

Deep Learning

what is Deep Learning ?

It is a process of training a neural network.

so what is Neural network ?

credit : Andrew NG,

Lets say we are computing the price of house based on size of house, then the non linear function or (linear function) can be considered as a neuron.

The collection of this neurons is called a neural network, Each and every neural network takes into all the features, but uses which ones it thinks as necessary.

WHAT IS A NEURAL NETWORK?

Notes

Housing Price Prediction

The diagram illustrates a neural network architecture for predicting house prices. It consists of three layers: an input layer with four features (x_1 through x_4), a hidden layer with three neurons, and an output layer with one neuron labeled "y price". The input features are "size", "#bedrooms", "zip code", and "wealth". The output "y" is also labeled "price". Above the network, a separate diagram shows how these features might be processed by a function f to produce "family size", "walkability", and "school quality", which then contribute to the final price prediction. The interface includes a play button, a progress bar at 6:37 / 7:16, and a download button.

All the neurons in the middle is coined as hidden layer.

Supervised Learning with neural networks

Deep Learning models works very well, but most of the hype around DL is bought out by Supervised learning.

They got well in interpreting structured as well as unstructured data, for example identifying cats from a group of dog pictures, which is unstructured data. rolling out specific ads to target audience (structured data)

for Real estate and Online advertising we use standard neural networks.

Supervised Learning with Neural Networks

Notes

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 |

Save note

Download



Why is Deep Learning Taking off ?

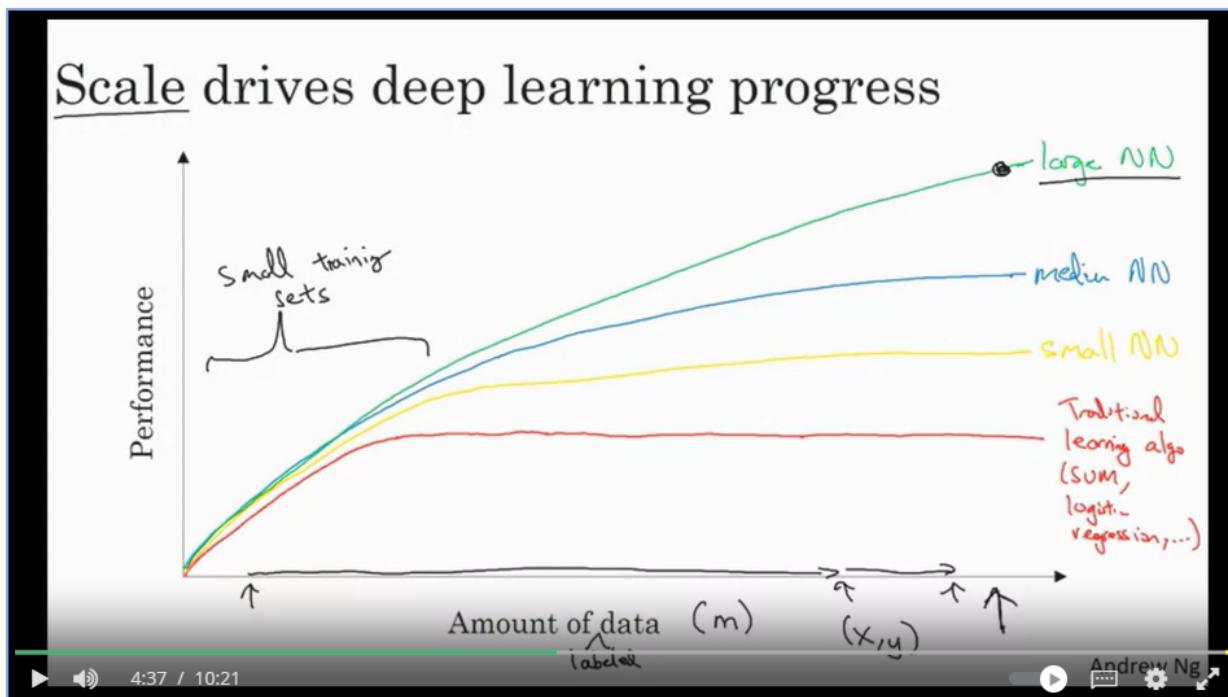
For all the traditional machine learning algorithms, the accuracy is almost a plateau, after some point in time, but the accuracy of neural networks explode, which makes the neural networks interesting.



All traditional ML algos, Neural networks perform similarly for small amounts of training data and it boils down to hyper parameter tuning. However, when we have large amounts of data, there is a significant increase in accuracy when neural networks are used.

for large data

traditional ml algos < small neural networks < medium NN < large NN



Scale is driving Deep Learning :

1. Algorithms :

one of the notable break through is using RELU function ?? instead of Sigmoid function ?? (It makes gradient Descent work much faster)

2. Computation :

It impacts how faster we can get the experimental result back

3. Data

Logistic Regression as a Neural Network :

To build a classifier that says if the image is cat or dog, we need logistic regression.

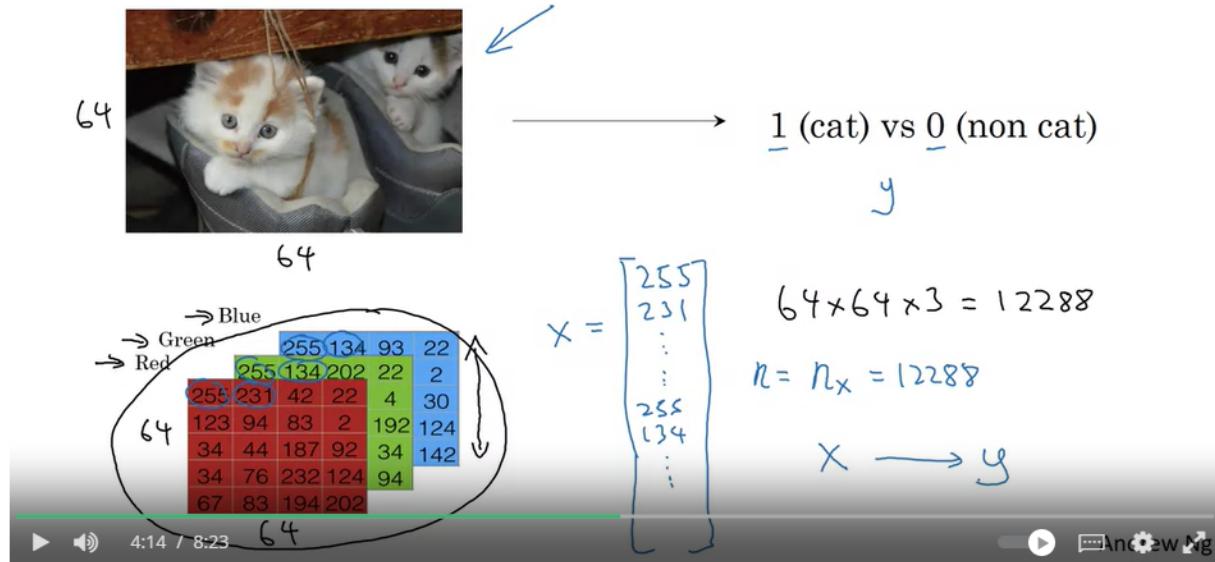
The images are represented using 3 matrixes they are red, blue and green, (have dimensions each equal to the image). Every element has some value associated with it.

Now the X feature is aggregated as the each element of red, green and blue matrix placed one below another, this is the X feature, N_x is the number of rows.

Binary Classification

[Notes](#)

Binary Classification



[Save note](#) [Download](#)

If there are m such training sets, then the training examples are stacked as below diagram named as X .

which also has Y (consisting of values 0 or 1).

Notation

$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$

m training examples : $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

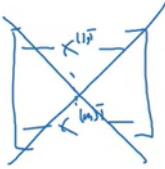
$M = M_{\text{train}}$

$M_{\text{test}} = \# \text{test examples.}$

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \quad \begin{matrix} \uparrow \\ n_x \\ \downarrow \end{matrix}$$

$X \in \mathbb{R}^{n_x \times m}$

$X.\text{shape} = (n_x, m)$



$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$

7:54 / 8:23

▶ Ancestry

Logistic Regression

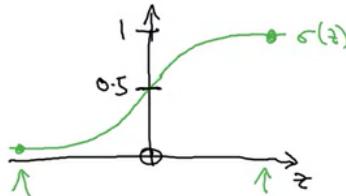
To build a model, we need a function.

Logistic Regression

Given x , want $\hat{y} = \frac{P(y=1|x)}{0 \leq \hat{y} \leq 1}$
 $x \in \mathbb{R}^{n_x}$

Parameters : $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$.

Output $\hat{y} = \sigma(w^T x + b)$



$$x_0 = 1, x \in \mathbb{R}^{n_x+1}$$

$$\hat{y} = \sigma(\theta^T x)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \quad \left. \begin{array}{l} \text{b} \leftarrow \\ \text{w} \leftarrow \end{array} \right\}$$

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

If z large $\sigma(z) \approx \frac{1}{1+0} = 1$

If z large negative number

$$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+\text{Bignum}} \approx 0$$

Andrew Ng

the function is outlined in the blue line beside output.

Since the output $W^T X + B$, would be a continuous variable, which we do not want, We require it to be between 0 and 1, for that purpose we apply sigmoid to the value, which converts the values between 0 and 1.

We could use the notation (θ), but that notation is computationally expensive (in RED LINES), the only difference is that the b intercept is modelled into the matrix, with (θ_0) representing B in our output.

Logistic Regression Cost Function :

Logistic Regression cost function

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

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$.

$x^{(i)}$
 $y^{(i)}$
 $z^{(i)}$

i-th example.

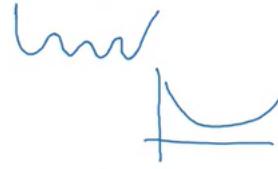
Loss (error) function: $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

$$L(\hat{y}, y) = -(\underbrace{y \log \hat{y}}_{\text{Want } \log \hat{y} \text{ large}} + \underbrace{(1-y) \log (1-\hat{y})}_{\text{Want } 1-\hat{y} \text{ large}}) \leftarrow$$

If $y=1$: $L(\hat{y}, y) = -\log \hat{y} \leftarrow$ Want $\log \hat{y}$ large, want \hat{y} large.

If $y=0$: $L(\hat{y}, y) = -\log (1-\hat{y}) \leftarrow$ Want $\log (1-\hat{y})$ large ... want \hat{y} small

Cost function: $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$



▶ 🔍 7:32 / 8:12 ⏪ Andrew Ng

Use of traditional squared error, would give us a curve, for which finding gradient descent is not that easy. So we use another loss function (here applying gradient descent is much easier).

Gradient Descent :

We have to minimize the cost function, the above cost function is observed to look like a convex curve, which has only one global minima, where ever we initialize the point, we would converge to a global minima.

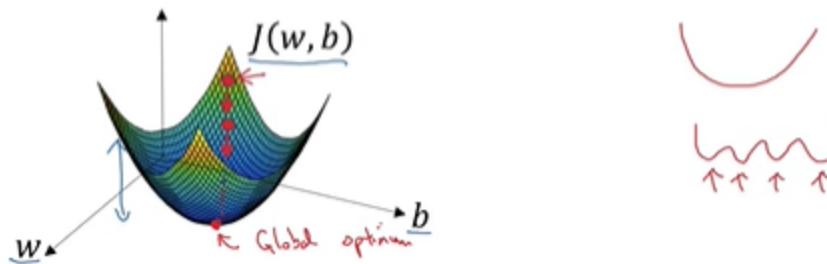
Intuition : we plot the cost function, in hopes of finding the variables x, y and (all other remaining dimensions), which gives us the global minima, as it can be considered as a best fit for our equation.

Gradient Descent

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

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

Want to find w, b that minimize $J(w, b)$

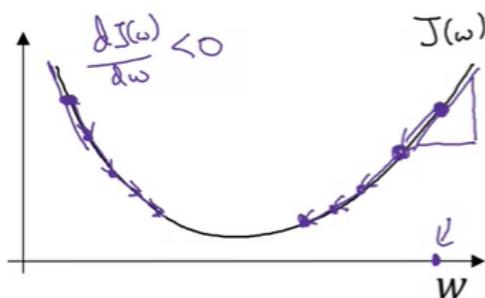


Andrew Ng

Note : There might be a lot of dimensions associated (w axis, b axis, so on), for this example, lets assume we only have one dimension.

We initialize at a particular point, and then converge down to a global minima, at each iteration we remove (learning rate * (derivative at that point))

Gradient Descent



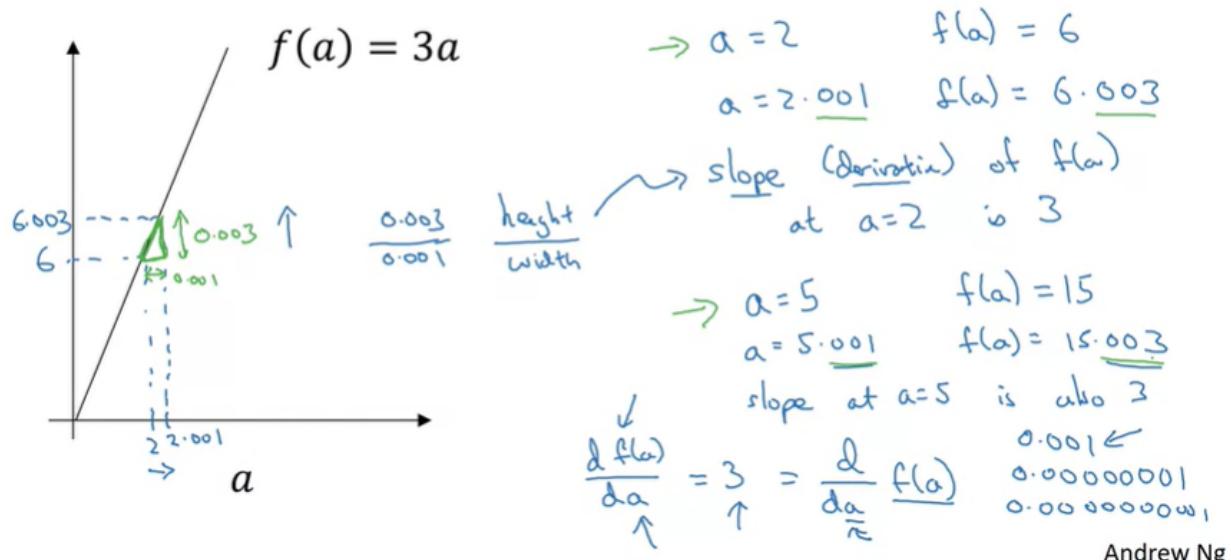
Repeat {
 $w := w - \alpha \frac{dJ(w)}{dw}$
 $w := w - \alpha \frac{d^2J(w)}{dw^2}$ }
 learning rate

In the above figure, we are updating the value of "w", by subtracting the learning rate * (the slope of the curve with respect to variable "w") and vice versa for "b".

Derivatives :

Derivative can be described as change in $f(a)$, for a change in a .

Intuition about derivatives

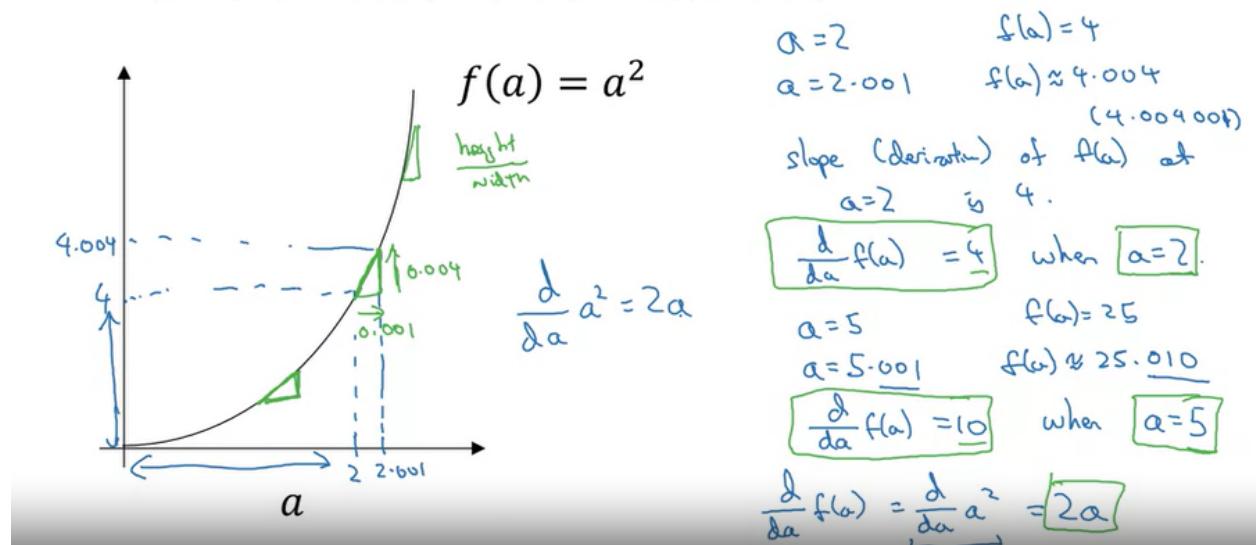


Note : the slope of line is same all around the points, However it is different when it comes to non-linear graphs

Now, let's jump to some more examples.

Note : The derivative of a function at a specific point represents the rate of change of the function at that point. It is a measure of how much the function changes as you make an infinitesimally small change in the input variable.

Intuition about derivatives



When you have a function ($f(a) = a^2$), you can calculate its derivative by looking at how the function changes as the input variable (a) changes. The derivative tells you the rate of change of the function.

1. When (a) is 2, ($f(a)$) is 4, and when (a) is 2.001 (**in reality it is much smaller than 0.001, (close to 0)**), ($f(a)$) is approximately 4.004. To calculate the derivative, you divide the change in $f(a)$ (0.004) by the change in (a) (0.001). This ratio, $0.004/0.001$, gives you approximately 4, which matches the derivative $2a$ (according to the formula $d(a^2) = 2*a$), where (a) is 2, and $2a$ is 4. This shows that the derivative accurately reflects the rate of change at ($a = 2$).

2. Now, when (a) is 2, (f(a)) is 4, and when (a) is 4, (f(a)) is 16. Again, you can calculate the derivative by dividing the change in (f(a)) (16 - 4 = 12) by the change in (a) (4 - 2 = 2). This ratio, 12/2, gives you 6, which **does not** match the derivative 2a, where |(a)| is 2, and 2a is 4. Because derivatives **define how much the function changes as you make an infinitesimally small change in the input variable, Here number increase of 2 units is really large, so the derivative formula does not work here.**

some more examples :

More derivative examples

$$f(a) = a^2$$

$$\frac{d}{da} f(a) = \underbrace{2a}_{4}$$

$$\begin{array}{l} a=2 \\ a=2.001 \end{array}$$

$$\begin{array}{l} f(a)=4 \\ f(a) \approx 4.004 \end{array}$$

$$f(a) = a^3$$

$$\frac{d}{da} f(a) = \underbrace{3a^2}_{3+2^2=12}$$

$$\begin{array}{l} a=2 \\ a=2.001 \end{array}$$

$$\begin{array}{l} f(a)=8 \\ f(a) \approx 8.012 \end{array}$$

$$\begin{array}{l} f(a) = \log_e(a) \\ \ln(a) \end{array}$$

$$\begin{array}{l} \frac{d}{da} f(a) = \frac{1}{a} \\ \downarrow \\ \frac{d}{da} f(a) = \boxed{\frac{1}{2}} \end{array}$$

$$\begin{array}{l} \downarrow \\ a=2 \\ a=2.001 \end{array} \quad \begin{array}{l} f(a) \approx 0.69315 \\ f(a) \approx 0.69265 \end{array}$$

Andrew Ng

Computation graph :

This diagram below states that

1. blue arrow depicts how the computation is done from left to right

2. Updation of weights (using derivatives is done from right to left)

Computation Graph

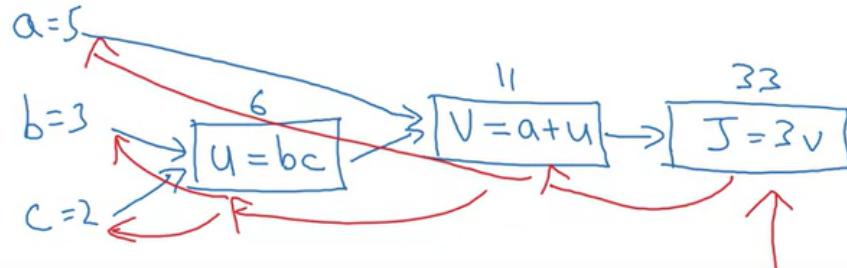
$$J(a, b, c) = 3(a + bc) = 3(5 + 3 \cdot 2) = 33$$

$\underbrace{}_v$
 $\underbrace{}_J$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



Derivatives with a Computation graph :

$$\begin{array}{c}
 a = 5 \\
 b = 3 \\
 c = 2
 \end{array}
 \xrightarrow{\quad b \quad}
 \boxed{u = bc} \xrightarrow{\quad v = a+u \quad}
 \boxed{v = 11} \xrightarrow{\quad J = 3v \quad}
 \boxed{J = 33}$$

[Effect of small change in
v on J]

$$\left[\frac{dJ}{dv} \right] = ? \\
 = \frac{0.003}{0.001} = 3.$$

(1)

$$\begin{cases}
 J = 3v \\
 v = 11 \rightarrow 11.001 \\
 J = 33 \rightarrow 33.003
 \end{cases}$$

This is basically called as

Because dJ is common for all.

$\frac{dJ}{da}$, $\frac{dJ}{db}$, $\frac{dJ}{du}$, etc., so we call
i.e. da , db , du , respectively.

i.e.

$$\text{Now } \left[\frac{dJ}{da} \right] = ? = \frac{0.003}{0.001} = 3.$$

$a = 5 \rightarrow 5.001$
 $v = 11 \rightarrow 11.001$
 $J = 33 \rightarrow 33.003$

This is written according to
rule as

chain
proof.

$$\begin{aligned}
 \frac{dJ}{da} &= \frac{dv}{da} \cdot \frac{dJ}{dv} \\
 &= 1 \cdot 3 \text{ from (1)} \\
 &= 3.
 \end{aligned}$$

a effects $\rightarrow v \rightarrow J$
(effects)

$$\frac{dv}{da} = ? \Rightarrow a \rightarrow 5 \rightarrow 5.001 \\
 \frac{0.001}{0.001} = 1$$

$v \rightarrow 11 \rightarrow 11.001$

Now

$$\frac{dJ}{db} = \frac{du}{db} \cdot \frac{dJ}{du}$$

$$\frac{dJ}{du} = \frac{dv}{du} \cdot \frac{dj}{dv}$$

$$= 1 \cdot 3$$

$$= 3.$$

$$u = 6 \rightarrow 6.001$$

$$v = 11 \rightarrow 11.00$$

$$\frac{dv}{du} = 1$$

$$db = \frac{dJ}{db} = \frac{du}{db} \cdot \frac{dv}{du} \cdot \frac{dj}{dv}$$

from (2)

$$b = 3 \rightarrow 3.001$$

$$u = 6 \rightarrow 6.002$$

$$= \frac{du}{db} \cdot \frac{dJ}{du}$$

$$= \frac{[0.002]}{0.001} \cdot 3 = 6.$$

dj/dv means : the effect of dj on small change in dv

similarly $dc = 9$

$$a = 5$$

$$\Delta a = 3$$

$$b = 3$$

$$\Delta b = 6$$

$$\begin{aligned} & \left. \begin{aligned} & b \\ & u = b - c \\ & \Delta u = 3 \end{aligned} \right\} \rightarrow \boxed{V = a + u} \\ & \rightarrow \boxed{J = 3V} \end{aligned}$$

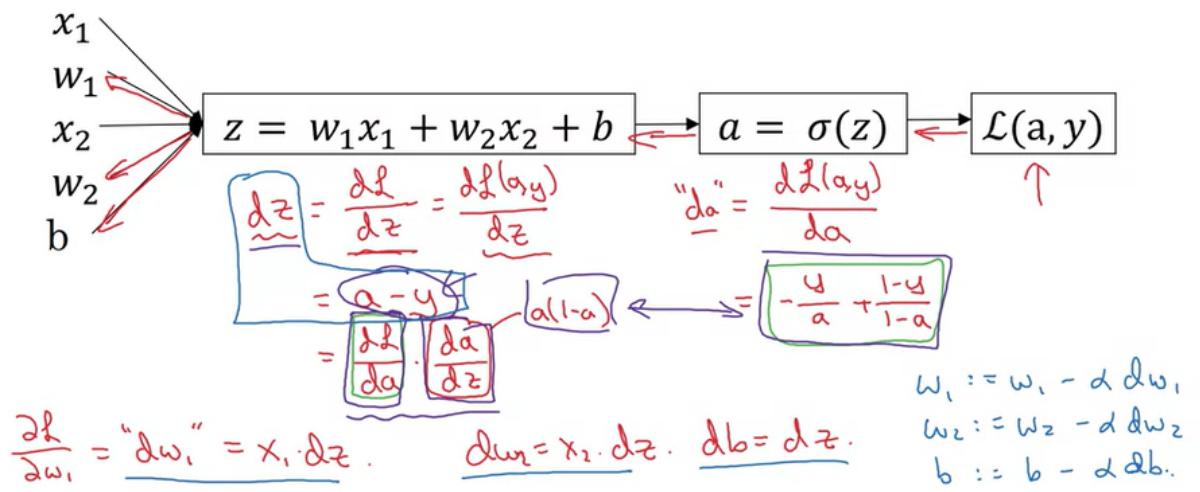
$$c = 2 \quad \Delta c = 9$$

effect of small change of

$$J = 3V$$

Finding the derivatives for Logistic Regression

Logistic regression derivatives



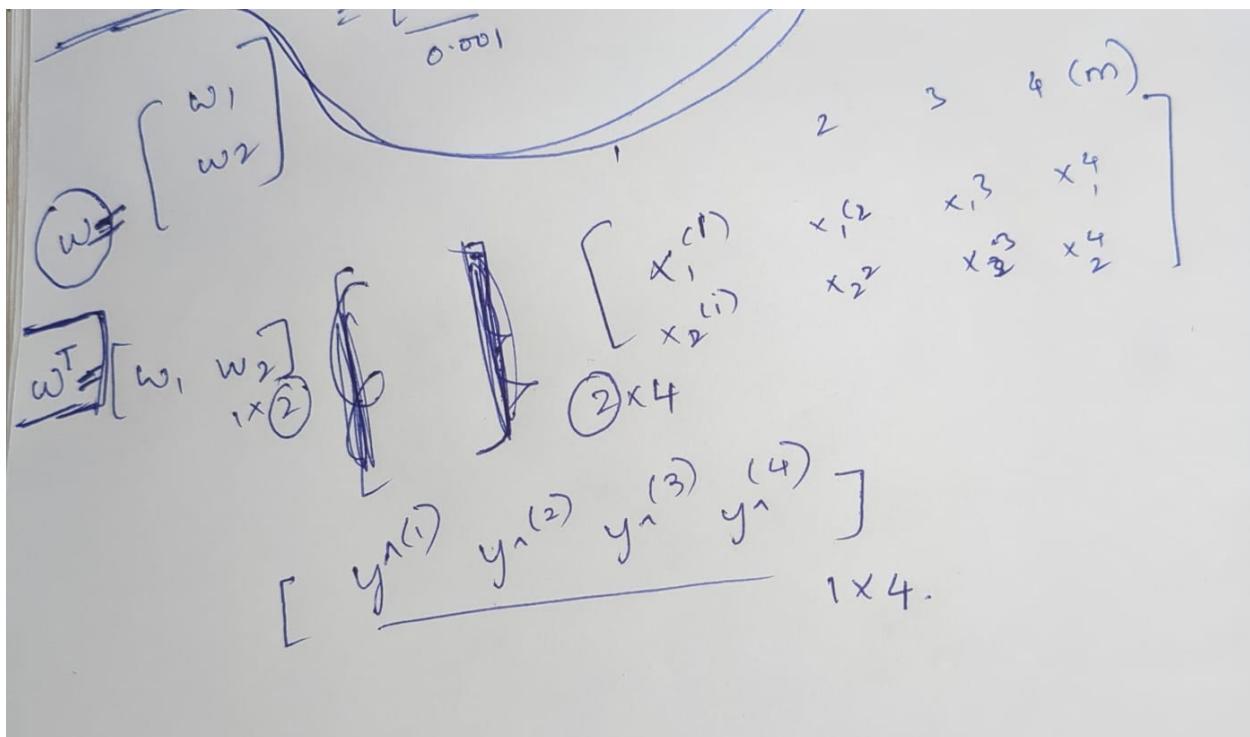
Andrew Ng

Note : The above example has only 2 features and only one training example.

$L(a,y)$ is the loss function, a is the predicted values, Loss function is same it is discussed above. apply derivative formulas lead to the above values.

Gradient Descent on m Examples

Note : here we have m examples, 2 features each



This is the sort of example considered

Logistic regression on m examples

$$\begin{aligned} J(w, b) &= \frac{1}{m} \sum_{i=1}^m \ell(a^{(i)}, y^{(i)}) \\ \rightarrow a^{(i)} &= \hat{y}^{(i)} = g(z^{(i)}) = g(w^T x^{(i)} + b) \end{aligned} \quad (\underline{x^{(i)}}, \underline{y^{(i)}})$$

$$\underline{dw_1^{(i)}}, \underline{dw_2^{(i)}}, \underline{db^{(i)}}$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_1} \ell(a^{(i)}, y^{(i)})}_{}$$

dw_1 value is computed as above. in the same way we calculate dw_2 and db

Algorithm to find gradient descent on M Examples

$$\begin{aligned}
 J &= 0; \quad \underline{\Delta w_1 = 0}; \quad \underline{\Delta w_2 = 0}; \quad \underline{\Delta b = 0} \\
 \text{For } i &= 1 \text{ to } m \\
 z^{(i)} &= \omega^T x^{(i)} + b \\
 a^{(i)} &= \sigma(z^{(i)}) \\
 J_t &= -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})] \\
 \underline{\Delta z^{(i)}} &= a^{(i)} - y^{(i)} \\
 \Delta w_1 &+= x_1^{(i)} \Delta z^{(i)} \quad \left. \right\} n=2 \\
 \Delta w_2 &+= x_2^{(i)} \Delta z^{(i)} \\
 \Delta b &+= \Delta z^{(i)} \\
 \Delta w_1/m &\leftarrow \Delta w_1 \\
 \Delta w_2/m &\leftarrow \Delta w_2 \\
 \Delta b/m &\leftarrow \Delta b
 \end{aligned}$$

$$\Delta w_1 = \frac{\partial J}{\partial w_1}$$

$$\begin{aligned}
 w_1 &:= w_1 - \alpha \underline{\Delta w_1} \\
 w_2 &:= w_2 - \alpha \underline{\Delta w_2} \\
 b &:= b - \alpha \underline{\Delta b}.
 \end{aligned}$$

Andrew Ng

note : we will compute dw1, dw2 and db for each example and take average.

note : all the formulas above are computed and explained in the Finding the derivatives for Logistic Regression section above.

and then update $w_1 := w_1 - (\text{learning rate}) * \Delta w_1$, vice versa for all.

we can observe that in the above example there is a for loop, in reality, if we have more features, then more associated (w's), which makes the use of another "for" loop for weights (w's). This makes our DL model work slow, so we use vectorization.

Vectorization

The use of for-loops drastically decreases the efficiency of the program, so while training the DL model, we might have to wait for a long amount of time just for it to train.

What is vectorization?

$$z = \underline{w^T x + b}$$

Non-vectorized:

```
z = 0  
for i in range(n - x):  
    z += w[i] * x[i]  
z += b
```

$$w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$w \in \mathbb{R}^{n_x}$$
$$x \in \mathbb{R}^{n_x}$$

Vectorized

$$z = \underbrace{\text{np.dot}(w, x)}_{w^T x} + b$$

GPU } SIMD - single instruction
CPU } multiple data.

But we can get over with for-loops **sometimes(not in all scenarios)** using vectorization technique.

For a glimpse see the time taken to execute code using for-loop and vectorization below:

```

In [1]: 1 import numpy as np

In [15]: 1 a = np.random.rand(1000000)
          2 b = np.random.rand(1000000)

In [16]: 1 import time
          2
          3 tic = time.time()
          4 x = np.dot(a,b)
          5 toc = time.time()
          6
          7 print("vectorized time", str(1000*(toc - tic)), "ms")
          8 print(x)
          9
         10 x = 0
         11 tic = time.time()
         12 for i in range(1000000):
         13     x += a[i]*b[i]
         14 toc = time.time()
         15
         16 print("Time using for loop", str(1000*(toc - tic)), "ms")
         17 print(x)

```

```

vectorized time 1.1518001556396484 ms
249930.35482235448
Time using for loop 337.20874786376953 ms
249930.35482235378

```

Vectorization was about 300 more faster than for loop.

credit **Andrew-NG (transcript_week2_course1)** :

But all the demos I did just now in the Jupiter notebook where actually on the CPU. And it turns out that both GPU and CPU have parallelization instructions. They're sometimes called SIMD instructions. This stands for a **single instruction multiple data**. But what this basically means is that, if you use built-in functions such as this np.function or other functions that don't require you explicitly implementing a for loop. It enables Python to take much better advantage of parallelism to do your computations much faster. And this is true both computations on CPUs and computations on GPUs. It's just that GPUs are remarkably good at these SIMD calculations but CPU is actually also not too bad at that. Maybe just not as good as GPUs. You're seeing how vectorization can significantly speed up your code. The rule of thumb to remember is whenever possible, avoid using

explicit for loops. Let's go onto the next video to see some more examples of vectorization and also start to vectorize logistic regression.

Neural Network programming guideline: Whenever possible, avoid explicit for-loops.

Vectors and matrix valued functions

Say you need to apply the exponential operation on every element of matrix/vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$\rightarrow u = np.zeros((n, 1))$
 $\rightarrow \boxed{\text{for } i \text{ in range}(n):} \leftarrow$
 $\rightarrow u[i] = \text{math.exp}(v[i])$

import numpy as np
 $u = np.exp(v) \leftarrow$
 \nearrow
 $np.log(u)$
 $np.abs(u)$
 $np.maximum(v, 0)$
 $v**2$
 v/v

So, in the above figure we can see that left side is non-vectorized, and right hand side is vectorized.

From this Algorithm(finding gradient descent on m examples)

Logistic regression derivatives

$$\begin{aligned}
 J &= 0, \quad \boxed{\cancel{dw1 = 0}, \cancel{dw2 = 0}}, \quad db = 0 \quad \cancel{dw} = np.zeros((n_x, 1)) \\
 \rightarrow \text{for } i &= 1 \text{ to } m: \\
 z^{(i)} &= w^T x^{(i)} + b \\
 a^{(i)} &= \sigma(z^{(i)}) \\
 J &+= -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})] \\
 \downarrow \text{for } j &= 1 \dots n_x \\
 dz^{(i)} &= a^{(i)} - y^{(i)} \\
 dw_1 &+= x_1^{(i)} dz^{(i)} \quad n_x = 2 \quad dw += x^{(i)} dz^{(i)} \\
 dw_2 &+= x_2^{(i)} dz^{(i)} \\
 db &+= dz^{(i)} \\
 J &= J/m, \quad \boxed{dw_1 = dw_1/m, \quad dw_2 = dw_2/m, \quad db = db/m} \quad \cancel{dw /= m}
 \end{aligned}$$

lets try to get rid of one for loop. In this figure dw1, dw2 are just two weights [**but in reality there could be a lot of them (equal to the number of features)**, (**Here, for simplicity we have just taken 2 features each of m training examples**)]. So this weights will require a for loop. Lets try to get rid of this specific one.

visualizing

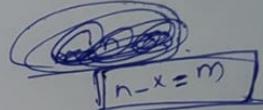
$$dw = x(i) \cdot dz(i).$$

w is a $(m \times 1)$ matrix

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \end{bmatrix} \quad m \times 1.$$

dw ~~has~~ has the same dimensions as w.

$$x = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \quad \leftarrow (x \text{ looks something like this})$$



so lets ~~recreate~~ recreate w's dimensions
 $dw = np.zeros(n - x, 1)$.

now $x^{(i)}$ looks like

$$\begin{bmatrix} \text{feature 1} \\ \text{feature 2} \end{bmatrix} \quad 2 \times 1$$

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

$$a^i \Rightarrow \sigma(z^{(i)}) \quad \& \quad z^{(i)} = w^T x^{(i)} + b$$

~~predicted value.~~

~~when we apply sigmoid function of $z^{(i)}$~~

we get \hat{a}^i .

\hat{a}^i is $(1, 1)$ shape matrix.

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

will result in $(1, 1)$ shape matrix.

so we can say

$$x^{(i)} \cdot dz^{(i)}$$

$$dw =$$

$$\begin{bmatrix} dw_1 \\ dw_2 \end{bmatrix} = \left[\begin{bmatrix} \text{feature-1 } x_1 \\ x_2 \end{bmatrix}_{2 \times 1} \right] \begin{bmatrix} dz^{(i)} \end{bmatrix}_{1 \times 1}$$

is nothing but $x^{(i)}$.

and then we are left with one thing, divide each element of dw with m.

for doubts refer again:

1.

2.

Vectorizing the entire Gradient descent :

Logistic regression derivatives

```

 $J = 0, \boxed{dw_1 = 0, dw_2 = 0}, db = 0 \quad dw = np.zeros((n_x, 1))$ 
 $\rightarrow \text{for } i = 1 \text{ to } m:$ 
 $z^{(i)} = w^T x^{(i)} + b$ 
 $a^{(i)} = \sigma(z^{(i)})$ 
 $J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$ 
 $\downarrow$ 
 $\text{for } j = 1 \dots n_x:$ 
 $\boxed{dw_1 += x_1^{(i)} dz^{(i)}, dw_2 += x_2^{(i)} dz^{(i)}, db += dz^{(i)}} \quad n_x = 2 \quad dw += x^{(i)} dz^{(i)}$ 
 $J = J/m, \boxed{dw_1 = dw_1/m, dw_2 = dw_2/m, db = db/m} \quad dw /= m.$ 

```

But from the above code snippet, we can see that, there is still a for loop, which effects our execution time, lets try to get a code without for-loop.

Vectorizing Logistic Regression

$$\begin{aligned}
&\rightarrow z^{(1)} = w^T x^{(1)} + b & z^{(2)} = w^T x^{(2)} + b & z^{(3)} = w^T x^{(3)} + b \\
&\rightarrow \boxed{a^{(1)}} = \sigma(z^{(1)}) & \boxed{a^{(2)}} = \sigma(z^{(2)}) & \boxed{a^{(3)}} = \sigma(z^{(3)}) \\
&\overline{X} = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} & \overline{w}^T \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} & \overline{b} \in \mathbb{R}^{1 \times m} \\
&\overline{Z} = \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \overline{w}^T \overline{X} + \begin{bmatrix} b & b & \dots & b \end{bmatrix}_{1 \times m} = \begin{bmatrix} w^T x^{(1)} + b \\ w^T x^{(2)} + b \\ \vdots \\ w^T x^{(m)} + b \end{bmatrix} & \rightarrow \overline{Z} = \underbrace{\overline{w} \cdot \overline{T} \cdot \overline{X}}_{\mathbb{R}^{1 \times 1}} + \underbrace{\overline{b}}_{\mathbb{R}^{1 \times 1}} & \text{"Broadcasting"} \\
&\overline{A} = \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(m)} \end{bmatrix} = \sigma(\overline{Z}) & \uparrow &
\end{aligned}$$

Andrew Ng

1. We are computing the predicted values, in each iteration ($\hat{z}^{(i)}$), this can be vectorized as in the above figure.
2. computing the sigmoid in each iteration can as well be vectorized, $\sigma(z)$

we still have to find vectorized calculation for dz , dw and db .

Vectorizing Logistic Regression

$$\begin{aligned}
 dz^{(1)} &= a^{(1)} - y^{(1)} & dz^{(2)} &= a^{(2)} - y^{(2)} & \dots \\
 dz &= [dz^{(1)} \ dz^{(2)} \dots dz^{(m)}] & & & \\
 A &= [a^{(1)} \dots a^{(m)}]. \quad Y = [y^{(1)} \dots y^{(m)}] \\
 \rightarrow dz &= A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots]
 \end{aligned}$$

$$\begin{aligned}
 \cancel{\frac{dw}{dw}} &\rightarrow dw = 0 \\
 dw &+= \cancel{x^{(1)} dz^{(1)}} \\
 dw &+= \cancel{x^{(2)} dz^{(2)}} \\
 &\vdots \\
 dw/m &= m
 \end{aligned}
 \quad
 \begin{aligned}
 db &= 0 \\
 db &+= \cancel{dz^{(1)}} \\
 db &+= \cancel{dz^{(2)}} \\
 &\vdots \\
 db &+= \cancel{dz^{(m)}} \\
 db/m &= m
 \end{aligned}$$

$$\begin{aligned}
 db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\
 &= \frac{1}{m} \text{np.sum}(dz) \\
 dw &= \frac{1}{m} X dz^T \\
 &= \frac{1}{m} \left[\begin{array}{c|c}
 x^{(1)} & \dots x^{(m)} \\
 \hline
 1 & 1
 \end{array} \right] \left[\begin{array}{c} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{array} \right] \\
 &= \frac{1}{m} \left[\underbrace{x^{(1)} dz^{(1)}}_{n \times 1} + \dots + \underbrace{x^{(m)} dz^{(m)}}_{n \times 1} \right]
 \end{aligned}$$

Andrew Ng

3. dz can be written as $A - Y$.
4. db as the average over m samples on the sum of dz .
5. for dw :

we are computing $dw += x(i)*dx(i)$, every time, this can be vectorized as

~~why can we write~~
for $i=1 \text{ to } m$ ~~we~~ $dw^+ = x^{(i)} \cdot dz^{(i)}$ as

$$dw = x \cdot dz^T$$

for 1st iteration

$$dw = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot [1] = \begin{bmatrix} 2 \\ 2 \\ 3 \end{bmatrix}$$

$$dw^+ = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \cdot [2] = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

$$dw = \begin{bmatrix} 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 2 \\ 4 \end{bmatrix} \\ = \begin{bmatrix} 4 \\ 7 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \end{bmatrix} = \underline{\underline{dz}}^T$$

$$dz^T = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Now:

$$\begin{bmatrix} 2 & 1 \\ 3 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}_{2 \times 1}$$

$$\begin{bmatrix} 2+2 & 3+4 \end{bmatrix}$$

$$\begin{bmatrix} 4 & 7 \end{bmatrix}. -\textcircled{2}$$

we can observe

that

$\textcircled{1}$ & $\textcircled{2}$ are equal

so we can write

$$dw = \alpha X \cdot dz^T$$

Implementing Logistic Regression

```

J = 0, dw1 = 0, dw2 = 0, db = 0
for i = 1 to m:
    z(i) = wTx(i) + b ←
    a(i) = σ(z(i)) ←
    J += -[y(i) log a(i) + (1 - y(i)) log(1 - a(i))]
    dz(i) = a(i) - y(i) ←
    dw1 += x1(i)dz(i) } dw1 = X(i) * dz(i)
    dw2 += x2(i)dz(i)
    db += dz(i)
J = J/m, dw1 = dw1/m, dw2 = dw2/m
db = db/m

```

```

for iter in range(1000): ←
    Z = wTX + b
    = np.dot(w.T, X) + b
    A = σ(Z)
    dZ = A - Y
    dw = 1/m * X * dZT
    db = 1/m * np.sum(dZ)
    w := w - α dw
    b := b - α db

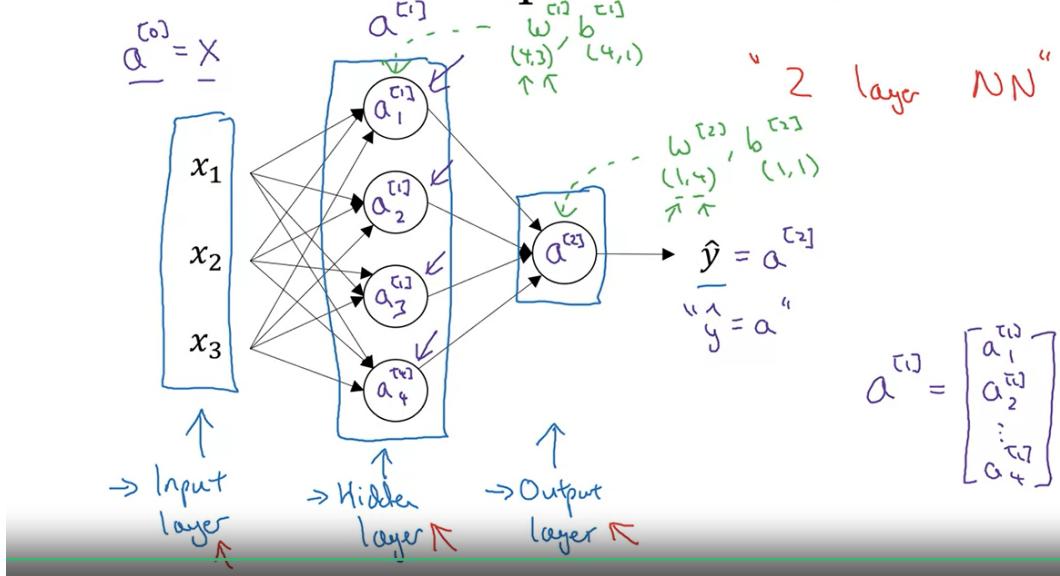
```

This is the highly optimized representation of the logistic regression to calculate derivatives in backward propagation.

Week 3 : Shallow Neural Networks

Neural Network Representation

Neural Network Representation

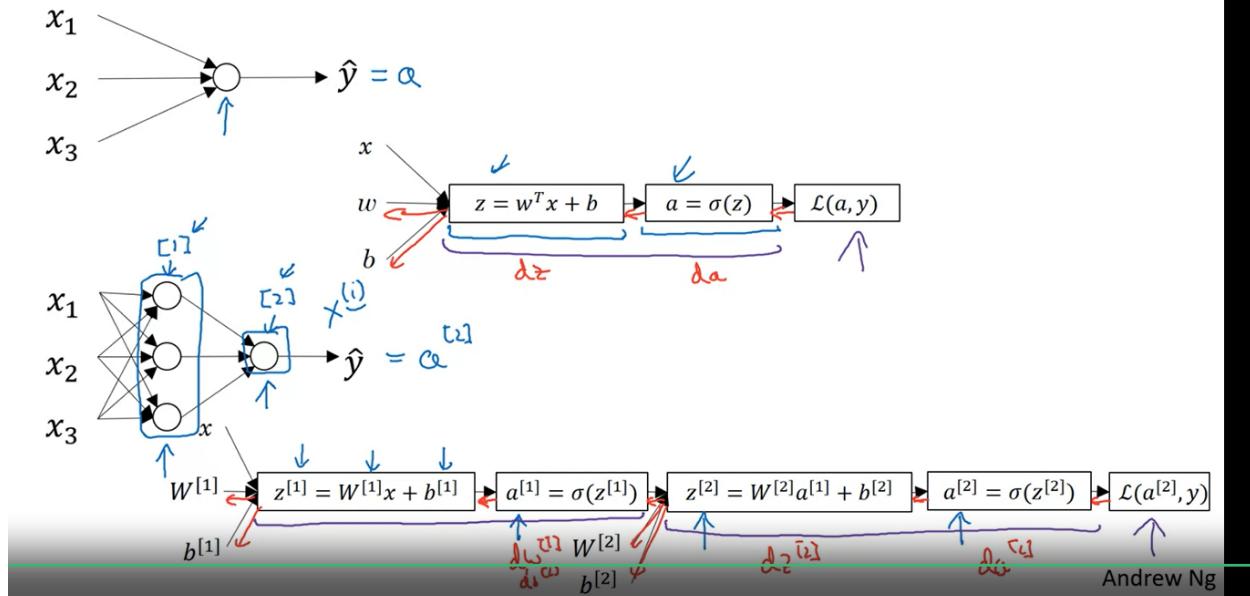


And ARMOURY CRATE

$a[0]$, represents that it is the input that will be passed on to the other layer, a is also known as activation.

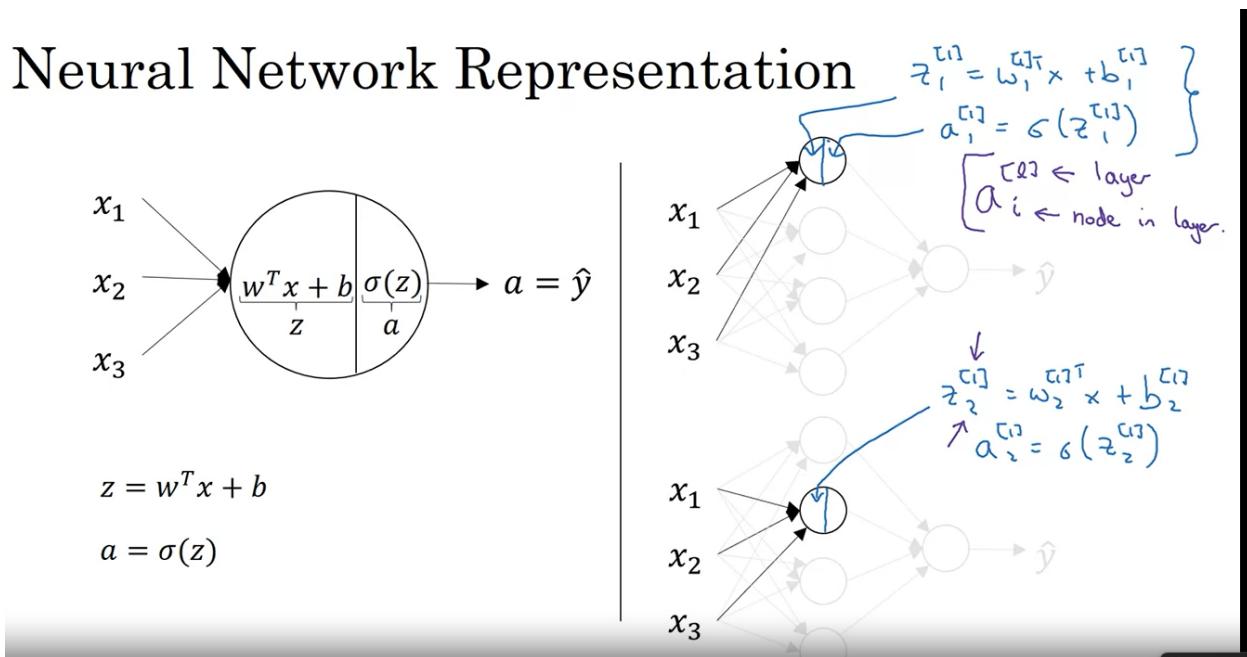
- In general Terms, a means that, this input will be passed on the next layer.
- The above image is a two layer neural network, (it should have been 3 layers), but by convention it is called 2 layered neural network.
- the first layer is called the input layer and the last one is obviously called the output layer.
- $a[1]$ and $a[2]$ are associated with parameters, w and b .

What is a Neural Network?



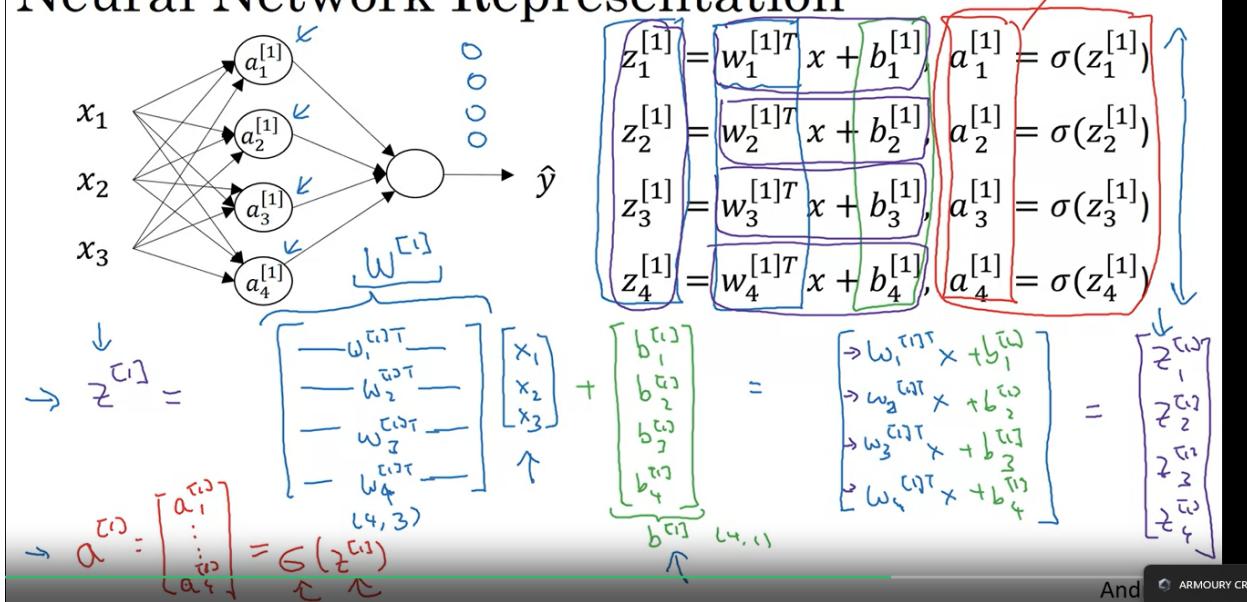
So, for each node in $a[1]$, we calculate z and corresponding sigma value.
as shown in the figure above.

Neural Network Representation



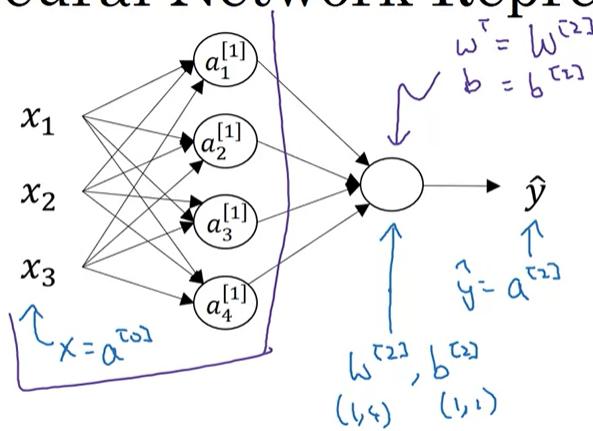
This gives a better and clear picture of what is going on.

Neural Network Representation



And ARMOURY CRA

Neural Network Representation learning



Given input x :

$$\rightarrow z^{[1]} = W^{[1]} \underset{(4,1)}{x} + b^{[1]} \underset{(4,1)}{a^{[0]}}$$

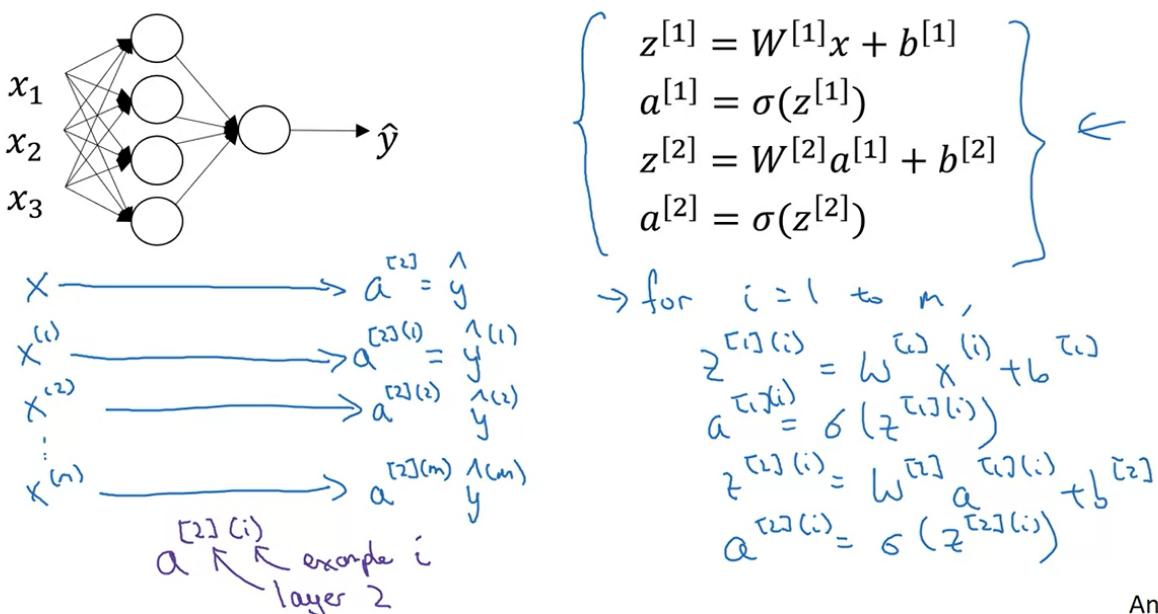
$$\rightarrow a^{[1]} = \sigma(z^{[1]}) \underset{(4,1)}{a^{[1]}}$$

$$\rightarrow z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \underset{(1,1)}{w^{[2]}} \underset{(4,1)}{a^{[1]}} \underset{(1,1)}{b^{[2]}}$$

$$\rightarrow a^{[2]} = \sigma(z^{[2]}) \underset{(1,1)}{a^{[2]}}$$

Until this point, we were focusing on features of one sample, now let's consider we have more than one examples

Vectorizing across multiple examples



And

This is how it looks like, when there are multiple examples, we can observe that $X(i)$ matrix represents features at sample i , and $a[2](i)$ represents predicted value at 2nd layer, that is the output, and i is that particular index.

and the **CODE** to the right, is how can we write code for multiple training examples.

Absolutely we need to vectorize it,

Vectorizing across multiple examples

for i = 1 to m:

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

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

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

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

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(n)} \\ | & | & & | \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

$$z^{[i]} = w^{[i]} x + b^{[i]}$$

$$= \mu^{(1)} = \epsilon(z^{(1)})$$

$$\Rightarrow A^{[2]} = \{1\}^{[2]}$$

$$\Rightarrow A^{(2)} = \zeta (z^{(2)})$$

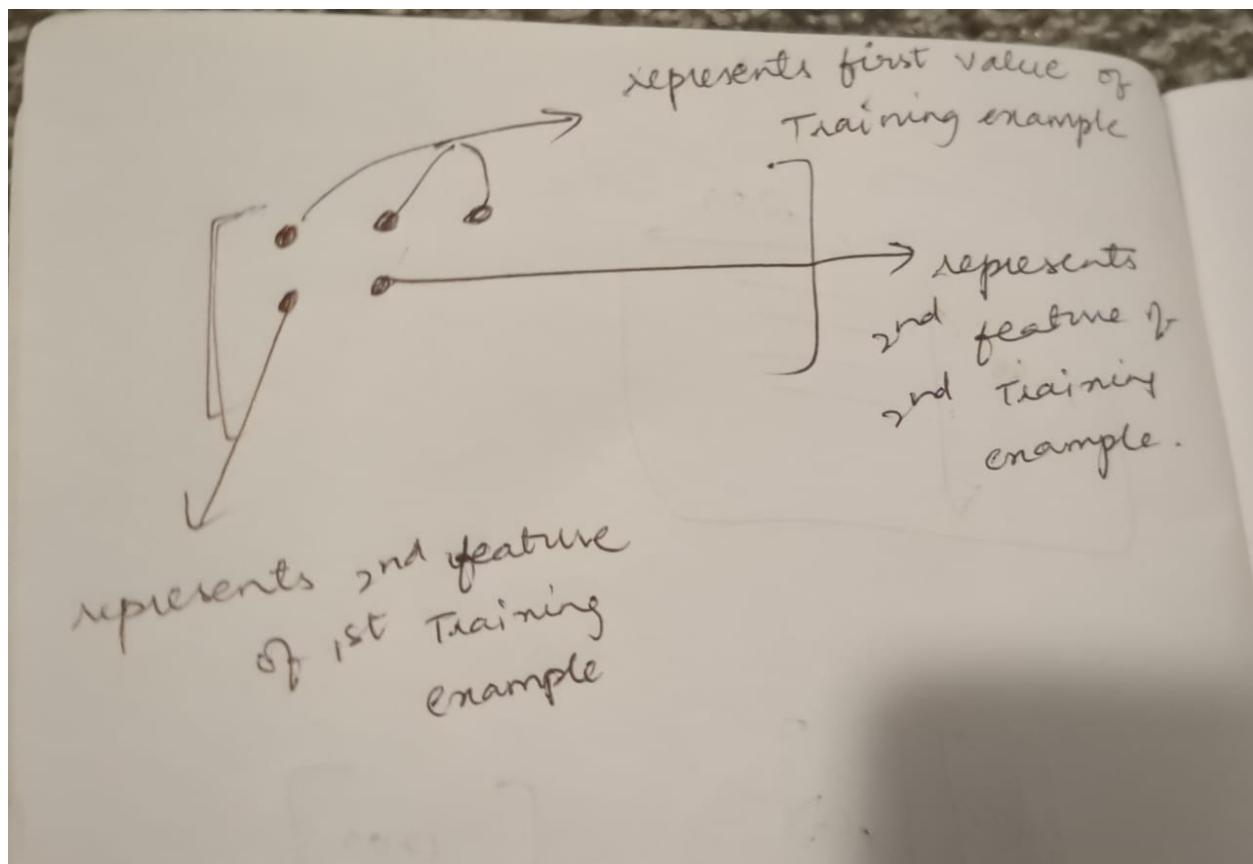
$$Z^{[1]} = \begin{bmatrix} 1 & z^{1} & z^{[1](2)} & \cdots & z^{[1](m)} \\ & 1 & 1 & & 1 \end{bmatrix}$$

$A^{[1]} = \begin{bmatrix} \bullet & \bullet & \bullet & \cdots & \bullet \\ z^{1} & a & z^{[1](m)} & \cdots & a \\ \bullet & a & a & \cdots & a \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}$

↑ hidden units

And

CODE to the right is how you vectorize it.



note :

$a[i]$, $z[i]$, in here $[] \rightarrow$ represents values at that particular layer

$a(i)$, $z(i)$, in here $() \rightarrow$ represents the value at that particular sample.

Justification for the Vectorized Implementation :

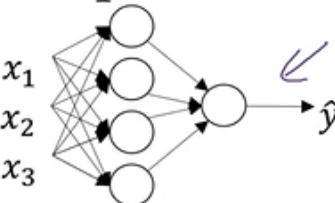
Justification for vectorized implementation

$$\begin{aligned}
 z^{(1)(1)} &= w^{(1)} x^{(1)} + b^{(1)}, \\
 z^{(1)(2)} &= w^{(1)} x^{(2)} + b^{(1)}, \\
 z^{(1)(3)} &= w^{(1)} x^{(3)} + b^{(1)}
 \end{aligned}$$

$w^{(1)} = \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix}$ $w^{(1)} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$ $w^{(1)} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$ $w^{(1)} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$

$w^{(1)} \begin{bmatrix} 1 & x^{(1)} & x^{(2)} & x^{(3)} \dots \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} z^{(1)(1)} \\ z^{(1)(2)} \\ z^{(1)(3)} \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} = z^{(1)}$

Recap of vectorizing across multiple examples



$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \end{bmatrix}$$

for i = 1 to m
 $\rightarrow z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$
 $\sim a^{[1](i)} = \sigma(z^{[1](i)})$
 $z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$
 $a^{[2](i)} = \sigma(z^{[2](i)})$

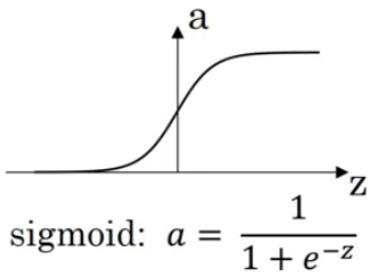
$Z^{[1]} = W^{[1]}X + b^{[1]}$
 $A^{[1]} = \sigma(Z^{[1]})$
 $Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$
 $A^{[2]} = \sigma(Z^{[2]})$

Andrew Ng

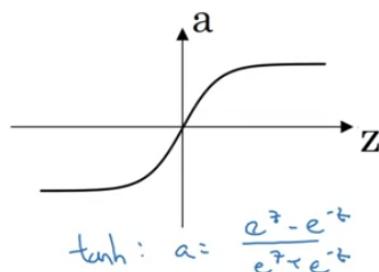
Thus can be written as above.

Activation Functions

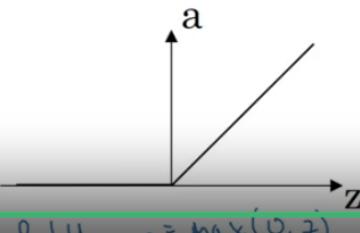
Pros and cons of activation functions



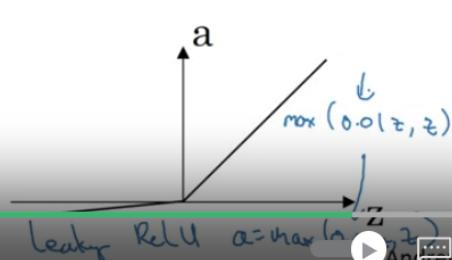
$$\text{sigmoid: } a = \frac{1}{1 + e^{-z}}$$



$$\tanh: a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



9:01 U 10:56



Leaky ReLU $a = \max(0.01z, z)$

these are the different activation functions available,

- do not use sigmoid except for output layer of the logistic regression.
- tanh almost always works well than sigmoid, so use it whenever possible. but we have to use the sigmoid in the output layer of logistic regression because, because we expect values from 0 to 1.'
- RELU is the best alternative possible, because for sigmoid, at extreme values the slope almost reaches 0, which makes the learning a lot slower, but for Rectified Linear Unit (RELU), the slope is 1 for all values greater than 0, which makes the learning faster.

If you observe clearly the slope at point zero is infinity, and the probability of output getting 0 is really low, hence we trust RELU.

The curve for sigmoid activation and tanh activation reaches near zero at some point, making learning slower, but why ?

NOTE : just in case, if you are wondering what slope has to do with activation functions, remember in backward propagation we calculate how the values of output changes by the small change in the input, which is nothing but slope right, if the curve we choose has low slope then definitely the gradient descent will be slower.

Why do you need non-linear activation functions ?

The below link explains why we need non-linear activation functions ?

<https://www.v7labs.com/blog/neural-networks-activation-functions#3-types-of-neural-networks-activation-functions>

when we use the linear activation function

$$g(z) = z$$

$g(z)$ is also called as identity activation function, provides the output as it is.

NOTE : Honestly, not intuitively clear why to use activation function.

- they are used to notice important patterns
- if you don't use activation functions, basically at the end, what we are seeing is just a linear regression model of machine learning.

Linear Activation Function:

- In the case of a purely linear activation function, the derivative is a constant, and it's always 1. When backpropagating through multiple layers, this constant gradient can cause the gradients to either vanish or explode, making it difficult for the network to effectively learn and adapt.

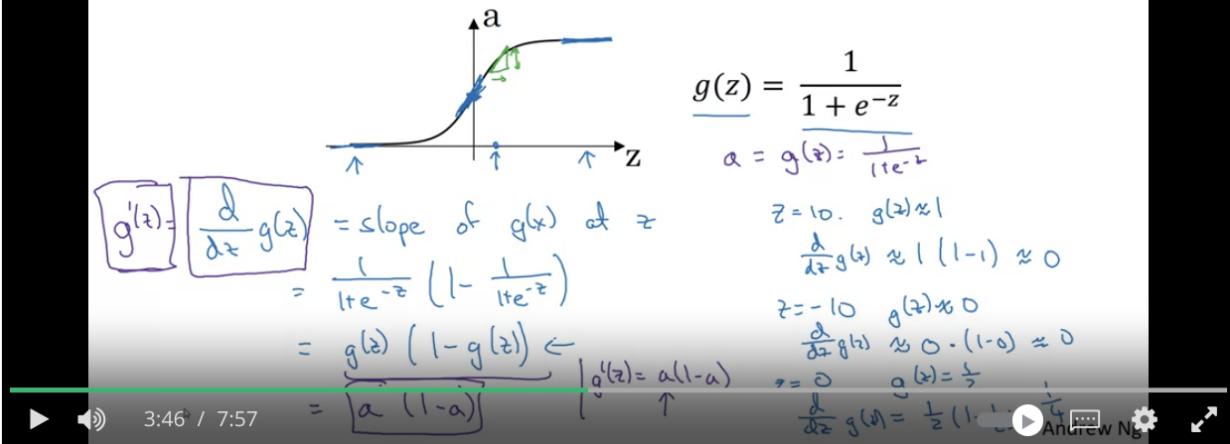
ReLU to the Rescue:

- ReLU, by setting negative values to 0, introduces non-linearity and helps mitigate the vanishing gradient problem. For positive values, the gradient is a constant 1, allowing for effective gradient propagation and weight updates. This makes ReLU well-suited for deep networks, where the ability to effectively learn and adapt through many layers is crucial.

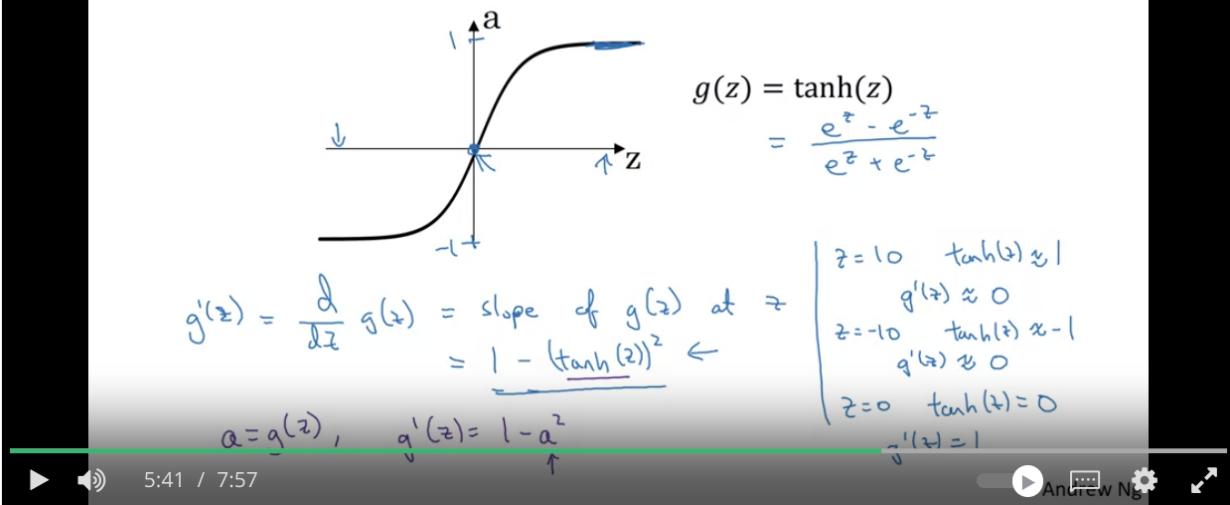
While ReLU helps address the vanishing gradient problem, it's not without its challenges, such as the potential for neurons to become "dead" (always outputting zero for any input less than or equal to zero). There are variations of ReLU, like Leaky ReLU and Parametric ReLU, designed to handle some of these issues. The choice of activation function is an essential consideration in designing and training neural networks.

Derivatives for Activation Functions

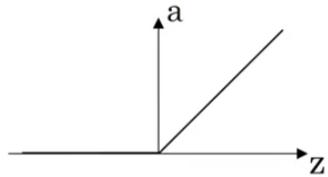
Sigmoid activation function



Tanh activation function



ReLU and Leaky ReLU



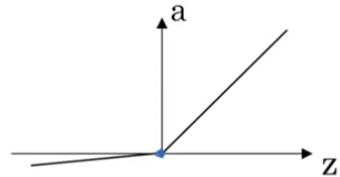
ReLU

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

$$\rightarrow g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

~~undefined if $z=0$~~

$\tilde{z} = 0.0000\cdots 0$



Leaky ReLU

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Andrew Ng

Deep Layer neural Network

Forward Propagation in a Deep Network

x_1, x_2, x_3

A^0

$A^1: z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$
 $a^{[1]} = g^{[1]}(z^{[1]})$
 $(z^{[1]}, W^{[1]}, a^{[1]}, b^{[1]})$
 $g^{[1]}(z^{[1]})$

$A^2: z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$
 $a^{[2]} = g^{[2]}(z^{[2]})$
 $(z^{[2]}, W^{[2]}, a^{[2]}, b^{[2]})$
 $g^{[2]}(z^{[2]})$

$A^3: z^{[3]} = W^{[3]} a^{[2]} + b^{[3]}$
 $a^{[3]} = g^{[3]}(z^{[3]})$
 $(z^{[3]}, W^{[3]}, a^{[3]}, b^{[3]})$
 $g^{[3]}(z^{[3]})$

$A^4: z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}$
 $a^{[4]} = g^{[4]}(z^{[4]})$
 $(z^{[4]}, W^{[4]}, a^{[4]}, b^{[4]})$
 $g^{[4]}(z^{[4]})$

$\hat{Y} = g^{[4]}(z^{[4]}) = A^{[4]}$

Vertical:
 $z^{[l+1]} = W^{[l+1]} A^{[l]} + b^{[l+1]}$
 $A^{[l+1]} = g^{[l+1]}(z^{[l+1]})$
 $(z^{[l+1]}, W^{[l+1]}, A^{[l+1]}, b^{[l+1]})$
 $g^{[l+1]}(z^{[l+1]})$

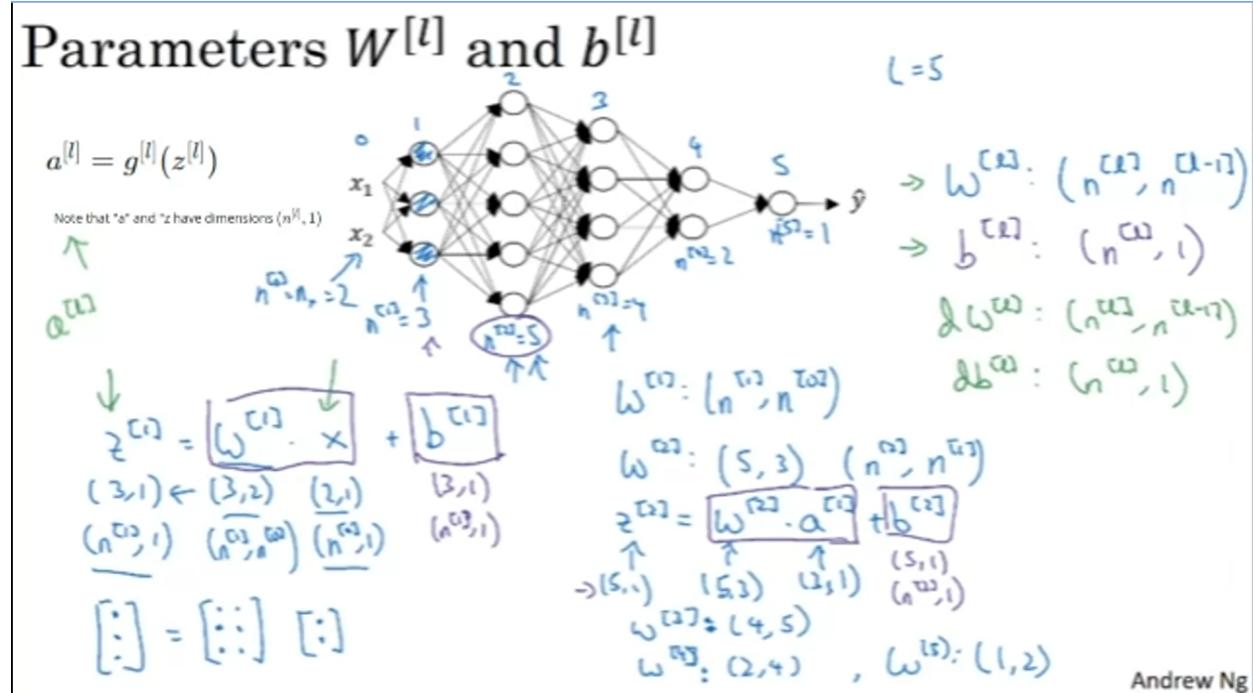
for $l=1..4$

Andrew Ng

left part of the image, shows implementing the neural network for single example, and the right part shows the way to implement neural network for multiple examples.

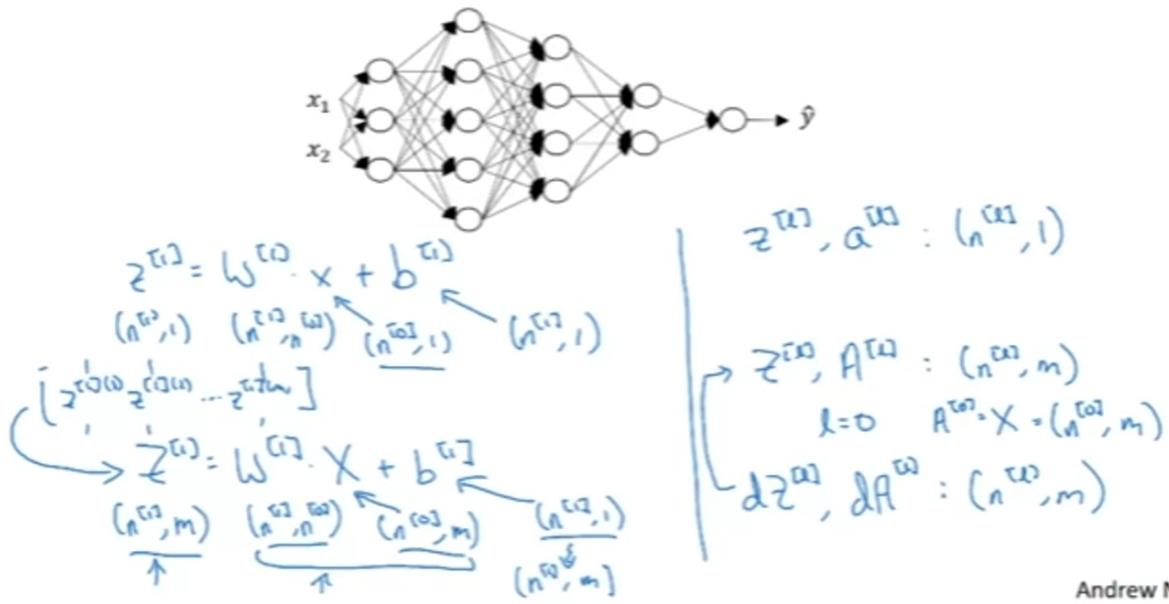
There is a need for using for loop, to calculate the Z , and A using for loop, for each and every layer. There is no known way to vectorize it currently.

Getting the Dimensions Right



These are the dimensions for the single example

Vectorized implementation

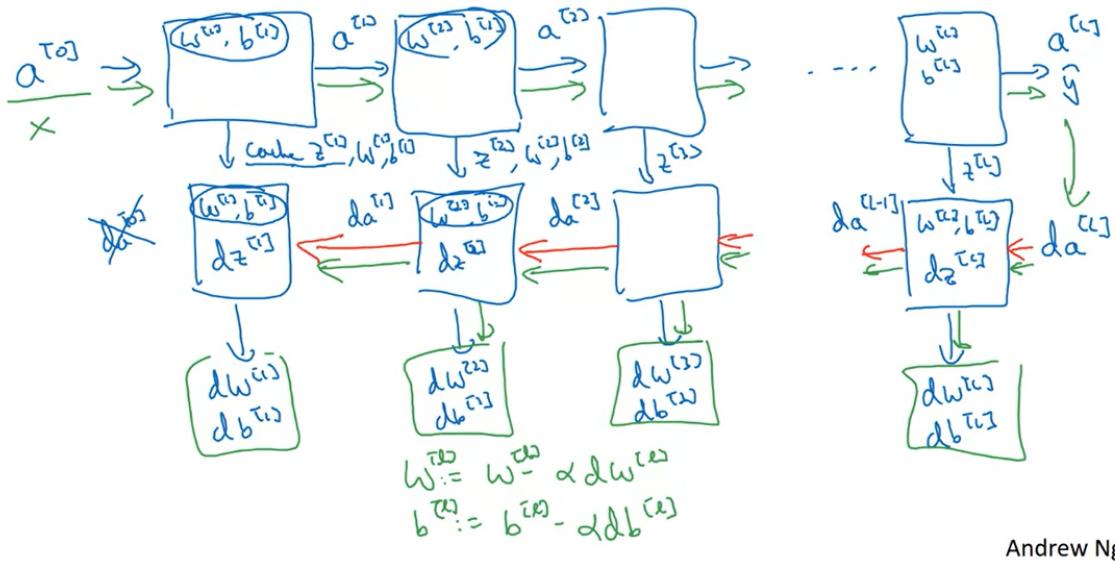


Andrew Ng

Dimensions for multiple examples

Building blocks of neural networks

Forward and backward functions



In each and every step, there is a forward propagation and backward propagation step going on, this is the visual representation.