

Deep Learning

Big Data & Machine Learning Bootcamp - Keep Coding



Outline

1. Intro to Neural Networks
2. Binary Classification and Logistic Regression
3. Gradient Descent
4. Vectorization
5. Activation functions
6. Backpropagation
7. Random initialization



ARTIFICIAL INTELLIGENCE

IS NOT NEW

ARTIFICIAL INTELLIGENCE

Any technique which enables computers to mimic human behavior



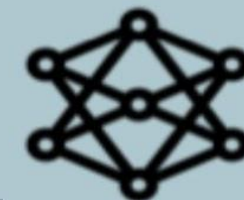
MACHINE LEARNING

AI techniques that give computers the ability to learn without being explicitly programmed to do so



DEEP LEARNING

A subset of ML which make the computation of multi-layer neural networks feasible



1950's

1960's

1970's

1980's

1990's

2000's

2010's

ORACLE

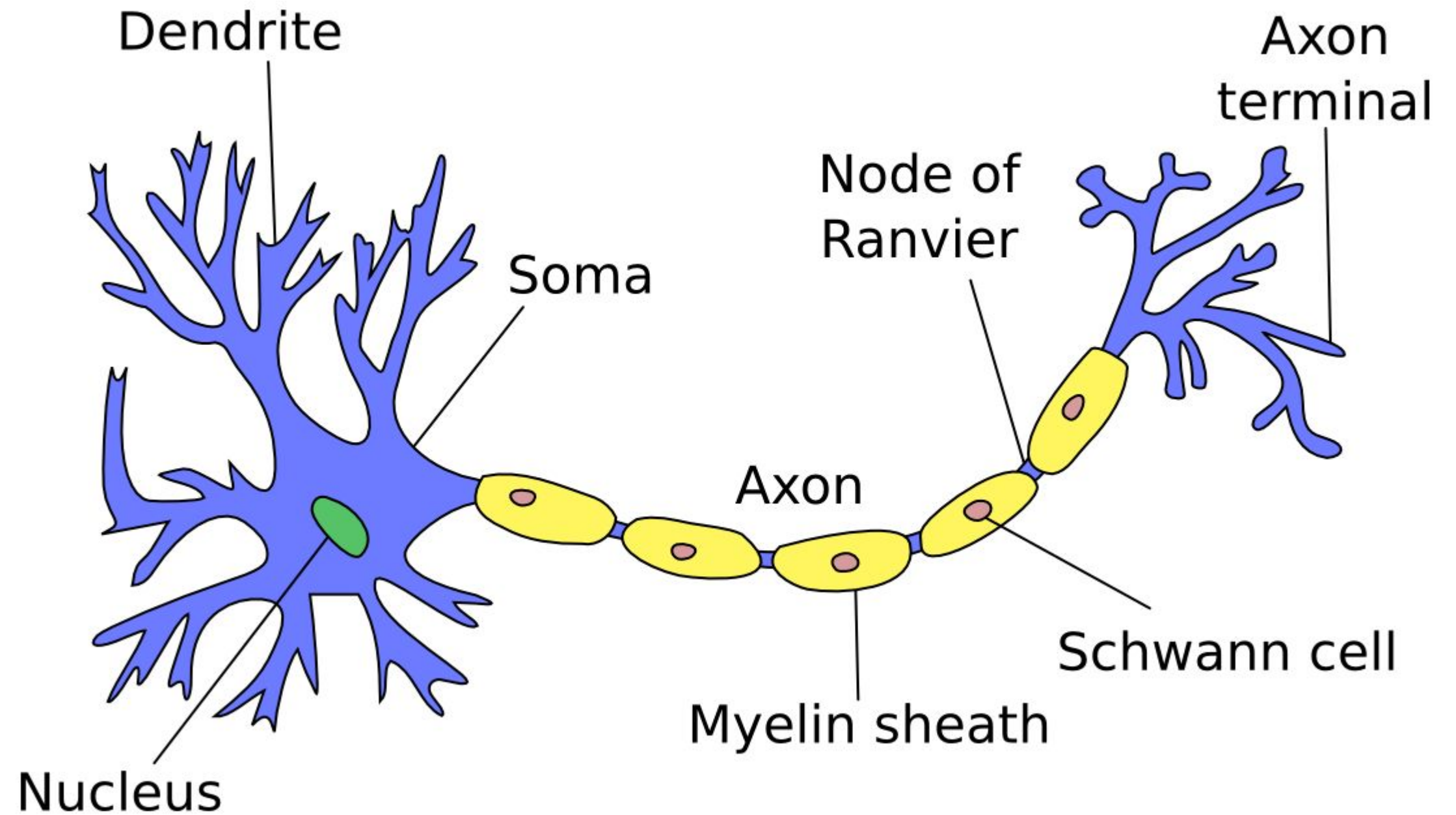
Copyright © 2018, Oracle and/or its affiliates. All rights reserved. |

Fuente: Oracle



1. Intro to Neural Networks

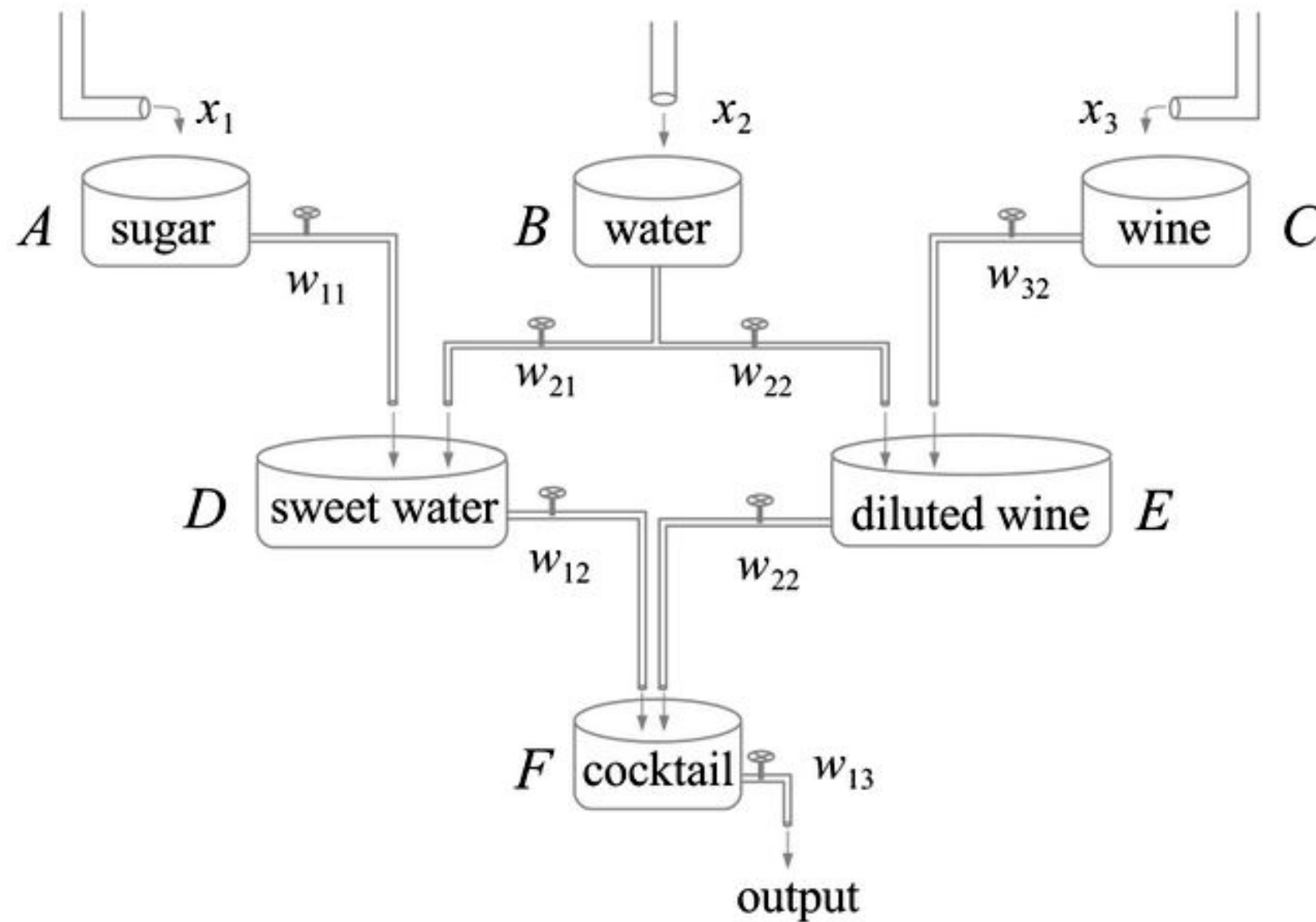
Source of inspiration



Source: <https://www.marekrei.com/blog/neural-networks-part-2-the-neuron/>

1. Intro to Neural Networks

Examples of neural networks

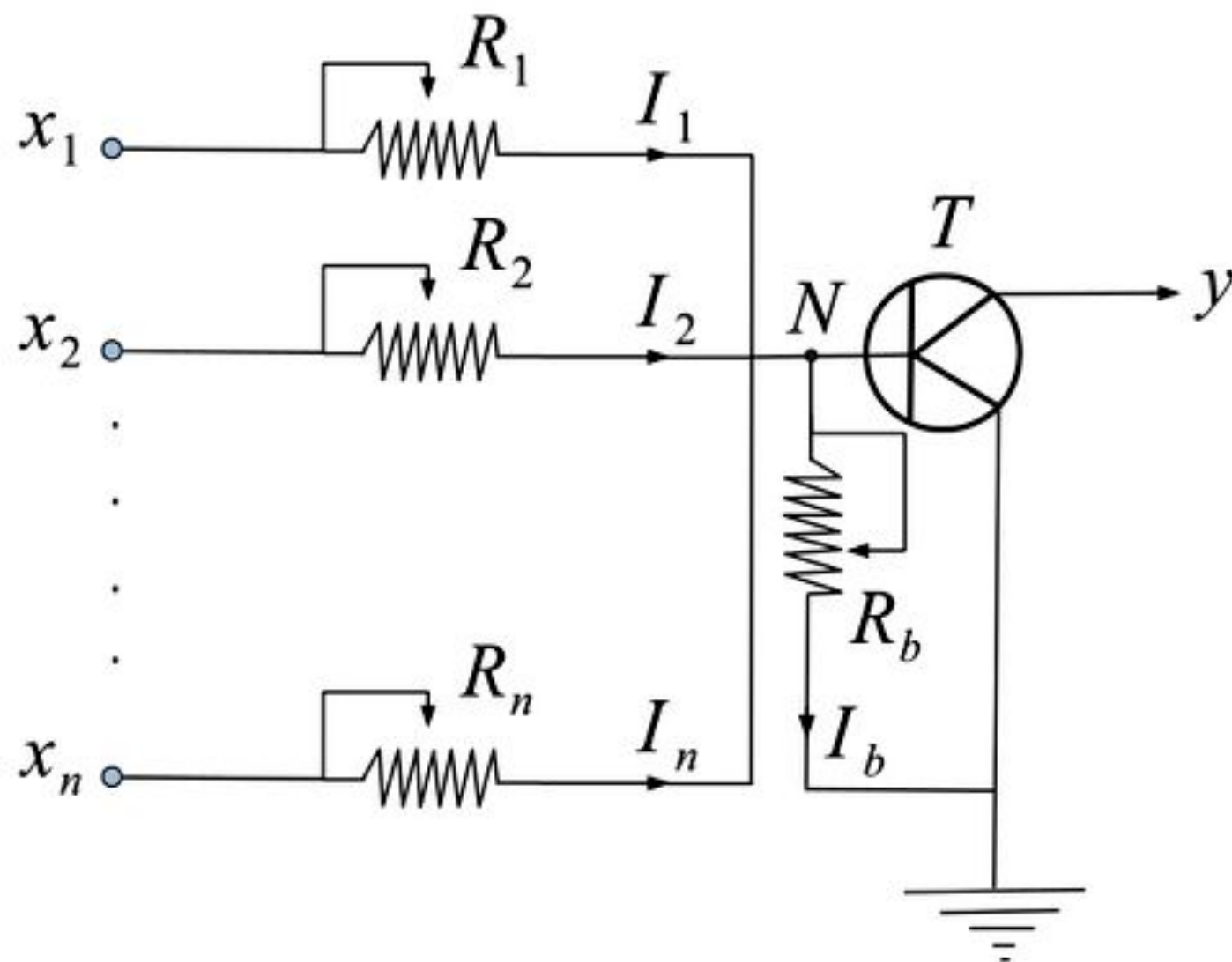


A cocktail factory as a neural network with two hidden layers



1. Intro to Neural Networks

Examples of neural networks

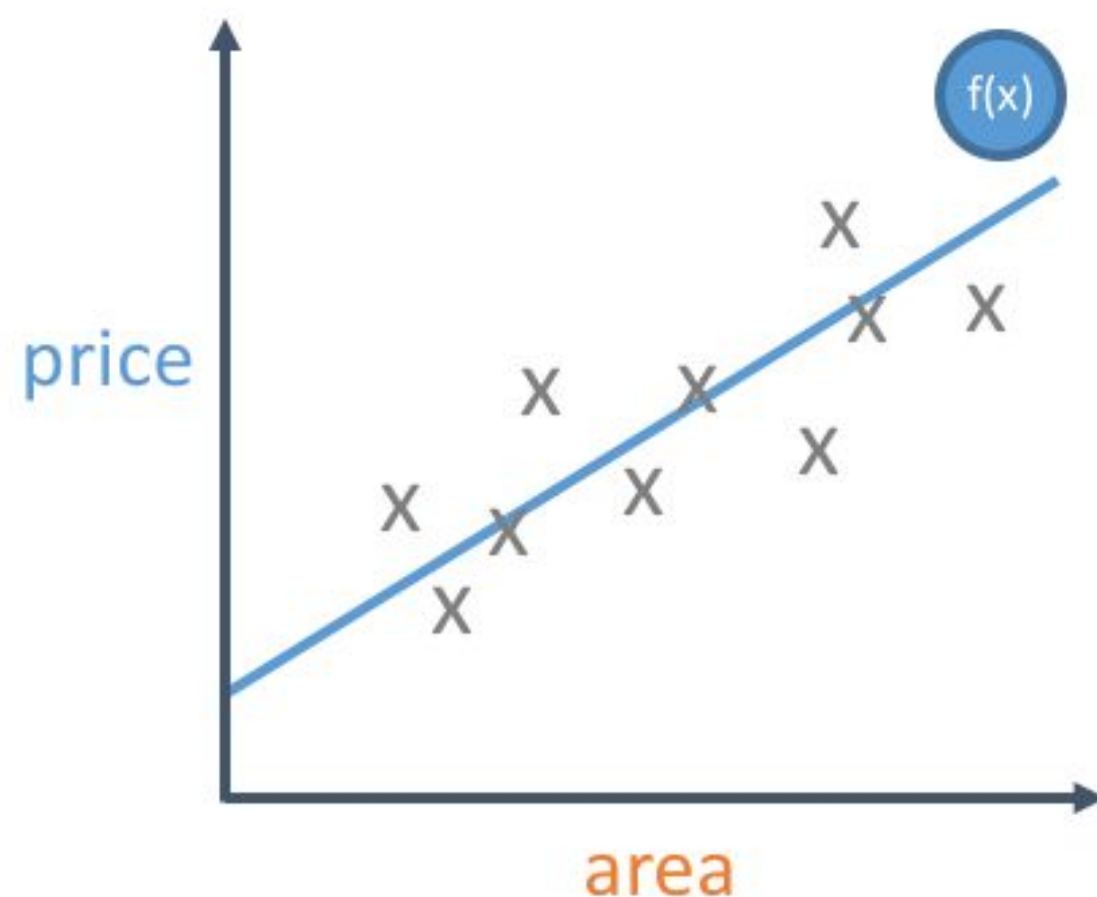


***The neuron implemented
as an electronic circuit***



1. Intro to Neural Networks

Let's continue by thinking on a simple linear model. We want to predict the house price based on the area.



This can be solve by drawing a line using linear regression.

Prices are non-negative!

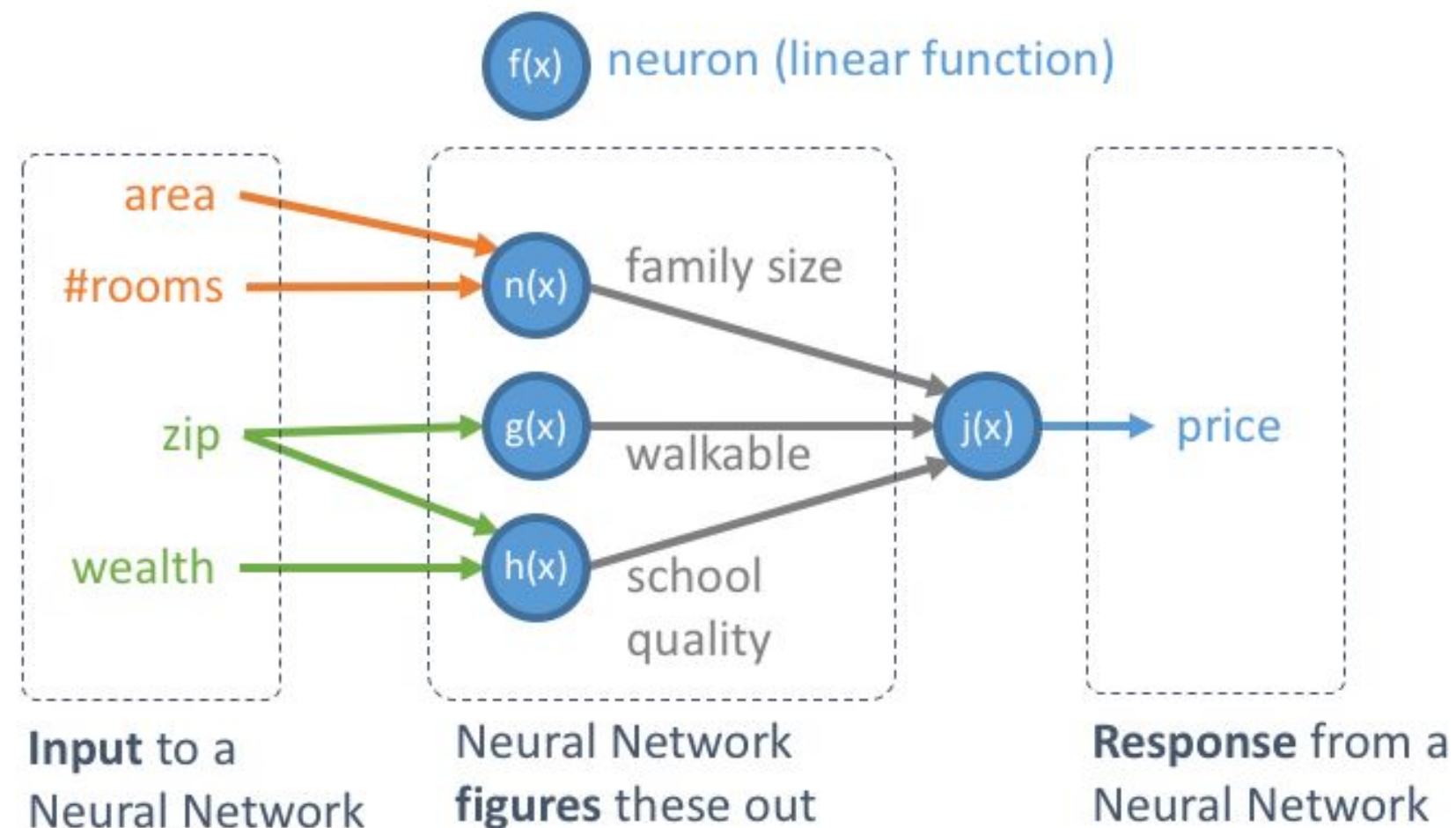


Sources:

- <https://manavseghal.com/lecture-notes-from-artificial-intelligence-is-the-new-electricity-by-andrew-ng-4712dcbf26e5>

1. Intro to Neural Networks

Now the problem is to predict house price for a given house with a number of input features like area, number of rooms, zip code, and wealth of the neighbourhood.

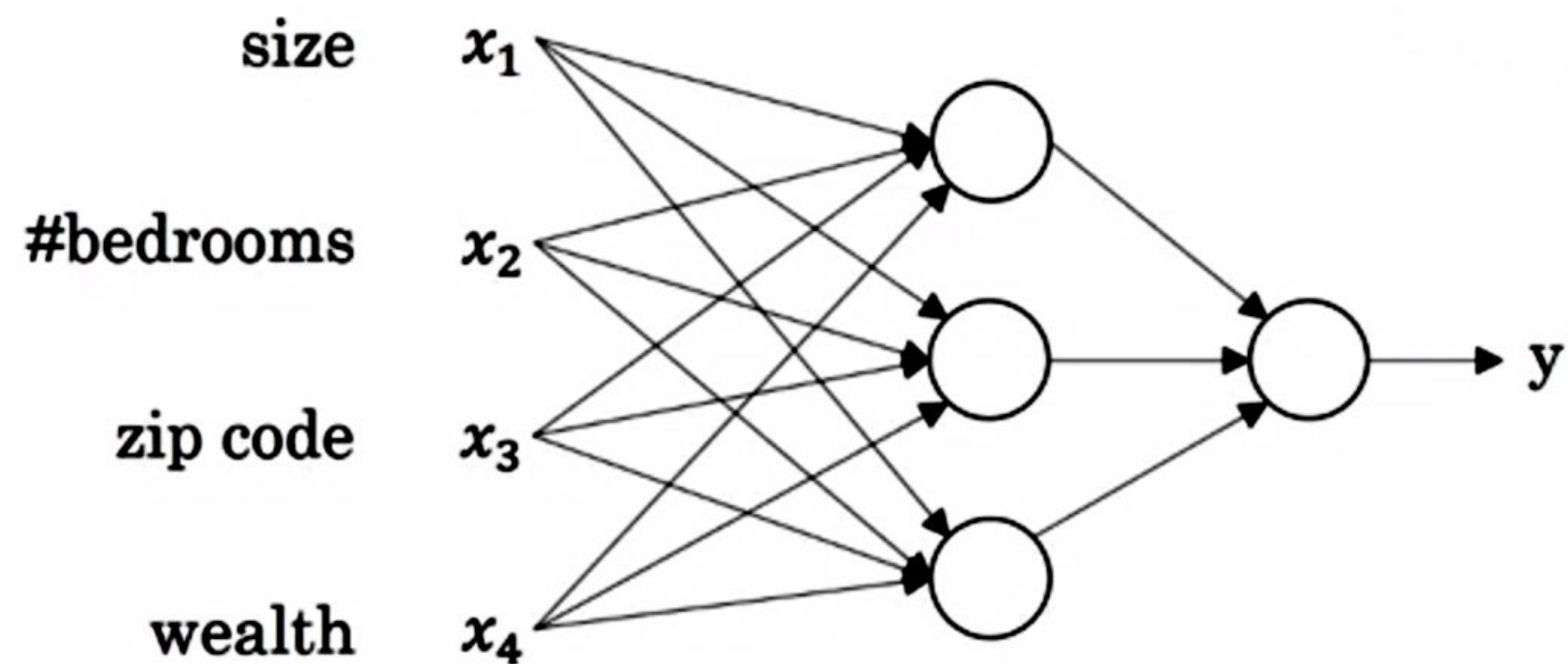


Sources:

- <https://manavseghal.com/lecture-notes-from-artificial-intelligence-is-the-new-electricity-by-andrew-ng-4712dcbf26e5>

1. Intro to Neural Networks

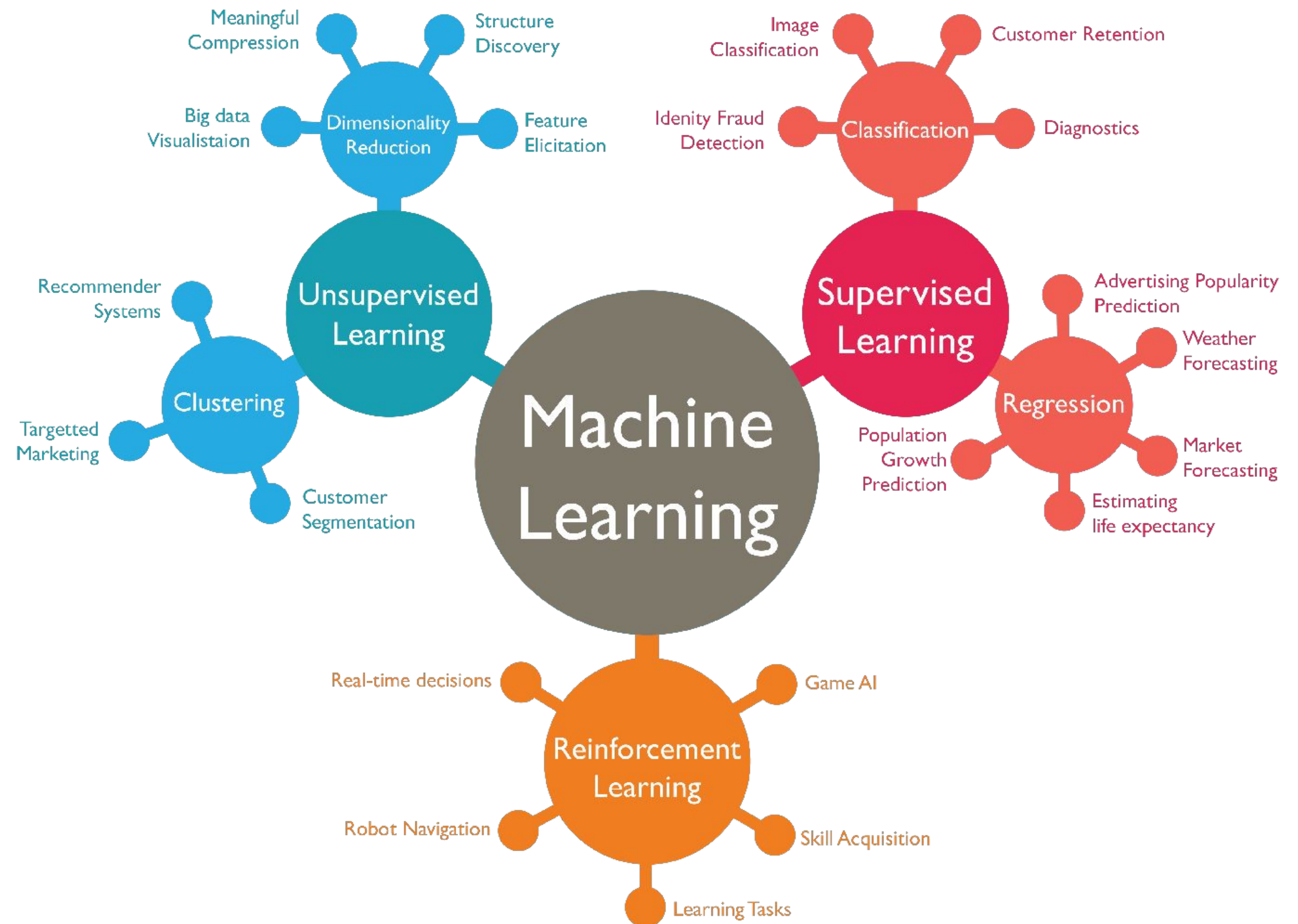
This is what we actually implement when using a neural network. We let the NN to learn whatever features/characteristics to predict the house price



Sources:
- Coursera

1. Intro to Neural Networks

Let's put this application in context. **Supervised learning!**



Sources:

- <http://www.favouriteblog.com/wp-content/uploads/2017/07/Types-of-Learning.png>

© All rights reserved. www.keepcoding.io

1. Intro to Neural Networks

Most of the hype in deep learning is given by their power on supervised learning applications

Supervised Learning

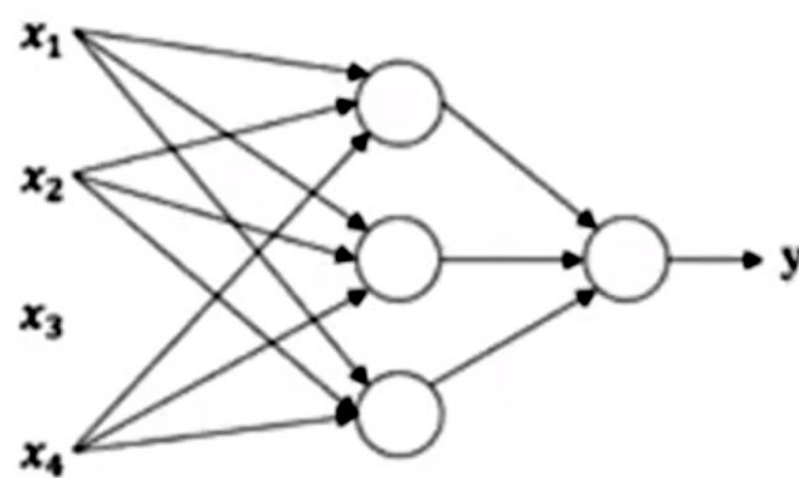
Input(x)	Output (y)	Application
Home features	Price	Real Estate
Ad, user info	Click on ad? (0/1)	Online Advertising
Image	Object (1,...,1000)	Photo tagging
Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Image, Radar info	Position of other cars	Autonomous driving



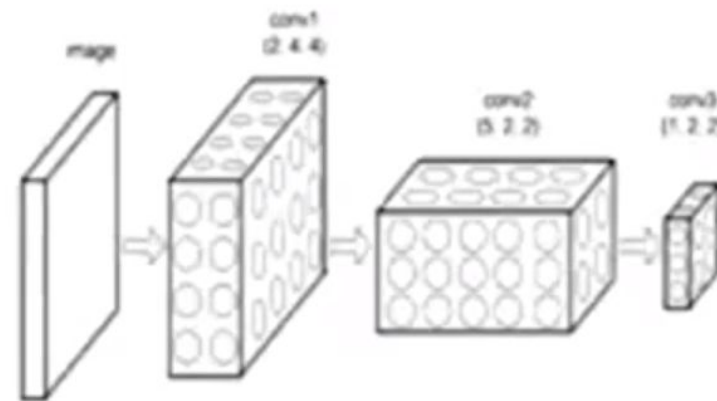
Sources:
- Coursera

1. Intro to Neural Networks

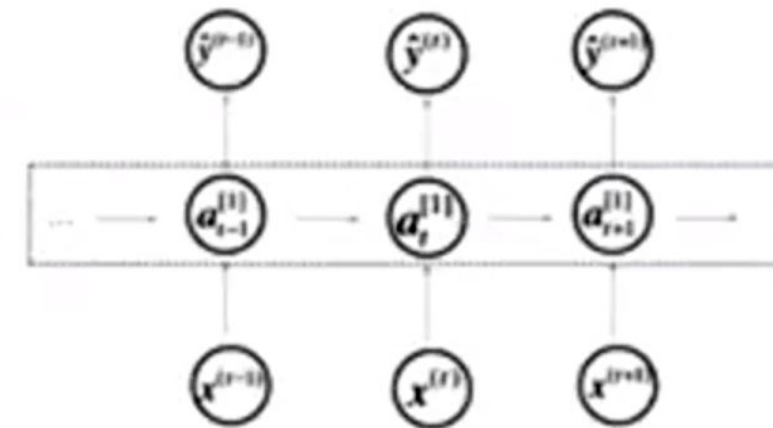
Depending on the application we use a standard NN, convolutional NN or a recurrent NN



Standard NN



Convolutional NN

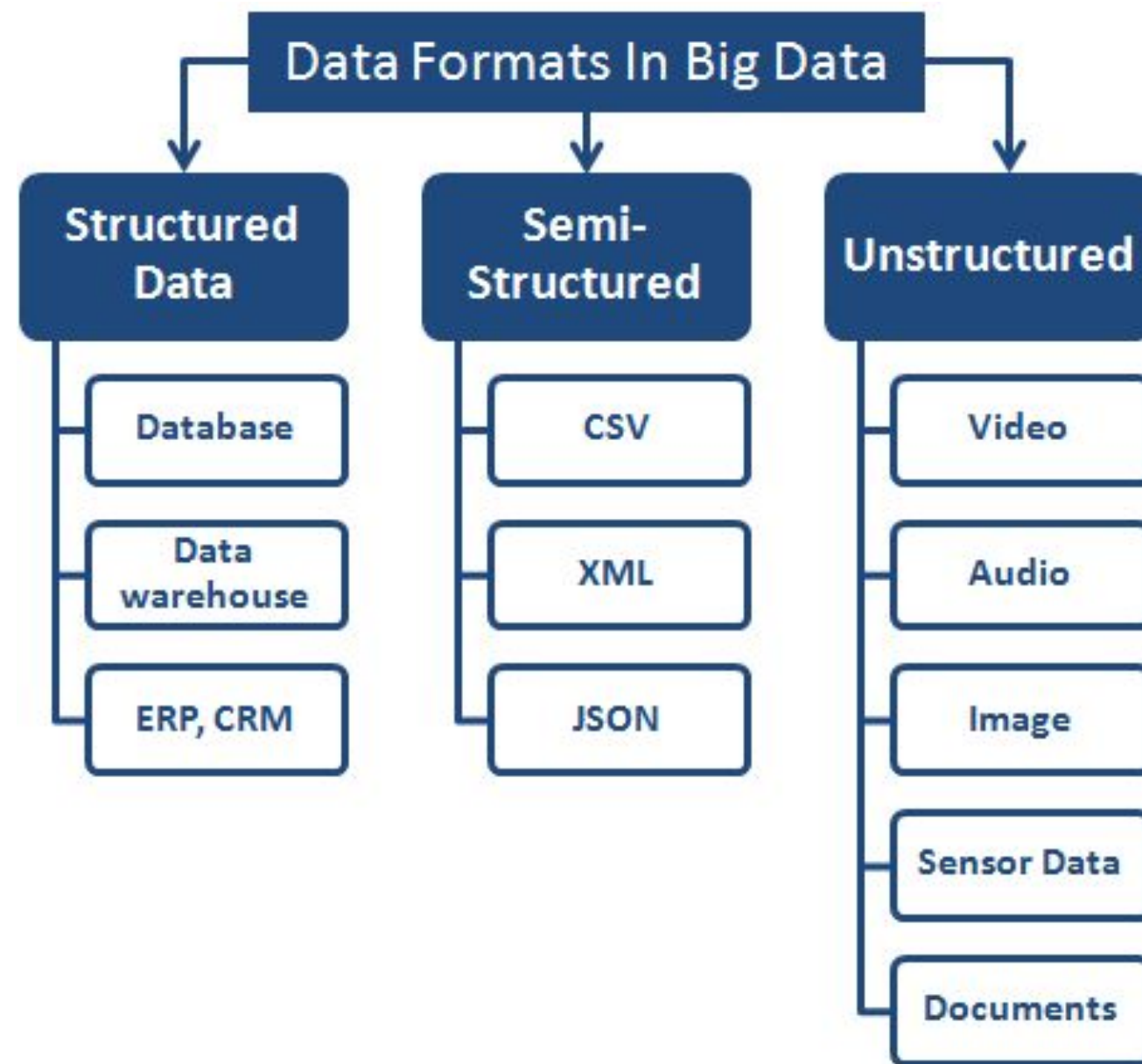


Recurrent NN



1. Intro to Neural Networks

There are also different data formats used to train/test NN depending on the application:

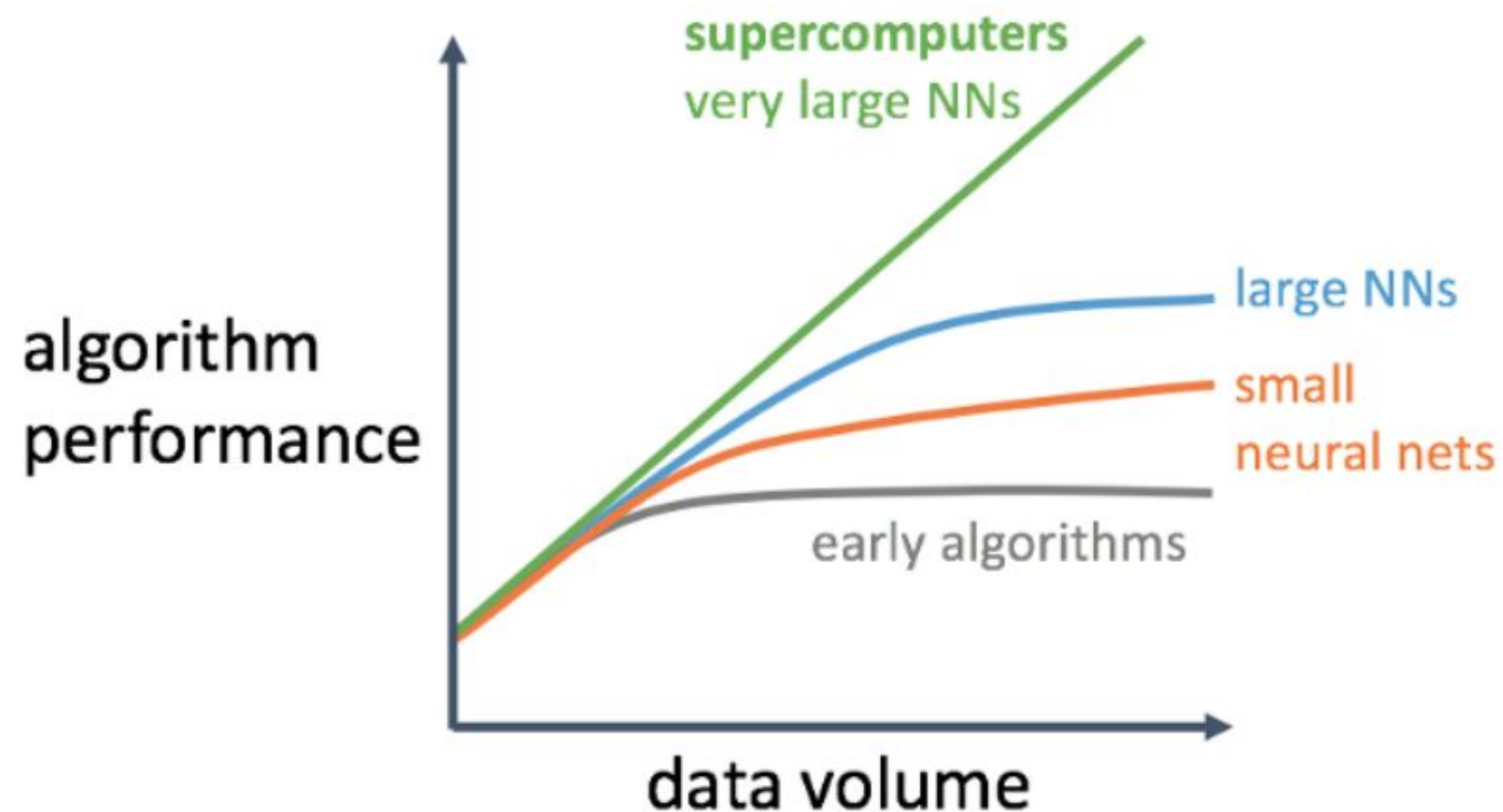


Sources:

- http://tryqa.com/wp-content/uploads/2017/07/Big_Data_Testing-Data_formats_tutorial.png

1. Intro to Neural Networks

The main concepts of deep learning have been around for many decades. But why now they that powerful and useful?



When having small amount of data algorithm performance is similar.



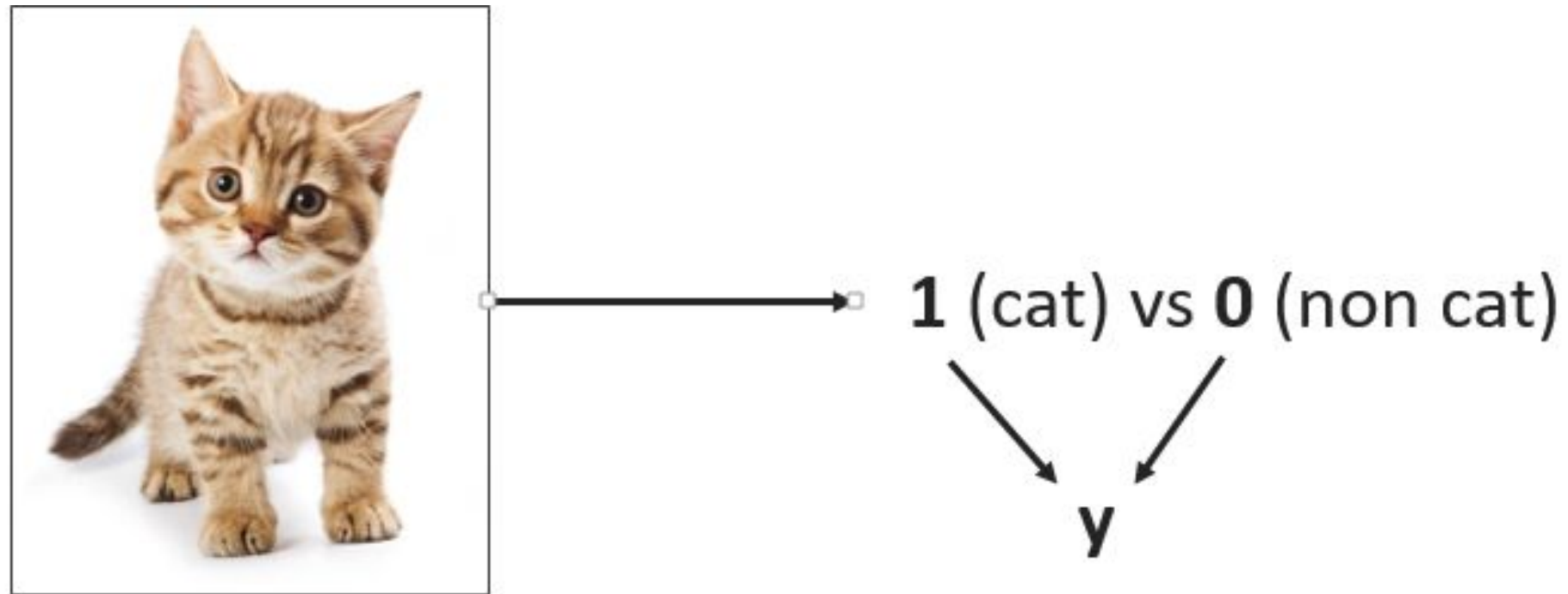
Sources:

- <https://manavseghal.com/lecture-notes-from-artificial-intelligence-is-the-new-electricity-by-andrew-ng-4712dcbf26e5>

2. Binary Classification and Logistic Regression

To clarify: *Logistic regression is an algorithm for binary classification!*

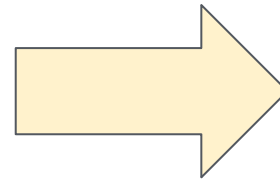
It is a simple way of explaining what a neural network can do.



Sources:
- <http://media5.datahacker.rs/2018/06/word-image-58.png>

2. Binary Classification and Logistic Regression

What is actually an image? Three 2D matrices that represent the Red, Green, and Blue colors



					165	187	209	58	7
					14	125	233	201	98
253	144	120	251	41	147	204			
67	100	32	241	23	165	30			
209	118	124	27	59	201	79			
210	236	105	169	19	218	156			
35	178	199	197	4	14	218			
115	104	34	111	19	196				
32	69	231	203	74					

In this case (for logistic regression), the labels/classes/Categories/Ground-Truth are represented by a vector of zeros or ones



Sources:

- <http://media5.datahacker.rs/2018/06/word-image-58.png>
- <https://yangwangresearch.com/2019/06/11/identifying-scene-changes-in-videos/>

2. Binary Classification and Logistic Regression

We use all the pixel values (RGB) to classify the images into a cat or non-cat class.
To do that the equation for logistic regression:

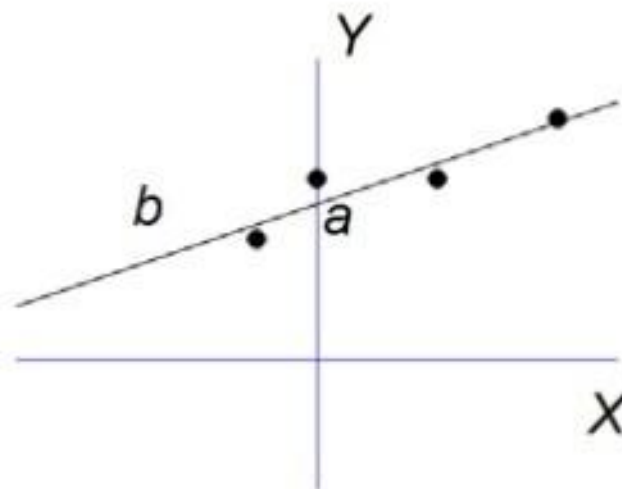
$$\hat{Y} = bX + a$$

predicted values of Y

b = slope = rate of predicted \uparrow/\downarrow for Y scores for each unit increase in X

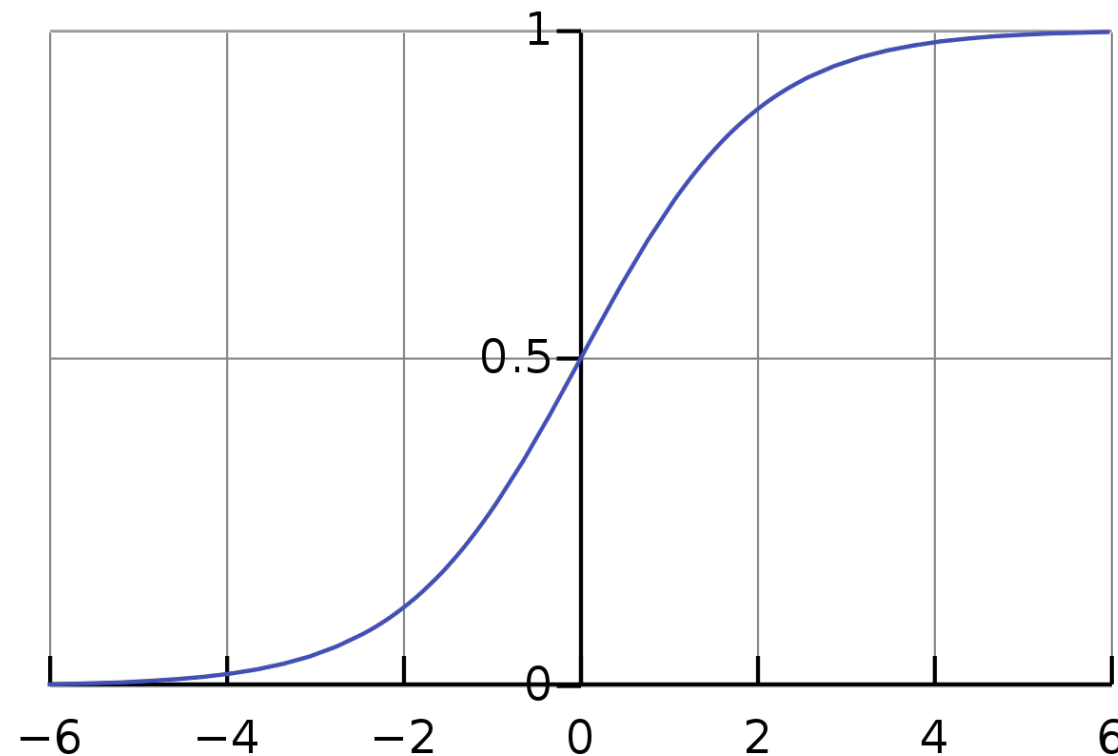
Y-intercept = level of Y when X is 0

Also called bias (b)



2. Binary Classification and Logistic Regression

What we need from the logistic regression equation are values between 0 and 1 representing the **probability of being a cat or a non-cat**. However, the output of the logistic regression could be bigger than 1 and smaller than 0. A sigmoid function will help us to deal with this problem.



Values in the X-axis that are smaller than zero will output values close to zero in the Y-axis.

Similar to values bigger than 1 in the X-axis. The output will be close to 1



Sources:
- https://en.wikipedia.org/wiki/Sigmoid_function

2. Binary Classification and Logistic Regression

So, what we have after using the sigmoid function and the logistic regression equation is the following:

$$\hat{y} = \sigma(w^T x + b), \text{ where } \sigma(z) = \frac{1}{1+e^{-z}}$$

How can we evaluate that the model is doing well? We use the so-called cost function:

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

The gradient descent doesn't work well for this equation.



2. Binary Classification and Logistic Regression

So, we defined a function that works for the gradient descent:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m \boxed{y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})}$$

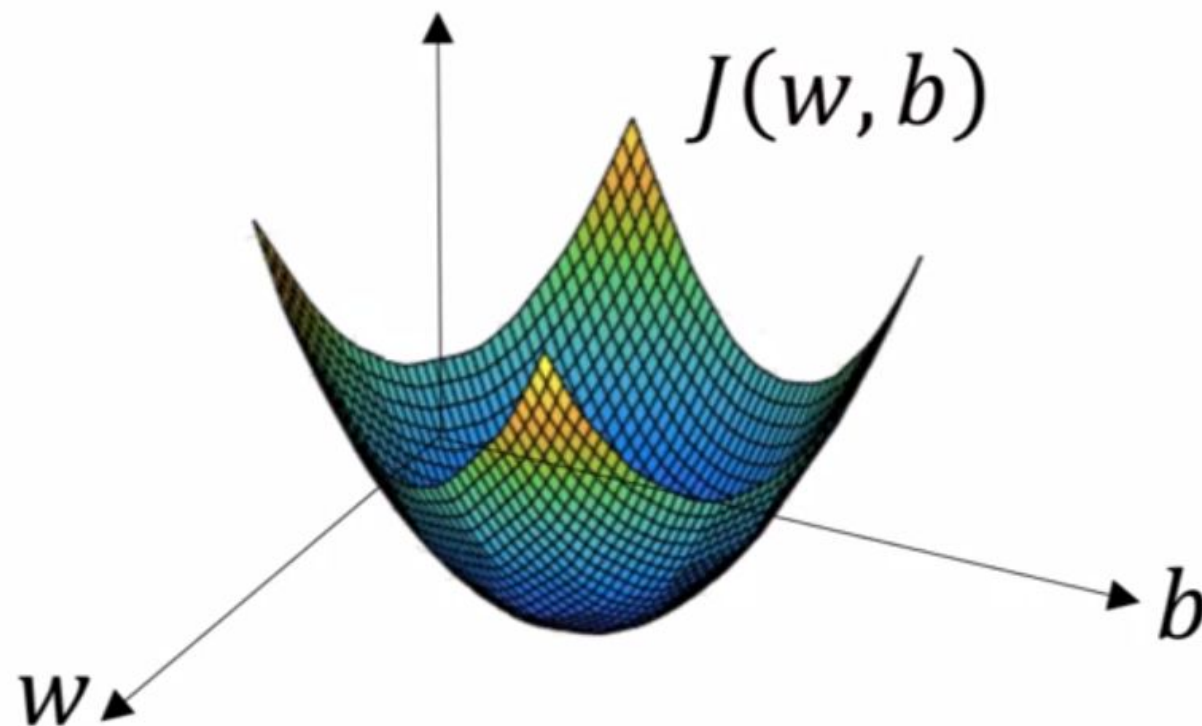
To clarify: A loss function is the error of a single training sample and the cost function is the average error of all the training samples.



3. Gradient Descent

So, we defined a function that works for the gradient descent:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m \boxed{y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})}$$



For intuition purposes, \mathbf{W} is unidimensional

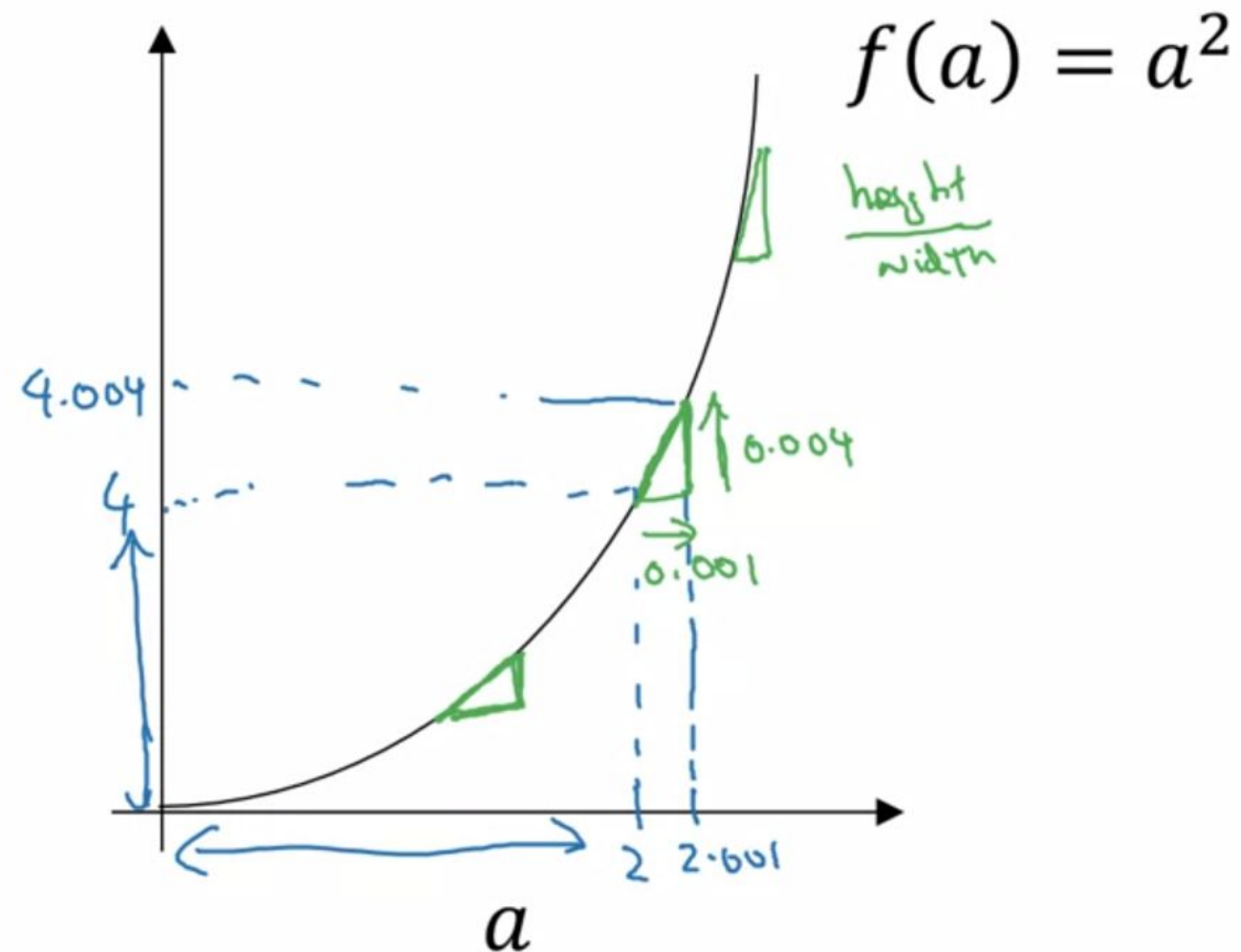
Gradient descent to the rescue! Derivatives! We compute the derivatives to find out the minimum of a function.

Think the derivative as the slope of a function



3. Gradient Descent

What is the main intuition about derivatives? *That the derivative is the slope of a function!*



$$\begin{aligned} a &= 2 & f(a) &= 4 \\ a &= 2.001 & f(a) &\approx 4.004 \\ & & & (4.004004) \end{aligned}$$

slope (derivative) of $f(a)$ at $a=2$ is 4 .

$$\frac{d}{da} f(a) = \underline{4} \quad \text{when } a = \underline{2}.$$

$$\begin{aligned} a &= 5 & f(a) &= 25 \\ a &= 5.001 & f(a) &\approx 25.010 \end{aligned}$$

$$\frac{d}{da} f(a) = \underline{10} \quad \text{when } a = \underline{5}$$

$$\frac{d}{da} f(a) = \frac{d}{da} a^2 = 2a$$



3. Gradient Descent

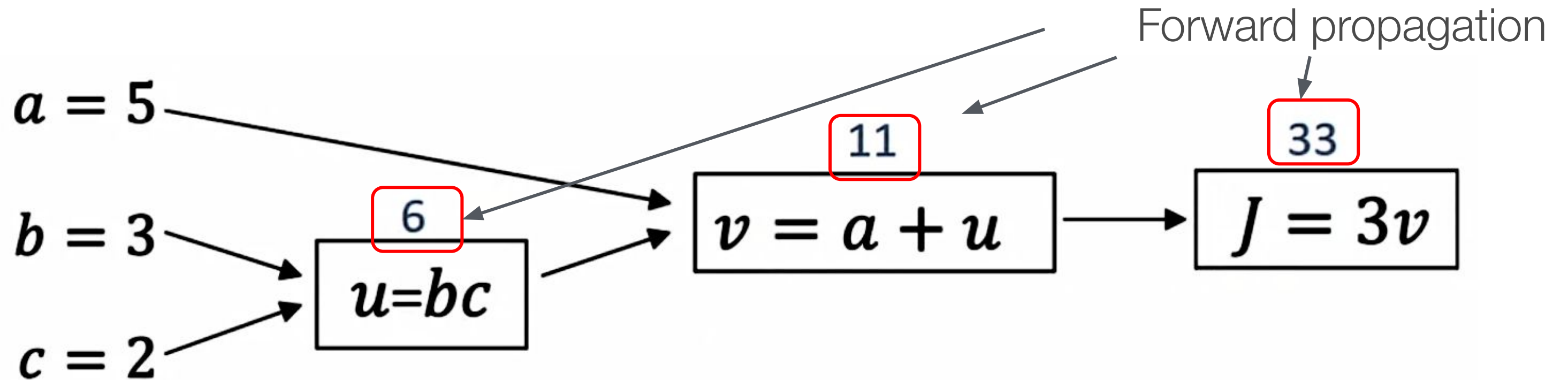
Take messages from the computation of derivatives:

- The derivative of a function is the slope of a function
- The slope of a function can be different at different points of a function
- The computation of a derivative of a function can be done analytically or using graphs



3. Gradient Descent

How do we apply this concept to the neural networks?

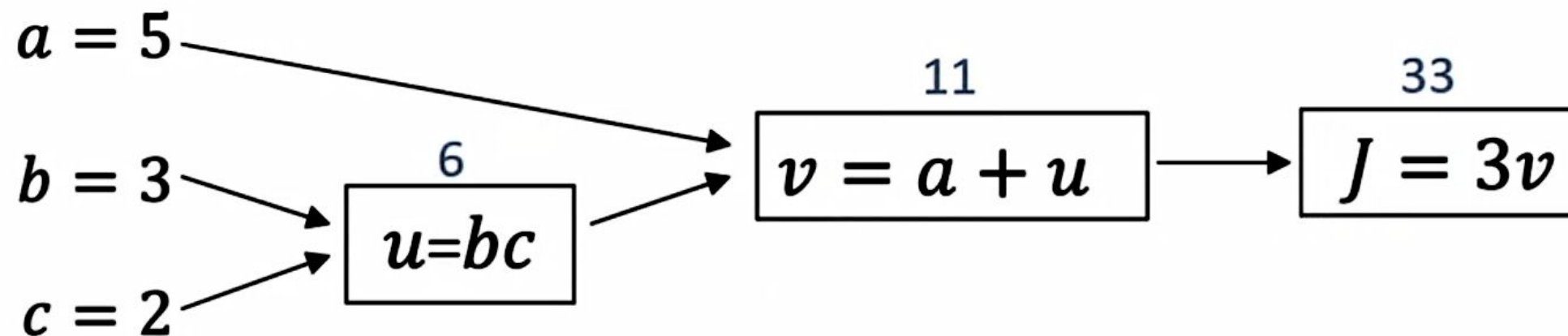


Applying the **chain rule** we can know the derivative w.r.t. a specific variable in the graph



3. Gradient Descent

Applying the chain rule we obtained that:

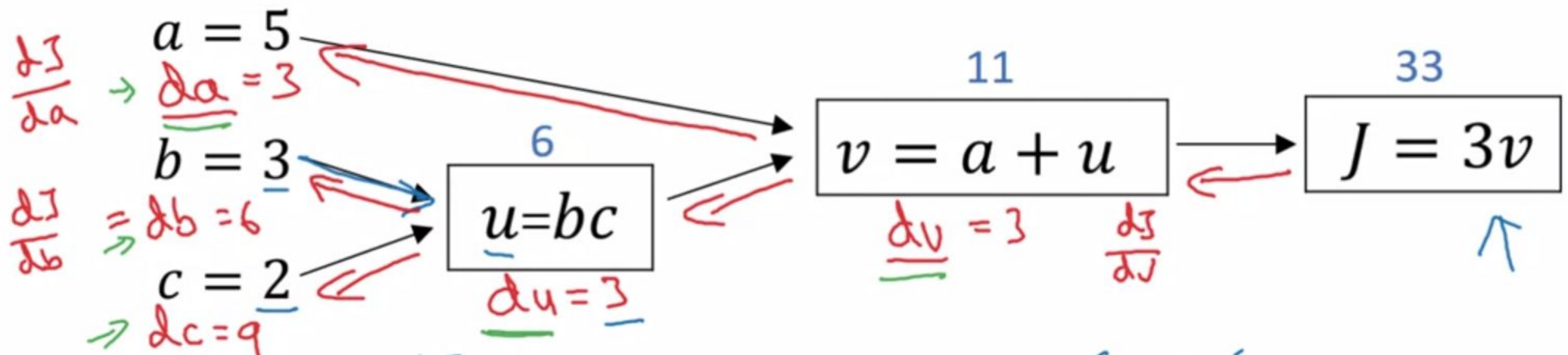


$$\begin{aligned}\frac{dJ}{du} &= 3 = \frac{dJ}{dv} \cdot \frac{dv}{du} \\ \frac{dJ}{db} &= \frac{dJ}{du} \cdot \frac{du}{db} = 6 \\ \frac{dJ}{dc} &= \frac{dJ}{du} \cdot \frac{du}{dc} = 9\end{aligned}$$



3. Gradient Descent

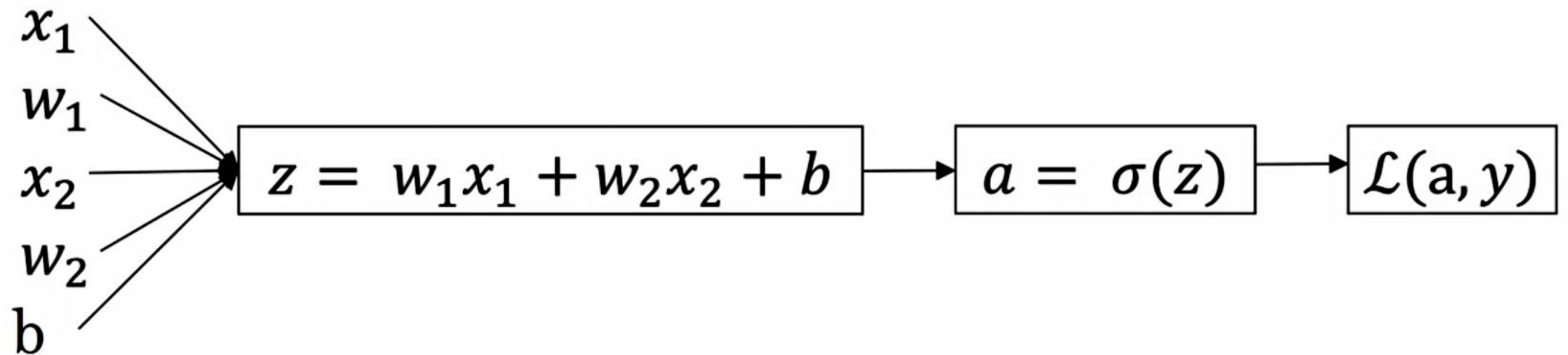
The key message from this is that if we want to compute derivatives, the most efficient way to do this is from right to left



3. Gradient Descent

Graph for gradient descent for logistic regression on one sample

$$\begin{aligned} z &= w^T x + b \\ \hat{y} &= a = \sigma(z) \\ \mathcal{L}(a, y) &= -(y \log(a) + (1 - y) \log(1 - a)) \end{aligned} \quad \left. \vphantom{\begin{aligned} z &= w^T x + b \\ \hat{y} &= a = \sigma(z) \\ \mathcal{L}(a, y) &= -(y \log(a) + (1 - y) \log(1 - a)) \end{aligned}} \right\} \textit{Forward pass}$$



3. Gradient Descent

Gradient descent for logistic regression on all training samples m :

$$J = 0, \quad dw_1 = 0, \quad dw_2 = 0, \quad db = 0$$

for $i = 1$ to m :

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

Forward pass

Backward pass

$$J = J/m, \quad dw_1 = dw_1/m, \quad dw_2 = dw_2/m$$

$$db = db/m$$



4. Vectorization

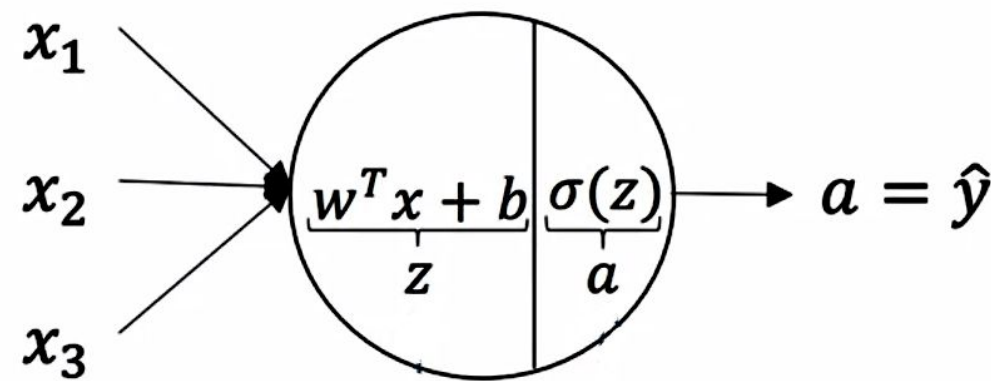
How we can do to take advantage of GPUs to compute gradient descent?

A neural network (NN) with two layers is basically one logistic regression follow by another logistic regression! A NN with three layers is a logistic regression repeated 3 times ...



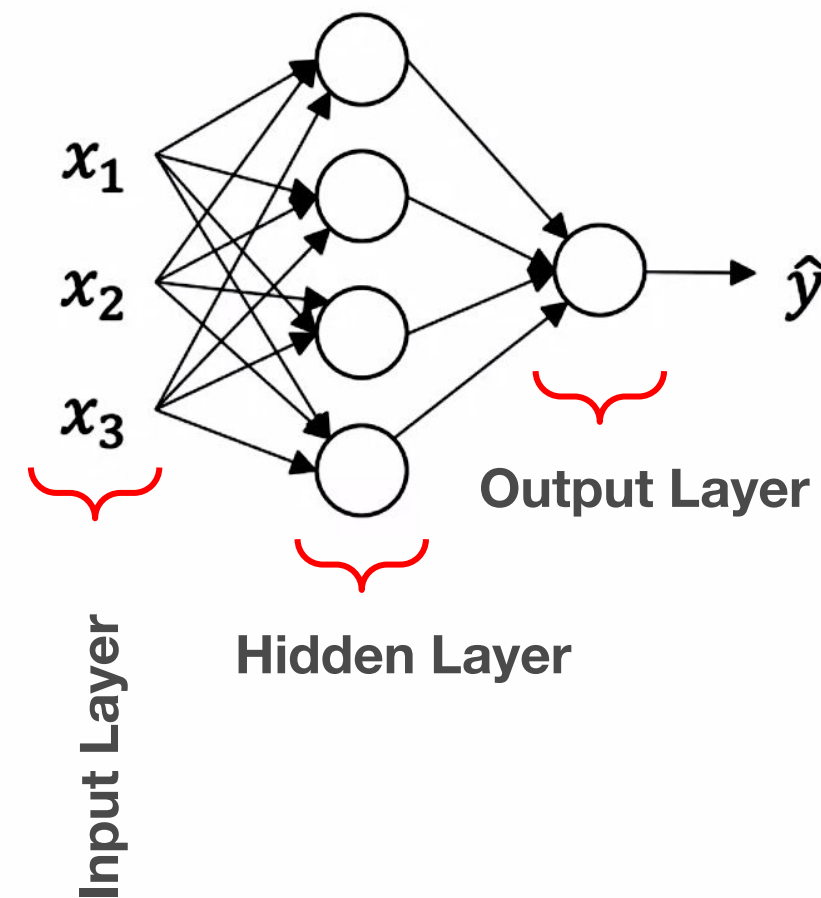
4. Vectorization

From the logistic regression to a neural network (NN)



$$z = w^T x + b$$

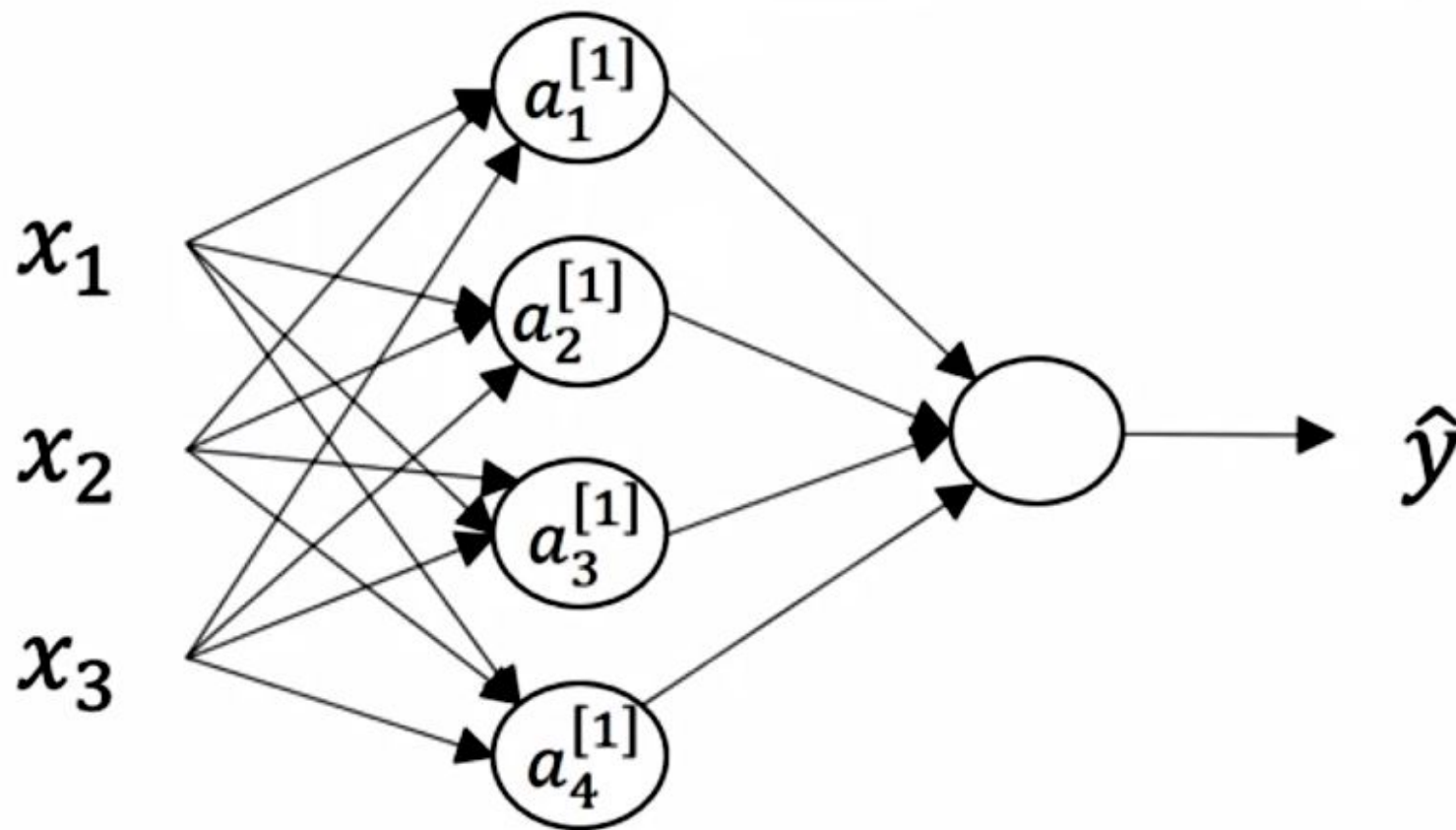
$$a = \sigma(z)$$



4. Vectorization

From the logistic regression to a neural network (NN).

In these equations, the index in square parenthesis represents the layer



$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$



4. Vectorization

From the logistic regression to a neural network (NN).

In these equations, the index in square parenthesis represents the layer

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

Computing these equations using a for loop is really inefficient!

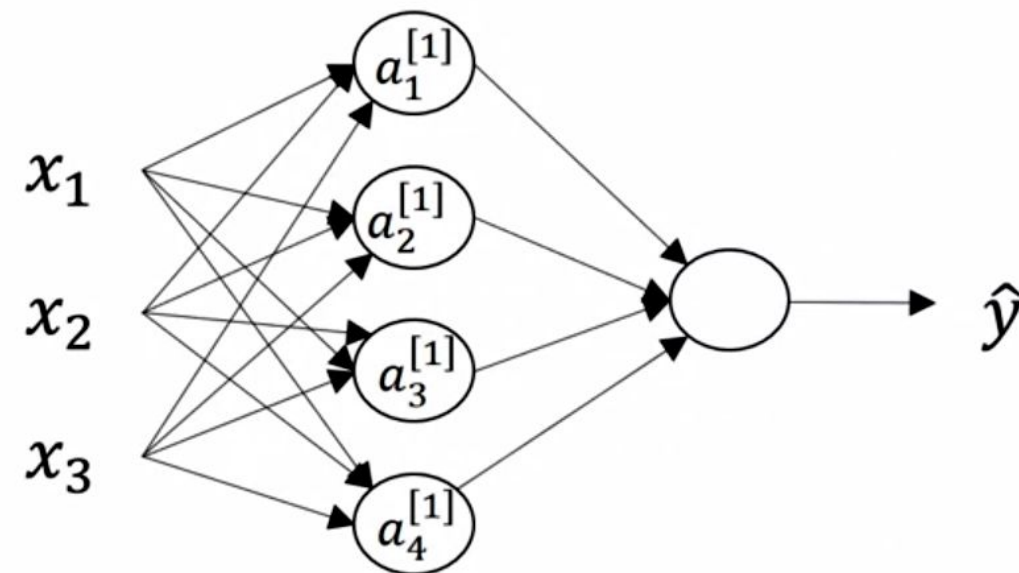
The more logistic units (neurons) we have the bigger is the for loop

For that reason vectorizing the weights, biases and Z values is important to increase performance



4. Vectorization

Vectorizing previous equations resulted on similar logistic regression equations. The only difference is that the **weights of the hidden layer are represented by a matrix of size (4,3)**, the biases by a matrix of size (4,1) and the z s by a matrix of (4,1)



Given input x :

Hidden Layer

$$\left\{ \begin{array}{l} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \end{array} \right.$$

Output Layer

$$\left\{ \begin{array}{l} z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right.$$

Hidden layer

Weights: 4 neurons, 3 inputs each

Biases: 4 values. One for each neuron

Output layer

Weights: 1 neurons, 4 inputs each

Biases: 1 value



4. Vectorization

So how can we compute **z** and **a** values for m samples in the training set?

Vectorizing across multiple samples!

```
for i = 1 to m:
```

$$z^{[1]}(i) = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1]}(i) = \sigma(z^{[1]}(i))$$

$$z^{[2]}(i) = W^{[2]}a^{[1]}(i) + b^{[2]}$$

$$a^{[2]}(i) = \sigma(z^{[2]}(i))$$

$z^{[1]}(i)$ -> Columns represent training samples and rows represent hidden units



5. Activation functions

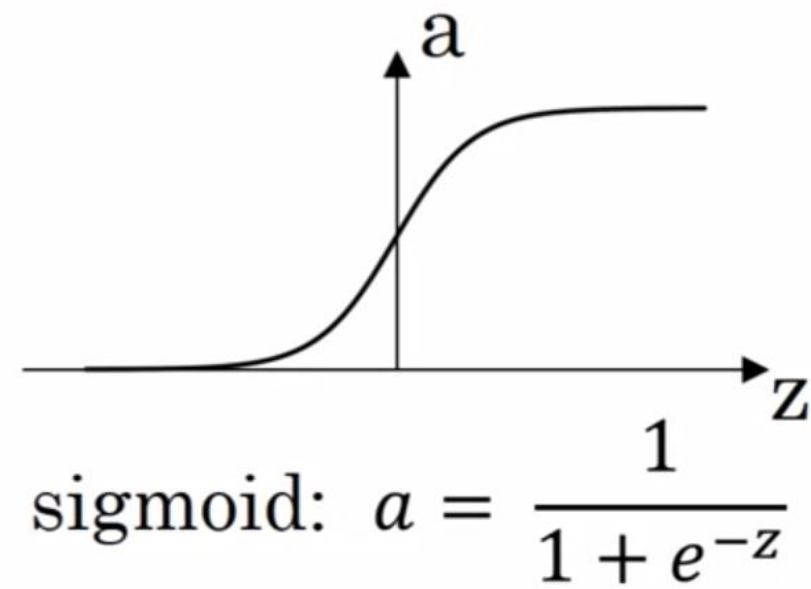
When building a neural network, one of the important parts that you need to decide is the activation function.

So far we have been using the sigmoid activation function, but it's been demonstrated that it is not the best option for hidden layers. **Tanh or ReLU (rectifier linear unit)** are more used and have better performance than the sigmoid in the hidden layers.

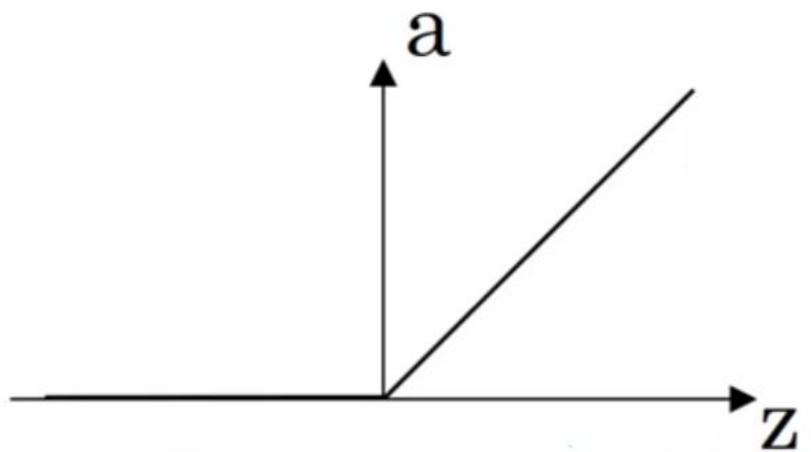
Sigmoid function is used in the output layer when classifying two classes (Binary classification)



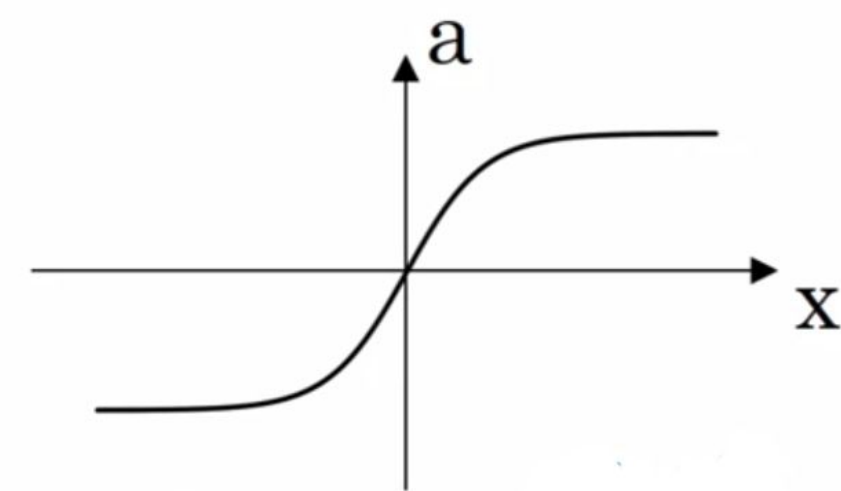
5. Activation functions



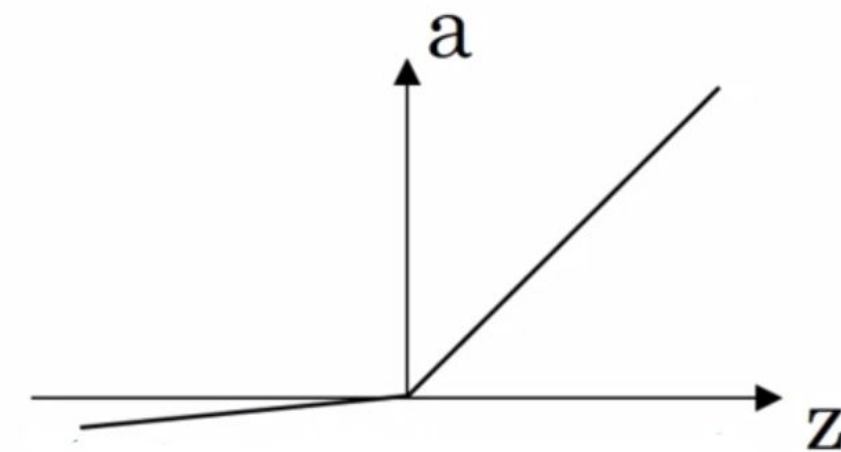
ReLU



Tanh



Leaky ReLU



5. Activation functions

Why neural networks need non-linear activation functions?

The most important reason for using non-linear activation functions is that they allow the NN to learn **non-linear mapping/functions between the input and the output!**

Otherwise, using linear activation function will only make the NN to learn **linear functions between input and output**



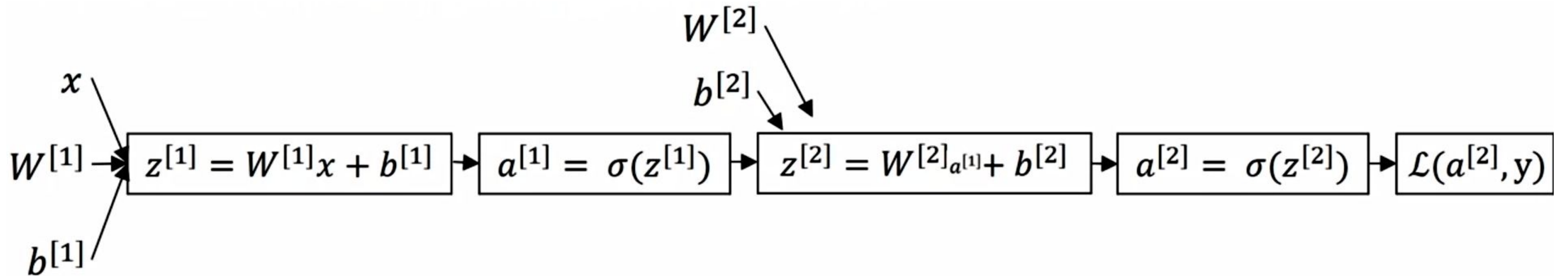
5. Activation functions

Remember that in order to implement gradient descent, we need **the derivatives of the activation functions**. Here are some of them:

Name	Equation	Derivative
Sigmoid	$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh	$f(x) = \tanh(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$	$f'(x) = 1 - f(x)^2$
Rectified Linear Unit (relu)	$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Leaky Rectified Linear Unit (Leaky relu)	$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$



6. Backpropagation



Output Layer

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

Hidden Layer

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

In these equations, dz actually means:

$$dz = \frac{\delta L}{\delta z} = \frac{\delta L}{\delta a} \frac{\delta a}{\delta z}$$



7. Random Initialization

Randomly initializing the weights in a neural network is highly beneficial for the gradient descent computation.

There are different methods developed in the literature presenting different ways of randomly initialized the weights.

Easy to do and highly beneficial!



Let's move to Google Colab!

Notebooks:

- *0_Numpy_Basic_Functions.ipynb*
- *1_First_Neural_Network.ipynb*
- *2_Neural_Network_with_One_Hidden_Layer.ipynb*

