

Week #3:- Hyperparameter Tuning

Training Process:-

So, far we have seen that, training a neural network involves setting up different hyperparameters. So, what's the setting of these hyperparameters that might lead you to optimized model. Let's find out about that.

→ One of the painful things is to find out about optimal setting of hyperparameters leaving from learning rate to many other independent variables like.

- $\alpha \Rightarrow$ learning rate , $\beta \Rightarrow$ momentum term
- $\beta_1, \beta_2, \epsilon \Rightarrow$ Adam's parameters

- # of layers in neural network
- # of hidden units
- Learning rate decay
- Mini-batch size

→ Some of them are more important than the others. Learning rate is the most important one.

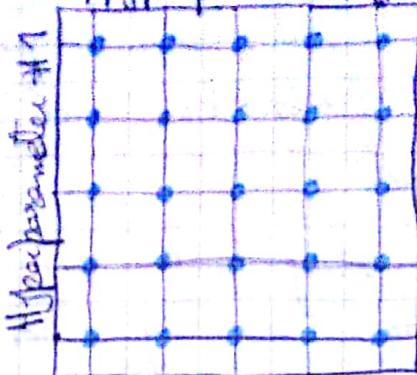
→ Now if you are tuning a set of values for these hyperparameters, then how do you go about choosing the values

→ In earlier generations of deep learning algorithms if you had two hyperparameters, it was common practice to sample the values in grid and systematically check all of them and pick

which hyperparameter works best

Hyperparameter #2

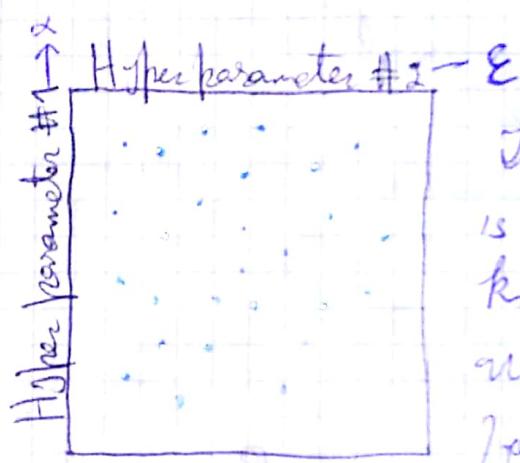
Scheme #1



This practice works
OK when no of
parameters are small.

What is recommended in practice is that, you build your list for points and choose points at random. Like previously we choose to test 25 points and choose them at random.

Scheme #2



The reason you do that is because you don't know which hyperparameters are going to be important for your problems.

- To take an example, let's say hyperparameter #1 turns out to be alpha the learning rate. To take the extreme, let's say that hyperparameter #2 is the ϵ from Adam optimization algorithm.

Following the scheme #1, you try out 25 different values for α and ϵ . You try 25 different models and find out that all of them give you the same cost. Here you try 5 different values of α .

unreasonable. In contrast, if you follow scheme #2 then you'll have checked 25 different values of α and you'll be more likely to find a value that works really well.

→ I have explained this using just 2 examples, in practice, you might be searching over many hyperparameters. So, if you are dealing with 3 hyperparameters, instead of searching over a square, you might be searching over a cube

→ Another technique in hyperparameter search is coarse to fine sampling scheme.



What you do here is that, you check this coarse grid and find that the bold points tend to work really best. What you'll do then is that, you'll select the area (legibly) around it and by optimizing your model there (selecting points randomly of course)

Using appropriate scale to pick hyperparameters:-

Lastly we saw that, how sampling the parameters for hyperparameters' values could allow you to pick your hyperparameters' values appropriately. But it turns out that, sampling at random doesn't mean sampling through random values uniformly over the range of values. Instead it is important to pick the appropriate scale & the hyperparameters' values.

Let's see how to do that.

→ Let's say you are trying to choose the value of number of hidden units and let's say you think that a good range is between 50 - 100.

$$n^{[l]} = 50, \dots, 100$$

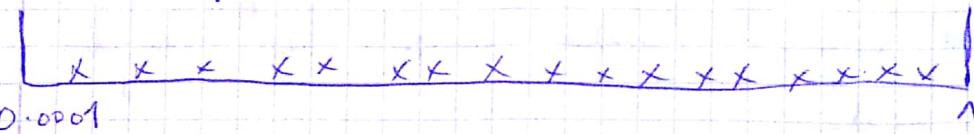
So picking a number for no. of hidden layers randomly between 50 - 100 might be OK for evaluation and see the results. But this is not OK for all hyperparameters.

Let's look at another example:-

Say You are searching for the hyperparameter "alpha"

You think that 0.0001 would be on the low end and it be as high as 1.

Now when you draw the number line



and sample values at random uniformly on this number line.

~~What's~~ You sample 90% values between 0.1 and 1
and only 10% values between 0.001 and 0.1

- That doesn't seem right. Instead it would be more reasonable to search the hyperparameters on the log scale. \uparrow Search uniformly b/w these partitions.

$$10^a \rightarrow 0.001 \quad 0.01 \quad 0.1 \quad 1 \quad \leftarrow 10^b \Rightarrow b = \log_{10}(1) = 0$$
$$a = \log_{10}(\frac{0.001}{0.01}) = -4$$

This gives you more resources to search for your hyperparameters.

In python the way you implement this is:-

Let : $r = -4 + np.random.rand() \rightarrow r \in [-4, 0]$

$$\alpha = 10^r \quad \leftarrow 10^{-4} \dots 10^0$$

$10^a \dots 10^b$ Sample r uniformly at random
between a and b $r \in [a, b]$

Hyperparameters for exponentially weighted averages:-

Let's say you suspect $\beta = 0.9 \dots 0.999$

Using $\beta = 0.9$ is like averaging over last 10 values.

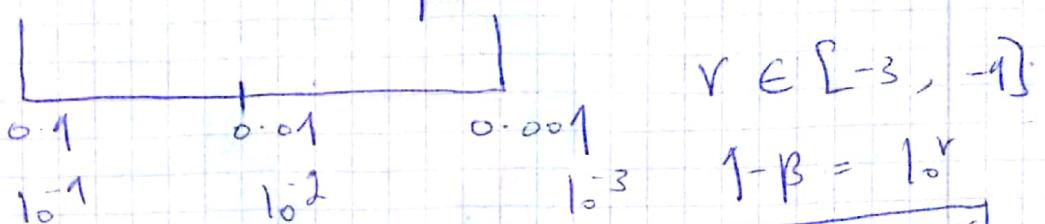
Remember $10^{-1} = 1/e$ and $\beta = 0.999$ is like
averaging over last 1000 values

\rightarrow or we can if you want to sample b/w 0.9 and 0.999, it doesn't make sense to sample on a linear scale.

Better way to think of it is to sample values

for $(1-\beta)$ i.e between 0.1 and 0.001

So maybe we will sample on the scale



Thus may you would have much more resources to explore the value of β .

Now question is why is this choosing of linear scale is a bad idea? What is the mathematical intuition behind it?

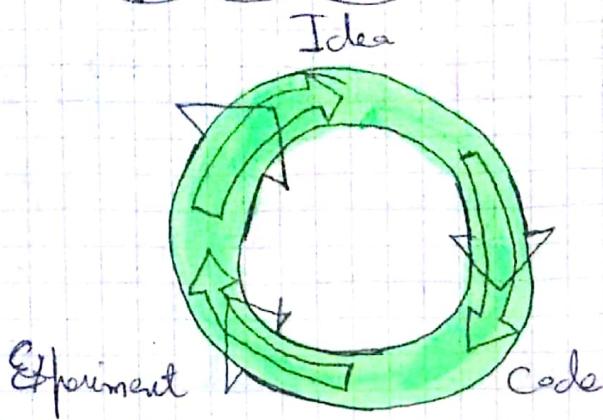
→ Now when β is close to 1. The sensitivity of your results changes even with a very small change in the input - i.e. β

→ e.g if β goes from 0.9000 \rightarrow 0.9005 it is fine no big deal.

→ But if $\beta: 0.999 \rightarrow 0.9995$ then this will have a huge impact on your results.

Because in the 1st case you are averaging over ≈ 10 examples but in the second case you go $1000 \rightarrow 2000$ examples to average over.

Hyperparameter tuning in Practice: Pandas Vs Cavier :-



Deep learning is usually applied to many fields like NLP, Vision, Speech, Ads, logistics and knowledge from one field regarding optimized hyperparameters may or may not transfer to the other field.

- But some fields' ideas can be applicable amongst each other, like Vision's can be applied to speech, NLP and vice versa.

In terms of hyperparameters, Intuitions do get stale.

- Even if you work on just one problem say logistics, you might find a good setting for your hyperparameters and kept on developing over the month on the updated data. So setting of hyperparameters can get stale, so I recommend to re-evaluate the setting of your hyper-parameters once at-least every few months to make sure you are happy with the values you have.

There are two major school of thoughts for the setting of hyper-parameters. One is Babysitting One model.

- This can be the case if you have not a lot of powerful computational resources but have huge amount of data so you can afford to train only one model.

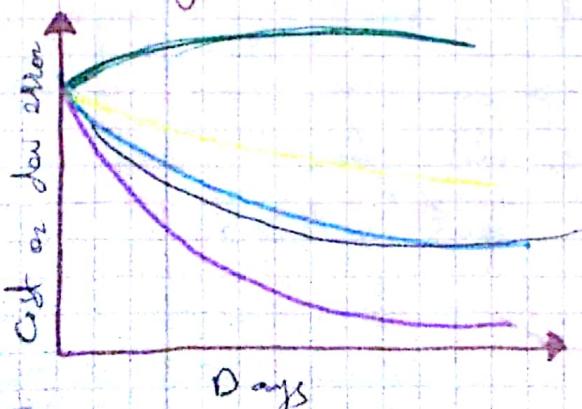
Babysitting One Model:-



The approach is that, you train the model and see its performance after fixed intervals, like after Day 1 you see your model doing pretty good so you think of increasing the learning rate. It's still doing fine. You tune other parameters and so one. Work like this until you find optimal setting.

- Other approach is that, you train many models in parallel.

Training many models in Parallel:-



Different colors
Dipict different models

Here you let the models run in parallel and analyze their performance with different hyperparameter settings.

→ In order to make an analogy, let's call the 1st approach, The Panda approach.

When pandas have children and few children, so they make sure and keep all their effort on one child at a time so that, baby panda survives. That's more like baby sitting.

→ 2nd approach is what a fish do. Let's call it Caviar strategy.

Some of the fish lay 100 million eggs in one mating season and keep their attention on a bunch of them that they survive.

→ The way to choose b/w them is that, it depends upon how much computation power you have.

Batch Normalization :-

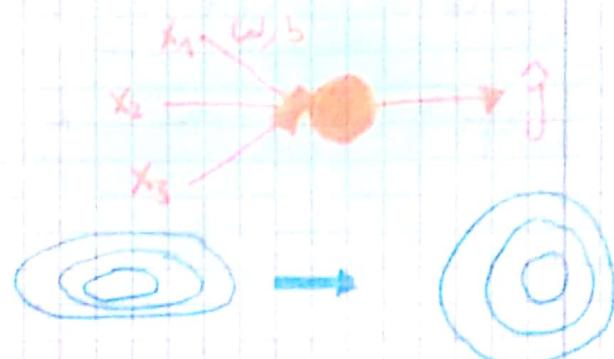
→ Normalizing activations in Network:-

One of the better ideas when it comes to deep learning is batch normalization. It makes your hyperparameter search problem much easier, makes your neural network much more robust. Choice of hyperparameters is a much bigger range of hyperparameters that work well, and will also enable you to train much more easily very deep neural networks.

Let's see how that works.

→ When training a network, you know that, normalizing the input features can speed up your learning.

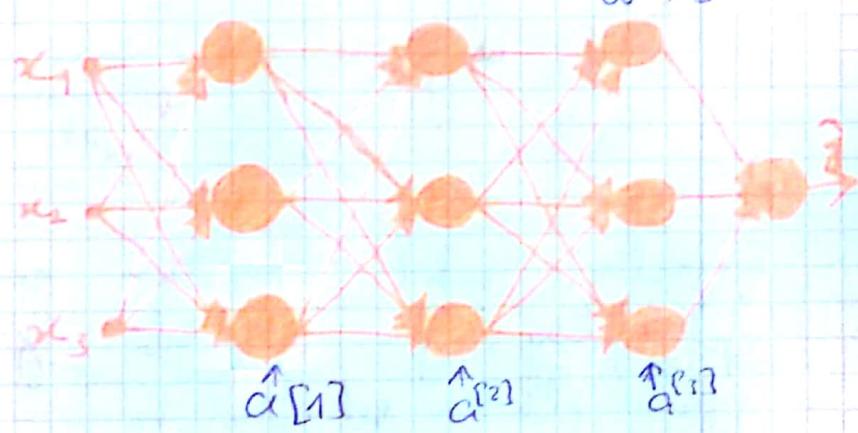
$$\begin{aligned} M &= \frac{1}{m} \sum_i x^{(i)} \\ X &= X - M \\ \sigma^2 &= \frac{1}{m} \sum_i x^{(i)2} \\ X &= X / \sigma^2 \end{aligned}$$



This works in normalizing the i/p feature values to a neural network or a logistic regression

- How about a deeper model?

Here you ~~do~~ not only have i/p features but also 3 more layers. Suppose you want to train the



- parameters $w^{[3]}, b^{[3]}$. Wouldn't it be nice to train them by normalizing the input ($a^{[2]}$)?

Here you'll normalize $z^{[2]}$ not $a^{[2]}$.

→ There is a debate regarding whether to normalize activations of z 's. But in practice more often z 's are normalized.

Here is, how you implement batch normalization.

Given some intermediate values in neural network

$$z^{[1]}, \dots, z^{[l]}$$

Compute mean $\Rightarrow \bar{y} = \frac{1}{m} \sum z^{[1]}$

$$\sigma^2 = \frac{1}{m} \sum (z_i - \bar{y})^2$$

$$z_{\text{norm}}^{[1]} = \frac{z^{[1]} - \bar{y}}{\sqrt{\sigma^2 + \epsilon}}, \text{ for numerical stability.}$$

This will give you mean zero and s.d 1 for every layer and we don't want that.

So what we are gonna do :-

$$z^{[1]} \xrightarrow{\gamma z_{\text{norm}} + \beta} \text{learnable parameters of model.}$$

Regardless what algorithm you are using for optimization, you update γ and β just like you do w and b .

→ Effect of γ and β is that, you can set the mean of $z^{[1]}$ whatever you want it to be.

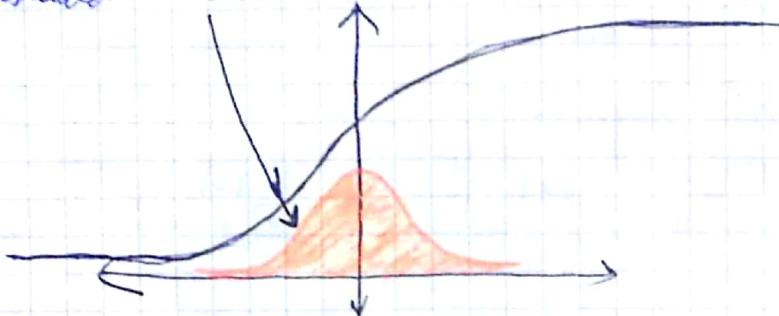
e.g. if $\gamma = \sqrt{\sigma^2 + \epsilon}$ and $\beta = \bar{y}$, then

$$z_{\text{norm}}^{[1]} = z^{[1]}$$

→ Now you'll use $z^{[1]}$ instead of $z_{\text{norm}}^{[1]}$ in your computations of neural network.

One difference though between input values and hidden units is that, you might not want your hidden unit values to be forced to have mean zero and variance 1.

e.g. if you are using Sigmoid activation functions, you might not want all of your values to be clustered like this:-



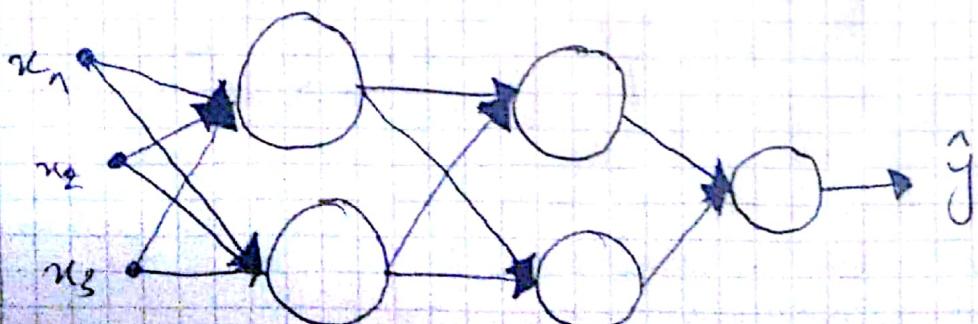
You might want it to have larger variance or mean that is different than zero, in order to take advantage of the non-linearity of sigmoid function in a better way.

That's why with parameters γ and β , you can make sure that, your network has the range of values you want.

Fitting batch norm into a neural network :-

Let's see how to fit a batch norm to a whole neural network.

Let's consider the following network



Each of the unit compute 2 things :-

$$\begin{aligned} z^{[i]} &= w^{[i]} a^{[i-1]} + b^{[i]} \\ a^{[i]} &= g(z^{[i]}) \end{aligned} \quad \text{Normally}$$

But what we do in batch norm is that :-

$$X \xrightarrow{w^{[i]}, b^{[i]}} z^{[i]} \xrightarrow[\substack{\text{batch} \\ \text{norm} \\ (\text{BN})}]{} \tilde{z}^{[i]} \xrightarrow{g(\tilde{z}^{[i]})} a^{[i]}$$

Now for next layer :-

$$a^{[1]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\substack{\text{BN} \\ \dots}]{} \tilde{z}^{[1]} \xrightarrow{g(\tilde{z}^{[1]})} a^{[2]} \xrightarrow{\dots}$$

So in short for firing of neurons, instead of using unnormalized z , we are going to use normalized values \tilde{z} .

→ Now we have parameters, $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$
 $\gamma^{[1]}, \beta^{[1]}, \gamma^{[2]}, \beta^{[2]}, \dots, \gamma^{[L]}, \beta^{[L]}$

→ Now that, you have other parameter too in your model rather than w and b , you will use optimization algorithm to update them too.

Like compute $d\beta^{[l]}$ and $d\gamma^{[l]}$ and

then update them usually like

$$\beta_{\text{new}}^{[l]} = \beta^{[l]} - \alpha d\beta^{[l]}$$

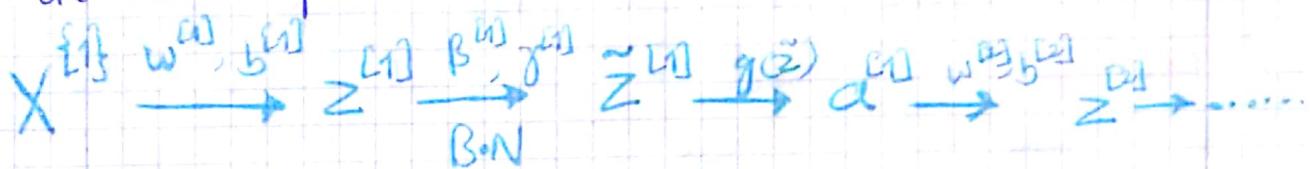
$$\gamma_{\text{new}}^{[l]} = \gamma^{[l]} - \alpha d\gamma^{[l]}$$

In programming frameworks for Machine learning like

Tensorflow or Pytorch you don't have to implement them by yourself from scratch, you can use pre-defined functions from those libraries.

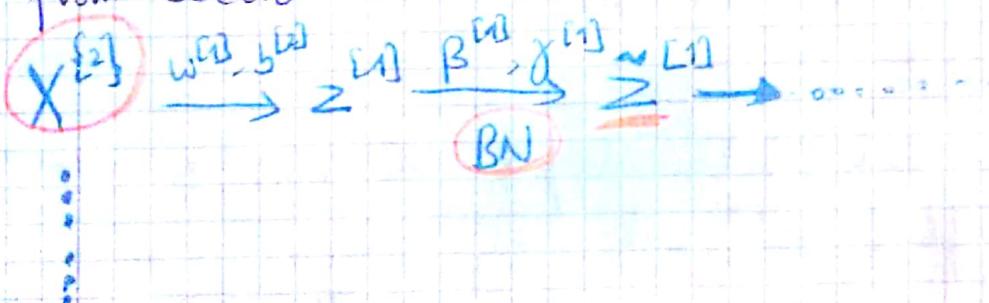
Till now we have talked about batch normalization like we are using batch gradient descent. But in **practice** we apply this technique with mini-batches.

→ Now what you do is that, you take first mini-batch and apply batch-norm on that mini-batch only and follow the rest of the trend as it then repeat the whole process with other mini-batches.



Now you do the same with other mini-batches.

→ One thing to be noted here is that, your values of μ and σ here which are calculated from just using the data from second mini-batch.



One other detail regarding the parameterization of the neural network now is that as follows:-

We established that we have following parameters:-
 $w^{[1]}, b^{[1]}, \mu^{[1]}, \sigma^{[1]}$

Notice that $Z^{[2]}$ is computed like this

$$z^{[l]} = \frac{x^{[l]} - \mu}{\sigma} + b$$

then normalize $z^{[l]}$ to first have mean zero and standard variance, then rescale them by β and γ .

\rightarrow what that means is that, whatever is the value of $b^{[l]}$ is just gonna get subtracted out because in that batch normalization, you are going to compute the mean of $z^{[l]}$ and subtract out the mean. So any constant you add will get canceled out by the "mean subtraction step".

So, what you can do is eliminate b terms.

\rightarrow As batch norm zeros out the mean so there is no point in having " b " term there, instead it is kind of replaced by β which now control the mean and variance:

$$\beta^{[l]}, \gamma^{[l]} \in \mathbb{R}^{(n^{[l]} \times 1)}$$

Implementing Gradient Descent Using minibatch

for $t = 1 \dots$ num batches:

Compute forward prop on $X^{[t]}$

In each hidden layers use BN to replace $z^{[l]}$ with $\tilde{z}^{[l]}$

Use back prop and compute $dW^{[l]}, d\beta^{[l]}, d\gamma^{[l]}$

update parameters

$$w^{[l+1]} = w^{[l]} - dW^{[l]}$$

$$\beta^{[l+1]} = \beta^{[l]} - d\beta^{[l]}$$

$$\gamma^{[l+1]} = \gamma^{[l]} - d\gamma^{[l]}$$

works with momentum, RMS as well

Why does Batch Norm Work?

You have already seen that by normalizing input features you can speed up learning. One intuition is that, you are doing the same thing in the hidden units.

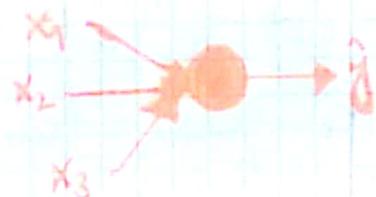
Now this is the shallow picture; there are couple of other things which can help you to gain better understanding of what batch norm is doing.

- One reason is that, batch norm makes weights later in the neural network more robust to changes earlier in the neural network.

To explain what does it mean, let's look at this example:-

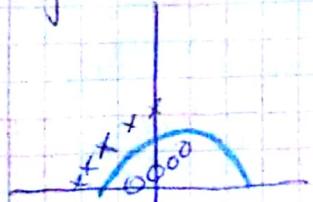
Let's say this shallow network,

may be a logistic regression.



- Let's say you have trained your neural network on all black cats. If you now apply your neural network to data with colored cats then your classifier might not do very well

e.g. if you have an algorithm where you use training data like this:- (the one on left)



then your model
can't perform well
on the ~~new~~ data



This idea of data distribution changing goes by the name of "**Covariate shift**".

The Idea is that -

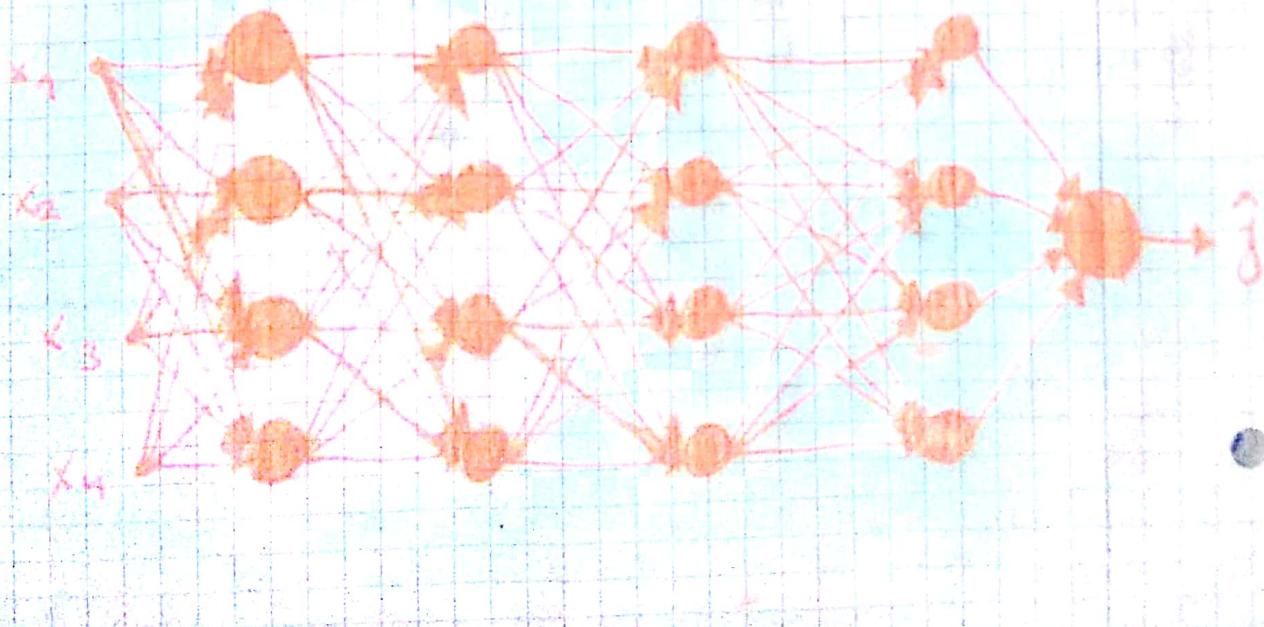
if you are using some data to train model that shows $X \rightarrow Y$ mapping and after training if your data distribution changes then you will need to ~~stop~~ re-train your model.

→ This is true even if ground true function mapping from X to Y remains unchanged. Like in the earlier problems that we tackled in assignments ask the ground true function uses, whether this is bad or not.

If the ground true function changes then the need to retrain your function becomes even worse.

Now:- how does this problem of Covariate Shift apply to neural network?

Consider a deep neural network like this.



Let's look at the learning process from the perspective of 3rd hidden layer:

So, this network, has to learn parameter $w^{[3]}, b^{[3]}$

So, this layer got parameters $a_1^{[2]}, a_2^{[2]}, a_3^{[2]}, a_4^{[2]}$ from previous layer and the job of 3rd hidden layer is to find a mapping that takes these incoming values and map them to g .

Network is also adapting parameters $w^{[2]}, b^{[2]}$ and $w^{[1]}, b^{[1]}$ so, as these values change, $a^{[2]}$ will also change and this is the problem of covariant shift.

What batch norm does is that, it reduces the amount of ~~the~~ the distribution of ~~these~~ which these hidden units' values shift around with.

If on one iteration, if $a_1^{[2]}, a_2^{[2]}$ change then on the next iteration it will again change. Batch norm ensures that, the mean and variance remains the same even if the values change.

→ Batch norm reduces the problem of input values changing. It really causes these values to become more stable so that later values have much more stable grounds to stand on.

Even if the earlier distribution changes, it changes less, and what it does is it forces less the ~~other~~ later layers to change as a result of change.

in earlier layers. In loose terms, it loosens the coupling of earlier layers with the later ones, and enable each layer to learn sort of independently. As a result this part in effect the speeding of neural network training.

→ It turns out batch norm has slight regularization effect too.

Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
- This adds some noise to the values $z^{[l]}$ within the mini-batch. So similar to dropout, it adds some noise to each hidden layer's activation.
- This has a slight regularization effect.
 - Because by adding noise to the hidden units, it is forcing the downstream of hidden units to rely too much on any one hidden unit.
 - Because the added noise is small, this is not a huge regularization effect and you might wanna use batch norm together with dropout.
 - If you increase the batch size, you reduce the noise and hence the regularization effect gets damped out.

Batch Norm at test time :-

- \rightarrow Batch norm processes your data, one mini-batch at a time.

But at test time, you might need to process your examples one at a time.

Let's see how you adapt your network to do that.

Recall the equations you use for training of examples:-

$$\rightarrow \bar{y} = \frac{1}{m} \sum_i z^{[i]} \quad \boxed{\text{Computed on entire mini-batch.}}$$

$$\rightarrow \sigma^2 = \frac{1}{m} \sum_i (z^{[i]} - \bar{y})^2$$

$$\rightarrow z_{\text{norm}}^{[i]} = \frac{z^{[i]} - \bar{y}}{\sqrt{\sigma^2 + \epsilon}}$$

$$\rightarrow \tilde{z}^{[i]} = \gamma z_{\text{norm}}^{[i]} + \beta$$

- \rightarrow So, in order to apply it at test time, you need to come up with different estimate of \bar{y} and σ^2 .

One estimate used is that, you estimate using exponentially weighted average (across mini-batches).

Let's pick some layer L and going through mini batches

$$x^{[1]}, x^{[2]}, x^{[3]}, \dots$$



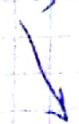
$$m^{[1]L}$$

θ_1



$$s^{[1]L}$$

θ_2



$$\sigma^{[1]L}$$

θ_3

You saw earlier that how to use weighted average to compute the means. You could do that to keep track of what's the exponentially weighted

average becomes of this mean vector that is the latest. So that exponentially weighted average becomes your estimate for what the mean of Z_s is for that hidden layer and similarly you use this exponentially weighted average to keep track of these values of sigma squared that you see on the first mini-batch in that layer, in the second mini batch and so on.

→ At test time you would use that exponentially weighted average for M and σ^2 to compute

Z_{norm}

$$Z_{\text{norm}}^{(c)} = \frac{Z^{(c)} - M}{\sqrt{\sigma^2 + \epsilon}}$$

Then you would compute $\tilde{Z}^{(i)}$ using γ and β that you learned during your learning of neural network.

$$\tilde{Z}^{(i)} = \gamma Z_{\text{norm}}^{(c)} + \beta$$

→ You can also run your whole network (final) through training data and use those computed final M and σ^2 .

But in practice people use exponentially weight average of M and σ^2 .

Multi-class Classification:-

Softmax Regression:-

- Till now we have considered only one class for classification, like is it a cat or not?

What if you have multiple classes.

There is a generalization of logistic regression called "Softmax regression", where you try to recognize one of multiple classes, rather than one.

- Let's say instead of recognizing cats only, you need to recognize cats, dogs and baby chicks

Let's say Dogs \Rightarrow class 2

Cats \Rightarrow class 1

Baby chicks \Rightarrow class 3

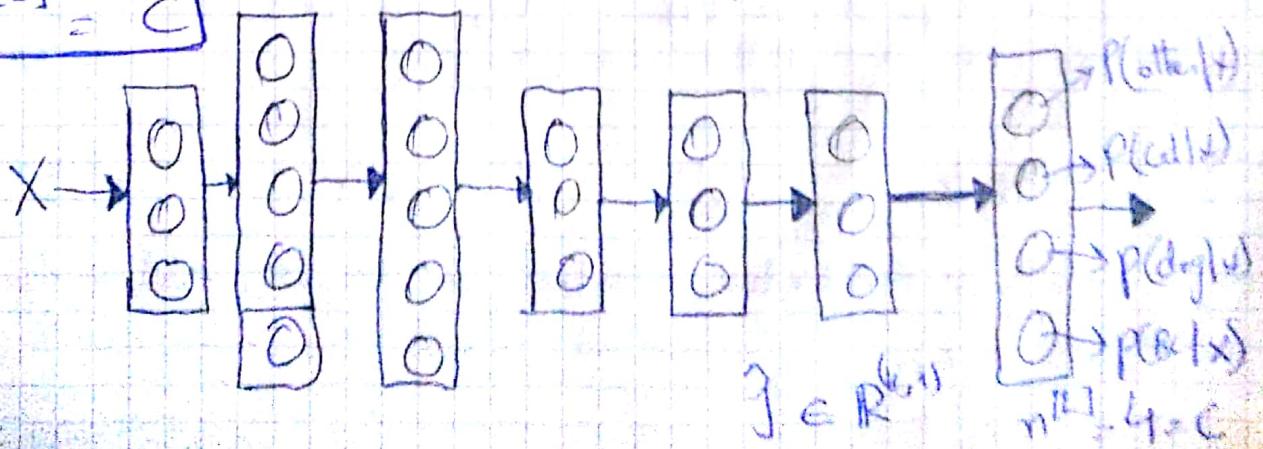
and if none of the above then class 0 or other

Notations:-

$C = \# \text{ of classes} \Rightarrow$ In this class we have 4
In code it is from (0, ..., 3)

In this example we are going to build
a neural network with C output units

$$n^{[L]} = C$$



→ Should sum to 1.

- The standard model that a neural network uses to perform this task is called softmax layer in the output layer in order to generate these outputs.
- So, in the final layer: you are going to compute:

$$Z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$$

Now having computed $Z^{[L]}$, you now need to apply softmax activation function.

Activation function:-

$$t \in \mathbb{R}^{(4 \times 1)}$$

for our example

$$t = Z^{[L]} \Rightarrow \text{elementwise}$$
$$a^{[L]} = \frac{e^{z_i^{[L]}}}{\sum_{j=1}^4 t_j}$$

$$a_{(1)}^{[L]} = \frac{t_1}{\sum_{j=1}^4 t_j}$$

e.g.: $Z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \Rightarrow t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 168.4 \\ 7.4 \\ 0.4 \\ 20.4 \end{bmatrix}$

So $\sum_{j=1}^4 t_j = 176.3$

$$a^{[L]} = \frac{t}{176.3}$$

So $a_{(1)}^{[L]} = 0.842$

$$a_{(2)}^{[L]} = 0.042$$

$$a_{(3)}^{[L]} = 0.002$$

$$a_{(4)}^{[L]} = 0.114$$

Summarize

$$a^{[L]} = g^{[L]}(Z^{[L]})$$

Softmax activation, takes a vector and outputs a vector.

Training a Softmax Classifier:-

Let's try to dig deeper into softmax classifier.

Let's say we have $C=4$ and recall our original example - $\mathbf{z}^{(1)} \in \mathbb{R}^{(4,1)}$

$$\mathbf{z}^{(1)} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \Rightarrow t \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \Rightarrow \mathbf{a}^{(1)} = \begin{bmatrix} 0.842 \\ 0.142 \\ 0.002 \\ 0.914 \end{bmatrix}$$

The name softmax comes from in contrast to what comes as "hardmax".

What hardmax does is takes a vector and puts 1 in place of the biggest element of the vector.

like for this example it will be $[1 \ 0 \ 0 \ 0]^T$

→ whereas softmax is a gentler mapping from \mathbf{x} to outputs.

Softmax regression generalizes logistic regression to "C" classes

if $C=2$ then softmax essentially reduces to logistic regression. So, C should be more than 2.

→ Let's see how we train a neural network with a softmax output layer.

In particular let's define the "loss function".

You used to train your neural network.

Let's say you have an example with a target output of ground true label of $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$

So, following previous example, this would be an image of a cat.

You get your output vector let's say,

$$a^{(3)} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

so clearly our model is not doing very well in this example

In softmax classification, the loss we typically use is :-

$$L(\hat{y}, y) = -\sum_{j=1}^4 y_j \log \hat{y}_j$$

$$\hat{y}_1 = \hat{y}_3 = \hat{y}_4 = 0$$

so only at $\hat{y}=2$ this function will give some output.

→ As we are trying to minimize the loss so, $-\log(\hat{y}_2)$ should be bigger and for it to be bigger correspondingly y value should be bigger also

→ This is the loss on one training example, how about the cost on the whole data.

$$J(w^{(3)}, b^{(3)}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

Output vector:- $Y = [L(\hat{y}^{(1)}, y^{(1)}), \dots, L(\hat{y}^{(m)}, y^{(m)})]$

$$= \begin{bmatrix} 1 & \dots & 1 \\ \hat{y}^{(1)} & \dots & \hat{y}^{(m)} \end{bmatrix} \in \mathbb{R}^{(k, m)}$$

Some dimensions will be of Y .

- Final detail, let's take a look about how do you implement gradient descent when you have a softmax output layer.

Key point is to start backprop and you do it like this

$$dZ^{(l)} = \hat{y} - y$$

(h, 1)

Rest is the same.

Introduction to Programming frameworks:-

- Till now we have learned to implement deep neural networks from scratch. But when we train very big models or heavy applications, it is not practical to do that, rather there are very sophisticated programming frameworks that allow you to do the aforementioned job
- It's just like multiply 2 matrices by yourself via using explicit loops and using numerical linear algebra functions.

Different deep learning frameworks:-

Caffe/Caffe2, CNTK, DL4J, Keras, Lasagne
mxnet, PaddlePaddle, TensorFlow, Theano, Torch

Criterion for Choosing framework:-

- Ease of programming (development and deployment)
- Running speed
- Truly open (open source with good governance)

Tensor Flow

As a motivating problem, let's say you have some function

$$J(w) = w^2 - 10w + 25 \quad \textcircled{1}$$

$$J(w) = (w-5)^2$$

The value that minimizes this is $w=5$

Let's say we didn't know that and we have just
① function and let's see how we can implement
tensor flow to minimize this.

Import tensorflow as tf

Import numpy as np

$w = tf.Variable(0, dtype=tf.float32)$

Cost = tf.add(tf.add(w**2, tf.multiply(-10, w)), 25)

train = tf.train.GradientDescentOptimizer(0.01).minimize(Cost)

↑
learning rate

init = tf.global_variables_initializer()

Session = tf.Session() OR with tf.Session() as session

Session.run(init)

Print(session.run(w))

⇒ 0.0

Session.run(train)

Print(session.run(w))

0.1

for i in range(1000):

Session.run(train)

Print(session.run(w))

⇒ 4.999

This was a fixed function, when learning a neural network, training data can change. How to load that data.

X = tf.placeholder(tf.float32, [3, 1]) → tells the program that, you'll provide the value for this variable

$$\text{Cost} = x[0][0] * w^{**2} + x[1][0] * w + x[2][0]$$

Coefficients = np.array([1.1, -1.0, 25.3])

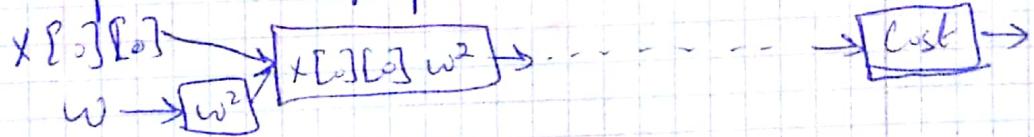
→ Now we have to get coefficients into X
for that, we write follows:-

Session.run(train, feed_dict={x: coefficients})

Print(session.run(w))

→ You just have to give tensorflow the function and it will it self figure out the derivatives and all.

→ what highlighted line is doing is allowing tensorflow to build a computational graph



One nice thing about tensorflow is that you give it the functionality for forward propagation and it has already built-in functionality for backward propagation.

→ By using the built-in functions for forward function, it can automatically do the backward function as well to compute backpropagation.

Writing a program in the tensorflow has the following steps:-

- 1- Create tensor (variables) that are not yet executed/evaluated.
- 2- Write operations b/w those tensors.
- 3- Initialize your tensors.
- 4- Create a session.
- 5- Run the session. This will run the operations you'd written above.

E.g:-

$\hat{J} = \text{tf.constant}(3b, \text{name} = 'J-hat')$

$J = \text{tf.constant}(39, \text{name} = 'J')$

$\text{loss} = \text{tf.Variable}((\hat{J} - J)^2, \text{name} = 'loss')$

$\text{init} = \text{tf.global_variables_initializer()}$

With tf.Session() as sessions

Session.run(init)

$\text{Print(Session.run(loss))}$