

COURSE 2: IMPROVING DEEP NEURAL NETWORKS

- Applied ML is **highly iterative process**.
- Intuition from one application to another are **not transferable**; it depends on data, CPUs (GPUs) used & many other factors.

* Train / Dev / test sets :

Data	-			
	- Train set	- Hold-out cross-validation set	- Test Set	
		- Development set 'Dev' set		

Previous Practice : train test train dev test
 70 / 30 80 / 20 / 20

However, with the amount of data (upto 1,000,000 points) the **size** of dev & test set is **decreasing rapidly**.

Eg: 1,000,000 : $\frac{10,000}{98/1/1\%}$, $\frac{10,000}{99.5/0.4/0.1\%}$.

Thumble

Thumble Rule: Make sure dev & test sets come from same distribution.

- No test set might also be okay, (only dev test).
- Dev Test set is set up to evaluate data on different models and find the best one for your dev set (and thus train set).
- Test set then runs on the chosen model to determine the accuracy of it & gives us an unbiased estimate of performance.
- In many places companies go use only 2 sets, train/test. usually Here, their test set actually works as hold-out cross validation set (wrong use of terminology) as we're overfitting the test train set data & not computing an unbiased estimate.
- However,

Dev Set:TRAIN SET

Max. ratio of our data set as we're training the model here.

DEV SET

Evaluate this data on different models to find perfect one for data.

Estimate is not

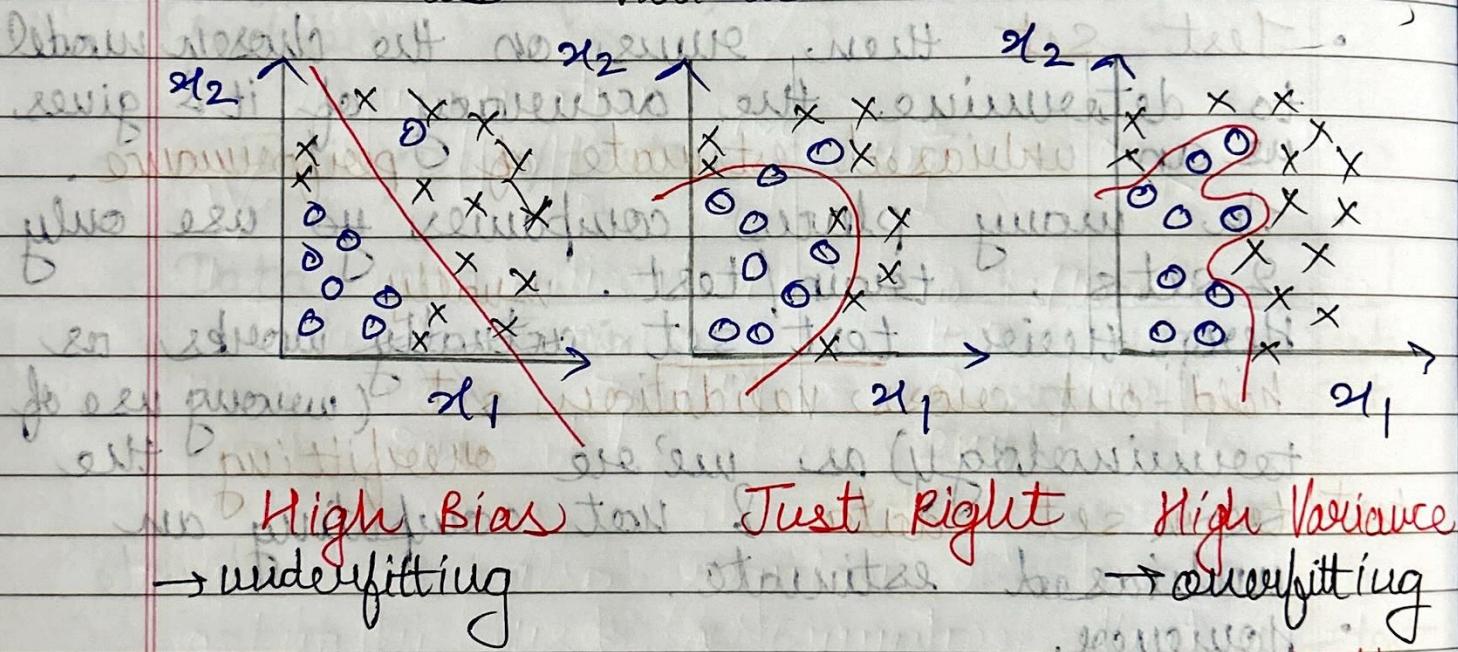
unbiased as we iterate through the diff. models to best fit the data.

TEST SET

Evaluate this data on the chosen model & get unbiased estimate.

- Thus, places where train/test is used in place of train/dev, it is seen as an **okay practice** when we **don't** need a completely unbiased estimate.
- Because we fit data to the dev set, the dev set does **NOT** give a completely unbiased estimate.

* ~~Best~~ Bias - Variance:



Train Set Error 1% 15% 15% 0.5%

Dev Set Error 11% 18% 80% 1%

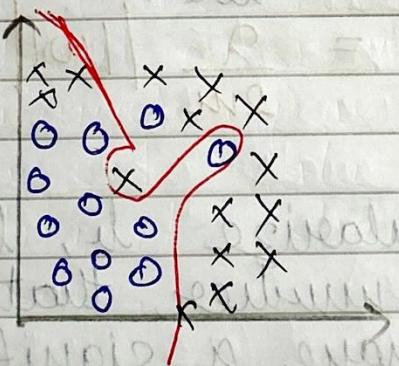
→ **high variance**, **high bias**, **high variance**, **high bias**

→ **keep the model**, **bias**, **high variance**
is **overfitting** the **train** set

as it is not even fitting the training set properly

The below analysis is based off the fact that the human error in differentiating cat/dog is 0. Thus, optimal (Bayes) error (here) : 0

- With the change in Bayes error the analysis may change.
- Train Test Error: in comparison to Bayes error tells us if we how well we're fitting training set & gives an idea about bias
- Dev Set Error: depending on how higher it is from train set error we understand the variance of the model.



high bias,
high variance.

* Basic Recipe for ML :

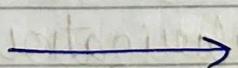
After training the model:

~~High bias?~~ (train set perf.)

~~go for~~
~~if not managed~~

~~Bigger network~~
~~train longer~~

~~Variance?~~
(dev set perf.)
done.



~~more data~~ regularization
NN architecture search

High bias \rightarrow more data won't help.

- As we've more options for managing bias & variation, we do not need to perform a bias-variation trade-off.

* Regularisation in Logistic Regression:
Applying regularisation to Logistic Reg.:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m I(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|^2$$

L₂ regularisation:

$$\|w\|_2^2 = \sum_{j=1}^{n_{\text{xc}}} w_j^2 = \alpha w^T w$$

L₁ regularisation \rightarrow w is/makes sparse (here).

$$\frac{\lambda}{2m} \sum_{i=1}^{n_{\text{xc}}} |w_i| = \frac{\lambda}{2m} \|w\|_1$$

- We can also regularize b, but it's actually just a number that is added. Thus, does not have a significant role in overfitting.

- On the other hand, w parameter is being multiplied & affects all the data-points ~~in~~ significantly.

λ = regularization parameter
 \rightarrow set according to the dev set
 \rightarrow another hyperparameter

* Regularization in a NN

$$J(\omega^{[1]}, \omega^{[2]}, \dots, \omega^{[L]}, \omega^{[L]})$$

$$= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\omega^{[l]}\|_F^2$$

$$\|\omega^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l-1]}} \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2$$

Frobenius Norm
(sum of square of elements of a matrix) dimension.
Gradient Descent : $\omega^{[l]} (n^{[l-1]}, n^{[l]})$

$$\underbrace{d\omega^{[l]}}_{\text{from backprop.}} + \frac{\lambda}{m} \omega^{[l]}$$

$$\omega^{[l]} = \omega^{[l]} - \alpha \underbrace{d\omega^{[l]}}_{\text{from backprop.}}$$

$$= \omega^{[l]} - \alpha \left[\text{(from backprop.)} + \frac{\lambda}{m} \omega^{[l]} \right]$$

$$\underbrace{d\omega^{[l]}}_{\text{from backprop.}} = \omega^{[l]} \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \text{(from backprop.)}$$

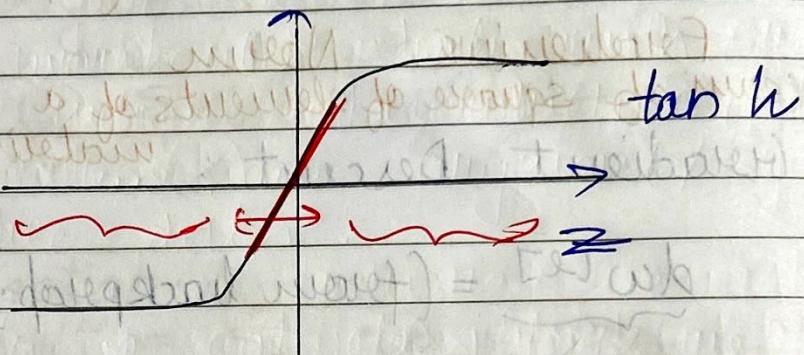
Thus, L_2 regularization \equiv weight decay
owing to $\underbrace{\left(1 - \alpha \frac{\lambda}{m} \right)}_{\text{weight decay}} \omega^{[l]}$.

- If $\lambda \gg \text{high}$, $\omega^{[l]} \approx 0$, that is, few of the features lose significance, making the data set simpler network, which is less prone to overfitting.

$\Delta \uparrow$ $w^{[l-1]} \downarrow$

$$z^{[l-1]} = \underbrace{w^{[l-1]} * a^{[l-1]}}_{\boxed{z^{[l-1]}}} + b^{[l]}$$

→ thus activation function will be relatively linear



* Dropout Regularization:

- Inverted Dropout

We realised that in d3 matrix we find hidden features to zero-out (randomly) at each pass.

keep-prob = 0.8 (here)

$d_3 = np.random.rand(a_3.shape[0], a_3.shape[1]) < keep_prob$

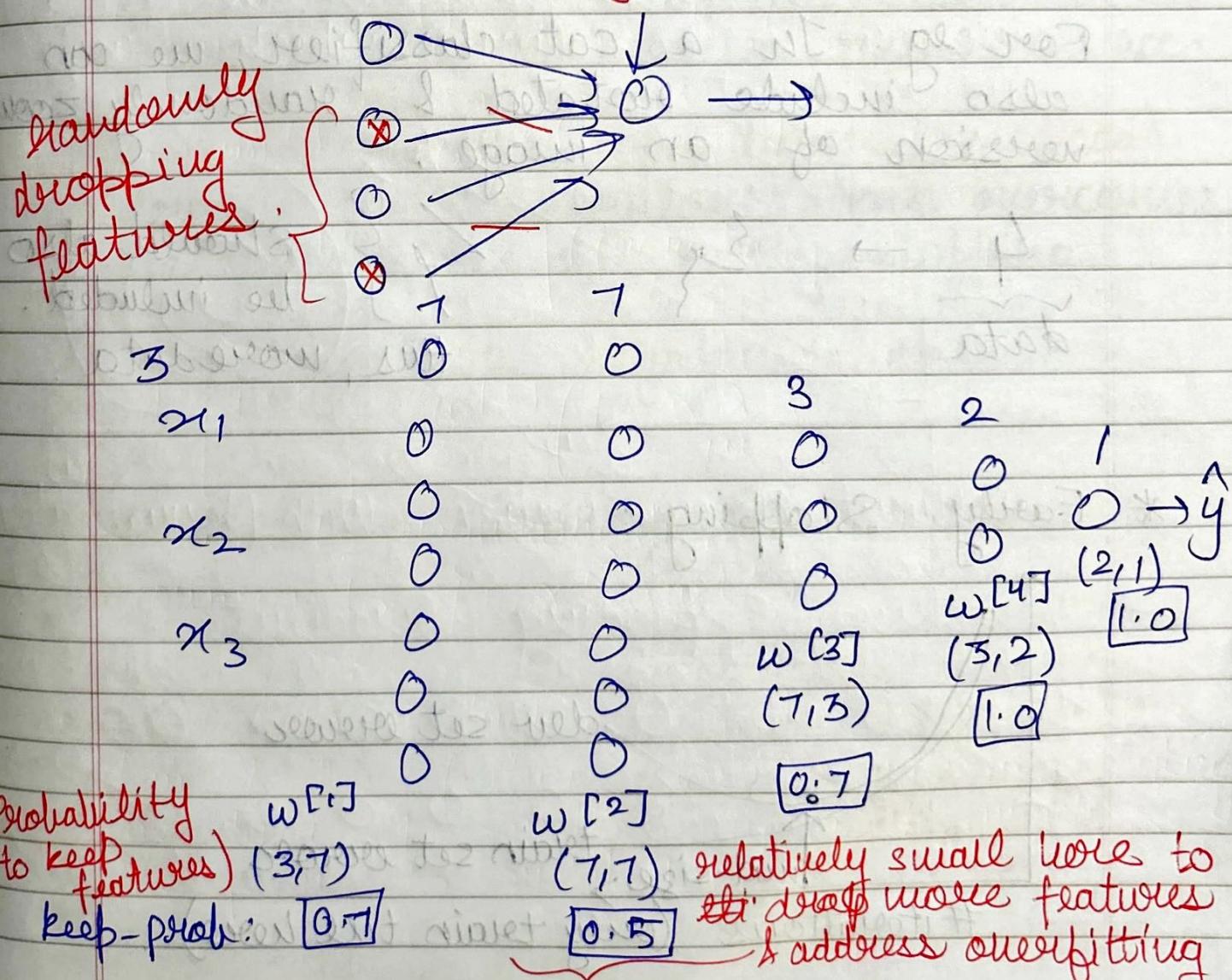
$a_3 = np.multiply(a_3, d_3)$

$a_3 / = \text{keep-prob}$ | ↳ ensures to not add extra scaling during test time.

to standardise/give the desired output despite zeroing out 20% (0.2) features.

We do not use dropout at test time as that would give us a random output & introduce noise.

- * Why does drop-out work?
(as randomly any can be dropped)
- Intuition: Can't rely on any one feature, so have to spread out weights.
 - shrinking weights
 - dealing with overfitting
- Drop-out was considered an adaptive form of L2 regularisation.



Downside of dropout is that the cost function J is not well-defined thus with each iteration of gradient descent we can get different loss.

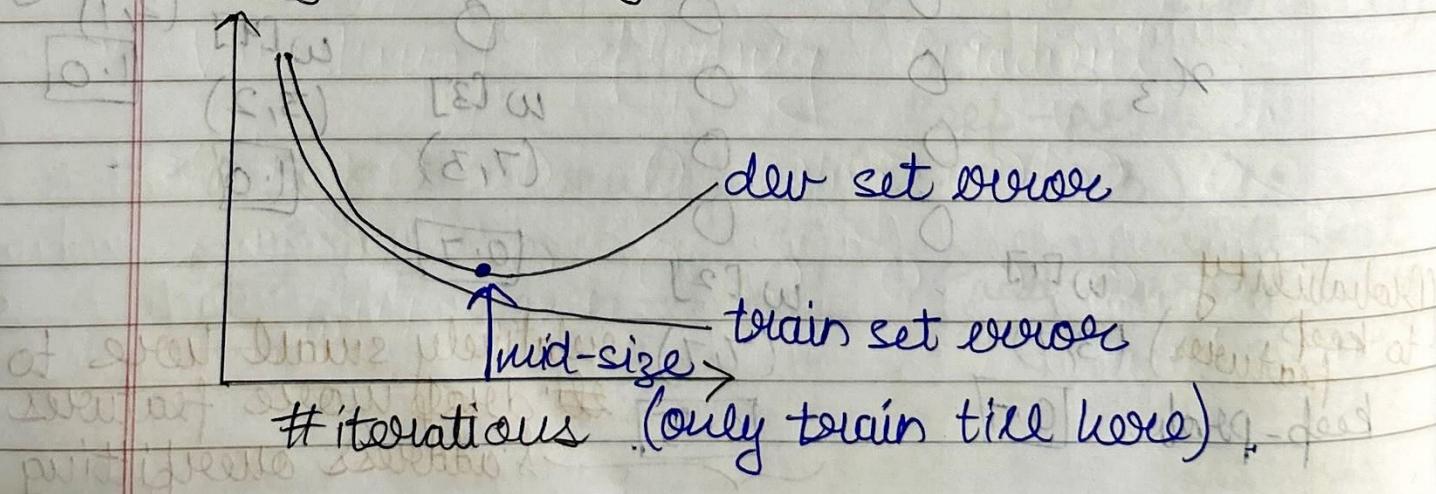
* Data Augmentation:

- Technique for regularisation.
- If we want to deal with overfitting, we need more data, which can be expensive. Thus, we augment/synthesise the available data.

For eg: In a cat classifier, we can also include rotated & randomly zoomed version of an image.

$\underbrace{4}_{\text{data}} \rightarrow \{x_j\}$ $\{x_j\}$ should also be included. thus, more data.

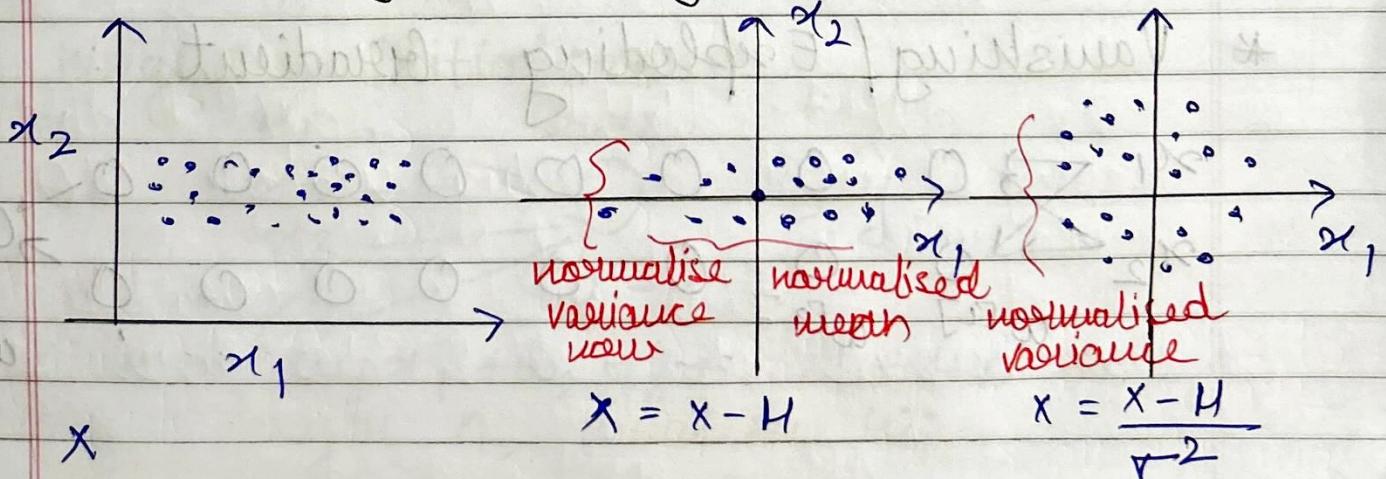
* Early Stopping:



* Orthogonalisation:

- Taking a particular set to solve a given problem, one problem at a time.
- . Downside of early stopping as it tries to solve the problem of overfitting but affects the cost funcⁿ (reduces its iterations) and gradient descent in process.
- . Advantage of early stopping is that it tries all range of values for w in a single gradient descent, which in case of L₂ regularisation we'll have to iterate many times as acc. to adjusting lambda. (which increases computation time/flood).
- . However, L₂ regularisation overcomes the downside of early stopping.
- . Preferred: L₂ regularisation

* Normalising Training Sets:



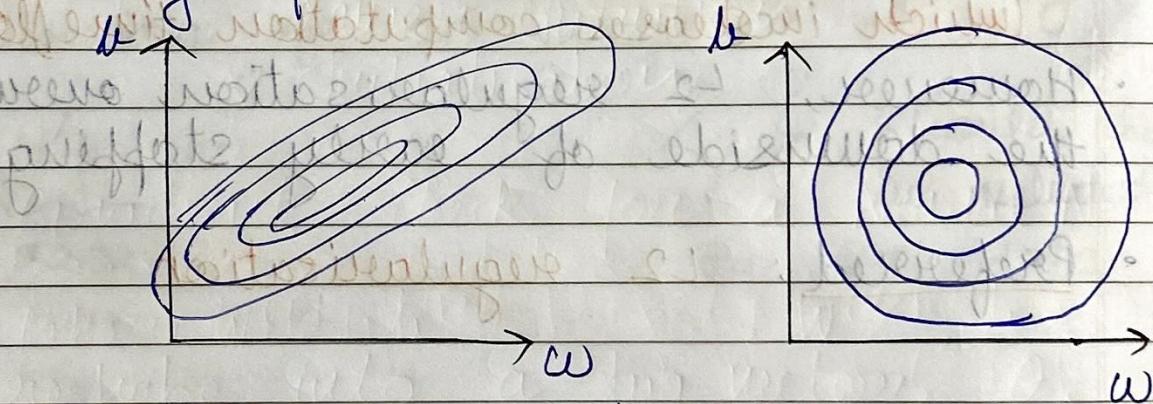
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} * * 2$$

→ use same μ, σ^2 for test set

* Why Normalise?

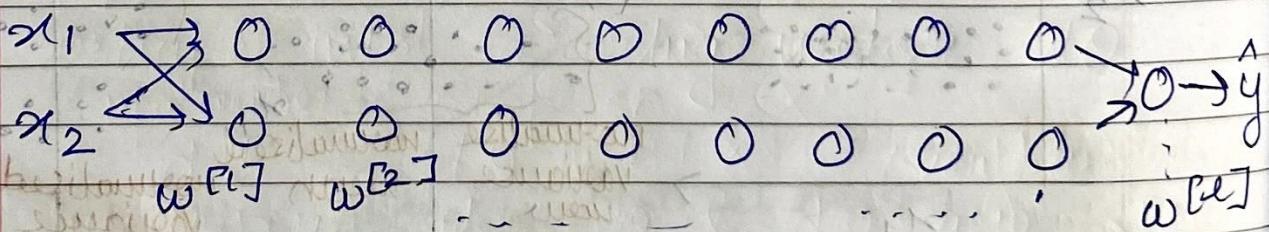
- Cost function (J) is very varied & oscillating without normalisation, thus then we'll have to use a **slow learning rate (α)**.

- On normalisation, however, we get a **symmetric -like- func'** which will converge faster



without Normalisation With Normalisation

* Vanishing / Exploding Gradient :



$$4-x = x$$

$$4-x = x$$

$$\text{let } \hat{y} = w^{[l]} = 0, g(z) = z \\ = w^{[l]} \cdot w^{[l-1]}, \dots, w^{[2]} w^{[1]} X$$

$$g^{(1)}(z) = a^{[1]}$$

$$a^{[2]} = w^{[2]} a^{[1]}$$

$$a^{[2]} = w^{[2]} w^{[1]} X$$

thus, $\hat{y} = a^{[l]}$
 $= \underline{w^{[l]}}$

$$\text{let } w^{[n]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \quad (\text{a little greater than identity matrix})$$

$$\hat{y} = w^{[l]} \begin{bmatrix} 0.5 \\ 1.5 & 0 \\ 0 & 1.5 \\ 0.5 \end{bmatrix}^{l-1} X$$

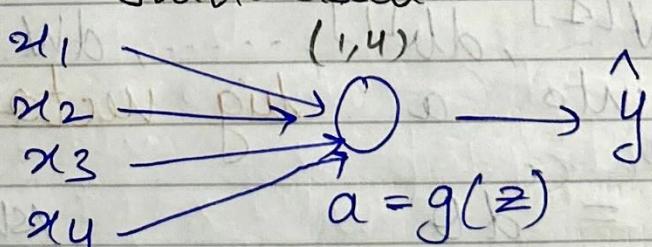
\Leftarrow very large value, & gradient
 very small value, & gradient

$w^{[n]} > I \Rightarrow$ exploding gradient
 $w^{[n]} < I \Rightarrow$ vanishing gradient

I = identity matrix.

We observe that deep NN can have the problem of exploding & vanishing gradient.

* Weight Initialisation:



$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

larger n , smaller w_i

$$\text{Var}(w_i) = \frac{2}{n} \quad \left. \right\} \text{ReLU}$$

$$= \frac{1}{n} \quad \left. \right\} \begin{array}{l} \tanh \text{ & few} \\ \text{other} \end{array}$$

$w^{[l]} = \text{np.random.random(shape)} *$

$$\text{np.sqrt} \sqrt{\frac{2}{n^{[l-1]}}}$$

for ReLU
(most used)

Other Variants:

$$\tanh = \underbrace{\sqrt{\frac{1}{n^{[l-1]}}}}_{\text{Xavier Intuition}} \quad \text{or} \quad \sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

Xavier Intuition.

Initialisation

* Gradient Check for a Neural Network:

1. Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape into a big vector θ .

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$$

2. Take $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$
concatenate

$$dJ = d\theta \quad (\text{should be})$$

• dimensions of θ & $d\theta$ are same.

Checking : for each i :

$$d\theta_{approx}^{(i)} = J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots, \theta_n) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots, \theta_n)$$

$$d\theta_{approx} \approx d\theta \approx \frac{dJ}{d\theta_i} \quad (\frac{\partial J}{\partial \theta_i})$$

Check: $\| d\theta_{approx} - d\theta \|_2 \leq \text{euclidean distance}$
w/o $d\theta_{approx}$ & $d\theta$.

$$\| d\theta_{approx} \|_2 + \| d\theta \|_2$$

euclidean length of $d\theta_{approx}$

$$\approx 10^{-1} \rightarrow \text{great}$$

$10^{-5} \rightarrow \text{may get better, check}$

$$10^{-3} \rightarrow \text{worsey}$$

* Implementation Notes :

- Don't use in training \rightarrow only to debug (as it is too slow).
- If algo fails grad check, look at components to try to identify bug.
- Remember regularisation term in J , if using regularisation in train set.
- Doesn't work with dropout, as it randomly changes the features in every iteration. If using dropout:
 - \rightarrow turn off dropout
 - \rightarrow use grad-check to ensure algo is correct
 - \rightarrow turn on drop out.
- Run at random initialization, perhaps again after some training, so that w & b wander away from 0.