

So multiplying by 0.01 would be reasonable to try.

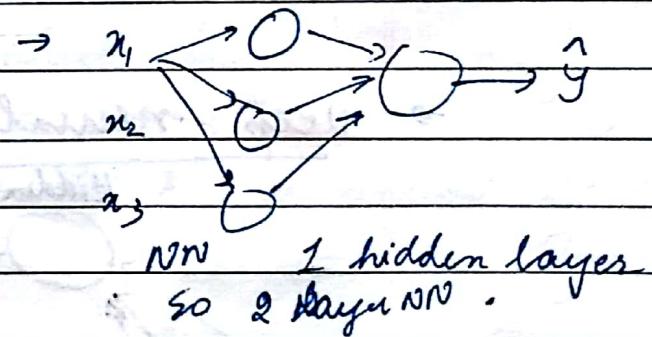
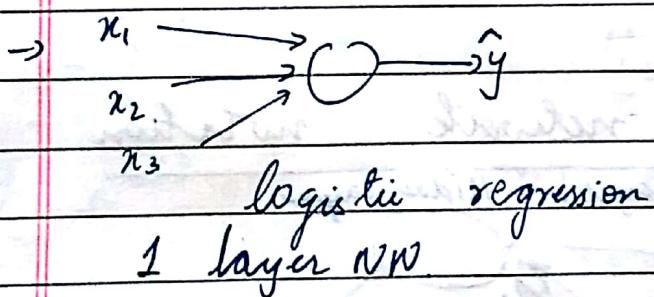
And same for w_2 .

→ When we have small NN then 0.01 will work ok.

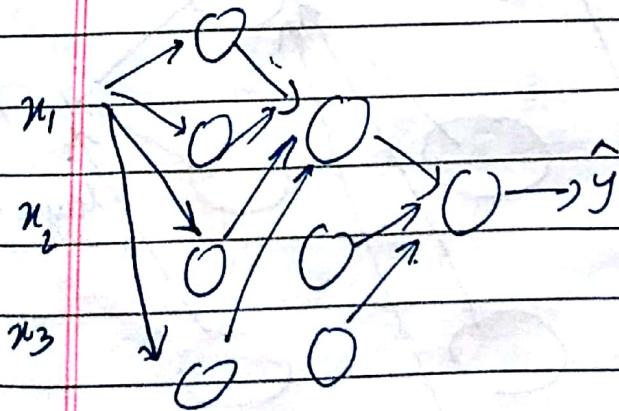
→ But when training a very deep NN, then you might need to pick up a diff constant than 0.01.

Week 4

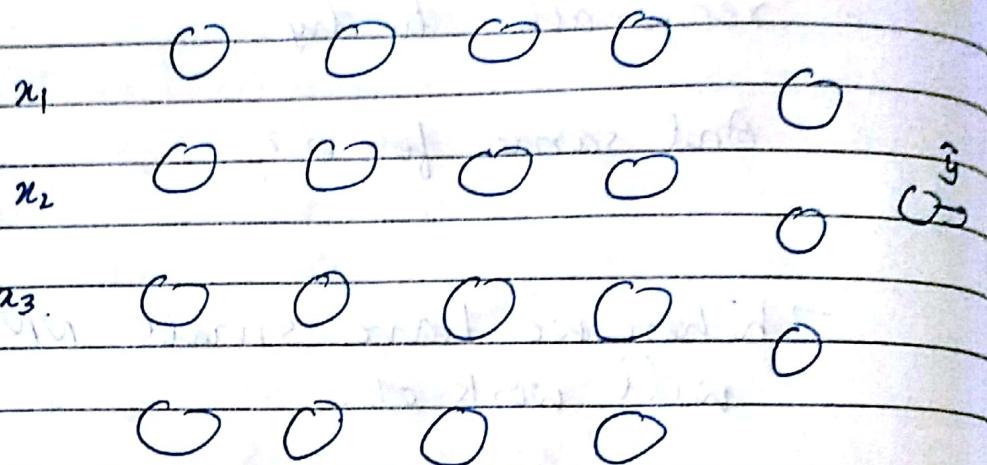
What is a Deep Neural Network



→ NN with 2 hidden layers



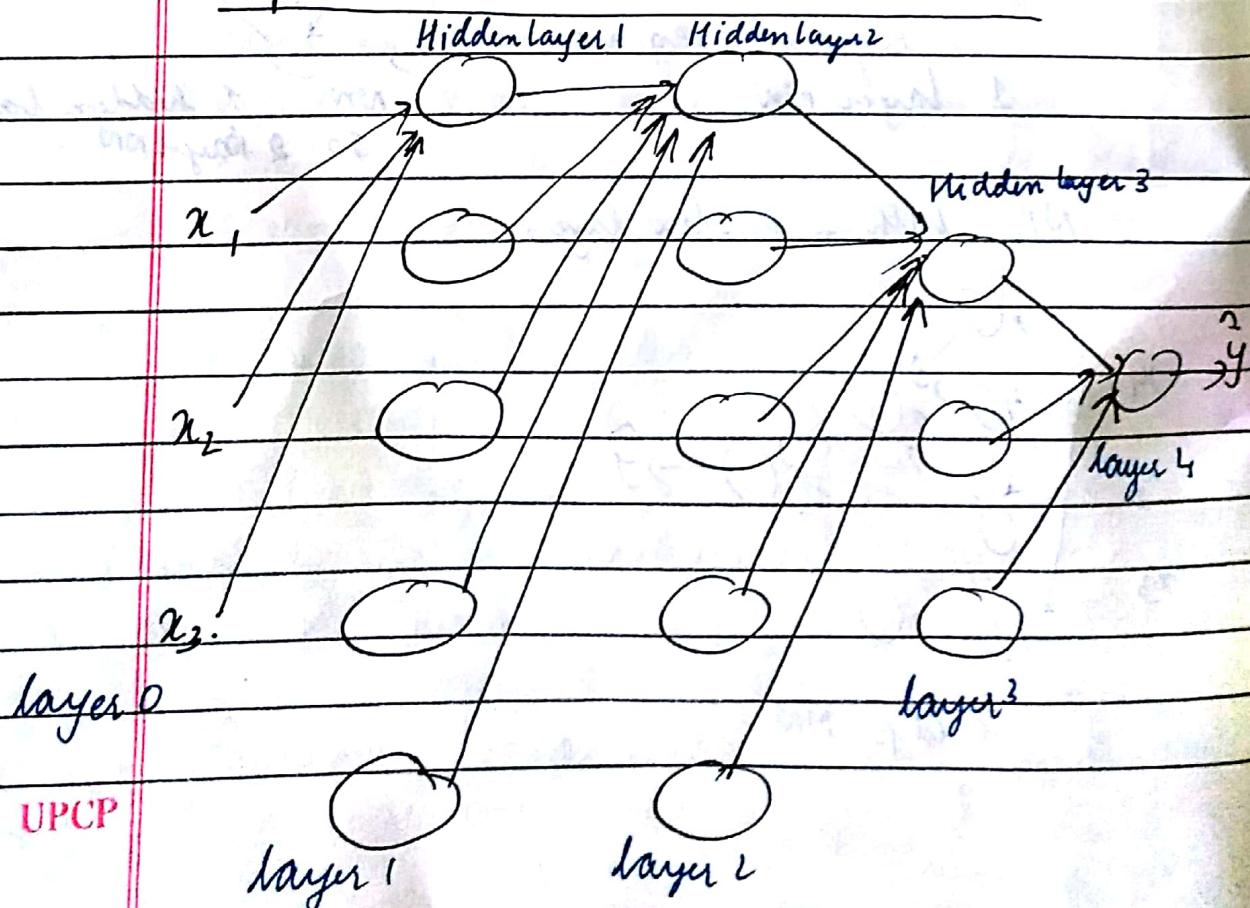
→ NN with 5 hidden layers.



* logistic regression is a very shallow model whereas as NN with 5 hidden layers is a deep model.

→ It is difficult to tell how NN with how many layers is used for a problem.

→ Deep neural network notation



→ 4 layer neural network.

→ It is a NN with 3 hidden layers. Hidden layer 1 has 5 hidden units. Hidden layer 2 has 5 hidden units. Hidden layer 3 has 3 hidden units.

$L = 4$ // % denote the no of layers in the net.

* $n^{[l]}$ = // % denote no of units in layer l .

$$\begin{array}{ll} n^{[1]} = 5 & \text{Hidden layer 1 has 5 units} \\ n^{[2]} = 5 & " " " 2 " \\ n^{[3]} = 3 & " " 3 " 3 " \\ n^{[4]} = n^{[L]} = 1 & \end{array}$$

$$n^{[0]} = n_x = 3 \quad // \text{Input layer.}$$

* $a^{[l]}$ = // % denote the activations in layer l .

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

$$\begin{aligned} w^{[l]} &= \text{weights for computing values in } z^{[l]} \\ b^{[l]} &= \text{biases used} \quad - \quad - \end{aligned}$$

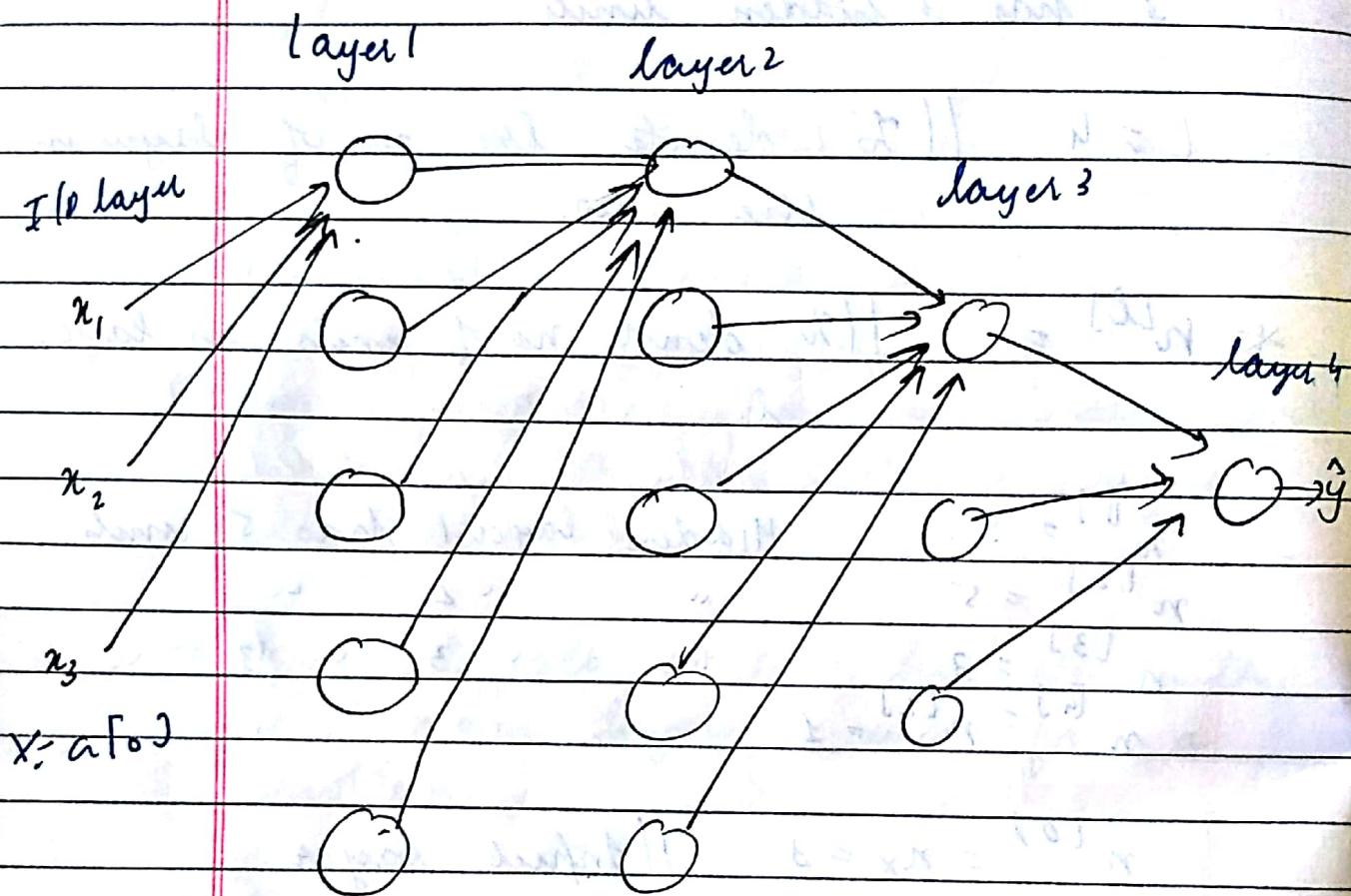
* The IP features are called x . x is the activation of layer 0.

$$\text{So } x = a^{[0]}$$

$$\hat{y} = a^{[L]}$$

FORWARD PROPAGATION IN A

Deep Network



Let's see forward propagation for 1 training example.

Later on will see vectorized version where we will carry out forward propn for entire training set at same time.

For training eg, x we compute forward propagation as follows.

For layer 1:

$x: z^{[1]} = w^{[1]} x + b^{[1]}$ or $w^{[1]} a^{[0]} + b^{[1]}$.
where $w^{[1]}$ and $b^{[1]}$ are the parameters
that affect activation in layer 1.

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

depends on what layer you are at.
So $g^{[1]}$ represents activⁿ funⁿ for layer 1.

layer 2:

$$\begin{aligned} z^{[2]} &= w^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} &= g^{[2]}(z^{[2]}) \end{aligned}$$

.

:

layer 4

$$\begin{aligned} z^{[4]} &= w^{[4]} a^{[3]} + b^{[4]} \\ a^{[4]} &= g^{[4]}(z^{[4]}) \end{aligned}$$

General rule of eqⁿ

$$\begin{aligned} z^{[l]} &= w^{[l]} a^{[l-1]} \\ a^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

Vectorized Version \rightarrow for entire training data
 $z^{[1]} = w^{[1]} X + b^{[1]}$ or Replace X with $n^{[0]}$.

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

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

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

$$\hat{y} = g(z^{[4]}) = a^{[4]}$$

$$\bar{Z}^{(2)} = \begin{bmatrix} \bar{z}^{(2)(1)} & \bar{z}^{(2)(2)} \\ | & | \\ \vdots & \vdots \end{bmatrix}$$

$A^{(2)}$ = Similarly.

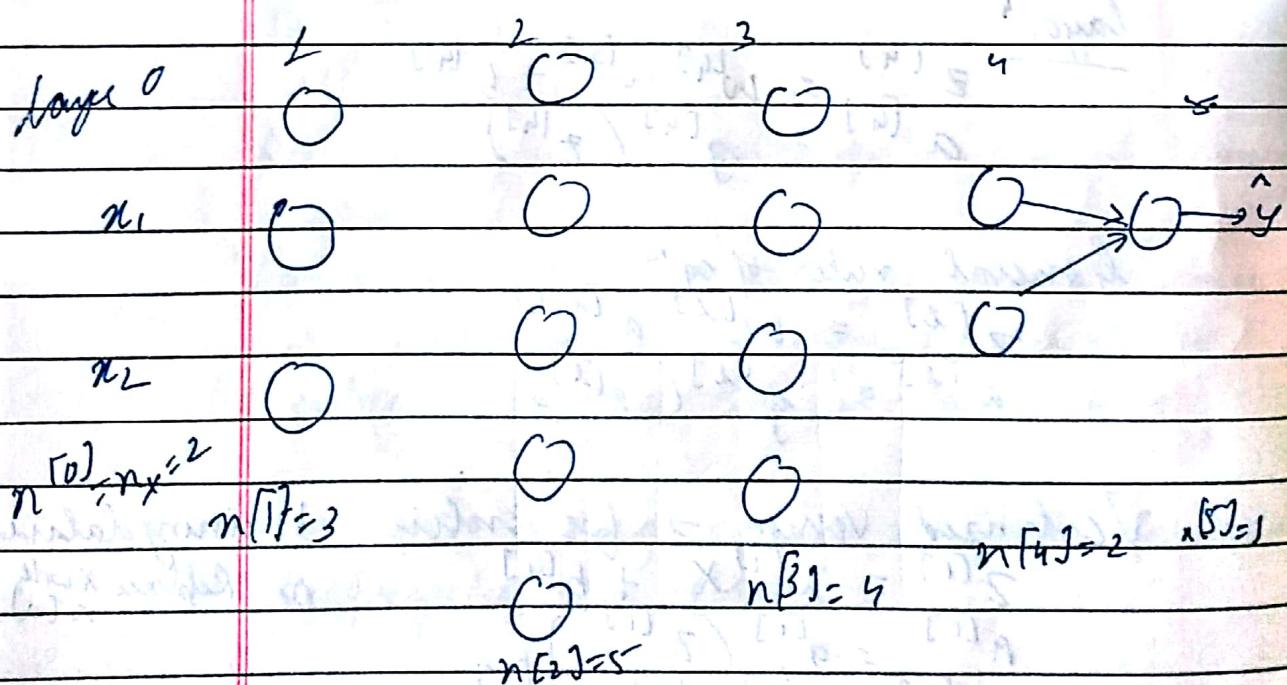
i.e all training egs are stacked horizontally.

General vectorized form:

$$\begin{aligned} Z^{[x]} &= W^{[x]} A^{[x-1]} + b^{[x]} \\ A^{[x]} &= g^{[x]}(\bar{Z}^{[x]}) \end{aligned}$$

Getting your matrix dimensions right.

Parameters $W^{[x]}$ and $b^{[x]}$.



L=5 bcz there are 5 layers
UPCP 1 I/P layer which is not counted among layers of neural net.

1 O/P layer.

4 hidden layers

$$z^{[1]} = w^{[1]}x + b^{[1]}$$

$$\text{dimensions of } w^{[1]} : (n^{[1]}, n^{[0]}) \\ w^{[2]} : (5, 3) = (n^{[2]}, n^{[1]})$$

Generalized dimension:

$$w^{[x]} : (n^{[x]}, n^{[x-1]})$$

$$b^{[x]} : (n^{[x]}, 1)$$

In backpropagation,

dimension of dw should be same as w
L " " obs " B

$$\alpha^{(1)} = g^{(1)}(z^{(1)})$$

So α & a should have same dimensions

* In virtualized implementation,

$$\begin{aligned} \hat{y}^{(1)} &= w^{(1)} x + b^{(1)} \\ (n^{(1)}, 1) &\quad (n^{(1)}, n^{(0)}) \\ &\quad (n^{(0)}, 1) \end{aligned}$$

$$z^{[1]} = w^{[1]} X + b^{[1]}. \quad // \text{Vectorized.}$$

$(n^{[1]}, n^{[0]})$ $(n^{[1]}, n^{[0]})$ \downarrow
 \downarrow no of training
egs $(n^{[0]}, m)$ $\underbrace{(n^{[1]}, b^{[1]})}_{\downarrow \text{Broadcasting}}$

dimensions of $Z^{[1]}$, $a^{[1]}$
 $dZ^{[1]}, da^{[1]} = (n^{[1]}, m)$

$Z^{[1]}, a^{[1]} = (n^{[1]}, 1)$

Why deep Representations?

Why "Nets" work well?

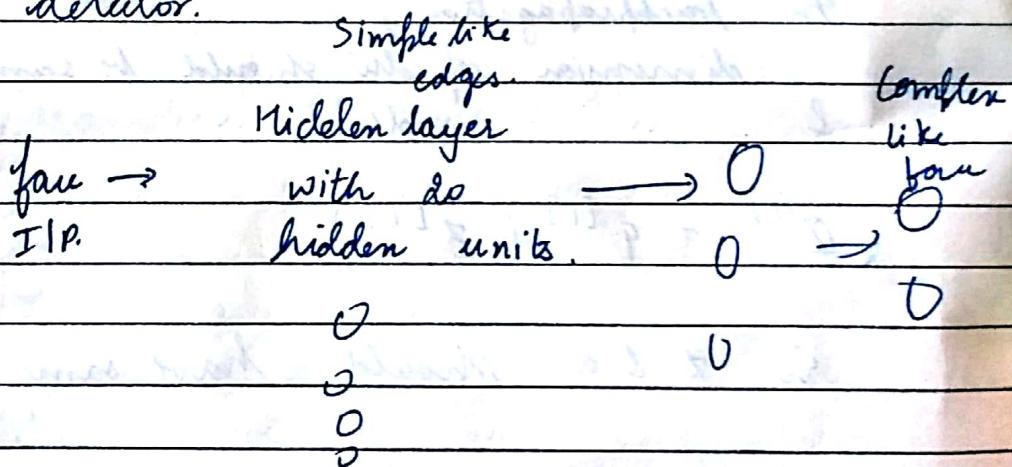
1. What is deep network computing.

When we are building a system for face recognition or face detection.

Here is what a deep neural network does.

I/P is the picture of face.

The first layer of Neural Network is a feature detector or an edge detector.



e.g. this
represents a
hidden unit that
is trying to
visualize the
orientation of
edges in image.

UPCP

There are neurons to find eyes. Some neurons try to find part of the nose.

- By putting together lots of edges it can start to detect diff parts of the faces.
And finally after putting together diff parts of the faces like eyes, nose, chin, it can try to recognize or detect diff types of faces.
- Earlier layer of neural network detect functions like edges and then composing them together in the later layers of a neural network so that it can learn more and more complex structure func's.
- * The edge detectors are looking in relatively small areas of an image.
Facial detectors look at much larger areas of image.

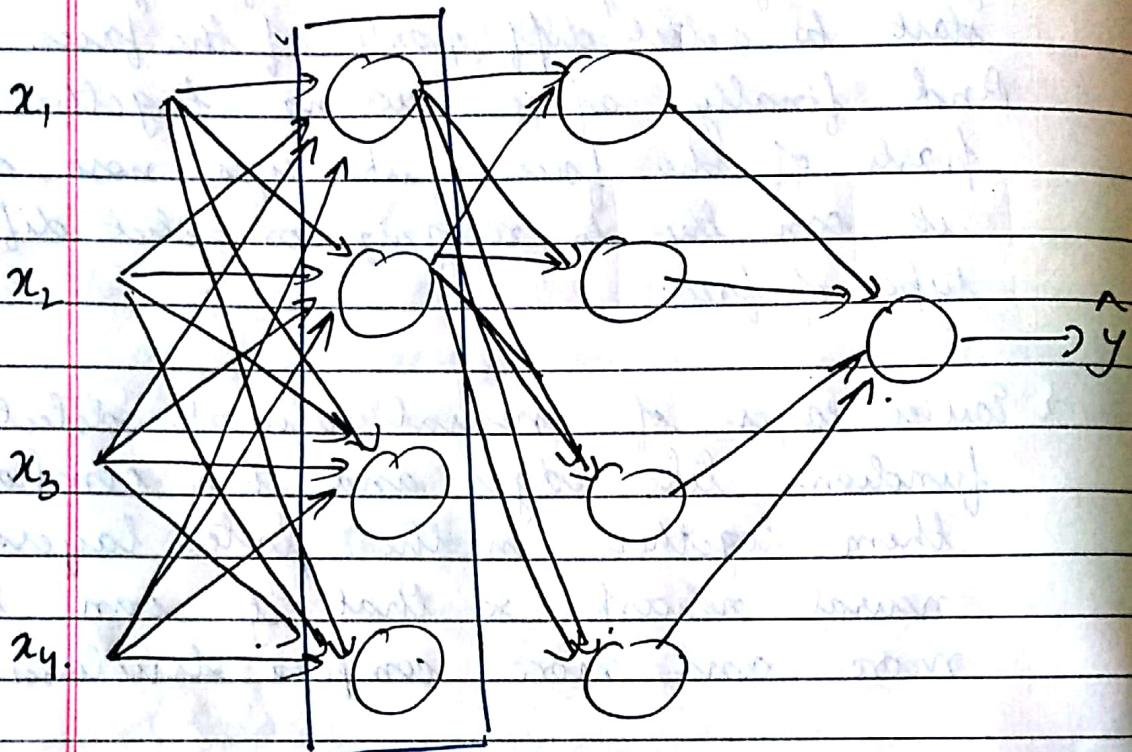
This type of simple to complex hierarchical representation applies in other types of data than images & face recognition as well.

Circuit theory and deep learning

There are functions you can compute with a "small" L-layer deep neural network that shallow networks require exponentially more hidden units to compute.

Building blocks of deep neural network

Forward and Backward functions :



Network of few layers is above.

Let's pick 1 layer & focus on computation on that layer for now

So for any layer, we have the parameters $w^{[l]}$, $b^{[l]}$

parameters : $w^{[l]}, b^{[l]}$

forward prop:

IIP \rightarrow activation $a^{[l-1]}$
OIP \rightarrow .. $a^{[l]}$.

for this we compute :

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

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

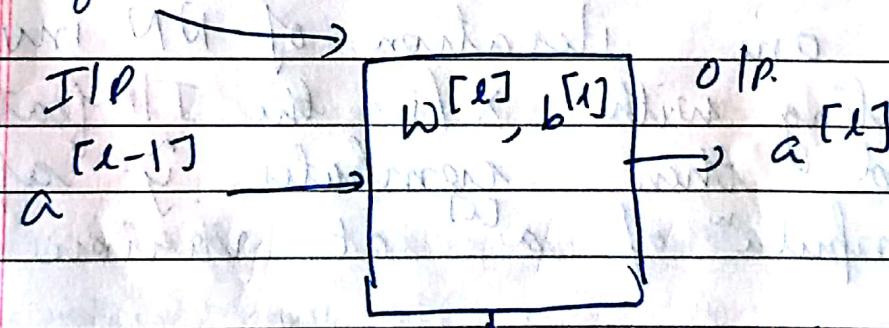
As we know $z^{[l]}$ is used later.
 So the value of $z^{[l]}$ is cached for
 later use. bcz storing the value of z
 will be useful for the backward prop^n
 step.

Backward Propagation \rightarrow

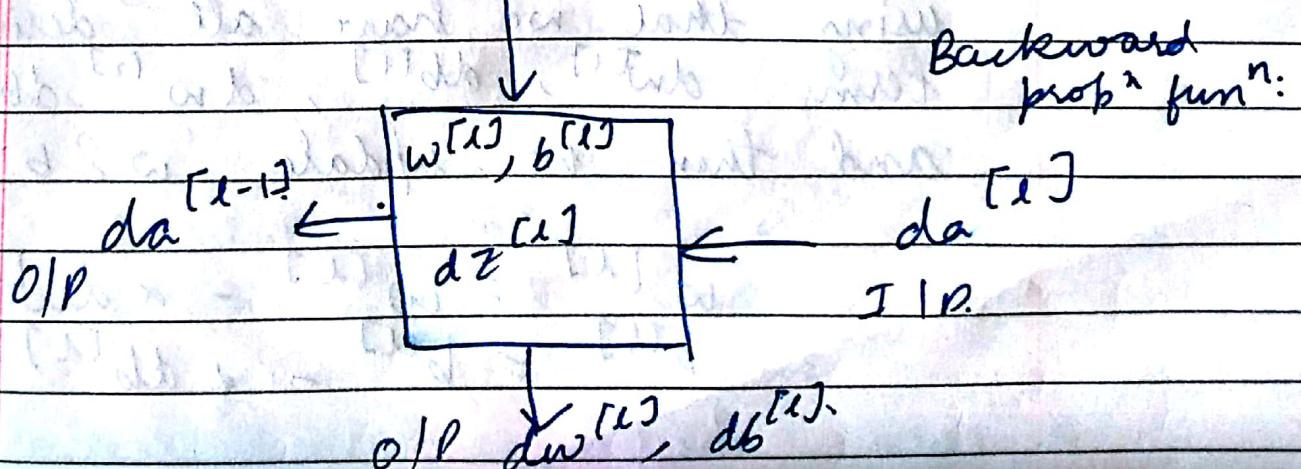
I/P $a^{[l]}$, & cache ($z^{[l]}$)
 O/P $da^{[l-1]}$, $dw^{[l]}$, $db^{[l]}$

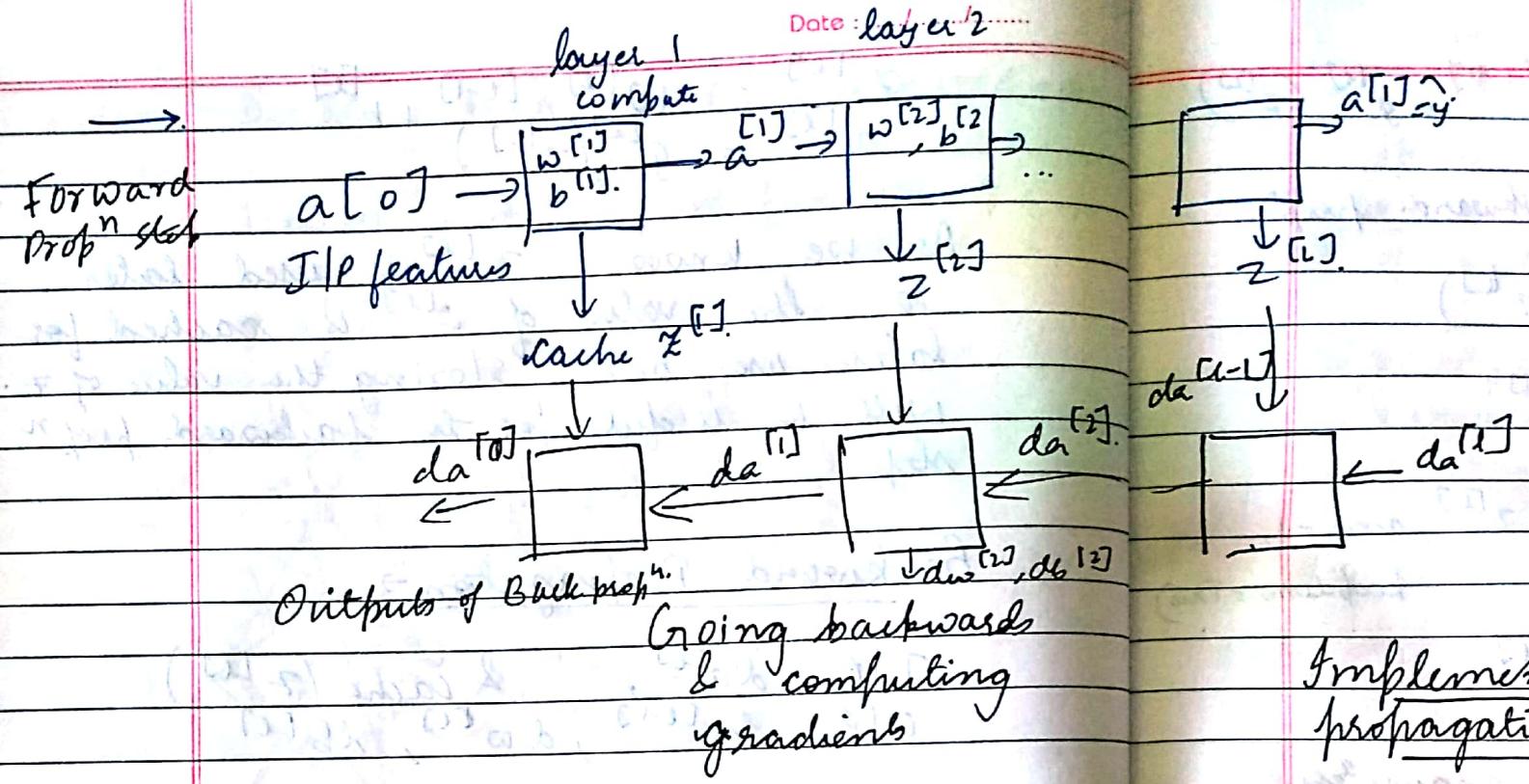
* Summary

In layer l , we have forward prop^n
 fun^n



O/P Cache $z^{[l]}$.





No need for derivative of I/P features in case of supervised NN

So one iteration of NN involves starting with x (i.e. the I/P features) and then computes \hat{y} after computation of a^l at various layers.

And then back prop & then using that we have all derivative terms $dw^{[1]}, db^{[1]}, dw^{[2]}, db^{[2]} \dots$ and then we update w & b .

$$w^{[1]} = w^{[1]} - \alpha dw^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha db^{[1]}$$

Cache is used to store the value of Z for the backward functions.

It may be convenient to store w & b also in cache to get them in the backward propagation.

Implementing forward & backward propagation steps:

Forward propagation for layer l :

- Input $a^{[l-1]}$
- Output cache ($Z^{[l]}$), $a^{[l]}$.

$$Z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$
$$a^{[l]} = g^{[l]}(Z^{[l]})$$

Backward propagation step for layer l .

- Input $da^{[l]}$
- O/p $da^{[l-1]}, dw^{[l]}, db^{[l]}$

$$dZ^{[l]} = da^{[l]} * g^{[l]}'(Z^{[l]})$$

$$dw^{[l]} = dZ^{[l]} * a^{[l-1]}$$

$$db^{[l]} = dZ^{[l]} * \mathbf{1}$$

$$da^{[l-1]} = W^{[l]}^T \cdot dZ^{[l]}$$

$$d\hat{z}^{[L]} = w^{[L+1]T} d\hat{z}^{[L+1]} \times g^{[L]'}(\hat{z}^{[L]})$$

Vectorized "tuple" of backward prop:

$$d\hat{z}^{[L]} = dA^{[L]} \times g^{[L]'}(\hat{z}^{[L]})$$

$$dw^{[L]} = \frac{1}{m} d\hat{z}^{[L]} A^{[L-1]T}$$

$$db^{[L]} = \frac{1}{m} \text{np.sum}(d\hat{z}^{[L]}), \text{axis=1}$$

keepdims=True

$$dA^{[L-1]} = w^{[L]T} d\hat{z}^{[L]}$$

Parameters v/s hyperparameters

What are hyperparameters?

Parameters : $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots$

Hyperpara \Rightarrow learning rate (α),
 \downarrow

determine how parameters
 evolve. \Rightarrow no of iterations of gradient
 descent you can carry out.

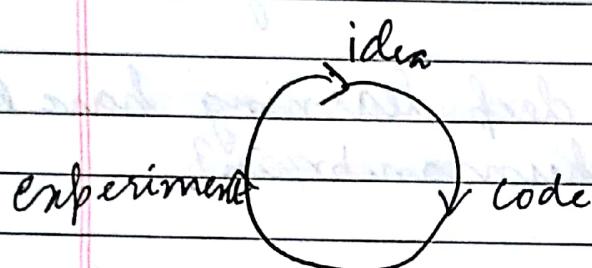
\rightarrow No of hidden layers [1]
 \rightarrow units

\rightarrow choice of activ' fun'
 like relu, sigmoid.

Here the hyperparameters control

Bcz they control values of W & b .
 Hence called hyperparameters. These hyperparameters determine the final value of W and b .

Applied deep learning is a very empirical process.



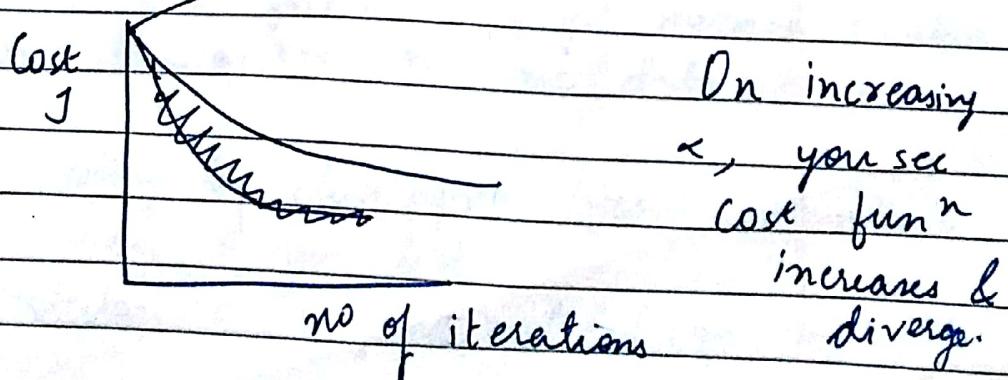
Eg \rightarrow You have an idea for the best value of learning rate. You say $\alpha = 0.01$.

Then you implement it. And see how it works.

Then based on outcome, you say that i changed my mind, you increase value to 0.05.

\rightarrow If you are not sure, what is the best value of α .

Then choose 1 value of α and build \Rightarrow cost & it goes down



Try other values & after trying certain values of α , it gives pretty fast learning & allows to converge to a lower cost

Empirical process bcz in deep learning we need to try a lot of things and see what exactly works.

What does deep learning have to do with human brain?

Classifying male or female given just their Body Measurements.

- * Data Science is the study of data.
 - " Scientist is someone who solves problem by studying data.

Dependency used in this gender classification program is scikit-learn.

↓
a machine learning package with pre built models for us to use.

```
from sklearn import tree
```

```
# [height, weight, shoe size]
```

```
X = [[181, 80, 44], [177, 70, 43],
```

```
[160, 60, 38], [154, 54, 37]] . . .
```

```
]
```

```
y = ['male', 'female', ' . . . ]
```

```
clf = tree.DecisionTreeClassifier() // stores decision tree classifier.
```

```
tree.DecisionTreeClassifier()
```

```
clf = clf.fit(X, y) // fit method trains the decision tree on one data set.
```

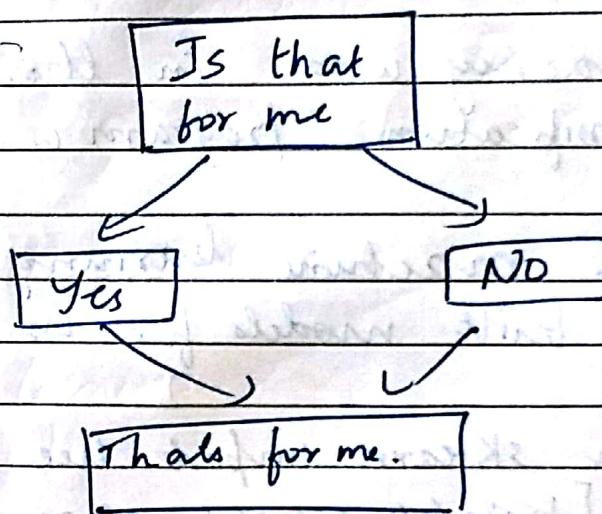
```
// Testing by classifying gender of someone given a new list
```

```
Prediction = clf.predict([[190, 70, 43]])
```

Tree \rightarrow It will help build a machine learning model called a decision tree.



It is like a flowchart that stores data.

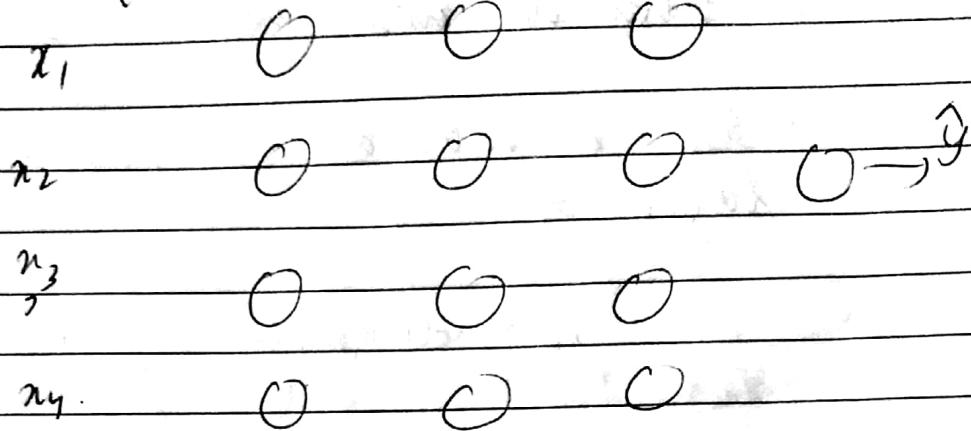


$$\frac{3}{b}$$

$$w = w - \alpha g_w$$

Dropout

→ Regularization Technique

overfitting

→ dropout → Go through network and set some probability of eliminating a node in your network.

e.g. At each node we set prob of 0.5 of keeping a node & 0.5 prob of eliminating it.

So after each epoch, we decide to eliminate those nodes & all the incoming & outgoing synapses from them. This results in smaller, diminished network

→ Train / eg on this much diminished network

So for each training eg you train using 1 of these new reduced eg.

→ Here for each eg training on a smaller neural network.

Implementing dropout (Inverted dropout)

Illustrate with layers $l = 3$.

→ Set a vector $d_3 = \text{dropout vector for layer 3}$

$d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape[1]})$

✓ < Keep prob
same size as a_3

Let $\text{Keep prob} = 0.8$.

// It generates random vector & compares with threshold value.

$a_3 = \text{activations computed}$.

$a_3 = \text{np.multiply}(a_3, d_3)$

$a_3 = a_3 * a_3$.

In this some nodes gives 0 activ.

→ $a_3 / = \text{keep prob}$ // This is used for scaling

// 50 neurons in hidden layer.

$50 \times 1 \rightarrow \text{dimension of } A_3$.

→ At test time, No drop out

$$a^{[0]} = X$$

$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

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

No need to eliminate hidden units.

If you implement dropout in test time,
it just adds noise in your predictions

During training neural networks is doing
well but during testing it never gets
correct O/P.

→ eg students who mug up for exams
score rubbish in a real exam.

→ Dropout → Creating a N/w where each
neuron has chance to be activated
or disabled

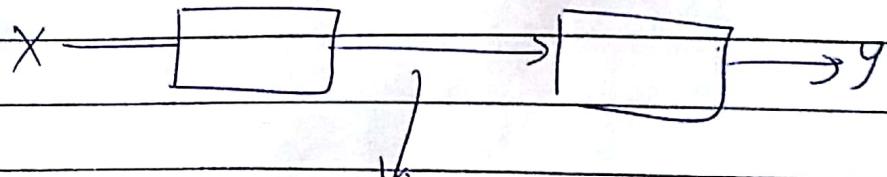
It helps to create N/w with reliable
units.

eg → It is more beneficial to consult
a committee of doctors than 1 single
doctor.

→ It gives the avg of a very large no of
possible neural networks, & we only have

we have to train / Network.

→ Dropout is slow in comparison to training / network.



Activations are values that go from 1 layer to other.

Some activations are somehow set to 0.

So in this way O/P does not rely on any given activation

So for the same I/P even if we don't have few hidden units, O/P remains same.

Everything remains fine at the end b/c even if few hidden layers are smashed, there are some other hidden layers that do the same task.

Dropout

→ Randomly eliminates hidden units. So for each iteration it seems as if

keras / optimization

Date : _____

We are training on a smaller neural network.

- 1 unit can't rely on any 1 I/O feature bcz any 1 of its I/O can go away at random.



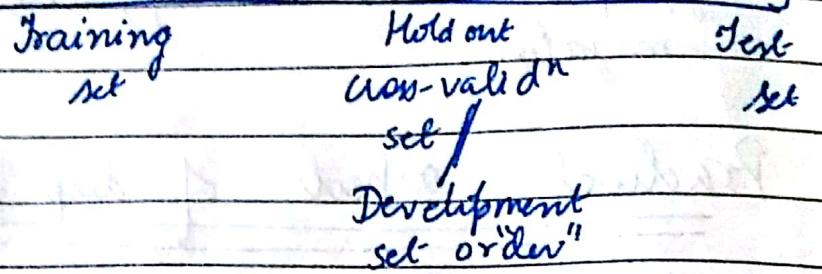
Improving deep Neural Networks:

Hyperparameter tuning, regularization & optimization.

Practical aspects of deep learning

- Making good choices in how you set up your training, development and test sets can make a difference in helping you find a good, high performance neural network.
- Some decisions for training a neural network
 - how many layers does your neural network have.
 - how many hidden units each hidden layer has.
 - What's the learning rate
 - What are the activation functions you want to use for diff. layers
- It is difficult to choose them in your first attempt.
- How to make progress quickly?
Progress is determined by how efficiently you can go around this cycle.
I also set up dataset idea
in training, development experiment →
and test sets.

Data



Workflow

- Keep on Training algo on training set.
- Dev set is used to check many different models and helps to find which model performs best.
- Take the best model found after doing step 1 & 2 for long time. This best model is evaluated on test set in order to get a unbiased estimate of how well your algo is doing.

*

When solving ML problem, data is divided into train/dev and test sets:

I] If we have relatively smaller data set, then ratio is like 80% training, 20% test, 20% dev.

When data set is much larger, it is fine to set up your dev and test dataset to be much smaller.

eg → Earlier in ML era, the data is divided in train, dev, test data in ratio of 60, 20, 20 or 80, 10, 10 or something like this. Here data is not very large.

In this era of Big data, the data is usually very large. So if total data is 1 billion then we need only about 10000 dev data that will be enough to help us decide which algo is best.

10000 test data is enough.

So here 98% training data, 1% dev, 1% test.

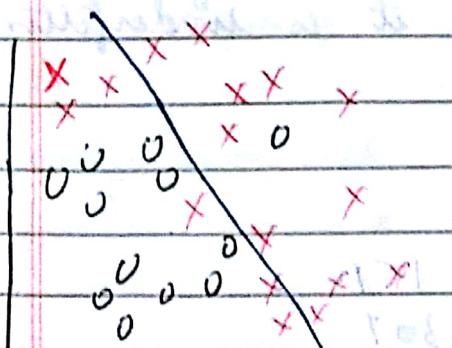
→ Mostly Models have 3

Training set
→ Cat pictures from webpages. (It is in high resolution)

Dev / test set.
→ Cat pictures from users.
(It may be blurred).

Make sure that dev & test data come from the same distribution. i.e. if test data is from random image from phone in casual manner by user. Then dev should also be from it.

Bias / Variance

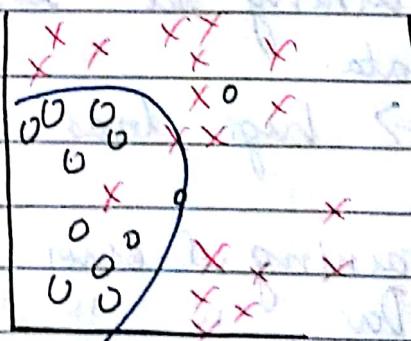


High bias

→ Maybe logistic reg fit on drawing straight line.

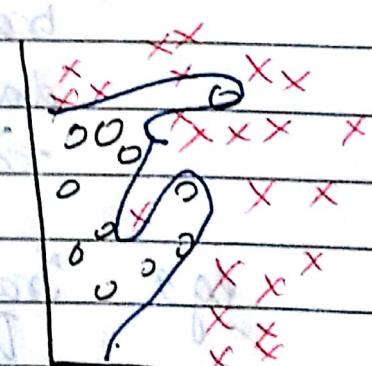
Not good fit to data.

Means high bias or underfitting



Just right

Using big NN with many hidden units, maybe you can fit the data perfectly



high variance

→ Overfitting data

Cat classifiⁿ

Model works $y = 1$ if cat

$y = 0$ if non-cat

To understand bias & variance, look at -

- eg
- Training set error = 1% (let's say)
 - Dev set error = 11% ("")

We are working good on training set but since large error ~~as 31/21~~ in dev set.

\Rightarrow Overfitting of data.

\Rightarrow So generalizⁿ to dev test is not done

In this eg, there is high variance.

eg. Training set error = 15%. (let's assume)
Dev " " = 16%.

Here algo is not even working good on training set then it is underfitting data.

\Rightarrow high bias.

eg Training set error = 15%.
Dev " " = 30%.

high bias. & high variance.

eg Training set error = 0.5%.

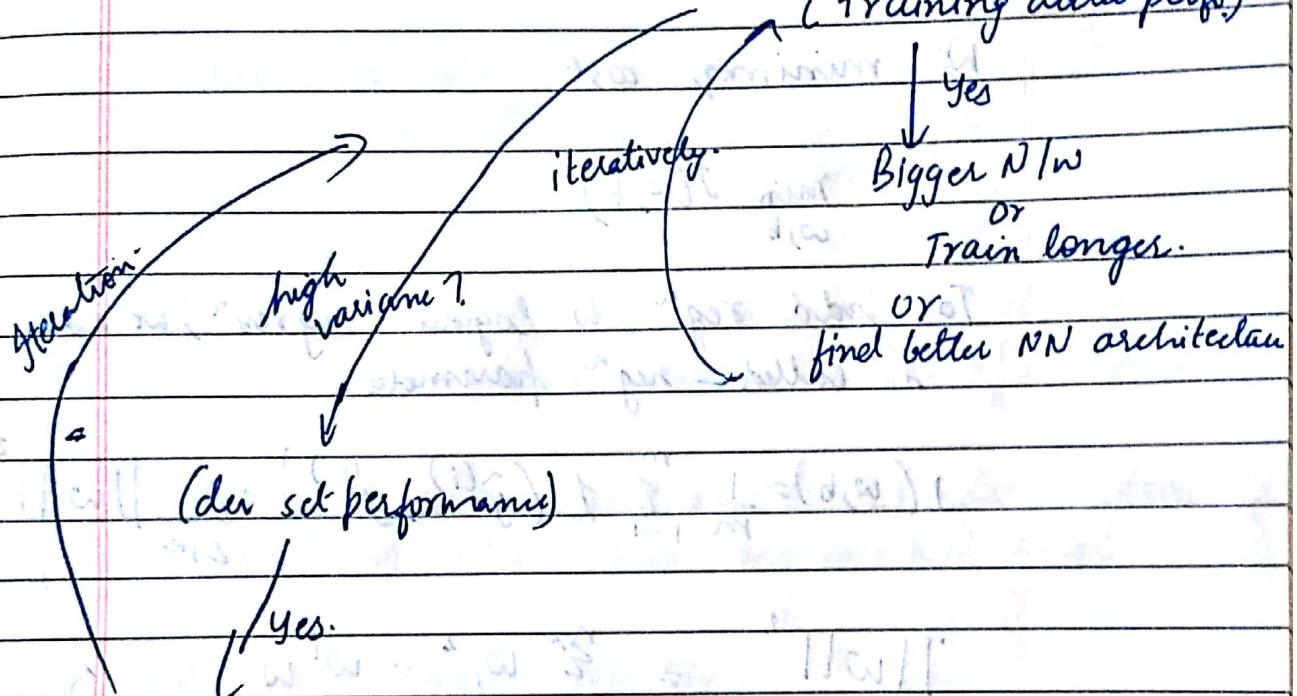
Dev " " = 1%.

low bias, low variance

Based on training error we can sense of bias
 " " Dev " variance

Basic Recipe for ML

- First train the model
- check whether algo has high bias?
 (training data perf.)



Get More data.

or

do Regularizⁿ

or

Change architecture

Regularizing Your NN

- If our NN is overfitting data i.e if it has high variance problem; & Methods to avoid this are:
 - Try regularizⁿ
 - Get more data → This is not possible always

How regularization works:

Logistic regⁿ:

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

$w \in \mathbb{R}^n$, $b \in \mathbb{R}$.

We minimize cost.

$$\min_{w, b} J(w, b)$$

To add regⁿ to logistic regressⁿ, we add a called regⁿ parameter.

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

$$\|w\|^2 = \sum_{j=1}^n w_j^2 = w^T w.$$

norm w^2

This is called L2 regularization because here we are using Euclidean norm or also called L2 norm with parameter vector w .

→ Now why we regularize just a parameter w .

Why not we add $\frac{\lambda}{2m} b^2$.

We can add. But since w is a high

B is a single no. So almost all parameters are in w . So B won't affect much.

$L_2 \text{ reg}^n$ is the most common type of regularization.

$L_1 \text{ reg}^n$

here we add

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

If we use $L_1 \text{ reg}^n$, w will be sparse i.e. w vector will have a lot of zeros in it.

$L_1 \text{ reg}^n$ may compress the model as has lots of zeros. Thus makes NN small.

* $L_2 \text{ reg}^n$ is used mostly.

λ is called regⁿ parameter. It is set over dev set

λ is hyperparameter.

Neural Network

→ $L_2 \text{ reg}^n$

$$\begin{aligned} \text{In NN, cost fun}^n \text{ is } J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) \\ &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \sum_{j=1}^L \|w_{ij}^{(l)}\|^2 \end{aligned}$$

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

$$W: (n^{[L-1]}, n^{[L]})$$

Called as Frobenius norm of a matrix denoted by f .
 \downarrow
means sum of squares of elements of matrix.

Q. How to find grad descent

→ First dW from backprop

$$dW^{[L]} = \frac{\partial J}{\partial W^{[L]}}$$

$$W^{[L]} = W^{[L]} - \alpha dW^{[L]}$$

After adding reg^n term to cost funⁿ, $dW^{[L]}$ changes

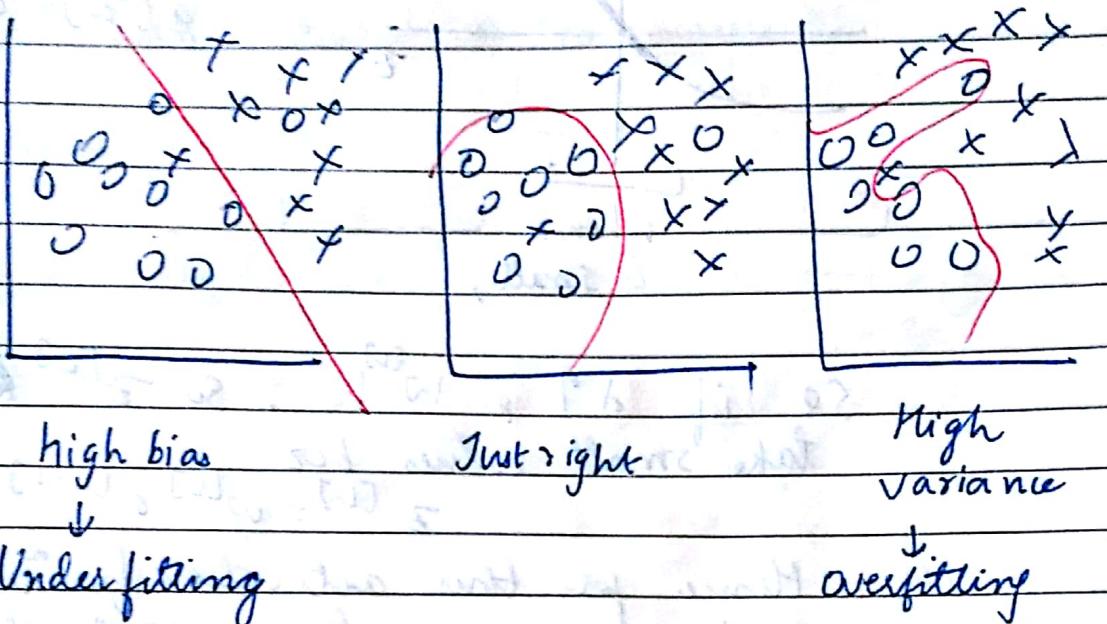
$$\text{so } dW^{[L]} = \frac{\partial J}{\partial W^{[L]}} + \frac{\lambda}{m} W^{[L]}$$

$$W^{[L]} = W^{[L]} - \alpha dW^{[L]}$$

For this reason, L2 regularization is sometimes called weight decay.

$$\begin{aligned} W^{[L]} &= W^{[L]} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{m} W^{[L]} \right] \\ &= W^{[L]} - \frac{\alpha \lambda}{m} W^{[L]} - \alpha (\text{from backprop}) \end{aligned}$$

How regularization prevent overfitting

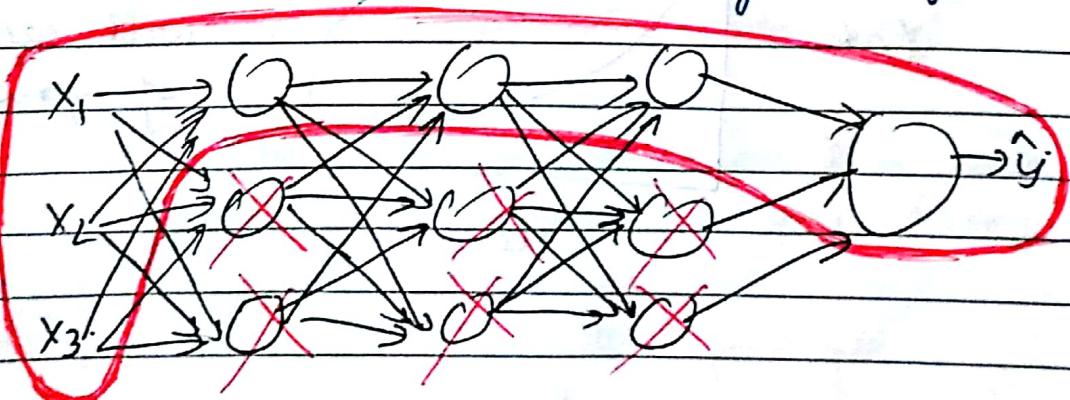


$$J(w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m d(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w_l\|_F^2$$

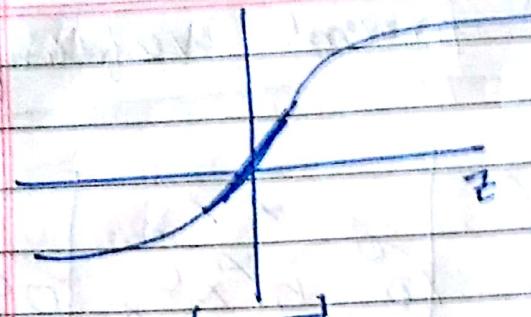
(with L2 regularization)

Now if λ is too large, $w^{[L]}$ for many hidden units become very small & they become close to zero.
 $w^{[L]} \approx 0$.

So our NN reduces to logistic regⁿ



It now becomes like logistic regⁿ
 So for high λ , high variance can reduce to low bias case. So intermediate value of λ should be chosen.

\rightarrow tanh activⁿ funⁿ

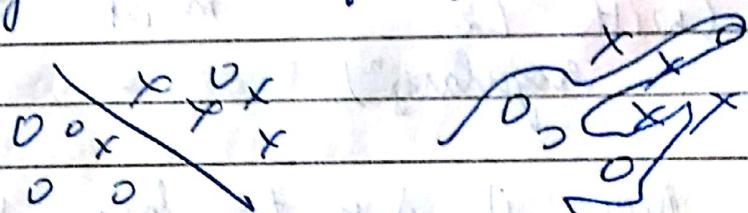
$$g(z) = \tanh(z)$$

when z
is small,

so if $\lambda \uparrow$, $w^{[l]} \downarrow$. So $z^{[l]}$ also take small values bcz

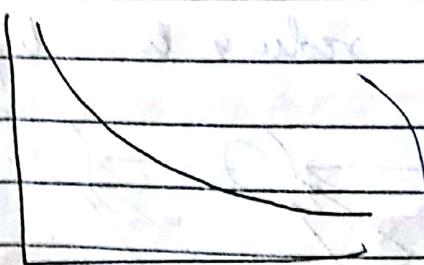
$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

Hence for this activation funⁿ results will be as in linear regn & so no curving and overfitting can occur.



not possible

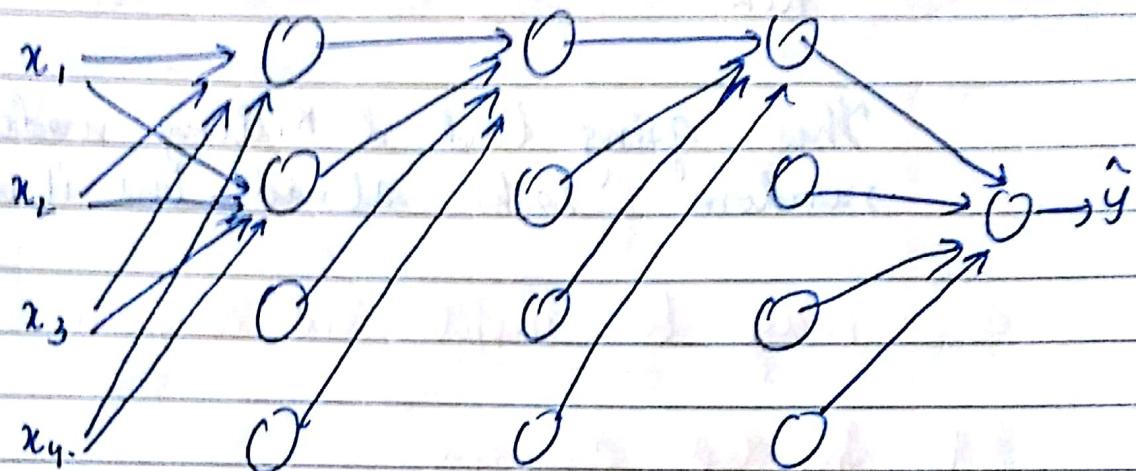
Cost



cost decreases
with each iterⁿ
of grad descent.

Dropout Regularizⁿ

→ Powerful regulⁿ technique.



When we train using this NN, it leads to overfitting.

To remove it overfitting, we do dropout
→ With dropout, we are going to go through each of the layers of the network & set some probability of eliminating a node in NN.

e.g. → for each layer, for each node, we toss a coin & have 0.5 chance of keeping the node & 0.5 chance of removing it.

So after coin tosses we eliminate those nodes. Also we remove incoming/out-going links from the node. Then we are left with a much smaller diminished network.

Then we do back propⁿ, training this e.g. on this diminished network.

Then for another eg, toss the coins & keep a diff set of nodes & for each training eg we train using new reduced netw.

This going back & killing nodes at random looks absurd but it works.

Ways to implement dropout

1. Inverted dropout

Let's assume we implement it with layer $l=3$.

d_3 = dropout vector for layer 3.

= $\text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1])$

< Keepprob

↓

Threshold

eg Keep-prob = 0.8 means 0.8 chance of keeping hidden unit & 0.2 chance of eliminating it.

a_3 = activations of layer 3

$a_3 = \text{np.multiply}(a_3, d_3)$

$a_3 / = \text{keep-prob}$ // This helps in scaling

$$z^{[4]} = w^{[4]} \cdot a^{[3]} + b^{[4]}$$

Making predictions at test time.

→ At test time, we do not dropout any hidden unit.

$$a^{[0]} = x$$

$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

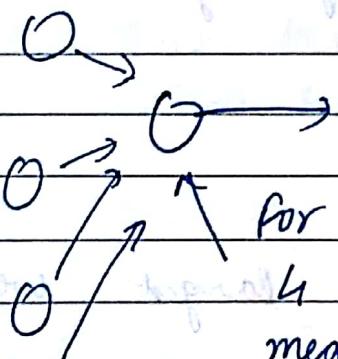
$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = .$$

.

.

\tilde{y}



For this hidden unit, it takes 4 I/Ps to generate some meaningful o/p.

Now with dropout, the I/Ps can get randomly eliminated. Everytime diff units may get eliminated.

So this unit can't rely on any 1 feature. bcz any feature can go at random.

| | | | |
|-------------------|---------|---------|-------------|
| | 0 | 0 | Date: _____ |
| $\rightarrow x_1$ | 0 | 0 | 0 |
| x_2 | 0 | 0 | 0 |
| x_3 | 0 | 0 | 0 |
| 3 I/P feature | 0 | 0 | 2 units |
| | 0 | 0 | 3 units |
| | 0 | 0 | |
| | 0 | 0 | |
| | 7 units | 7 units | |

\rightarrow sometimes it is feasible to vary keep-prob layer by layer.

$$W^{[1]} = 7 \times 3$$

$$W^{[2]} = 7 \times 7$$

so $W^{[2]}$ will be larger weight matrix of N.N.

So to reduce overfitting of this matrix $W^{[2]}$, we have to keep prob that is relatively low - say 0.5.

Layers where you worry less of overfitting you can keep higher keep-probs. eg 0.7

For layers you don't care of overfitting, you can keep keep-prob of 1 which means we are keeping every unit.