

Deep Learning Specialization

1. Neural Networks & Deep Learning

- Linear Activation

- It allows multiple outputs, not just binary like step functions.
- It has 2 major problems:
 - not possible to use backpropagation to train models
 - all layers of NN collapse into one.

- Non-Linear Activations

- They allow complex mappings between the network's inputs & outputs. They actually address the 2 major fallbacks of linear activations.

- Overfitting

- model does better on the training set than on test set.
i.e. model follow the noise, errors too closely.
- Most popular solutions for overfitting include:
 - cross-validation
 - train with more data
 - remove features
 - early stopping
 - regularization
 - ensembling (bagging - boosting)

- Formulas for computing derivatives

- **tanh** activation is zero centred, it usually works better than sigmoid for hidden layers because the mean of its output is closer to zero; so it centres the data better for the next layer.
- For **tanh** activation function, using `np.random.randn(.,.) * 1000` will cause the inputs of tanh to be also very large, thus cause gradients to be close to zero. The optimization algorithm will thus become slow.

- Initializing \mathbf{w} to zero will result a symmetric \mathbf{w} , which is typically useless having all rows being the same. That's why we use `np.random.rand(a,b) * 0.01` yet with `b = np.zeros(a,b)`
- Initializing both \mathbf{w} & \mathbf{b} to be zero will cause each neuron in the first layer to perform the same computation. So even after multiple iterations of gradient descent each neuron in the layer will be computing the same thing as other neurons.
- Logistic regression does not have a hidden layer.

* Formulas for computing derivatives

$$Z^{[1]} = W^{[1]} X + b^{[1]} \quad \text{"forward"}$$

$$A^{[1]} = g^{[1]}(Z^{[1]}) \rightarrow \text{Linear unit}$$

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]})$$

"Backward"

$$dZ^{[2]} = A^{[2]} - Y, \quad Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1, \text{keepdims}=\text{True})$$

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

$(n^{[1]}, m)$ \downarrow $(n^{[1]}, m)$
 element-wise

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims}=\text{True})$$

- Forward Propagation

- Forward propagation equations:

* Forward propagation

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

For the vectorized version, replace a with A
 z with Z

- getting matrix dimensions for w & b :

- Getting matrix dim

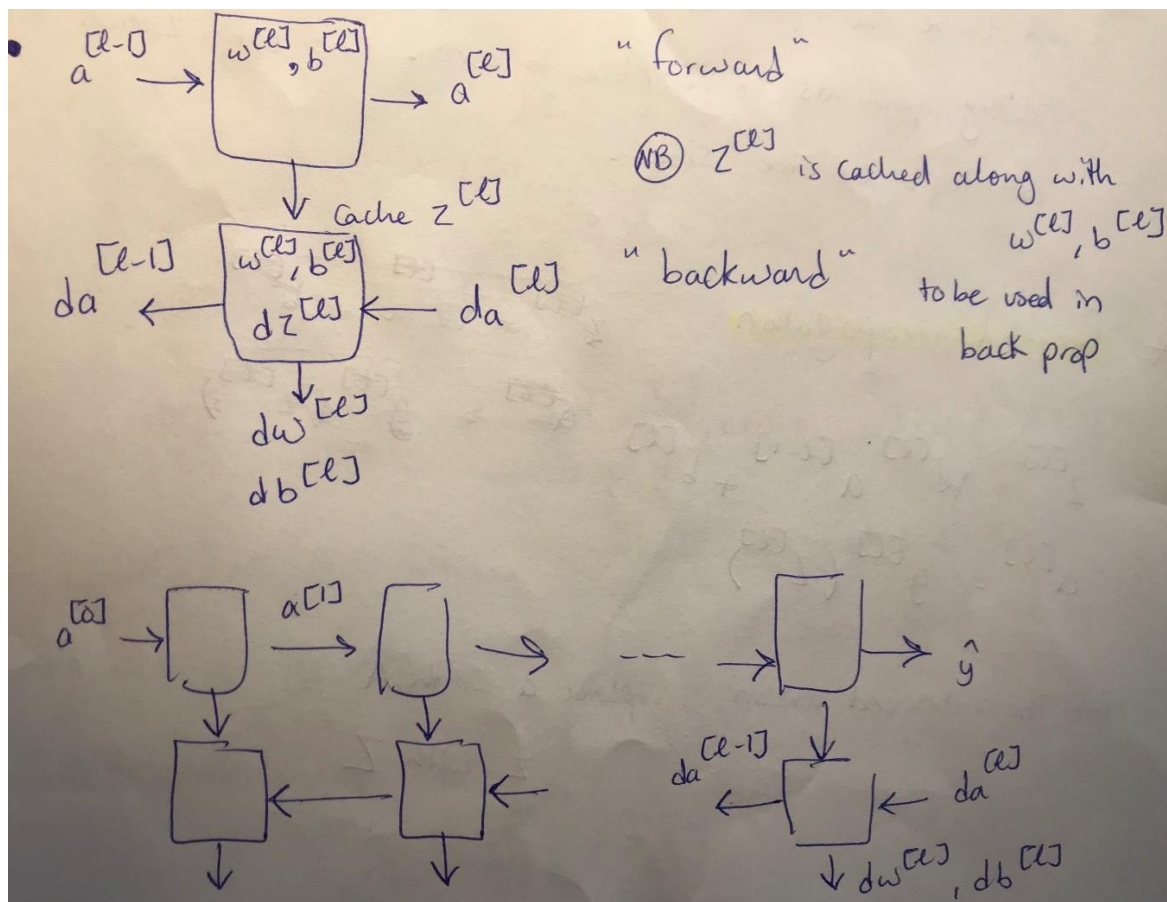
$$w^{[l]} : (n^{[l]}, n^{[l-1]}) \quad , \text{dim of } dw^{[l]} = \text{dim of } w^{[l]}$$

$$b^{[l]} : (n^{[l]}, 1) \quad , \text{dim of } db^{[l]} = \text{dim of } b^{[l]}$$

For the vectorized form of $z^{[l]}$ for example

$$\begin{bmatrix} | & | & & | \\ z^{[l](1)} & z^{[l](2)} & \dots & z^{[l](m)} \\ | & | & & | \end{bmatrix} \quad \text{so dim will be } (n^{[l]}, m)$$

- block diagram: how forward & backward props work



- summary of all equations:

• Again --

- Forward :
Prop input $a^{[l-1]}$
output $a^{[l]}$, cache $z^{[l]}$, $w^{[l]}$, $b^{[l]}$

- Backward :
Prop input $da^{[l]}$
output $da^{[l-1]}$, $dw^{[l]}$, $db^{[l]}$

$$dz^{[l]} = da^{[l]} * g^{[l]'}(z^{[l]})$$

$$dw^{[l]} = dz^{[l]} \cdot a^{[l-1]T}$$

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

$$da^{[l-1]} = w^{[l]T} \cdot dz^{[l]}$$

$$dz^{[l]} = w^{[l+1]T} \cdot dz^{[l+1]} * g^{[l+1]'}(z^{[l+1]})$$

Vectorized

$$dz^{[l]} = dA^{[l]} * g^{[l]'}(z^{[l]}) \rightarrow Y$$

$$dw^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{keepdims=True})$$

$$dA^{[l-1]} = w^{[l]T} \cdot dz^{[l]}$$

* For logistic reg $\ell(\hat{y}, y) \Rightarrow da^{[l]} = \frac{-y}{a} + \frac{1-y}{1-a}$

- **Hyperparameters vs Parameters**

- Hyperparameters: learning rate, # iterations, # hidden layers, choice of activation function.
- Parameters: w & b .
- The deeper layers of a NN are typically computing more complex features of the input than the earlier ones.
- Notation: $x^{(i)}$: i^{th} training example.
 $a_i^{[l]}$: i^{th} entry of the l^{th} layer's activation.
- To compute a function using a shallow network circuit, you will need a large network; where we measure size by # of logic gates in the network.
To compute it using a deep network circuit, you need only exponentially smaller network.