.2,-0.2]
    - x=[1,0,1],
      - Compute -0.4+0*0.2-1*0.2=-0.6<0, predict $\bar{y}$ = 0
      - The error is e = y- $\bar{y}$ = 0-0 = 0
      - Update w = w + 0.5*0 * [1,0,1] = [-0.4,0.2,-0.2]

# Perceptron Example: AND function

- Two input nodes (with bias term): x = [1,$x_1$, $x_2$]
- No hidden layer
- One output node: y
- The weights of this network: w=[b, $w_1$, $w_2$]

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



- Initialize: w=[0.1, 0.2, -0.2], set $\alpha = 0.5$
- Loop:
  - For each training instance (cont'd):
    - x=[1,1,0],
      - Compute -0.4+1*0.2-0*0.2=-0.2<0, predict $\bar{y}$=0
      - The error is e = y- $\bar{y}$ = 0-0 = 0
      - Update w = w + 0.5*0 * [1,0,0] = [-0.4,0.2,-0.2]
    - x=[1,1,1],
      - Compute -0.4+1*0.2-1*0.2=-0.4<0, predict $\bar{y}$ = 0
      - The error is e = y- $\bar{y}$ = 1-0 = 1
      - Update w = w + 0.5*1 * [1,1,1] = [0.1,0.7,0.3]

# Perceptron Example: AND function

- Two input nodes (with bias term): $x = [1, x_1, x_2]$
- No hidden layer
- One output node: y
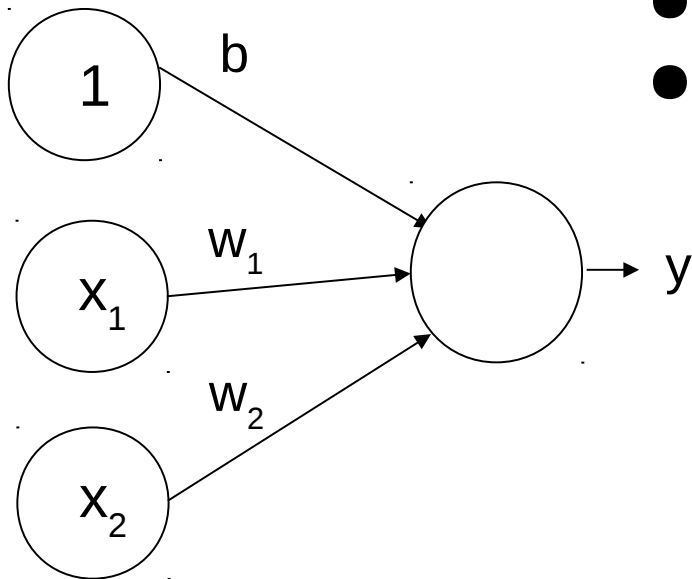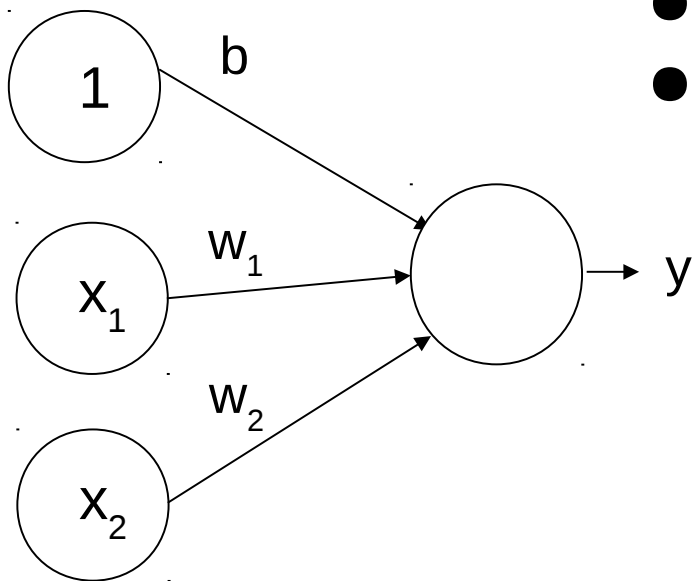- The weights of this network: $w = [b, w_1, w_2]$

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

- Initialize: $w = [0.1, 0.2, -0.2]$, set $\alpha = 0.5$
- Loop:
  - After the first epoch
    - $w = [0.1, 0.7, 0.3]$



1

b

$w_1$

$x_1$

$w_2$

$x_2$

y

# Perceptron: Applications?

- Spam filtering: given an email → **spam** or **not spam**



data

prediction

Spam
vs.
Not Spam

Bag of words

$$x = \begin{bmatrix} winner = 2 \\ prize = 1 \\ \$\$ = 3 \\ etc. \end{bmatrix}$$

Spam?

Not?

$b$

$w_1$

$w_2$

$w_3$

# The Perceptron: Limitations

- Perceptrons are incredibly limited in their abilities.
  *Can only solve linearly separable problems*



0,1             1,1

weight plane    output =1
output =0

0,0             1,0

The positive and negative cases cannot be separated by a plane

# The Perceptron: Limitations

We can improve things with a ***multi-layered perceptron*** (learning with hidden units)



Input layer   Hidden layer   Output layer

A network of many neurons. Each node is one **Linear Neuron**

- **Input Layer (layer 1):** neurons that receive the inputs
- **Hidden layers (layers 2,3, ...):** connected to neither the inputs nor the outputs of the network directly
- **Output layer (last layer):** neurons from which we read the results.

# Multi-layered perceptron

We can improve things with a ***multi-layered perceptron***



Input layer   Hidden layer   Output layer

- However, training is now much more complicated.

- With the simple perceptron, we could easily evaluate how to change the weights according to the error.

- Now there are so many different connections, each in a different layer of the network.

- **How does one know how much each neuron or connection contributed to the overall error of the network?**

# Multi-layered perceptron



Input layer   Hidden layer   Output layer

- **Input variables**: $x = [x_0, x_1, x_2]$ ($x_0$ is often defined as bias term: $x_0 = 1$)
- **Parameter matrix** $w^{(j)}$ controlling the function mapping from layer $j$ to layer $j+1$.
  - ↱ If layer $j$ has m units, and $j+1$ has n units, then $w^{(j)}$ is mxn matrix
  - ↱ e.g. $w^{(1)}$ is 3x4 matrix, $w^{(2)}$ is 4x2 matrix
- **Activation function** $\sigma_i^{(j)}$ at node i of the layer j (By activation, we mean the value which is computed and output by that node).
  - ↱ It receives the input to the node i, and transforms it.
  - ↱ Input layer could also use activation
- **Output variables**: $y = [y_0, y_1]$ (with activations $\sigma_0^{(3)}, \sigma_1^{(3)}$

Each node is a linear neuron with activations (e.g. binary threshold neuron, sigmoid neuron, etc.). Therefore the output after each node is the <span style="color:red">activation function</span> applied to the <span style="color:red">linear combination</span> of <span style="color:red">its inputs</span>

# Multi-layered perceptron



Input layer   Hidden layer   Output layer

- We calculate each of the layer-2 activations based on the input values.
- We then calculate the output (prediction) (i.e. two nodes in layer 3) using exactly the same logic, except in input is not x values, but the activation values from the preceding layer.
- The activation value on each hidden unit is equal to the activation function applied to the linear combination of its inputs
- Every input/activation goes to every node in following layer

Calculations of four hidden nodes:

Hidden node 1: $g_0 = \sigma_0^{(2)}\left(x_0 w_{00}^{(1)} + x_1 w_{10}^{(1)} + x_2 w_{20}^{(1)}\right)$

Hidden node 2: $g_1 = \sigma_1^{(2)}\left(x_0 w_{01}^{(1)} + x_1 w_{11}^{(1)} + x_2 w_{21}^{(1)}\right)$

Hidden node 3: $g_2 = \sigma_2^{(2)}\left(x_0 w_{02}^{(1)} + x_1 w_{12}^{(1)} + x_2 w_{22}^{(1)}\right)$

Hidden node 4: $g_3 = \sigma_3^{(2)}\left(x_0 w_{03}^{(1)} + x_1 w_{13}^{(1)} + x_2 w_{23}^{(1)}\right)$

Calculations of two output nodes:

$$y_0 = \sigma_0^{(3)}\left(g_0 w_{00}^{(2)} + g_1 w_{10}^{(2)} + g_2 w_{20}^{(2)} + g_3 w_{30}^{(2)}\right)$$

$$y_1 = \sigma_1^{(3)}\left(g_0 w_{01}^{(2)} + g_1 w_{11}^{(2)} + g_2 w_{21}^{(2)} + g_3 w_{31}^{(2)}\right)$$

# Multi-layered perceptron: Activations

- **Activation function** $\sigma_i^{(j)}$ at node i of the layer j (By activation, we mean the value which is computed and output by that node).



$$y = \sigma\left(\sum_i x_i w_i\right) = \begin{cases} 1: if \sum_i x_i w_i > 0 \\ 0: otherwise \end{cases}$$

Perceptron's activation function: $\sigma(z) = \begin{cases} 1: if\ z > 0 \\ 0: otherwise \end{cases}$

# Multi-layered perceptron: Activations

- Two common activation functions used in multi-layered perceptron
  - ✓ sigmoid and tanh functions

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

$$\sigma(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

# Multi-layered perceptron (MLP)



Input layer   Hidden layer   Output layer

Calculations of four hidden nodes:

$$g_0 = \sigma_0^{(2)}\left( x_0 w_{00}^{(1)} + x_1 w_{10}^{(1)} + x_2 w_{20}^{(1)} \right)$$

$$g_1 = \sigma_1^{(2)}\left( x_0 w_{01}^{(1)} + x_1 w_{11}^{(1)} + x_2 w_{21}^{(1)} \right)$$

$$g_2 = \sigma_2^{(2)}\left( x_0 w_{02}^{(1)} + x_1 w_{12}^{(1)} + x_2 w_{22}^{(1)} \right)$$

$$g_3 = \sigma_3^{(2)}\left( x_0 w_{03}^{(1)} + x_1 w_{13}^{(1)} + x_2 w_{23}^{(1)} \right)$$

- If the activations at last layer are sigmoid functions, this layer is just logistic regression
  - The only difference is, instead of input a feature vector, the features are just values calculated by the hidden layer
- NN is representation learning: Instead of being constrained by the original input features, a neural network can learn its own features to feed into logistic regression.
  - So we feed the hidden layers our input values, and let them learn whatever gives the best final result to feed into the final output layer.

Calculations of two output nodes:

$$y_0 = \sigma_0^{(3)}\left( g_0 w_{00}^{(2)} + g_1 w_{10}^{(2)} + g_2 w_{20}^{(2)} + g_3 w_{30}^{(2)} \right)$$

$$y_1 = \sigma_1^{(3)}\left( g_0 w_{01}^{(2)} + g_1 w_{11}^{(2)} + g_2 w_{21}^{(2)} + g_3 w_{31}^{(2)} \right)$$

# MLP: Vectorization
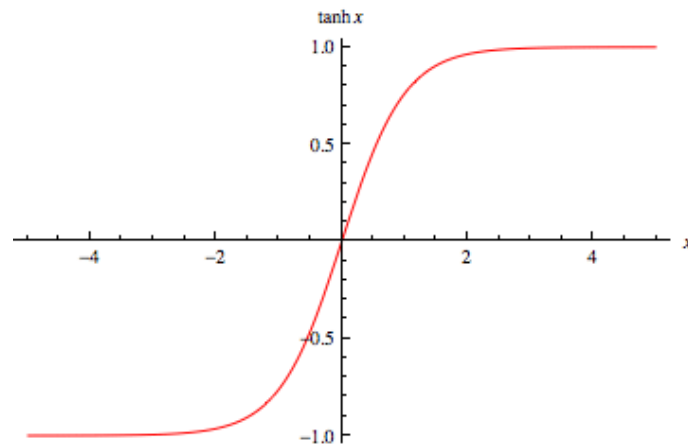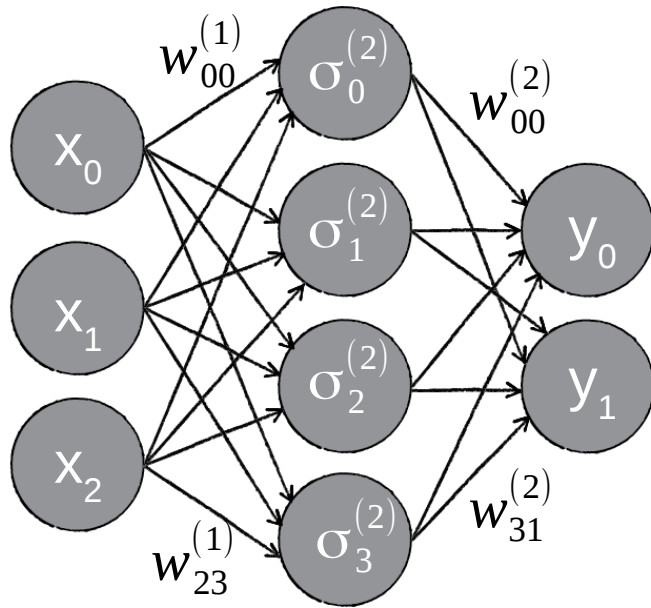


- Note that: $w^{(j)}$ is the paramter matrix mapping from layer j to layer j+1
  - e.g.: $w^{(1)}_{ij}$ is the parameter from note i (layer 1) to node j (layer 2)
  - Denote $w^{(j)}_i$ is the $i^{th}$ column of $w^{(j)}$
    - e.g.: $w^{(1)}_0 = [w^{(1)}_{00}, w^{(1)}_{10}, w^{(1)}_{20}]$
- Assume all nodes in the same layer using the same activations function, e.g. $\sigma^{(2)}_i = \sigma^{(2)}$

Input layer   Hidden layer   Output layer

$$\text{input } x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \quad \text{denote } z = \begin{bmatrix} x^T w^{(1)}_0 \\ x^T w^{(1)}_1 \\ x^T w^{(1)}_2 \\ x^T w^{(1)}_3 \end{bmatrix} = w^{(1)} x$$

Calculations of four hidden nodes:

$$g_0 = \sigma^{(2)}(x_0 w^{(1)}_{00} + x_1 w^{(1)}_{10} + x_2 w^{(1)}_{20})$$
$$g_1 = \sigma^{(2)}(x_0 w^{(1)}_{01} + x_1 w^{(1)}_{11} + x_2 w^{(1)}_{21})$$
$$g_2 = \sigma^{(2)}(x_0 w^{(1)}_{02} + x_1 w^{(1)}_{12} + x_2 w^{(1)}_{22})$$
$$g_3 = \sigma^{(2)}(x_0 w^{(1)}_{03} + x_1 w^{(1)}_{13} + x_2 w^{(1)}_{23})$$

vectorize

$$g_0 = \sigma^{(2)}(x^T w^{(1)}_0)$$
$$g_1 = \sigma^{(2)}(x^T w^{(1)}_1)$$
$$g_2 = \sigma^{(2)}(x^T w^{(1)}_2)$$
$$g_3 = \sigma^{(2)}(x^T w^{(1)}_3)$$

$$g = \begin{bmatrix} g_0 \\ g_1 \\ g_2 \\ g_3 \end{bmatrix} = \begin{bmatrix} \sigma^{(2)}(x^T w^{(1)}_0) \\ \sigma^{(2)}(x^T w^{(1)}_1) \\ \sigma^{(2)}(x^T w^{(1)}_2) \\ \sigma^{(2)}(x^T w^{(1)}_3) \end{bmatrix}$$

# MLP: Vectorization

Calculations of four hidden nodes:

$$g_0 = \sigma^{(2)}\left(x_0 w_{00}^{(1)} + x_1 w_{10}^{(1)} + x_2 w_{20}^{(1)}\right)$$
$$g_1 = \sigma^{(2)}\left(x_0 w_{01}^{(1)} + x_1 w_{11}^{(1)} + x_2 w_{21}^{(1)}\right)$$
$$g_2 = \sigma^{(2)}\left(x_0 w_{02}^{(1)} + x_1 w_{12}^{(1)} + x_2 w_{22}^{(1)}\right)$$
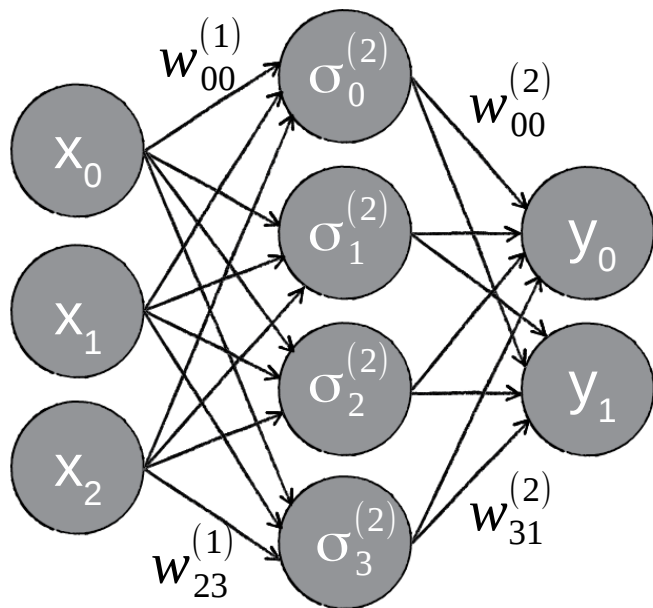$$g_3 = \sigma^{(2)}\left(x_0 w_{03}^{(1)} + x_1 w_{13}^{(1)} + x_2 w_{23}^{(1)}\right)$$
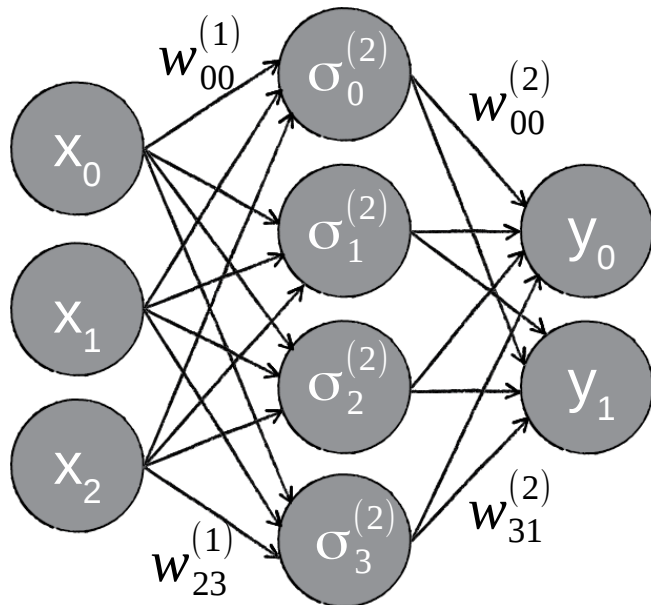
vectorize

$$\text{input} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \qquad \text{denote} \quad z = \begin{bmatrix} x^T w_0^{(1)} \\ x^T w_1^{(1)} \\ x^T w_2^{(1)} \\ x^T w_3^{(1)} \end{bmatrix} = w^{(1)} x$$

$$g_0 = \sigma^{(2)}\left(x^T w_0^{(1)}\right)$$
$$g_1 = \sigma^{(2)}\left(x^T w_1^{(1)}\right)$$
$$g_2 = \sigma^{(2)}\left(x^T w_2^{(1)}\right)$$
$$g_3 = \sigma^{(2)}\left(x^T w_3^{(1)}\right)$$

$$g = \begin{bmatrix} g_0 \\ g_1 \\ g_2 \\ g_3 \end{bmatrix} = \begin{bmatrix} \sigma^{(2)}\left(x^T w_0^{(1)}\right) \\ \sigma^{(2)}\left(x^T w_1^{(1)}\right) \\ \sigma^{(2)}\left(x^T w_2^{(1)}\right) \\ \sigma^{(2)}\left(x^T w_3^{(1)}\right) \end{bmatrix}$$

$$g = \sigma^{(2)} \begin{bmatrix} x^T w_0^{(1)} \\ x^T w_1^{(1)} \\ x^T w_2^{(1)} \\ x^T w_3^{(1)} \end{bmatrix} = \sigma^{(2)}(z) = \sigma^{(2)}\left(w^{(1)} x\right)$$



Input layer   Hidden layer   Output layer

# MLP: Vectorization

- Note that: $w^{(j)}$ is the paramter matrix mapping from layer j to layer j+1
  - e.g.: $w^{(1)}_{ij}$ is the parameter from note i (layer 1) to node j (layer 2)

input $x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$     values at hidden layer $g = \begin{bmatrix} g_0 \\ g_1 \\ g_2 \\ g_3 \end{bmatrix}$     output $y = \begin{bmatrix} y_0 \\ y_1 \end{bmatrix}$

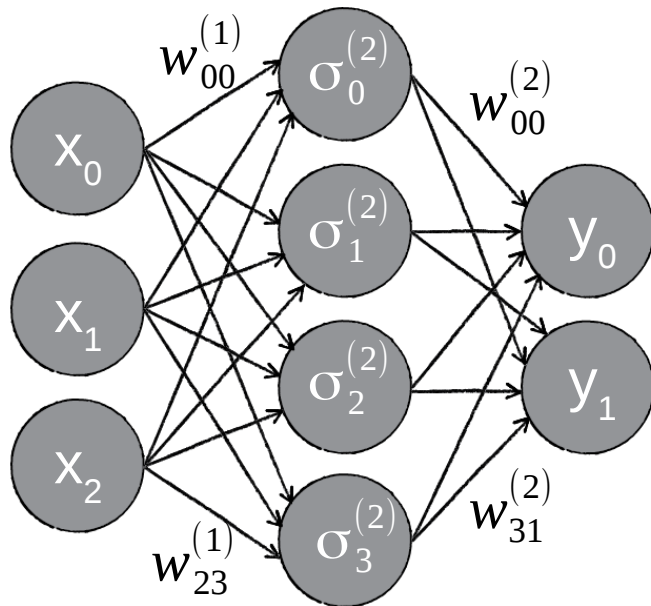Calculations of four hidden nodes (x is the input and $w^{(1)}$ are weights to this layer):

$$g = \sigma^{(2)}(w^{(1)} x)$$

Likewise, calculations of two output nodes (g is the input and $w^{(2)}$ are weights to this layer):

$$y = \sigma^{(3)}(w^{(2)} g)$$

A full forward pass (also called forward propagation):

$$y = \sigma^{(3)}\left(w^{(2)} \sigma^{(2)}(w^{(1)} x)\right)$$
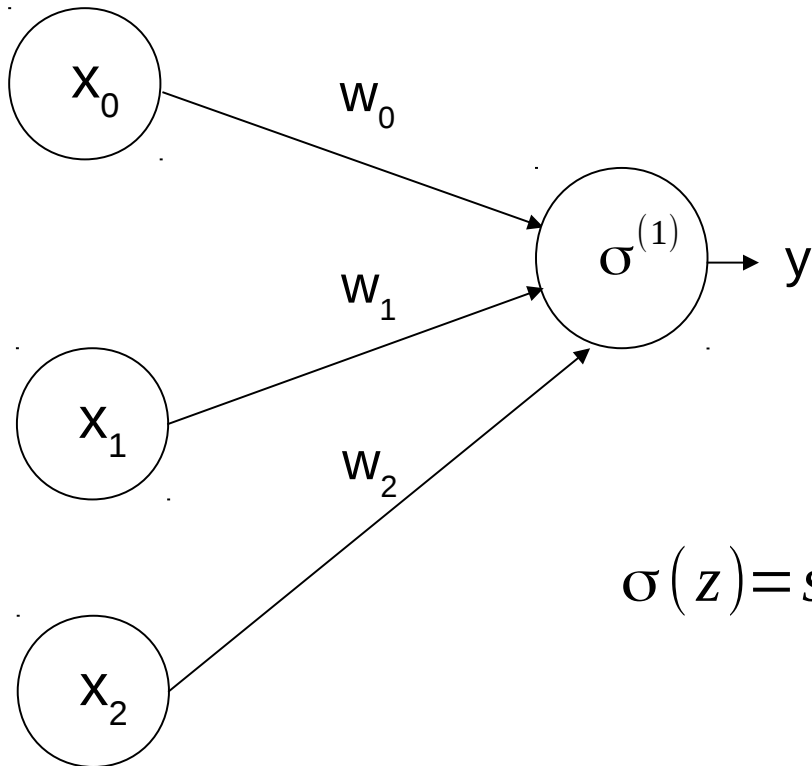
- Start off with activations of input unit

- Forward propagate and calculate the activation of each layer sequentially

- This is a vectorized version of this implementation



Input layer   Hidden layer   Output layer

# MLP Example: AND function

- A simple MLP network:
  - Two input nodes (with bias term at $x_0$): $x = [1, x_1, x_2]$
  - No hidden layer
  - One output node (**with sigmoid activations**): y
  - The weights of this network: $w^{(1)} = [w_0, w_1, w_2]$

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



$$y = \sigma^{(1)}\left(w^{(1)} x\right) = \frac{1}{1 + e^{-w^{(1)} x}} = \frac{1}{1 + e^{-w_0 - w_1 x_1 - w_2 x_2}}$$

$$\sigma(z) = sigmoid(z) = \frac{1}{1 + e^{-z}}$$

# MLP Example: AND function



$$y = \sigma^{(1)}(w^{(1)} x) = \frac{1}{1 + e^{-w^{(1)} x}} = \frac{1}{1 + e^{-w_0 - w_1 x_1 - w_2 x_2}}$$

| $x_1$ | $x_2$ | y (ground-truth) | predicted |
|-------|-------|------------------|-----------|
| 0 | 0 | 0 | Sigmoid(-30) = 0 |
| 0 | 1 | 0 | Sigmoid(-10) = 0 |
| 1 | 0 | 0 | Sigmoid(-10) = 0 |
| 1 | 1 | 1 | Sigmoid(10) = 1 |

$$\sigma(z) = sigmoid(z) = \frac{1}{1 + e^{-z}}$$

$$\bar{y}_1 = \sigma^{(1)}(w^{(1)}[1,0,0]) = \frac{1}{1 + e^{-30}}$$

$$\bar{y}_2 = \sigma^{(1)}(w^{(1)}[1,0,1]) = \frac{1}{1 + e^{-10}}$$

$$\bar{y}_3 = \sigma^{(1)}(w^{(1)}[1,1,0]) = \frac{1}{1 + e^{-10}}$$

$$\bar{y}_4 = \sigma^{(1)}(w^{(1)}[1,1,1]) = \frac{1}{1 + e^{10}}$$

# MLP Example: XOR function

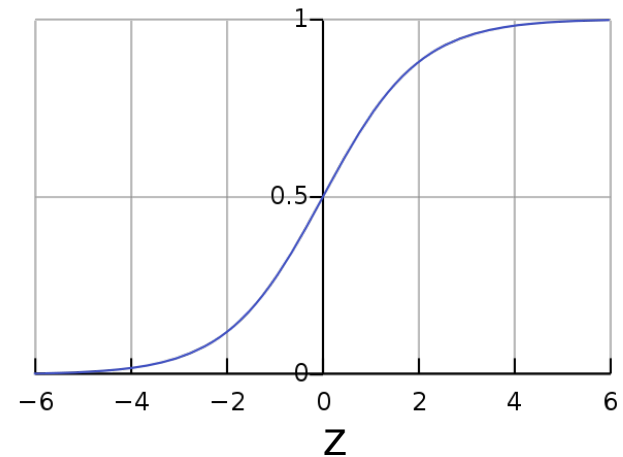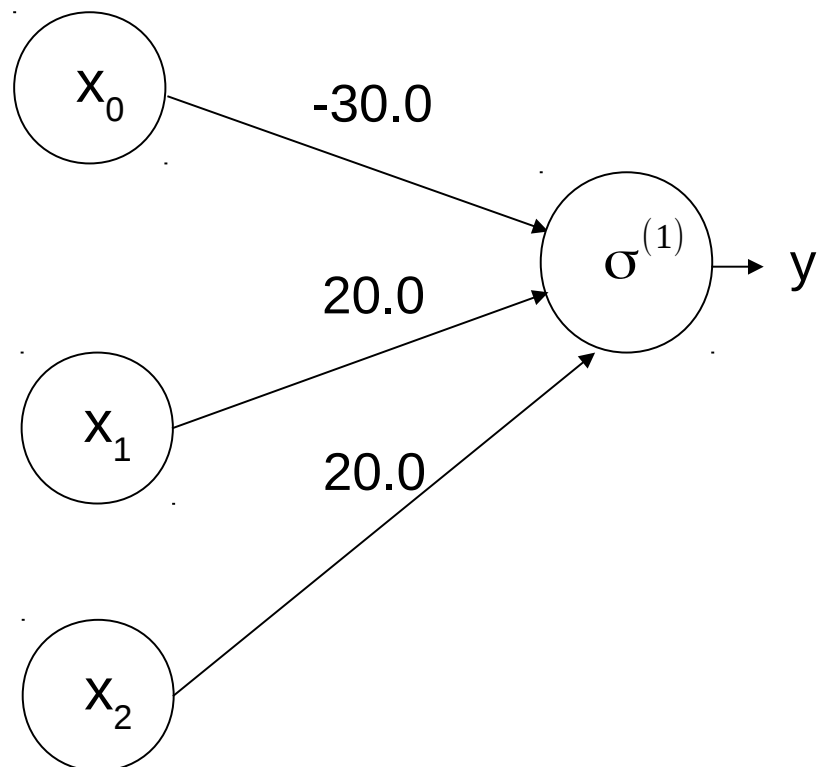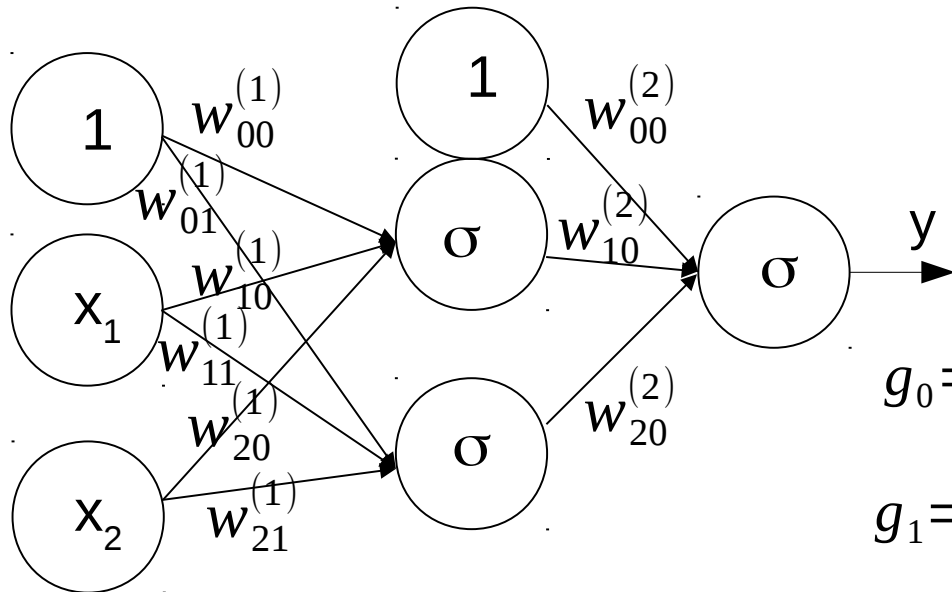| $x_1$ | $x_2$ | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

- The data is non-linearly separable, so a MLP network:
  - Two input nodes (with bias term at $x_0$): $x = [1, x_1, x_2]$
  - One hidden layer with two nodes (**with sigmoid activations**)
  - One output node (**with sigmoid activation**): y
  - The weights of this network: $w^{(2)} = [w^{(2)}_{00}, w^{(2)}_{10}, w^{(2)}_{20}]$

$$w^{(1)} = \begin{bmatrix} w^{(1)}_{00} & w^{(1)}_{01} \\ w^{(1)}_{10} & w^{(1)}_{11} \\ w^{(1)}_{20} & w^{(1)}_{21} \end{bmatrix}$$



Hidden node 1:
$$g_0 = sigmoid(w^{(1)}_{00} + w^{(1)}_{10} x_1 + w^{(1)}_{20} x_2) = sigmoid(x^T w^{(1)}_0)$$

Hidden node 2:
$$g_1 = sigmoid(w^{(1)}_{01} + w^{(1)}_{11} x_1 + w^{(1)}_{21} x_2) = sigmoid(x^T w^{(1)}_1)$$

Prediction:
$$y = sigmoid(w^{(2)}_{00} + w^{(2)}_{10} g_0 + w^{(2)}_{20} g_1) = sigmoid([1, g_0, g_1] w^{(2)})$$

# MLP Example: XOR function



$$w^{(1)} = \begin{bmatrix} -30 & 10 \\ 20 & -20 \\ 20 & -20 \end{bmatrix}$$

$$w^{(2)} = [10, -20, -20]$$

| $X_1$ | $X_2$ | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

**Prediction for x=[1,0,0]**

Hidden node 1:

$$g_0 = sigmoid\left([1,0,0]\begin{bmatrix} -30 \\ 20 \\ 20 \end{bmatrix}^T\right) = sigmoid(-30) \approx 0$$
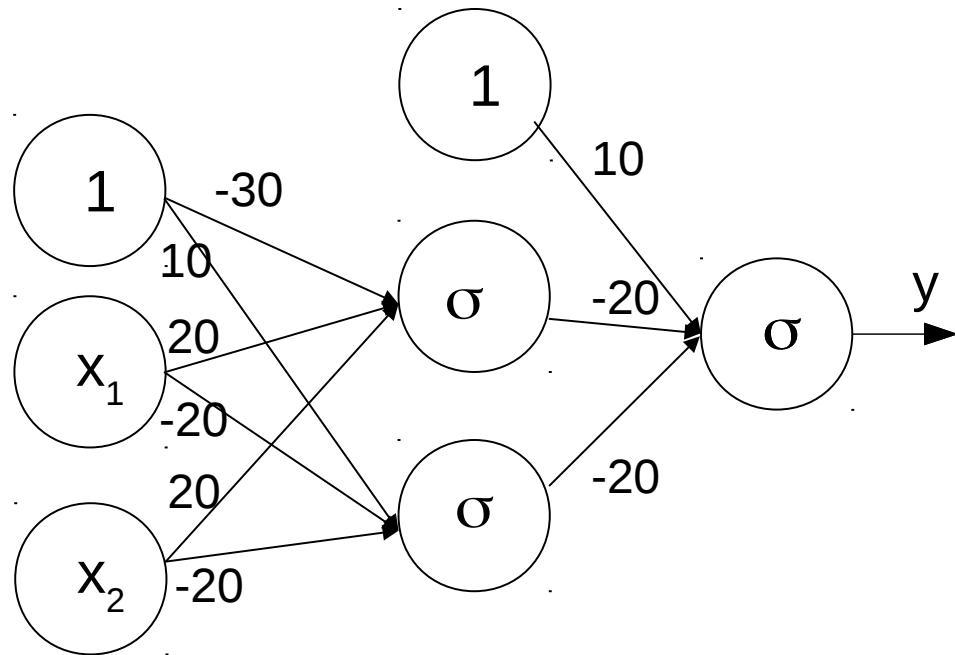
Hidden node 2:

$$g_1 = sigmoid\left([1,0,0]\begin{bmatrix} 10 \\ -20 \\ -20 \end{bmatrix}^T\right) = sigmoid(10) \approx 1$$

Prediction:

$$y = sigmoid\left(\begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}^T \begin{bmatrix} 10 \\ -20 \\ -20 \end{bmatrix}\right) = sigmoid(-10) \approx 0$$

# MLP Example: XOR function



$$w^{(1)} = \begin{bmatrix} -30 & 10 \\ 20 & -20 \\ 20 & -20 \end{bmatrix}$$

$$w^{(2)} = [10, -20, -20]$$

| X₁ | X₂ | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

**Prediction for x=[1,0,0]** $\quad \bar{y} = sigmoid\left(\begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}^T \begin{bmatrix} 10 \\ -20 \\ -20 \end{bmatrix}\right) = sigmoid(-10) \approx 0$

**Prediction for x=[1,0,1]** $\quad \bar{y} = sigmoid(10) \approx 1$

**Prediction for x=[1,1,0]** $\quad \bar{y} = sigmoid(10) \approx 1$

**Prediction for x=[1,1,1]** $\quad \bar{y} = sigmoid(-10) \approx 0$

# MLP for Multi-class Classification: 1 sv. All

- Example: Hand-written digits dataset
  - Each image input is 28x28 = 784 dims
  - One or many hidden layers
  - Output layer: 10 neurons (i.e. each correspond to one logistic regression classifier)
  - Using one-hot vector for ground-truth output

digit **2** has output

$$y = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

digit 8 has output

$$y = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

hidden layer
($n = 15$ neurons)

output layer

input layer
(784 neurons)

$\rightarrow 0$
$\rightarrow 1$
$\rightarrow 2$
$\rightarrow 3$
$\rightarrow 4$
$\rightarrow 5$
$\rightarrow 6$
$\rightarrow 7$
$\rightarrow 8$
$\rightarrow 9$

# MLP for Multi-class Classification: 1 sv. All
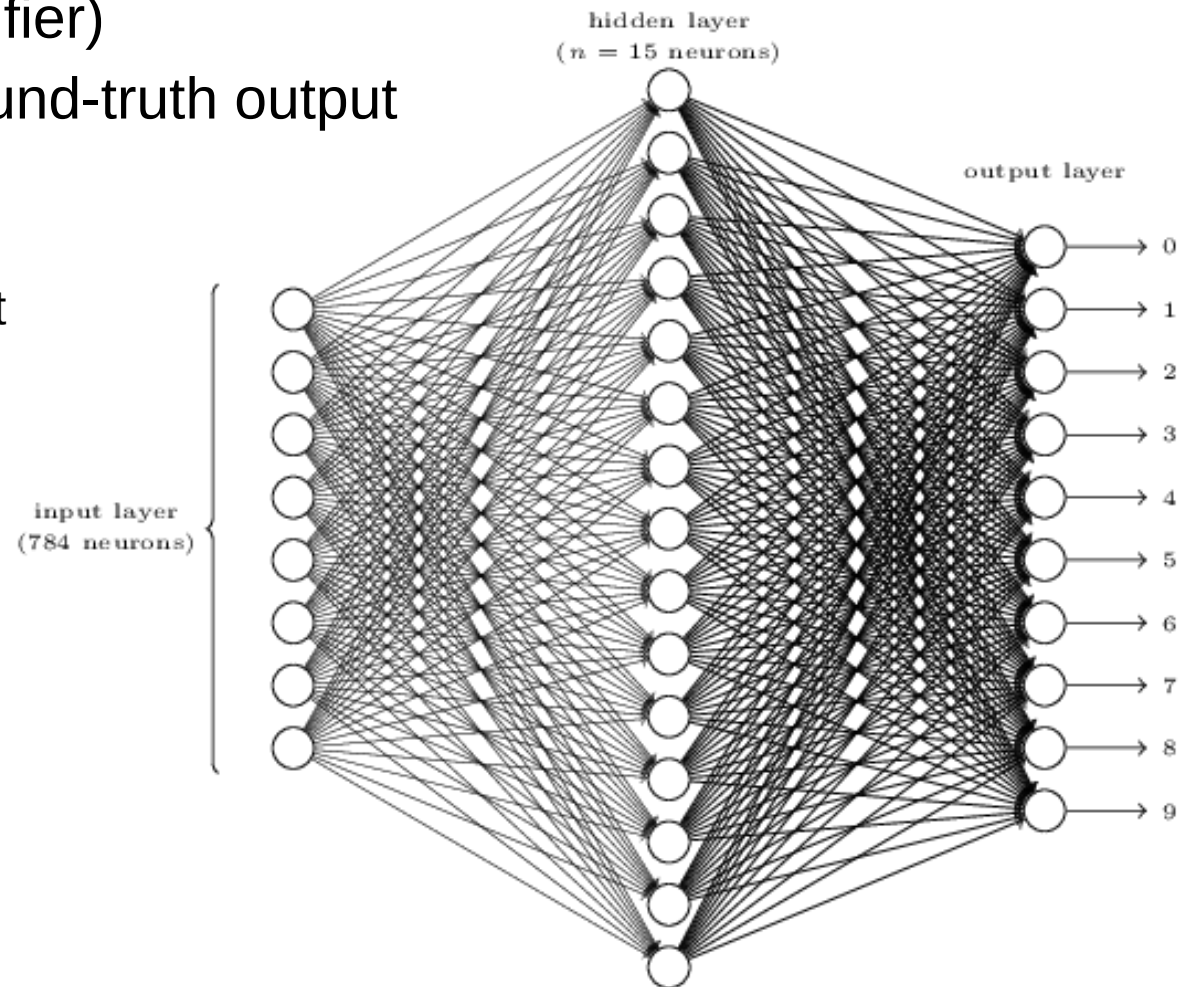
- Example: Hand-written digits dataset
  - Each image input is 28x28 = 784 dims
  - One or many hidden layers
  - Output layer: 10 neurons (i.e. each correspond to one logistic regression classifier if using sigmoid activation function)
  - Using one-hot vector for ground-truth output

digit **2** has output   **prediction**

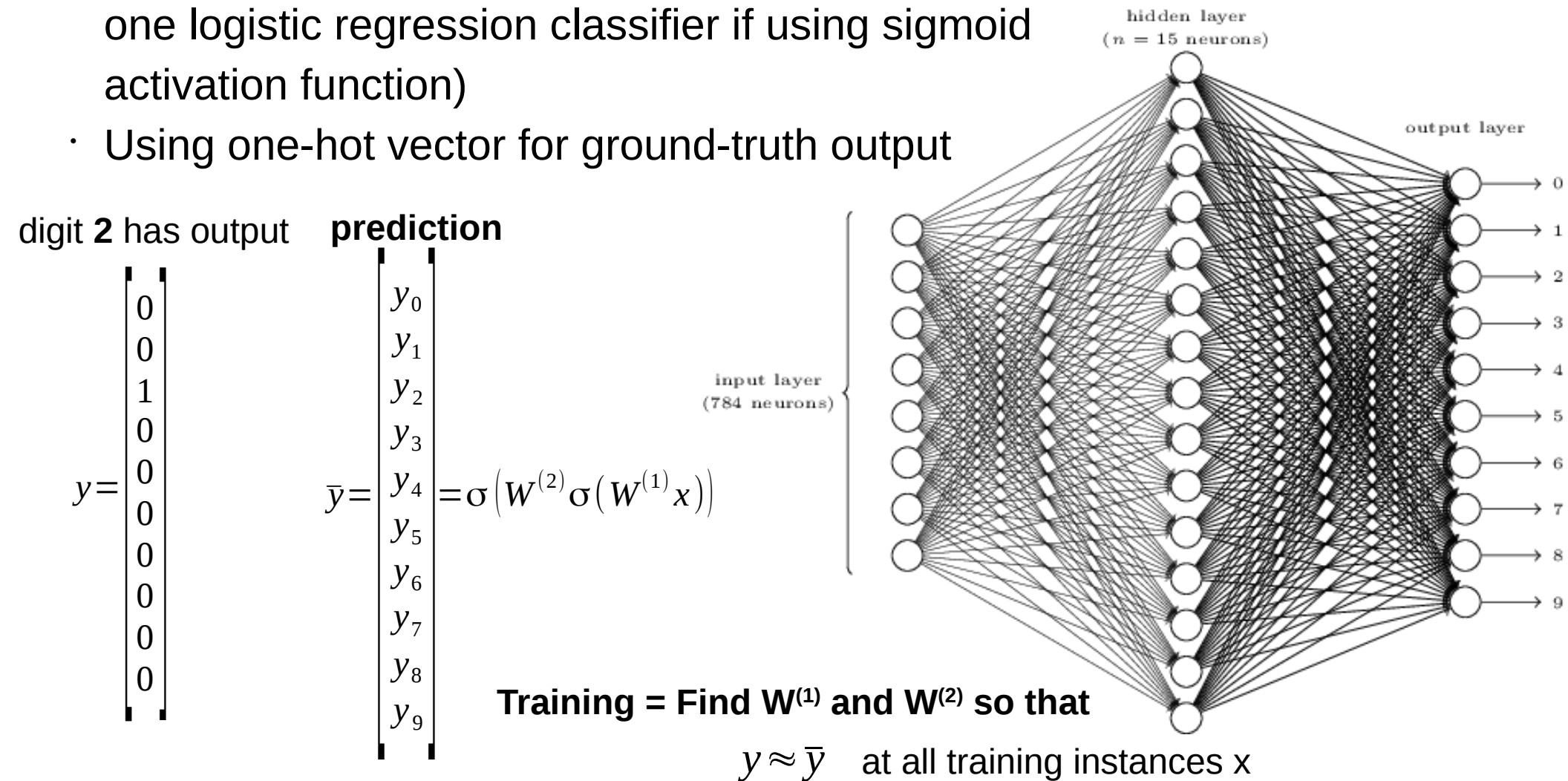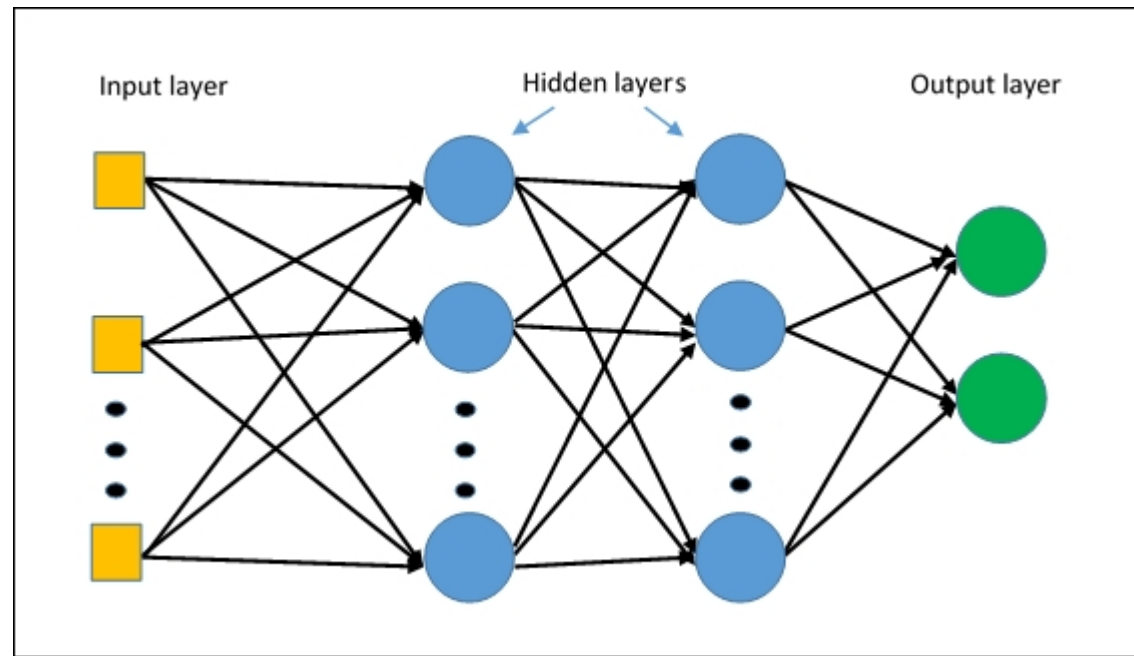$$y = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\bar{y} = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \\ y_7 \\ y_8 \\ y_9 \end{bmatrix} = \sigma\left(W^{(2)}\sigma\left(W^{(1)}x\right)\right)$$

hidden layer
($n = 15$ neurons)

output layer

input layer
(784 neurons)

→ 0
→ 1
→ 2
→ 3
→ 4
→ 5
→ 6
→ 7
→ 8
→ 9

**Training = Find $W^{(1)}$ and $W^{(2)}$ so that**

$$y \approx \bar{y} \quad \text{at all training instances x}$$

# MLP General Cases

- A multi-layer perceptron with K-1 hidden layers (deep neural network)
  - ➢ Activation function at layer j is $\sigma^{(j)}$
  - ➢ Matrix $W^{(j)}$ is the parameter matrix mapping from layer j to layer j+1
- Prediction is

$$y = \sigma^{(K)}\left(W^{(K-1)}\sigma^{(K-1)}\left(\ldots\sigma^{(2)}\left(W^{(1)}x\right)\ldots\right)\right)$$



Next lecture: Backpropagation method to train a neural network