

Feed-forward Neural Networks

Lecture 8 - DAMLF | ML1



Motivation: Non-Linear Hypotheses

Random Forests

Very good out-of-the-box performance.

However, random forests presume a set of inherent, well-defined features.

Reasonable in many business & finance applications, which is why they are popular.

What if you have a problem in which the inherent features are **not** well-defined?

Motivation: Computer Vision

What you see



Motivation: Computer Vision

What you see



What the computer sees[†]



Motivation: Computer Vision

What you see



What the computer sees[†]

243	232	230	223	221	33	42	68	78	88	28	43	21	19	23
245	229	223	227	210	40	31	59	41	74	23	39	24	17	19
241	235	239	229	225	39	38	53	48	63	31	28	19	18	21
243	230	229	231	240	51	36	38	71	84	20	27	16	21	25
238	232	230	240	231	43	70	29	78	86	27	32	17	18	29

[†] a 5×5 pixel patch of the red ball has $25 \times 3 = 75$ pixels (features).

Motivation: Computer Vision

What you see



What the computer sees

640×480 pixels = **921,600** Features

1200×1800 pixels = **6,480,000** Features % 10×15 cm image, 300 dpi

Dogs



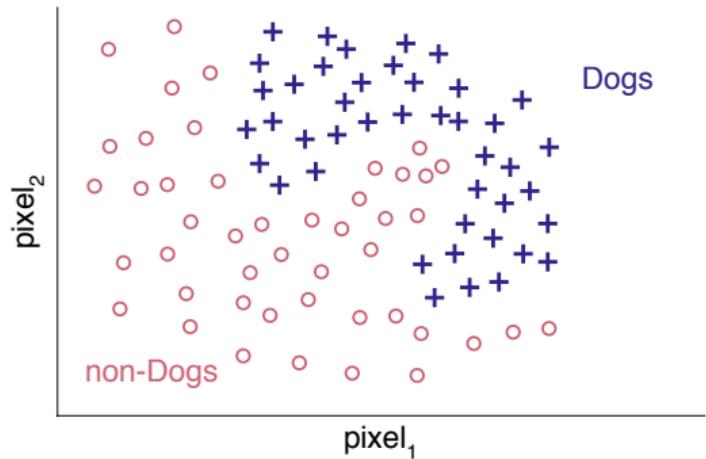
Dogs



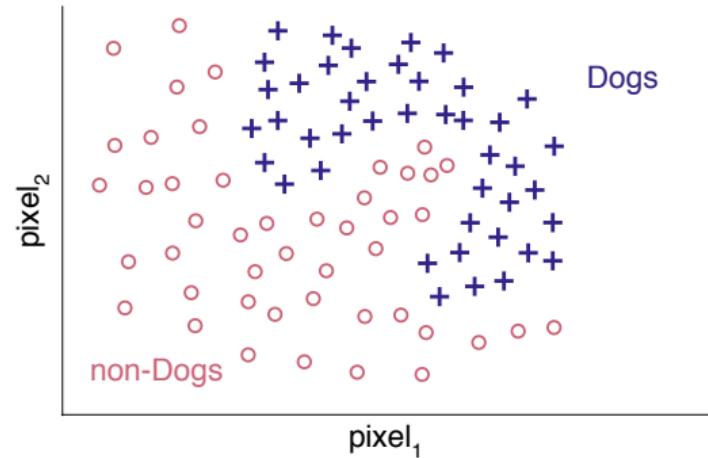
non-Dogs



Classification



Classification



- now simply add 1 million (or more) pixels
- and no idea about any inherent/general features

New Features

Natural language processing (NLP) is a straightforward example. ‘Reading’ all of the financial reports and extracting information automatically is an obvious application.

Image and Language Processing

New Features

Natural language processing (NLP) is a straightforward example. 'Reading' all of the financial reports and extracting information automatically is an obvious application.

Imagine processing may be less obvious.

Image and Language Processing



New Features

Natural language processing (NLP) is a straightforward example. 'Reading' all of the financial reports and extracting information automatically is an obvious application.

Imagine processing may be less obvious.

- Cars in retail parking lots

Image and Language Processing



Copyright © Birdi Ltd. - Image © 2018 DigitalGlobe, Inc.

New Features

Natural language processing (NLP) is a straightforward example. 'Reading' all of the financial reports and extracting information automatically is an obvious application.

Imagine processing may be less obvious.

- Cars in retail parking lots
- Shadows of oil storage tanks
(Ras Tunura refinery, Saudi Arabia)

Image and Language Processing



New Features

Natural language processing (NLP) is a straightforward example. 'Reading' all of the financial reports and extracting information automatically is an obvious application.

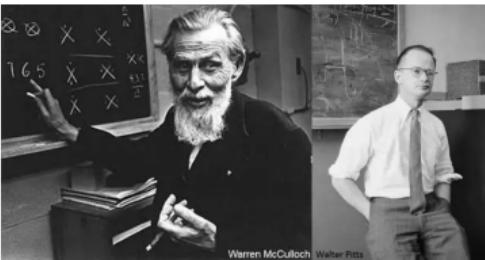
Imagine processing may be less obvious.

- Cars in retail parking lots
- Shadows of oil storage tanks
(Ras Tunura refinery, Saudi Arabia)

Applications will increase with IoT and 5G

Neural Networks

Some History

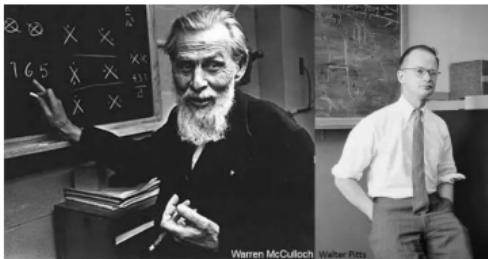


WARREN McCULLOCH & WALTER PITTS

A logical calculus of the ideas immanent in nervous activity, *Bul Math Biophys* 5:115–133, 1943

Neural networks refer to algorithms that **mimic the brain**.

Some History



WARREN McCULLOCH & WALTER PITTS

A logical calculus of the ideas immanent in nervous activity, *Bul Math Biophys* 5:115–133, 1943

Neural networks refer to algorithms that **mimic the brain**.

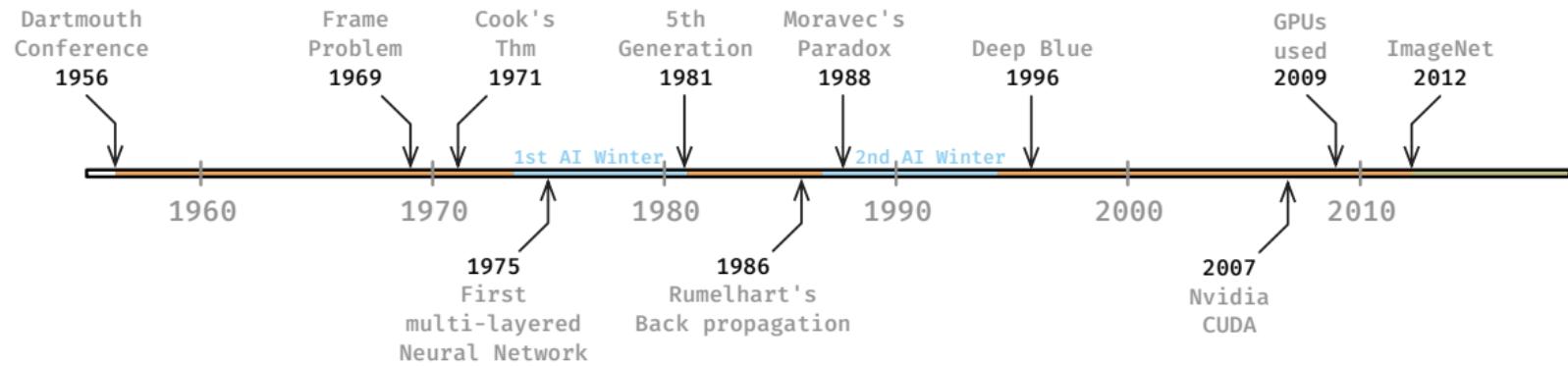
Two approaches:

- Biological process models
- Artificial neural networks (ANNs)

The One Algorithm Hypothesis

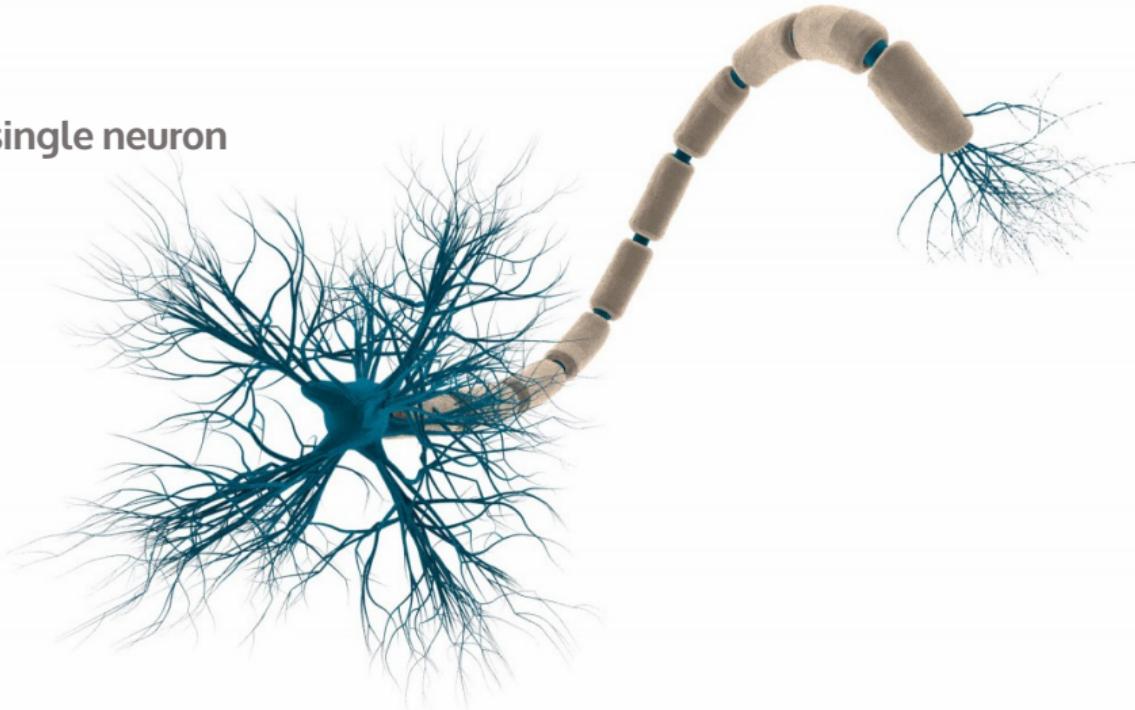
*Rather than 1000s of individual algorithms for different cognitive functions, perhaps the brain uses **just one**.*

Potted History of AI

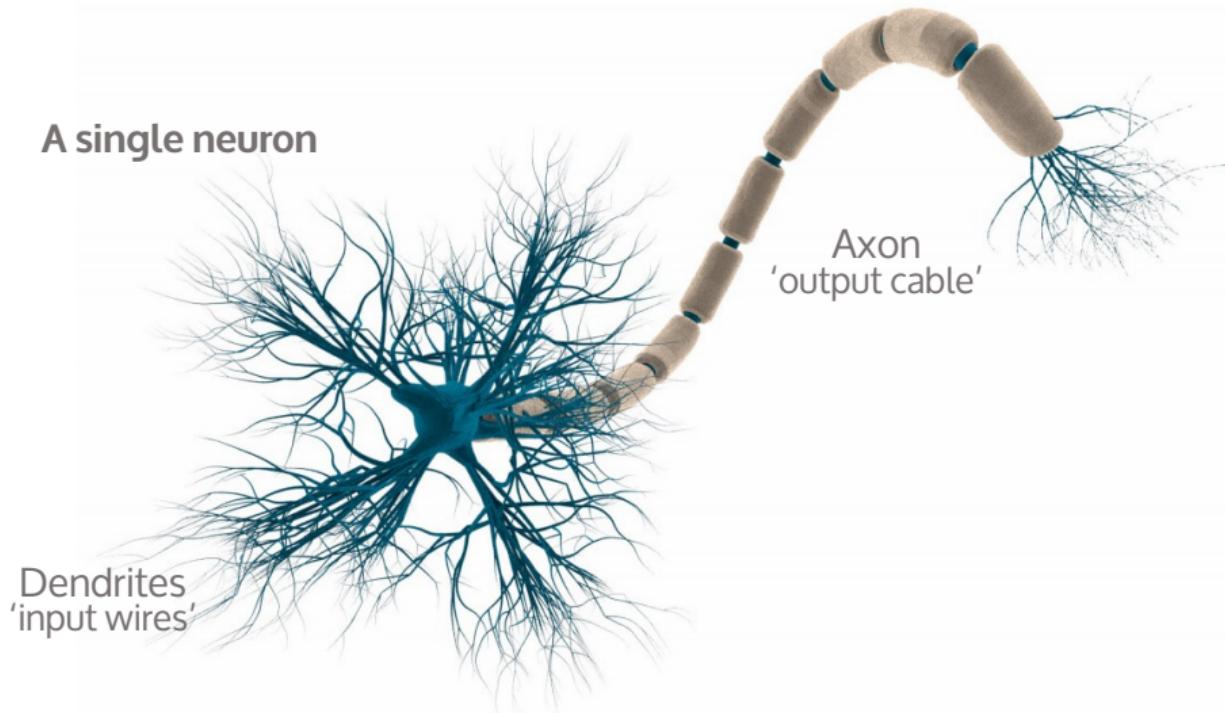


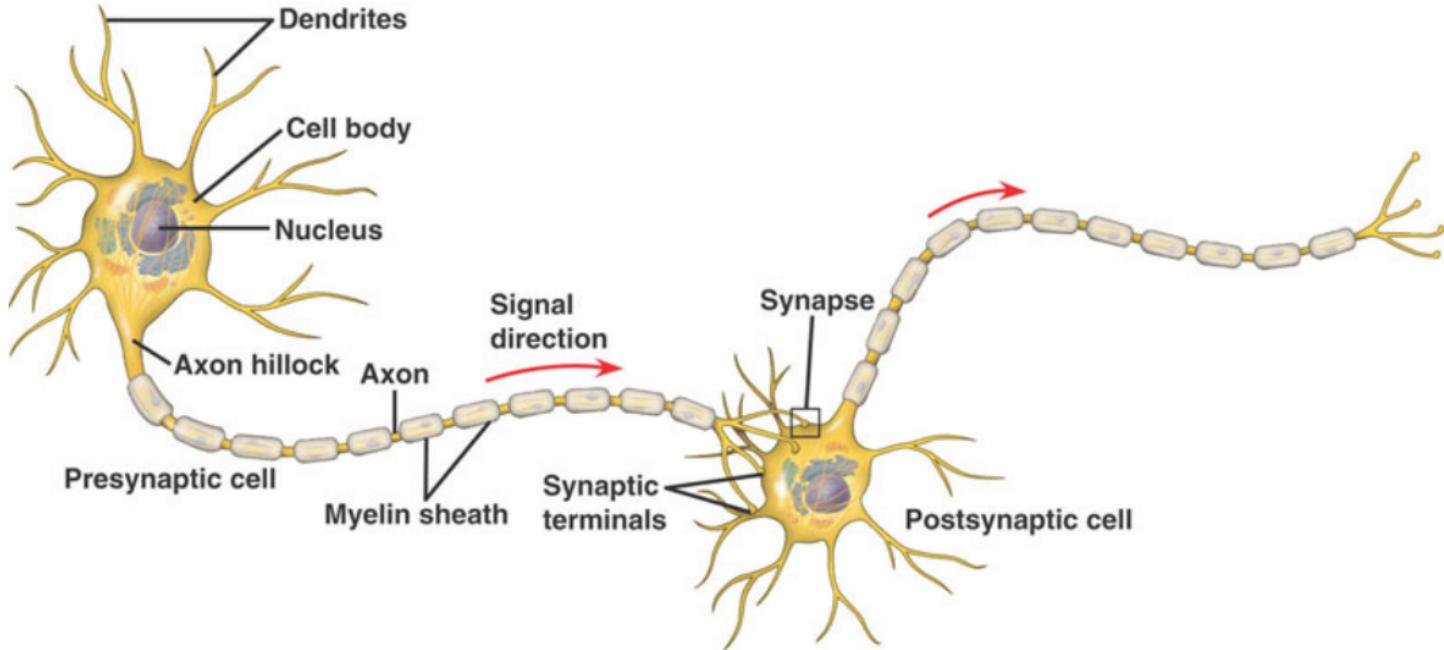
Neural Networks: Hypothesis Representation

A single neuron

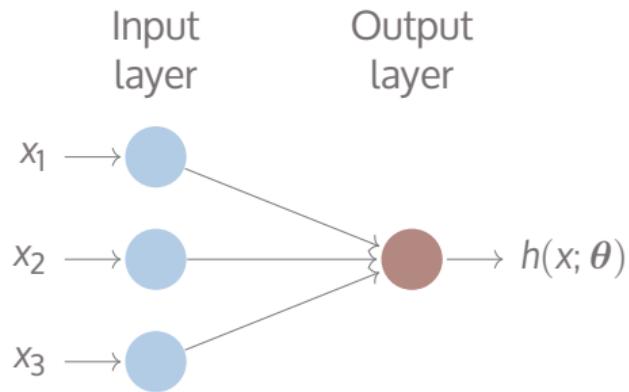


A single neuron

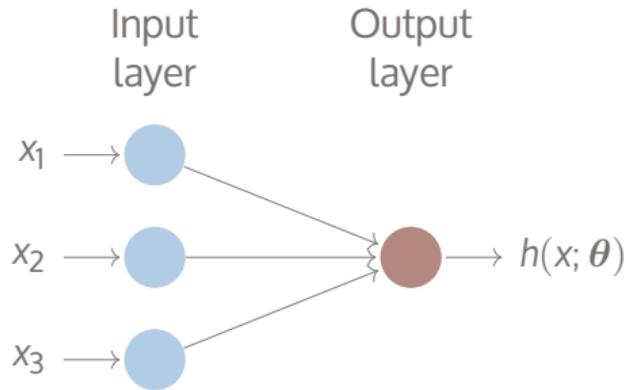




Neuron model: logistic unit



Neuron model: logistic unit

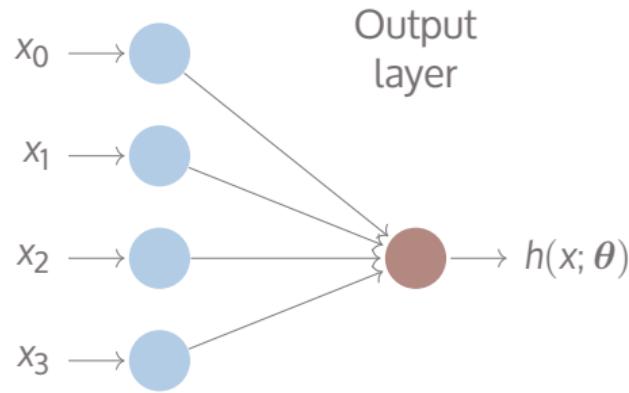


Recall:

$$h(x; \theta) = \frac{1}{1 + e^{-\theta^T x}}$$

where: $x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$, $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$

Neuron model: logistic unit



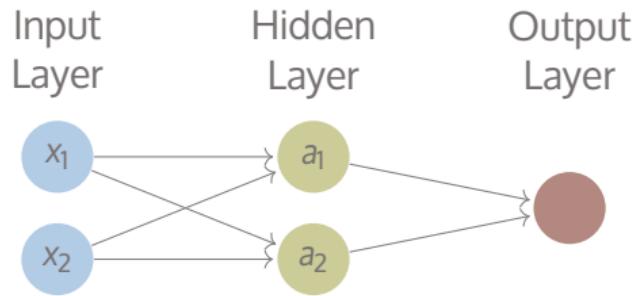
Note:

$$x_0 = 1$$

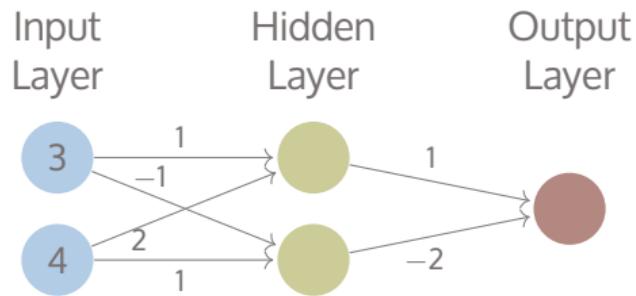
called the '**bias unit**'

where: $x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$

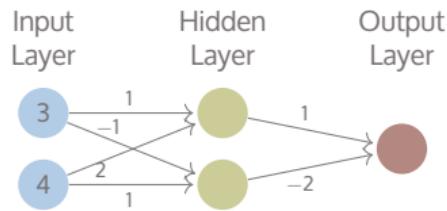
Adding Layers



Input Data, Output Prediction



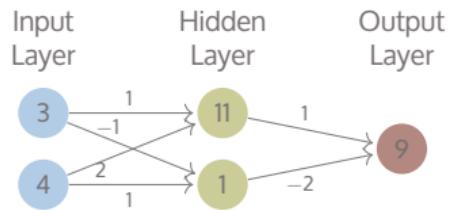
Forward Propagation Intuition



FORWARD PROPAGATION

- Input Data
- Dot product *multiply then add*
- One-data-point at a time
- Last output is the model's prediction

Forward Propagation Intuition



FORWARD PROPAGATION

```
In [1]: import numpy as np
```

```
In [2]: # Forward Propagation with
# Identity Activation Function
input_nodes = np.array([3,4])

Theta = {'node_a1': np.array([1, 2]),
         'node_a2': np.array([-1, 1]),
         'output': np.array([1, -2])}

node_a1_value = (input_nodes * Theta['node_a1']).sum()
node_a2_value = (input_nodes * Theta['node_a2']).sum()
```

```
In [3]: hidden_layer_values = np.array([node_a1_value, node_a2_value])
print(hidden_layer_values)

[11  1]
```

```
In [4]: output = (hidden_layer_values * Theta['output']).sum()
print(output)
```

Activation Function

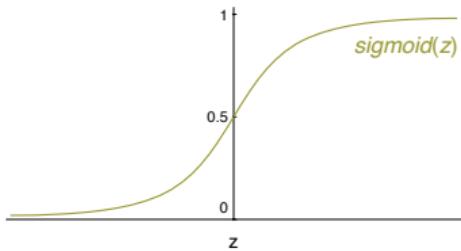
NON-LINEARITY

Activation functions:

- applied to **node inputs** to produce **node output**
- capture **non-linear relationships** in data

Examples:

Activation Function



NON-LINEARITY

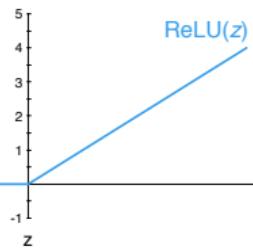
Activation functions:

- applied to **node inputs** to produce **node output**
- capture **non-linear relationships** in data

Examples:

Sigmoid Function

Activation Function



NON-LINEARITY

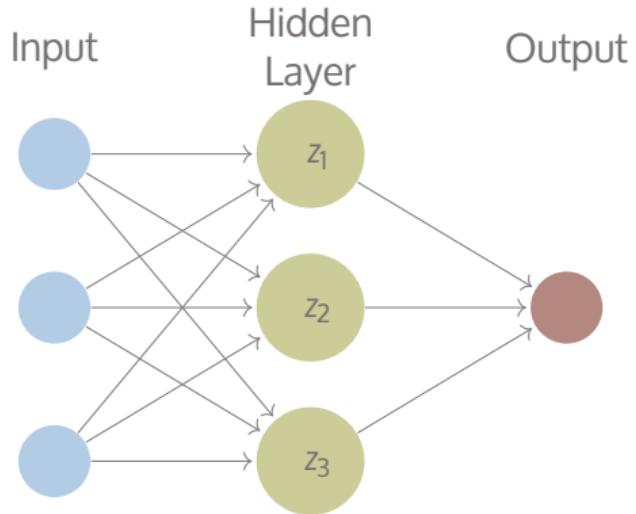
Activation functions:

- applied to **node inputs** to produce **node output**
- capture **non-linear relationships** in data

Examples:

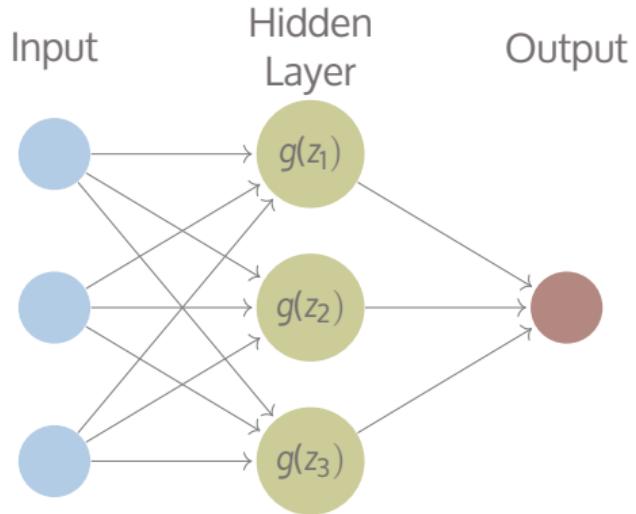
Rectified Linear Unit (ReLU)

Activation Function



$z_1 = \text{node_a1_value}$
 $z_2 = \text{node_a2_value}$
 $z_3 = \text{node_a3_value}$

Activation Function



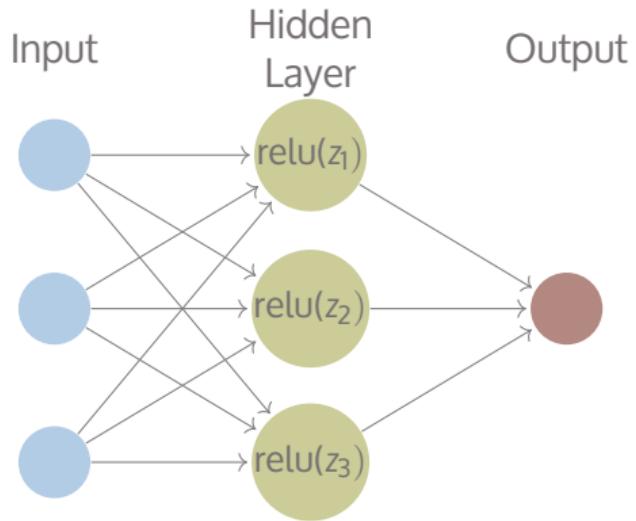
$z_1 = \text{node_a1_value}$

$z_2 = \text{node_a2_value}$

$z_3 = \text{node_a3_value}$

$g(z_i) = \text{sigmoid activation}$

Activation Function



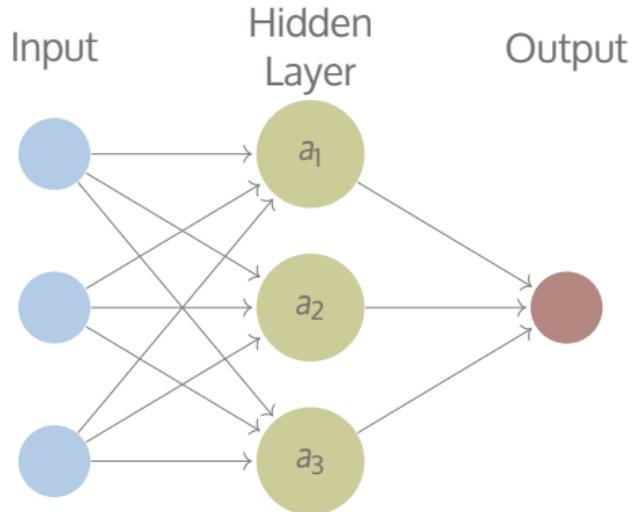
$$z_1 = \text{node_a1_value}$$

$$z_2 = \text{node_a2_value}$$

$$z_3 = \text{node_a3_value}$$

$$\text{relu}(z_i) = \text{ReLU}$$

Activation Function



$$z_1 = \text{node_}a_1\text{_value}$$

$$z_2 = \text{node_}a_2\text{_value}$$

$$z_3 = \text{node_}a_3\text{_value}$$

$$a_i = \text{activation_fun}(z_i)$$

Terminology

Features: x

Inputs: x

Parameters: θ

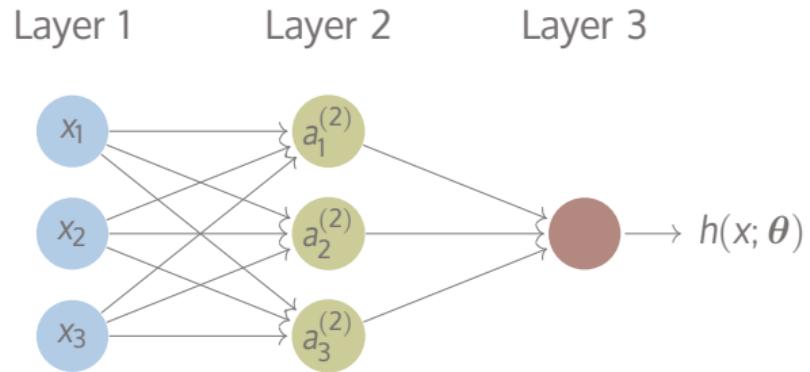
Weights: θ

Activation Function

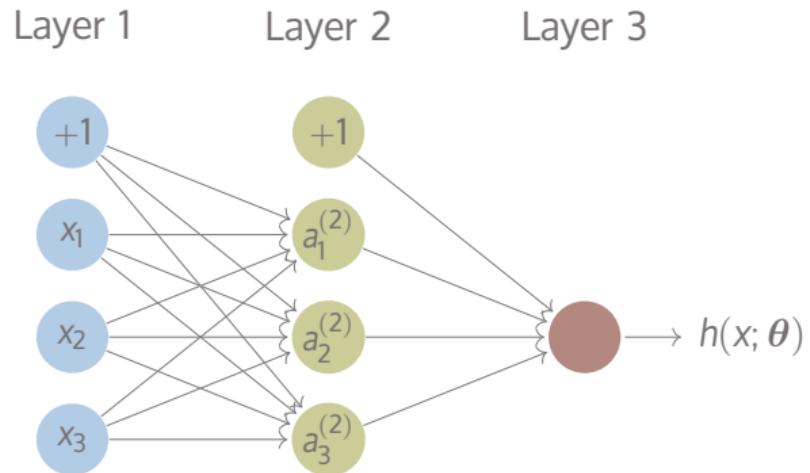
Some types of activation functions:

- Logistic regression hypothesis (Sigmoid)
- Hyperbolic Tangent (\tanh) *# rescaled sigmoid to (-1, +1)*
- Rectified Linear Unit (ReLU)
- Normalized Exponential Function (Softmax)

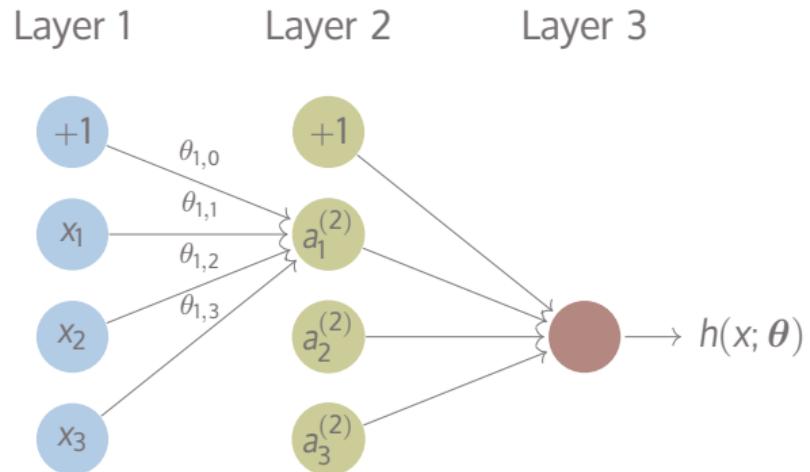
Neural Network Bias Nodes



Neural Network Bias Nodes



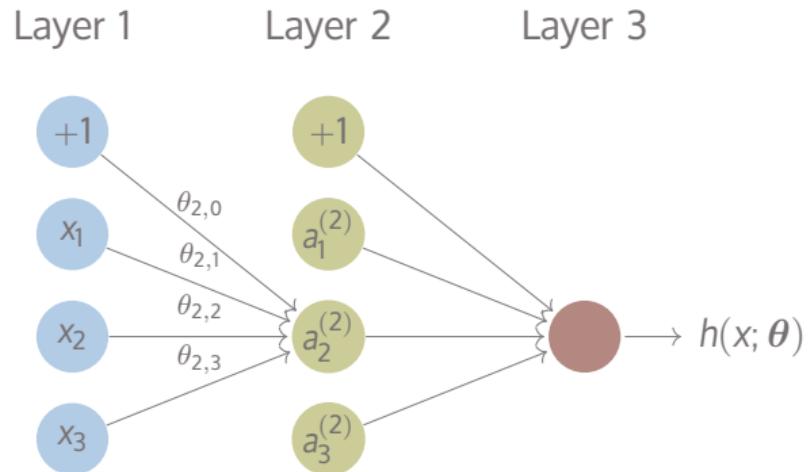
Neural Network Bias Nodes



Highlighted Step:

$$\theta_{1,0} + \theta_{1,1}x_1 + \theta_{1,2}x_2 + \theta_{1,3}x_3 = z_1^{(2)} \Rightarrow g(z_1^{(2)}) = a_1^{(2)}$$

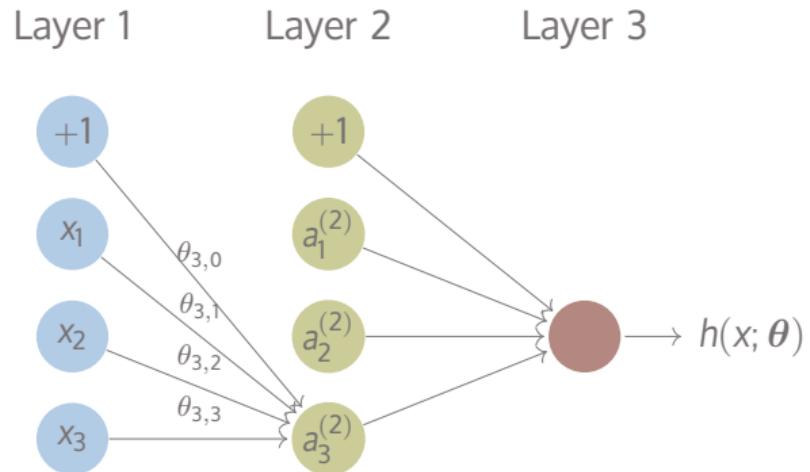
Neural Network Bias Nodes



Highlighted Step:

$$\theta_{2,0} + \theta_{2,1}x_1 + \theta_{2,2}x_2 + \theta_{2,3}x_3 = z_2^{(2)} \Rightarrow g(z_2^{(2)}) = a_2^{(2)}$$

Neural Network Bias Nodes



Highlighted Step:

$$\theta_{3,0} + \theta_{3,1}x_1 + \theta_{3,2}x_2 + \theta_{3,3}x_3 = z_3^{(2)} \Rightarrow g(z_3^{(2)}) = a_3^{(2)}$$

Comparison to Logistic Regression

Recall the form of the hypothesis, $h(x; \theta)$, for **logistic regression**.

For example,

$$\begin{aligned} h(x; \theta) &= g(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2) \\ &= g(\theta^T x) \end{aligned}$$

with **parameter** vector θ and **feature** vector x :

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Comparison to Logistic Regression

Recall the form of the hypothesis, $h(x; \theta)$, for **logistic regression**.

For example,

$$\begin{aligned} h(x; \theta) &= g(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2) \\ &= g(\theta^T x) \end{aligned}$$

with **parameter** vector θ and **feature** vector x :

$$\theta^T = [\theta_0 \quad \theta_1 \quad \theta_2] \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Comparison to Logistic Regression

$$\theta^T = [\theta_0 \quad \theta_1 \quad \theta_2] \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Observe that $\theta^T = [\theta_0 \quad \theta_1 \quad \theta_2]$ is a 1×3 **matrix**.

Comparison to Logistic Regression

$$\theta^T = [\theta_0 \quad \theta_1 \quad \theta_2] \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Observe that $\theta^T = [\theta_0 \quad \theta_1 \quad \theta_2]$ is a 1×3 **matrix**.

Let's call it $\Theta = [\theta_0 \quad \theta_1 \quad \theta_2]$.

Comparison to Logistic Regression

We can also relabel the elements of Θ using the convention $\Theta_{i,j}$ to refer to the element in the i th row and j th column of Θ .

$$\Theta = \begin{bmatrix} \Theta_{10} & \Theta_{11} & \Theta_{12} \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Comparison to Logistic Regression

We can also relabel the elements of Θ using the convention $\Theta_{i,j}$ to refer to the element in the i th row and j th column of Θ .

$$\Theta = \begin{bmatrix} \Theta_{10} & \Theta_{11} & \Theta_{12} \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

So, our familiar logistic regression hypothesis can be rewritten (baroquely) as

$$\begin{aligned} h(x; \theta) &= g(\Theta_{10}x_0 + \Theta_{11}x_1 + \Theta_{12}x_2) \\ &= g(\Theta x) \end{aligned}$$

Comparison to Logistic Regression

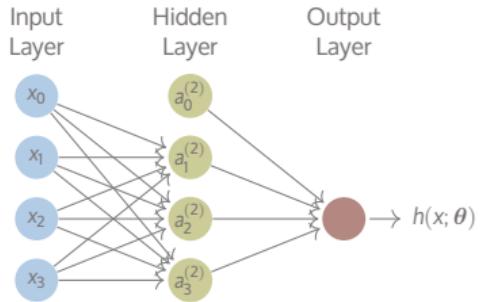
We can also relabel the elements of Θ using the convention $\Theta_{i,j}$ to refer to the element in the i th row and j th column of Θ .

$$\Theta^{(1)} = \begin{bmatrix} \Theta_{10}^{(1)} x_0 & \Theta_{11}^{(1)} & \Theta_{12}^{(1)} \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

So, our familiar logistic regression hypothesis can be rewritten (baroquely) as

$$\begin{aligned} h(x; \theta) &= g\left(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2\right) \\ &= g\left(\Theta^{(1)} x\right) \end{aligned}$$

A Neural Network

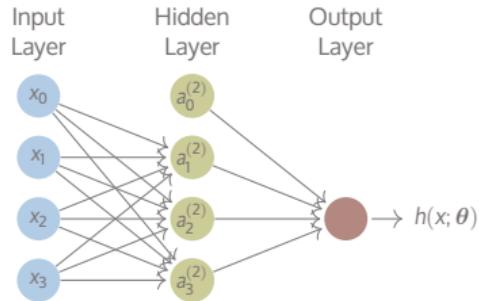


Notation:

$a_i^{(j)}$ = activation of unit i in layer j

$\Theta^{(j)}$ = matrix of weights applied to
the function mapping from
layer j to layer $j + 1$

A Neural Network



Notation:

$a_i^{(j)}$ = activation of unit i in layer j

$\Theta^{(j)}$ = matrix of weights applied to the function mapping from layer j to layer $j + 1$

Example:

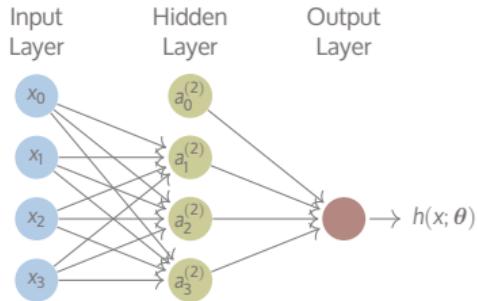
3 input units and 3 hidden units entails that $\Theta^{(1)}$ is a 3×4 matrix

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

A Neural Network



Notation:

$a_i^{(j)}$ = activation of unit i in layer j
 $\Theta^{(j)}$ = matrix of weights applied to the function mapping from layer j to layer $j + 1$

Example:

3 input units and 3 hidden units entails that $\Theta^{(1)}$ is a 3×4 matrix

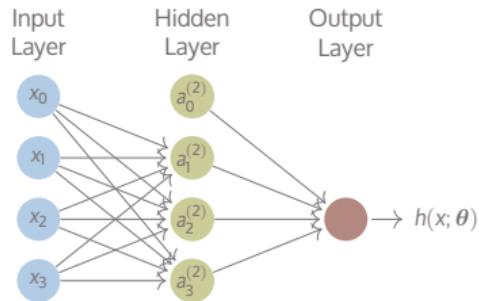
$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

In general, if layer j has s_j units and layer $j + 1$ has s_{j+1} units, then $\Theta^{(j)}$ has dimension $s_{j+1} \times (s_j + 1)$.

A Neural Network



Notation:

$a_i^{(j)}$ = activation of unit i in layer j
 $\Theta^{(j)}$ = matrix of weights applied to the function mapping from layer j to layer $j + 1$

Example:

3 hidden units and 1 output unit entails that $\Theta^{(2)}$ is a 1×4 matrix

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

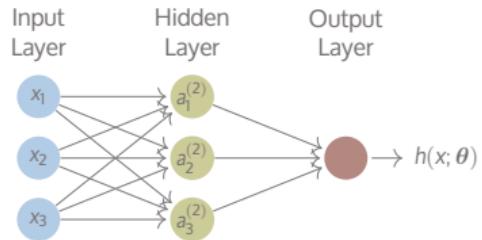
$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

$$a_1^{(3)} = g\left(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)}\right) = h(x; \theta)$$

The **vectorized** implementation of these calculations is called **forward propagation**.

Forward propagation computes the value of the output layer, $h(x; \theta)$, in terms of the values of inputs, $x_0 \dots x_n$, propagated through the layers of a neural network.

Vectorized Implementation



$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

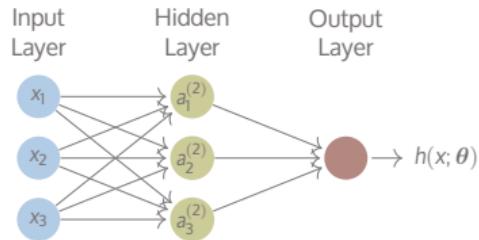
$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

$$a_1^{(2)} = g\left(z_1^{(2)}\right)$$

$$a_2^{(2)} = g\left(z_2^{(2)}\right)$$

$$a_3^{(2)} = g\left(z_3^{(2)}\right)$$

Vectorized Implementation



Notation:

$$z_1^{(2)} = \Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3$$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

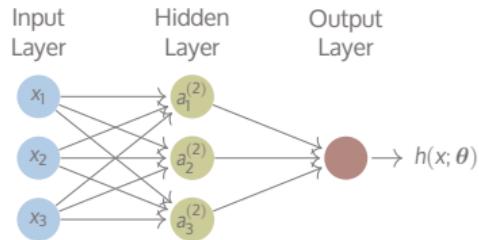
$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

$$a_1^{(2)} = g\left(z_1^{(2)}\right)$$

$$a_2^{(2)} = g\left(z_2^{(2)}\right)$$

$$a_3^{(2)} = g\left(z_3^{(2)}\right)$$

Vectorized Implementation



Notation:

$$z_1^{(2)} = \Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3$$

$$z_2^{(2)} = \Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3$$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

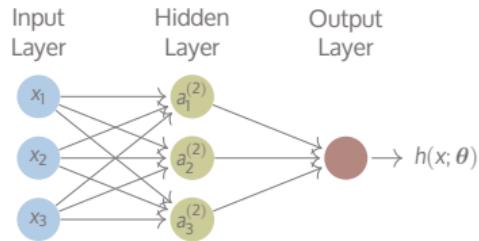
$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

$$a_1^{(2)} = g\left(z_1^{(2)}\right)$$

$$a_2^{(2)} = g\left(z_2^{(2)}\right)$$

$$a_3^{(2)} = g\left(z_3^{(2)}\right)$$

Vectorized Implementation



Notation:

$$z_1^{(2)} = \Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3$$

$$z_2^{(2)} = \Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3$$

$$z_3^{(2)} = \Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3$$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

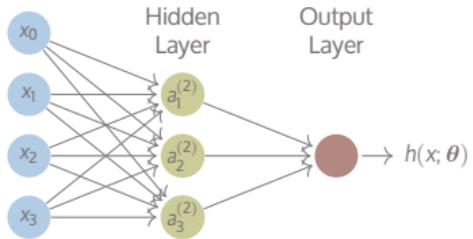
$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

$$a_1^{(2)} = g\left(z_1^{(2)}\right)$$

$$a_2^{(2)} = g\left(z_2^{(2)}\right)$$

$$a_3^{(2)} = g\left(z_3^{(2)}\right)$$

Vectorized Implementation



Notation:

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

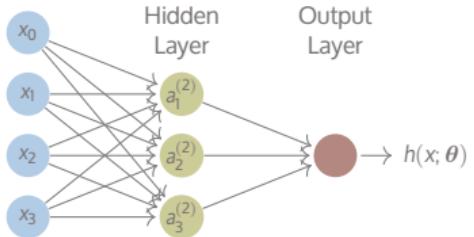
$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

$$a_1^{(2)} = g\left(z_1^{(2)}\right)$$

$$a_2^{(2)} = g\left(z_2^{(2)}\right)$$

$$a_3^{(2)} = g\left(z_3^{(2)}\right)$$

Vectorized Implementation



Notation:

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

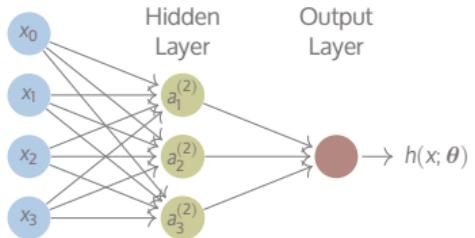
$$z^{(2)} = \Theta^{(1)}x$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

$$a^{(2)} = g(z^{(2)})$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

Vectorized Implementation



Notation:

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

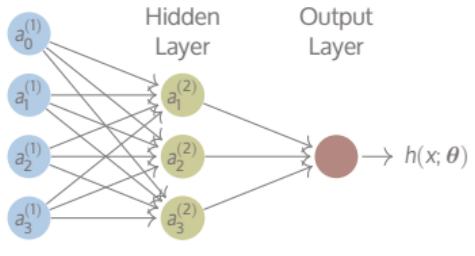
$$z^{(2)} = \Theta^{(1)}x$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

$$a^{(2)} = g(z^{(2)})$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

Vectorized Implementation



Notation:

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad a^{(1)} = \begin{bmatrix} a_0^{(1)} \\ a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \end{bmatrix}$$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

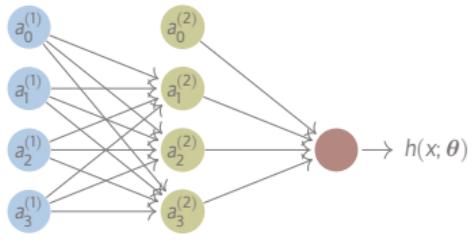
$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

Vectorized Implementation



Notation:

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad a^{(1)} = \begin{bmatrix} a_0^{(1)} \\ a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \end{bmatrix}$$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right)$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right)$$

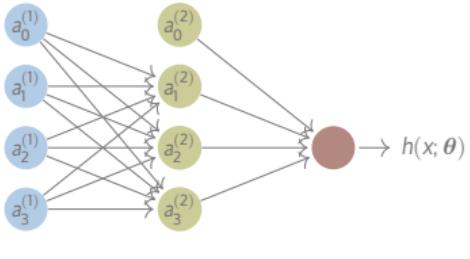
$$h(x; \theta) = g\left(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)}\right)$$

$$z^{(2)} = \Theta^{(1)}a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

$$a_0^{(2)} = 1 = a_0^{(1)}$$

Forward Propagation



Equations to plug in:

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad a^{(1)} = \begin{bmatrix} a_0^{(1)} \\ a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \end{bmatrix}$$

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)} a_0^{(1)} + \Theta_{11}^{(1)} a_1^{(1)} + \Theta_{12}^{(1)} a_2^{(1)} + \Theta_{13}^{(1)} a_3^{(1)}\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)} a_0^{(1)} + \Theta_{21}^{(1)} a_1^{(1)} + \Theta_{22}^{(1)} a_2^{(1)} + \Theta_{23}^{(1)} a_3^{(1)}\right)$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)} a_0^{(1)} + \Theta_{31}^{(1)} a_1^{(1)} + \Theta_{32}^{(1)} a_2^{(1)} + \Theta_{33}^{(1)} a_3^{(1)}\right)$$

$$h(x; \theta) = g\left(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)}\right)$$

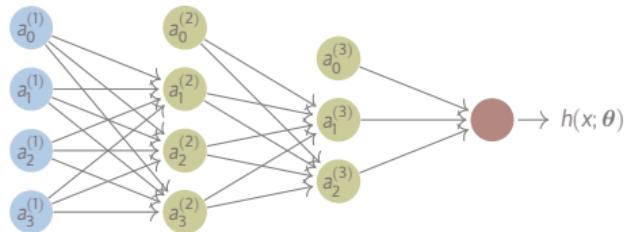
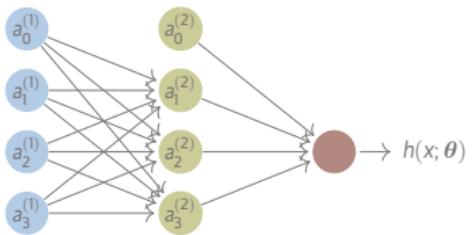
$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = h(x; \theta) = g(z^{(3)})$$

Adding Hidden Layers



References

Nielsen, M. (2019).
Neural Networks and Deep Learning.
Determination Press.