

26/7/18

Course-2

Improving DNN, Hyperparameter tuning, Regularization and Optimization.

@ Week 1:

Data:	60%.		
	train	cv	test

If dataset is relatively small, traditional ratios might be okay but in case of a much larger dataset, we set cv (devset) and test set to be much smaller than 20% or even 10% of entire data.

e.g. in 10,000 dataset, 6k, 2k, 2k distribution is fine i.e 60-20-20% but in 10 million " , just 2% or even 0.02% data is enough for cross validate and test set.

Dev set is to check which of the different models applied on train set is the best to give results. If for example:

train set

Cat pictures from webpages (high resolution images)

dev/test set

Cat pictures from users using the cat recognition app (blurry and comparatively low resolution images)

In such a case, ensure that test set and dev set come from the same distribution of data. Doing this, progress in ML would be much faster.

Note It is okay to not have a test set bcoz it is there just to have an unbiased estimate of performance. Many a times only train set and dev set are given, and people start calling this dev set the test set, which is not a very good use of terminology as they are then overfitting the test set, but it is okay if we don't need an unbiased estimate of the performance and all the data that we are concerned about is only that train and dev set.

This train, dev and test setup will allow a more efficient measure of the bias and variance of our algorithm and better ways can be thus selected to improve our algorithm.

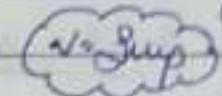
• Bias and Variance:

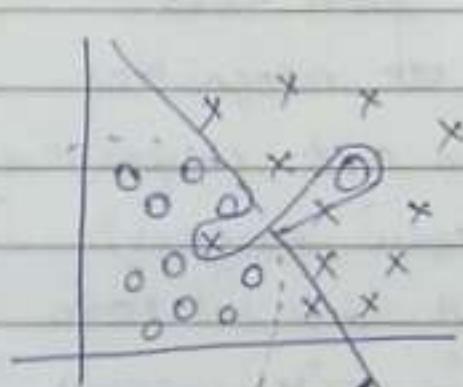
Train set error: 1%. 15%. 15%. 0.5%.

Dev set error 11%. 16%. 30%. 1%.

high variance	high bias	high bias & high variance	low bias & low variance
---------------	-----------	---------------------------	-------------------------

- assuming train & dev sets are drawn from the same distribution.
- assuming human error or Optimal bayes error = 10%

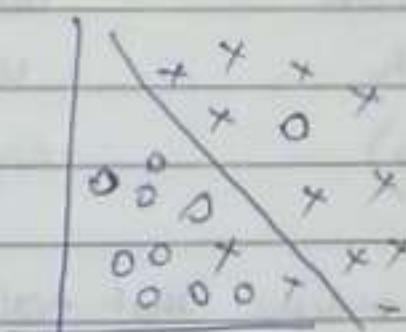
So, by looking at the training set error, one can get an idea of how well the model fits atleast the training set and tells us if we have a bias problem, and then seeing how much the error rises on going from train to dev set, we get a sense of how bad is the variance problem. 



high bias &
high variance
model.

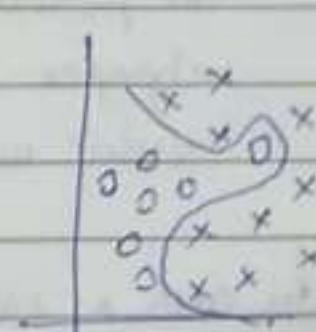
(bcz in some region it shows high bias and in some it shows high variance)

by taking into consideration even the mislabeled or outlier data)



High bias
to solve this:

- use bigger n/w
- train longer
- NN architecture search



high variance
to solve this:

- get more data
- do regularization
- NN architecture search

Low bias, low variance - desired.

Note

Bias variance trade off is an old concept when we didn't have tools which could just reduce one error without affecting the other.

steps like \uparrow up data or doing regularization just reduces variance without affecting bias, and getting a bigger ^{NN for} data just reduces bias without affecting variance much so long as regularization is done, which is a technique to reduce variance. Regularization \downarrow s the bias a bit but not much if we have a huge enough network.

→ Regularization: reduces overfitting and reduces errors in the network.

• in logistic regression -

$$\min_{w, b} J(w, b) \quad w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \left(\frac{1}{2m} \|w\|_2^2 \right)$$

$$\|w\|_2^2 = \sum_{j=1}^{n_y} w_j^2 = w^T w$$

(norm of w , squared)

penalizes the weight matrices from being too large.

[L₂ regularization]

bcz here euclidean norms are used also called L₂ norm with parameter vector w .

Note regularization term is added only for w and not for b like $\frac{1}{2m} b^2$ because w is a high dimensional parameter vector, especially with a high variance problem while b is just a single number unlike w which might not be fitting all the parameters. So, adding $\frac{1}{2m} b^2$ won't make much of a difference because b is just 1 parameter over a very large no. of parameters in w .

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

$$= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{1}{m} \|w\|_1, \quad [\text{L}_1 \text{ regularization}]$$

↳ just a scaling constant.

Since L₁ regularization, w vector will end up being sparse i.e. it will have a lot of zeroes in it, so takes less memory space. L₂ is more often used than L₁.

L called the regularization parameter is set using the cross validation set and needs to be tuned.

- in Neural Network: {L = no. of layers}

$$J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|^2$$

\downarrow
squared norm

$$\|w^{(l)}\|^2 = \sum_{i=1}^m \sum_{j=1}^{n^{(l)}} (w_j^{(l)})^2 \quad \left\{ w: (n^{(l)} \times n^{(l-1)}) \right\}$$

no. of hidden units in layer l-1 and l.

F: "Frobenius norm"

(not called L₂ norm bcoz of linear algebra technical reasons.)

Frobenius norm simply means sum of square of elements of a matrix.

Now, $d\omega^{(l)} = (\text{from backprop}) + \frac{\lambda}{m} \omega^{(l)} \quad \text{--- } ①$

$$\omega^{(l)} = \omega^{(l)} - \alpha d\omega^{(l)} \quad \text{--- } ② \quad \hookrightarrow \text{bcoz of regularization term}$$

we know, $d\omega^{(l)} = \frac{\partial J}{\partial \omega^{(l)}}$ (derivative of the cost fn w.r.t the parameters w)

Hence, L₂ regularization is also sometimes called "weight decay".

Putting ① in ②:

$$\omega^{(l)} = \underbrace{\omega^{(l)} - \frac{\lambda}{m} \omega^{(l)}}_{+} - \alpha (\text{from backprop})$$

i.e whatever the matrix $\omega^{(l)}$ is, we are going to make it a little bit smaller.

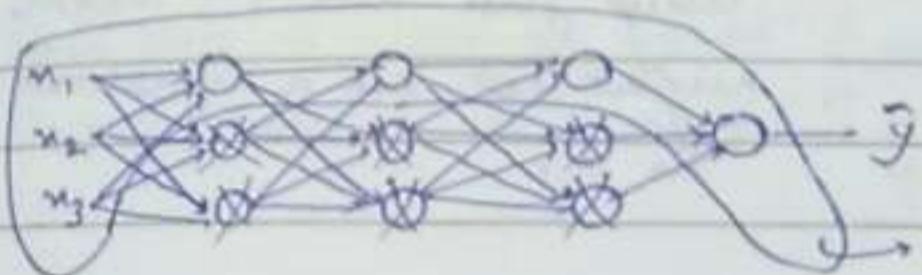
\rightarrow little less than 1
 $w^{(l)} = \left(1 - \frac{\alpha}{m}\right) w^{(l)} - \alpha$ (from backprop)

↓
 this is why L2 regularization is called "weight decay", bcoz α times the same gradient is subtracted but from a reduced value of $w^{(l)}$.

27/7/18

Q: Why regularization helps with overfitting?

Ans: Regularization shrinks the parameters and if made close to zero, it's basically zeroing out a lot of the impact of the hidden units.

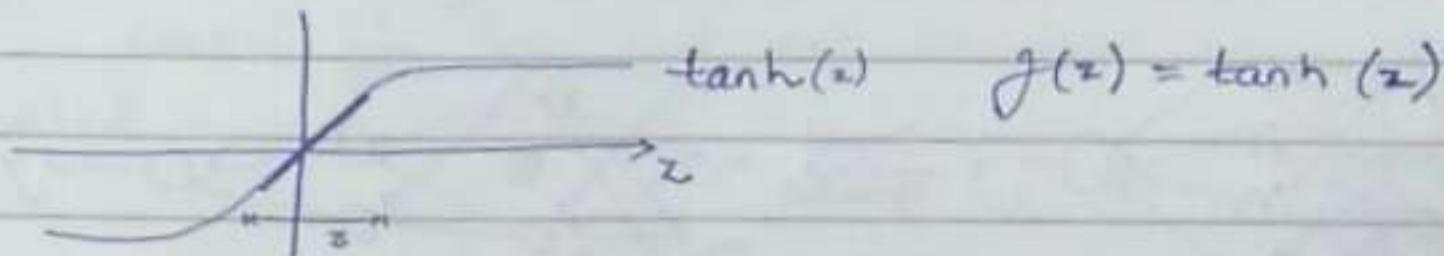


(NN reduced after regularization)

seems like logistic regression is being used.

removing hidden units this way, shift the model from overfitting to the underfitting side and there exists λ which will settle it midway giving just the right fit. Hidden units are not completely removed, it's just that their impact is reduced.

Also,



If z takes on only a small range of parameters, tanh is linear. Only if larger values are allowed for z , the activation function starts to become less linear.

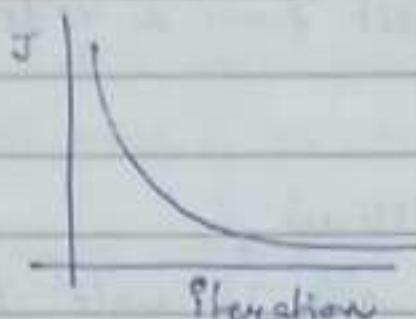
so, if λ is large, then our parameters will be relatively small bcoz they are penalized being large into a cos fn.

bcoz, $z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$

so, this makes the whole deep NN also a linear network, bcoz of the linear tendency of tanh for small z , this linearity is unable to fit very complex decisions thus allowing

"it to avoid overfitting."

So, large $\lambda \rightarrow$ small $w \rightarrow$ small $z \rightarrow$ linear model which doesn't overfit.



J goes monotonically after every "step" of gradient descent provided J is taken with regularization i.e.

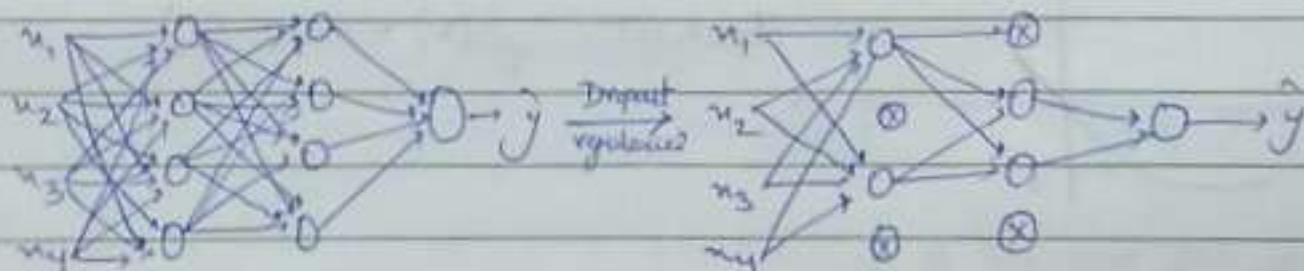
$$J = \sum l(\hat{y}, y) + \frac{\lambda}{2m} \|w\|_F^2$$

Plot J as fⁿ of no. of iterations of gradient descent to debug gradient descent.

Using just this 1st term might not see a monotonic ↓ in J.

- Dropout regularization:

In this regularization, we go through each layer and set some probability of eliminating a node in neural network. Then some neurons are completely dropped and the backpropagation is applied on the much diminished network.



In this manner, different set of nodes are dropped each time and each training example is trained using one of these NN with dropped layers, resulting into regularization.

Implementing dropout ("Inverted dropout")

Step: d₃ = np.random.rand(a₃.shape[0], a₃.shape[1]) < keep_prob

dropout vector

(probability that a given hidden unit will be kept.)

if keep-prob = 0.8 means that there is a 0.2 chance of eliminating any hidden unit.

Step 2: $a_3 = np \cdot \text{multiply}(a_3, d_3)$ # $a_3^* = d_3$
(element wise multiplication)

for every element of d_3 that's equal to zero (given there are 20% chance of it being 0), all the corresponding elements of a_3 end up being zero thus giving 0 activations.

Step 3: $a_3^* = \text{keep-prob}$ {inverted dropout technique'}

This step is done so as to keep the value of $z^{[L]}$ intact.

We know, $z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$
↑ reduced by 20% due to dropout regularization. So to compensate the loss 0.8 is divided to make up for the 20% loss.

So, inverted dropout is useful becz no matter what is the value of keep prob, it is ensured that the expected value of a_3 remains the same. When this step 3 is not used, error in the test set becomes more and more complicated.

'd' vector ensures that for different training ex.s, different hidden units are removed (bcz of the random fn used) and in fact for different iterations also, the same training example is trained by different network, with some different hidden units removed.

'd' decides which units to drop and which to keep during fwd and bwd propagation.

While doing predictions, no dropouts are done becz we don't want our output to be random. The complete original NN is used to make prediction having weight set without getting overfitted as dropouts were made at time of training.

Adding dropouts at test time just adds noise to the predictions, nothing else. Also that step 3 was done previously only so that it was not needed at test time to keep α_3 intact.

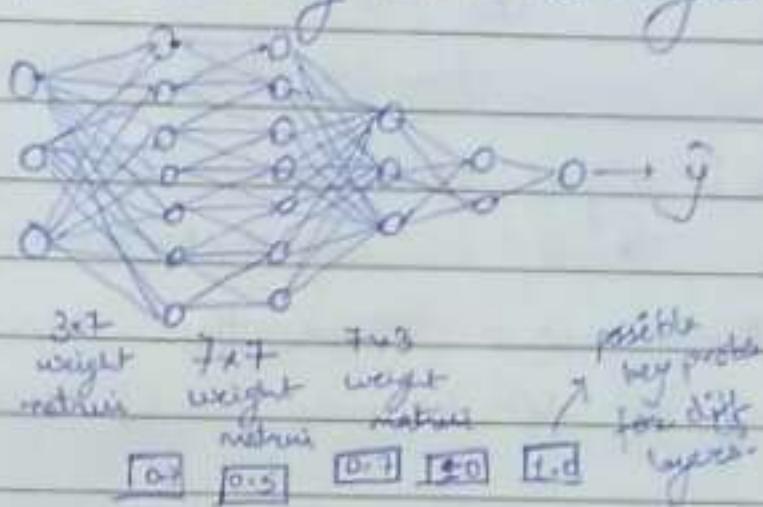
28/7/18

Q why does dropout work?

A As our network can't rely on any one feature, so the weights need to be spread out. Any one feature can go away at random and any one of its inputs could go away at random. So, it's unwise to put too much weight on any one input because it can go away. This dropout will shrink the squared norm of the weights and does some L2 regularization that helps everything prevent. Results of L2 regularization and dropout regularization are similar.

Note 'keep-prob' can vary with each layer also in the NN.

e.g.



In this case, 'keep-prob' for layer 2 can be relatively low bcoz of the big 7×7 matrix, so as to prevent overfitting.

'Keep prob' can be bigger like 0.7 or so when there is less fear of overfitting. 'Keep prob' 1 means no dropout for that layer.

This is like L2 only as there also different layers are having different regularization. The downside for this is that it increases the number of hyperparameters to be searched for using cross-validation, also J is not well defined.

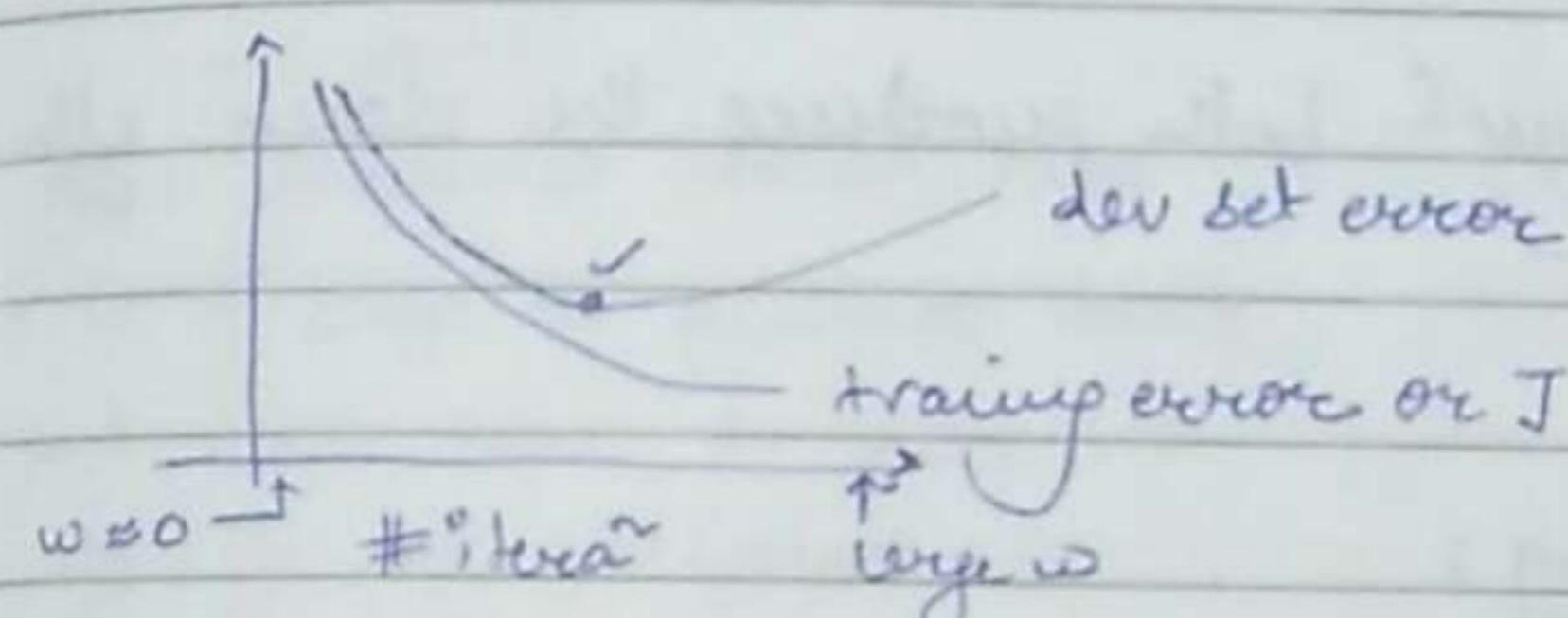
Sometimes, not very often keep-prob is applied to the input layer as well.

J can't be now used to double check the performance of gradient descent, as now J is certainly hard to calculate bcoz of being less defined with regularization.

- Other regularization methods:

1. Data augmentation - bcoz every time collecting more data can be infeasible and expensive. Though this won't be as good as taking new pictures or new dataset but is free of any expenses. The augmented data does not add that much info as a brand new image does but is an inexpensive way to give more data to the algo.

2. Early stopping - ,



Usually the dev set error goes down for a while and it will ↑ from there. Early stopping means we stop the iterations at the point where it looks like the NN is doing best and take whatever value achieved this best dev set error. As the training is done, w gets bigger and bigger unlike the small values (usually 0) with which it was initialized. So, stopping the process midway give us a mid-sized w , and hence giving a less overfitting NN.

Note

Downside of early stopping is that it couples the two tasks of optimizing cost fn and trying not to overfit thus hindering the orthogonality i.e. in early stopping we no longer can work on these two problems independently. Bcoz stopping midway is not letting the $J(w, b)$ minimize in order for the model to not overfit.

Optimizing cost fn needs Gradient descent, adagrad, etc.

Preventing overfitting needs regularization, more data, etc.

Q Difference b/w early stopping and L2 regularization?

Abs Since L2 regularizes, we can train the NN as long as we want and that makes the search space of hyperparameters easier to decompose and easier to search over but the downside is that we might have to try a lot of values of λ making it computationally expensive.

Adv. of early stopping is that running the gradient descent process just once, we get to try out values of small w , mid-size w and large w without needing to try a lot of values as in L_2 .

L_2 is still used more, though both produce the same effect more or less.

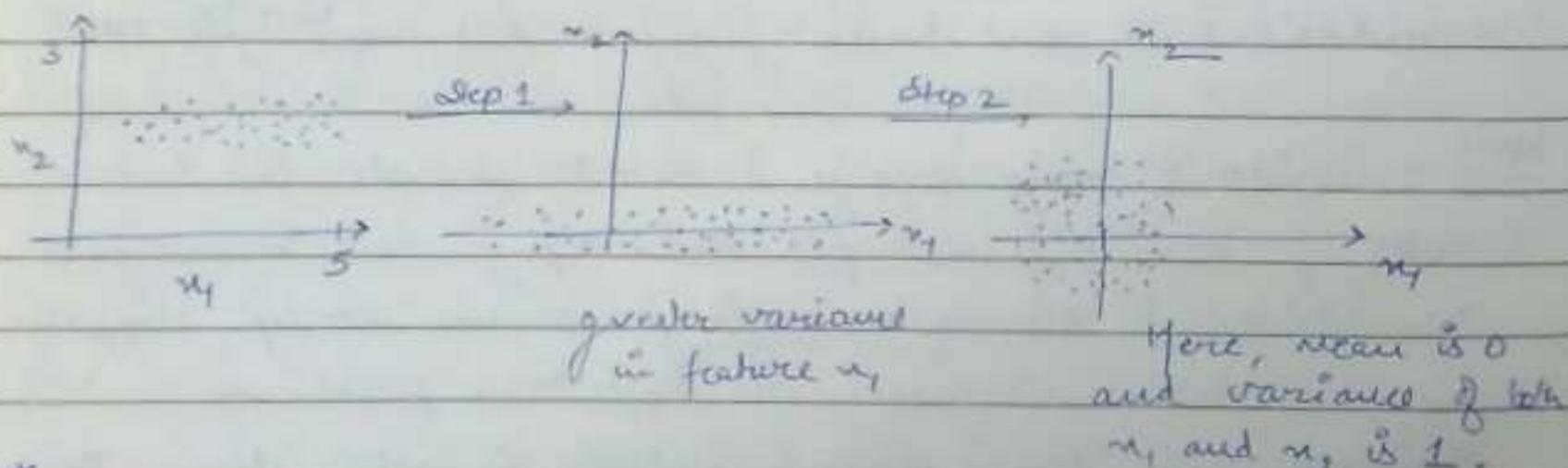
→ Setting up the optimization problem:

One of the techniques to speed up training is normalize your inputs.

Step 1: Subtract mean $\{u = \frac{1}{m} \sum_{i=1}^m x^{(i)}\} \rightarrow x = x - u$

Step 2: Normalize variance $\{v = \frac{1}{m} \sum_{i=1}^m x^{(i) \times 2}\} \rightarrow x' = v^{-2}$

\downarrow elementwise scaling
 vector with variances
 of all the features.



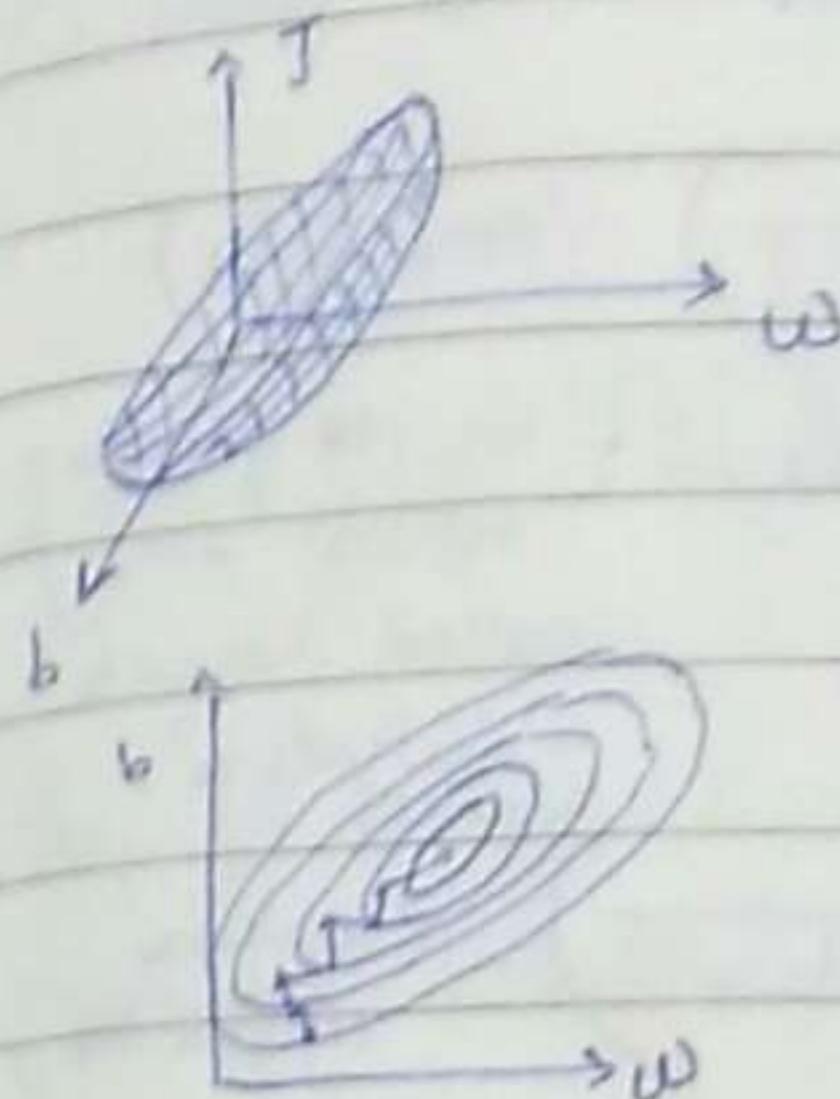
* Use same u and v for transforming the test set as well.

Q Why normalize the dataset?

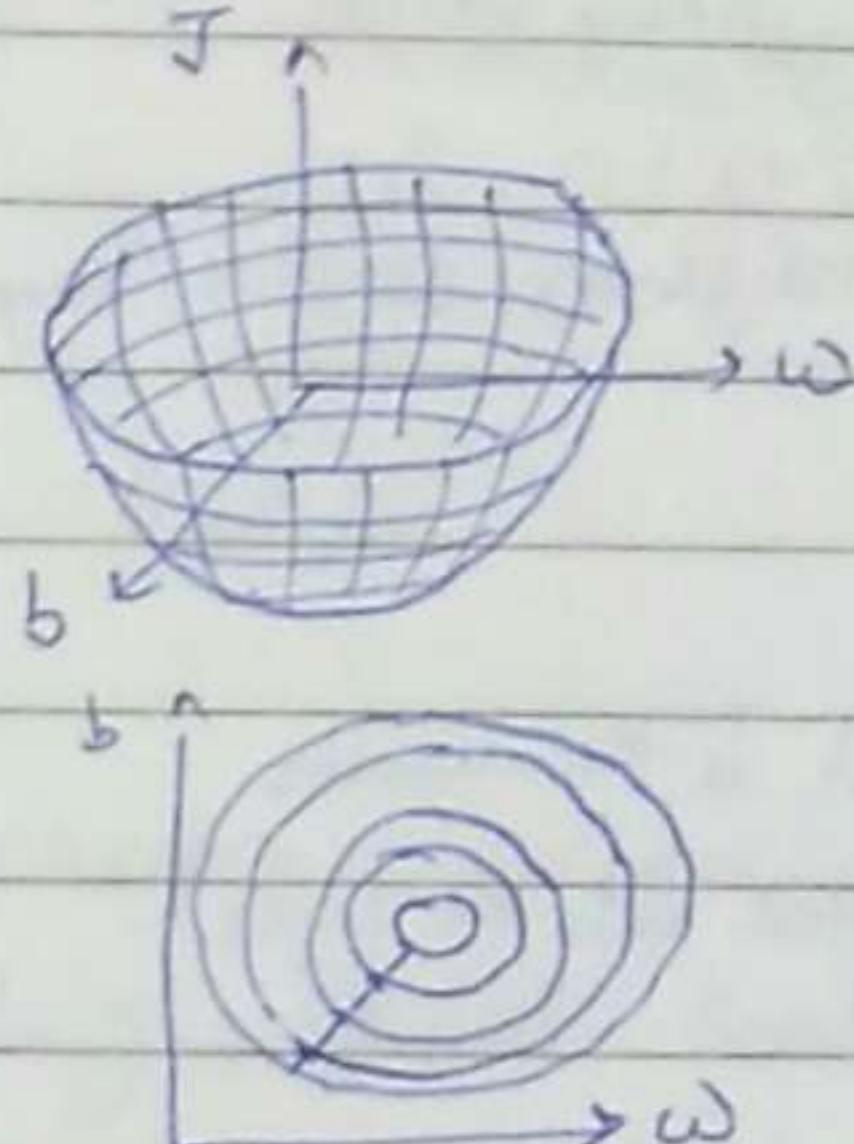
Abs If the range of the different features vary too much like $w_1: 1 \text{ to } 1000$ and $w_2: 0 \text{ to } 1$ then cost fn comes out to be quite elongated. So, to

have a more symmetric J , "normalize" is done. Note

Unnormalized:



Normalized: (makes the cost fw easier and faster to optimize.)



α : small

bcz gradient decent takes a lot of steps to oscillate back and forth here to find its way to the minimum.

α : large

bcz here gradient decent can take much larger steps to reach the minimum.

- Vanishing / Exploding Gradient Decent :

Problem with training a very deep network, is that the derivatives or slopes can sometimes get either very big or very small, thus making training difficult.

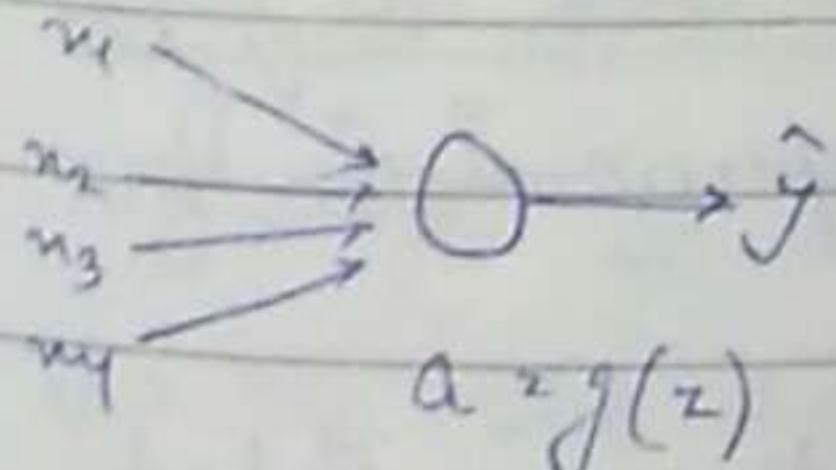
When $w^{[L]} > 1$, activations can explode. as $a^{[L]} = g(z^{[L]})$

when $w^{[L]} < 1$, activations can vanish as a f^o of L (no. of layers)

In a similar manner, the gradient and slopes vanish or explode, can be shown by using the same formulas for dw & all.

Note If gradients are smaller than 1, then gradient descent will take tiny little steps and it will take a long time for gradient descent to learn anything.

This problem can be solved by carefully initializing the weights during random weight initialization.



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

for z not to blow up & not become too small, larger the n is, smaller the w should be.

So, it can be done that $\text{var}(\omega) = \frac{1}{n}$ $\{n = \# \text{ input features}\}$
into each neuron

If we are using ReLU activation fn then $\text{var}(\omega) = \frac{3}{n}$ works better.

$$\omega^{(l)} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{2}{n^{(l-1)}}\right) \quad \{\text{for ReLU}\}$$

$$\rightarrow \text{np.sqrt}\left(\frac{1}{n^{(l-1)}}\right) \quad \{\text{in case of tanh, activation fn, called Xavier initialization.}\}$$

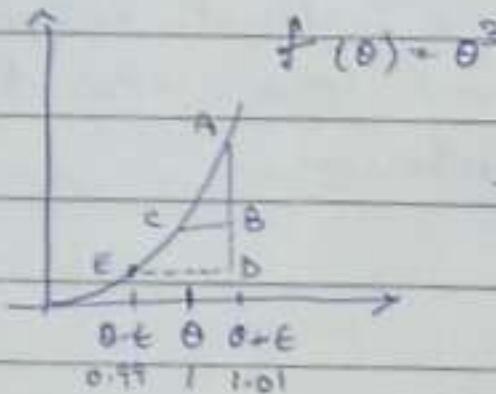
Note $n^{(l)}$ bcoz that is the number

of units being fed into the layer l .

$$= \text{np.sqrt}\left(\frac{2}{n^{(l-1)} + n^{(l)}}\right) \quad \{\text{another alternative.}\}$$

15/8/18

→ Numerical approximation of gradients:



$\frac{AB}{DE}$ over $\frac{AB}{BC}$ is used as it is a much better approximation to the derivative of θ .

Gradient checking is done to make sure that implementation of backprop is correct. So, for it we need to numerically approx the gradients first.

Now, $\frac{AD}{DE} = \frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \approx g(\theta)$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001$$

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.0001

when only one side was considered, error = 0.03

So, the two sided difference helped getting close to $g(\theta)$ and tells that $g(\theta)$ is a pretty good approximation of derivative of $f(\theta)$.

This runs twice as slow as compared to the one-sided

defuse.

Note $f'(0) = \lim_{\epsilon \rightarrow 0} \frac{f(0+\epsilon) - f(0-\epsilon)}{2\epsilon}$

2 order of error for this
is $O(\epsilon^2)$ and ϵ we
know is very very small.

but $\lim_{\epsilon \rightarrow 0} \frac{f(0+\epsilon) - f(0)}{\epsilon}$, error = $O(\epsilon)$.

So, by taking a two sided difference, we can numerically
verify whether or not a fn $g, g(0)$ that someone else
gives is a correct implementaⁿ of $f'(0)$. This can
be used to verify the correctness of backpropagation
implementation.

→ Gradient checking: to verify " " " "

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

Concatenate.

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

Now, reshape $d\theta^{[1]}, db^{[1]}, \dots d\theta^{[L]}, db^{[L]}$ into a big vector $d\theta$ of
same dimension as θ .

Q. Is $d\theta$ the slope or gradient of the cost fn $J(\theta)$?

Ans for each i :

$$d\theta_{\text{approx}}^{(i)} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\approx d\theta^{(i)} = \frac{\partial J}{\partial \theta_i} \text{ i.e. (partial derivative of } J \text{ w.r.t } \theta_i)$$

So, we need to check if $d\theta_{\text{approx}} \approx d\theta$ to see if $d\theta$ is
the slope. To see this we can compute Euclidean distance
b/w the two vectors.

$$\text{check } \frac{\|\theta_{\text{approx}} - \theta\|_2}{\|\theta_{\text{approx}}\|_2 + \|\theta\|_2} \quad \textcircled{1}$$

← (for normalize?)

→ Euclidean lengths

we take $\epsilon = 10^{-7}$

if $\textcircled{1} \approx 10^{-7}$, that's great as error is quite low.

If it comes out to be more like 10^{-5} or 10^{-3} , then there must be a bug somewhere and individual components of θ need to be checked if there is a specific value of i for which $d\theta_i$ is very different from $d\theta_{\text{approx}}^{[i]}$ and then that computation's error need to be tracked. If after debugging, error reduces to 10^{-7} , then probably you have a correct implementation. So, it must pass with a small grad check value.

Sup notes on gradient checking:

① Use grad check only for debugging, not training as the computation of gradcheck is very slow and so only during debugging its checked if $d\theta^{[i]}$ is close to $d\theta_{\text{approx}}^{[i]}$ and after that it is turned off.

② If algorithm fails grad check, look at individual components to try to identify bug.

③ Remember regularization term while doing grad check.

$$\text{i.e. if } J(\theta) = \frac{1}{m} \sum_i L(\hat{y}^{(i)}, y^{(i)}) + \frac{1}{2m} \sum_i \|w^{(i)}\|_F^2$$

then $d\theta = \text{gradient of } J \text{ w.r.t } \theta \text{ including the regularization term.}$

④ Gradcheck doesn't work with dropout becoz in dropout there isn't an easy to compute $J(\theta)$ that dropout is doing gradient descent on because of its property of randomly dropping parameters in every iteration. So implement grad check without dropout i.e. keep prob = 1.0.

Page

and then turn on dropout hoping its "implementation" to be correct
grad check can be used with dropout if the pattern of
dropout is fixed, but usually its recommended to turn
off dropout while grad check.

- ⑤ Run at random initialization; perhaps again after some
training as then w and b have some time to wander
away from 0 which tends to have correct implementation?
of backprop, larger values of w and b give more
inaccurate implementation.

Note You do not apply dropout and do not keep the \checkmark ~~keep prob~~
factor in the calculations used in training.

① Week 2:

18/8/18

→ Mini-batch gradient descent: splitting the training set into little
baby training sets give mini-batches
this helps get a faster algo as gradient descent starts to
make some progress even before it finishes processing the
entire training set, but just 1 mini batch.

$$\text{So, } X = [x^{(1)} \ x^{(2)} \ x^{(3)} \dots x^{(1000)} \ | \ x^{(1001)} \dots x^{(2000)} | \dots | \ x^{(n)}]_{(n \times m)}$$
$$x^{(1)} \quad x^{(2)} \quad x^{(3)} \dots \rightarrow y_{1000}$$

So, now we are having 5000 minibatches if $m = 5$ million.
The same thing is done for y as well.

So, Mini batch t : $x^{(t)}, y^{(t)}$

1000 training examples with the corresponding
1/p of p features/ pairs

Note $Z^{(t)}$; t^{th} mini batch , $Z^{(i)}$: i^{th} training e.g., $Z^{(1)}$: in 1st layer

code for mini batch gradient descent:

repeat L

for $t = 1, \dots, 5000 \{$ { considering 5K mini batches }

 fwd prop on X^{t+3}

$$Z^{[1]} = w^{[1]} X^{t+3} + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

:

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

 // as we are processing only the mini

 // batch t at a time.

 // vectorized implementation all over

 // processing 1K examples at a time

 // rather than 5 million examples.

$$\text{Compute cost } J = \frac{1}{1000} \sum_{i=1}^L L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2(1000)} \|w^{[1]}\|_F^2$$

examples from
the minibatch t.

↓
regularize term.

Backprop to compute gradients wrt J^{t+3} using (X^{t+3}, y^{t+3})

$$w^{[L]} = w^{[1]} - \alpha dw^{[1]}, b^{[L]} = b^{[1]} - \alpha db^{[1]}$$

} }

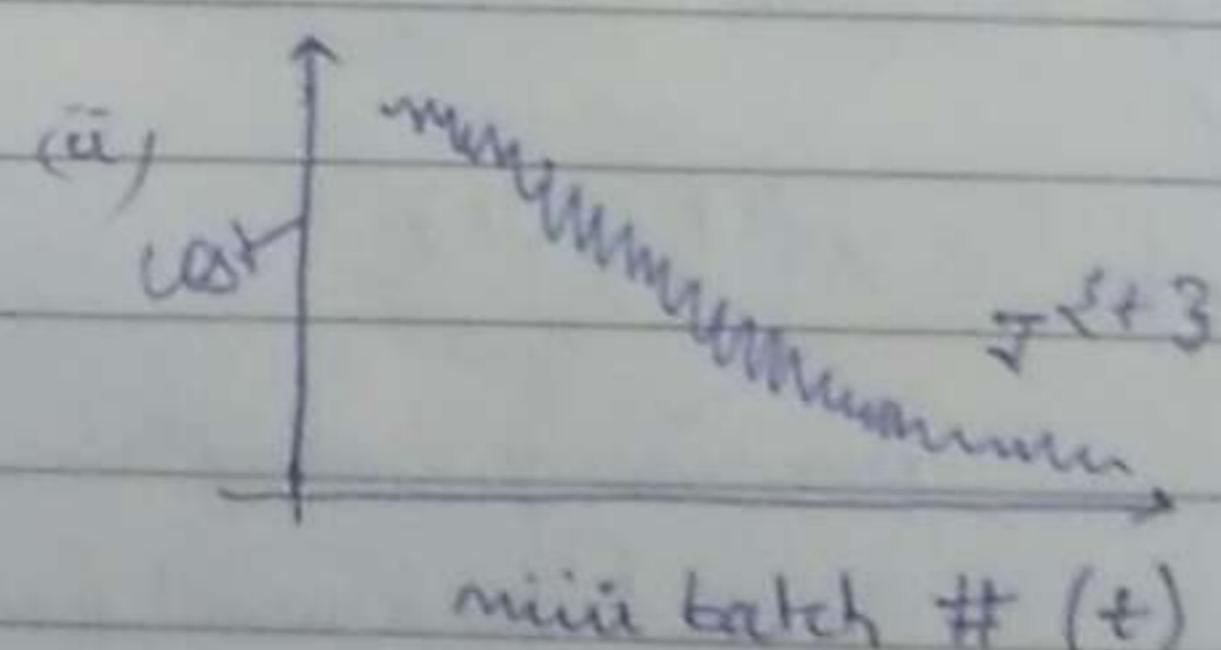
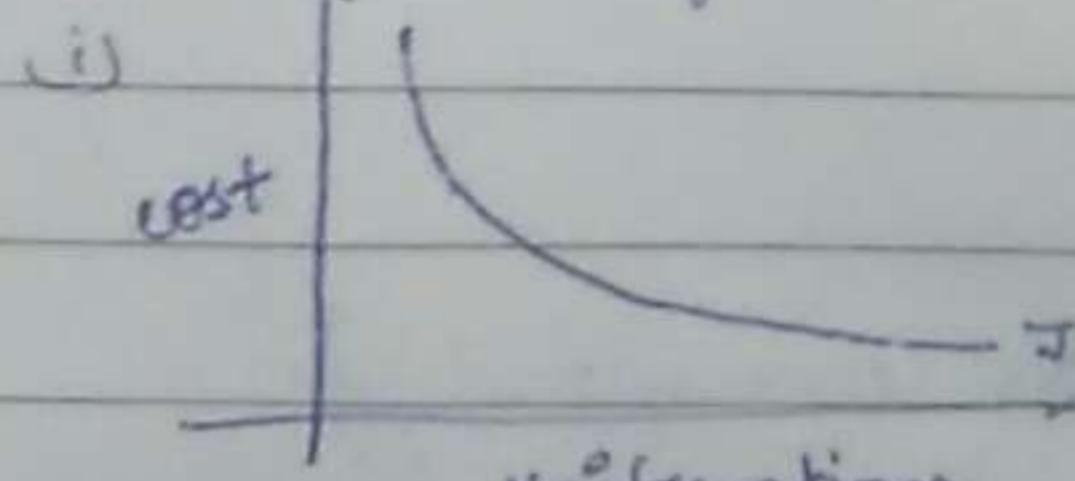
// this code is also called doing 1 epoch of training

// as it is = 1 pass through the training set using minibatch

// gradient descent.

Note Batch gradient descent, a single pass through the training set allows us to take only one gradient descent step while with mini batch gradient descent, single pass through training set, allows to take 5K (as in above e.g.) gradient descent steps. With large training set, mini-batch runs much faster than the batch gradient descent.

Batch gradient descent



{ goes through complete dataset on every iteration }

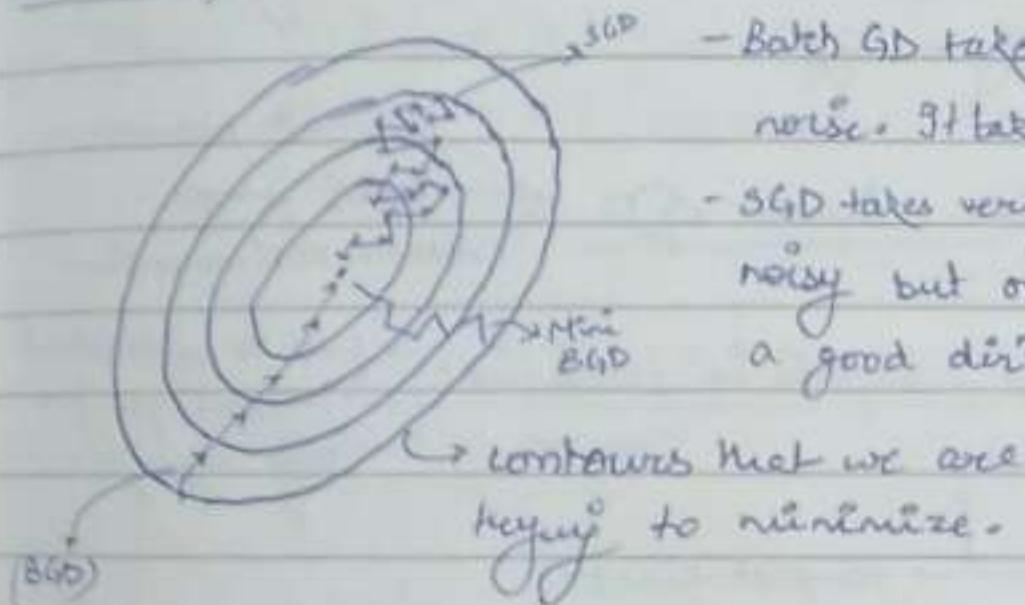
Plot J^{t+3} computed using X^{t+3}, y^{t+3} ,
{ it will go down but is noisy }

so (ii) is noisy bcoz may be the batch $x^{(1)}, y^{(1)}$ is just the scores of any mini batch, so cost might be little lower but maybe just by chance, $x^{(2)}, y^{(2)}$ is just a harder mini batch, as a result of which the cost will be a little higher.

Parameter: size of mini batch

- may deviate from m (= batch gradient descent) to 1 (stochastic gradient descent), in batch, there is just 1 batch while in stochastic, every example is a batch.

Now we will see, what these two extremes do on optimizing the cost fn.



- Batch GD takes greater steps, has low noise. It takes too long per step~.
- SGD takes very small steps and is highly noisy but on an average will take us on a good dir. SGD never converges to minimum, it keeps on oscillating nearby but never stays there.

Note: Noise in SGD can be overcome by using a small learning rate but the huge disadvantage is that we lose all the speedup that can be gained through vectorizeⁿ as we are processing a single example at a time. Its good that we are learning after every example but the time becomes a major issue.

BGD works fine with smaller datasets, both learning wise and time wise.

So, midway is best i.e. t somewhere b/w 1 and m, bcoz we to get a lot of vectorizeⁿ which would be much faster than processing the examples one at a time. Also, progress is

made without needing to wait till the processing of entire training set. In our previous eg. we had 5000 gradient steps per iteration or epoch.

- Guidelines to choose the mini batch size :

- i) in case of small training set, use SGD. (e.g. when m < 2k)
- ii) " " " big " " , typical mini batch sizes will be anything from 64 upto may be 512 (typical values) becz the way computer is laid down, somehow computers are faster when size is in power of 2. So, in previous example try using 1024 instead of 1000. this mini-batch size that is decided should fit the CPU, GPU memory otherwise performance will fall off a cliff.

* Optimized algorithm faster than GD

GD with momentum

RMS prop

Adam optimization

* Exponentially weighted moving averages in statistics: (to be understood before going to any algo)

$$v_t = \beta v_{t-1} + (1-\beta) \theta_t \quad \text{, temp at time t}$$

$\beta_1 = 0.9$: \approx averaging 10 days temp.

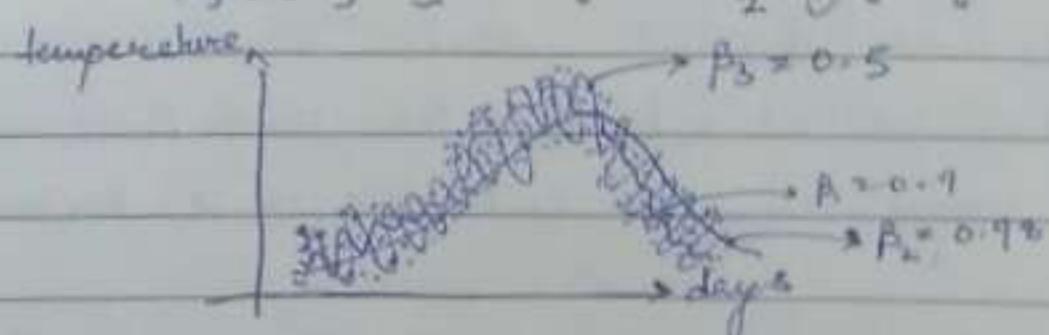
v_t approximately averages over

$\beta_2 = 0.98$: \approx 50 days temp

$\Rightarrow \frac{1}{\beta}$ days

$\beta_3 = 0.5 \approx 2$ " "

temperature.



Plot for β_2 is much smoother becz now we are averaging over more days of temperature but the curve has shifted further to right becz you're now averaging over a much larger window of temperatures and so it adapts more slowly with the temp changes.

exponentially weighted moving avg. formula

It is bias when β is large, much weight is given on the previous temperature values and only $(1-\beta)$ weight is given on the current temp. value. $\sum \beta^t v_{t-1} + (1-\beta) \theta_t$
 when $\beta = 0.5$, graph is much more noisy as only 2 days temp is being averaged over and so graph is more susceptible to outliers and adapts much more quickly as the temp. changes.

$$v_{100} = 0.9 v_{99} + 0.1 \theta_{100}$$

$$v_{99} = 0.9 v_{98} + 0.1 \theta_{99}$$

$$v_{98} = 0.9 v_{97} + 0.1 \theta_{98}$$

:

$$v_{100} = 0.1 \theta_{100} + 0.9 v_{99}$$

$$= 0.1 \theta_{100} + 0.9 (0.9 v_{98} + 0.1 \theta_{99})$$

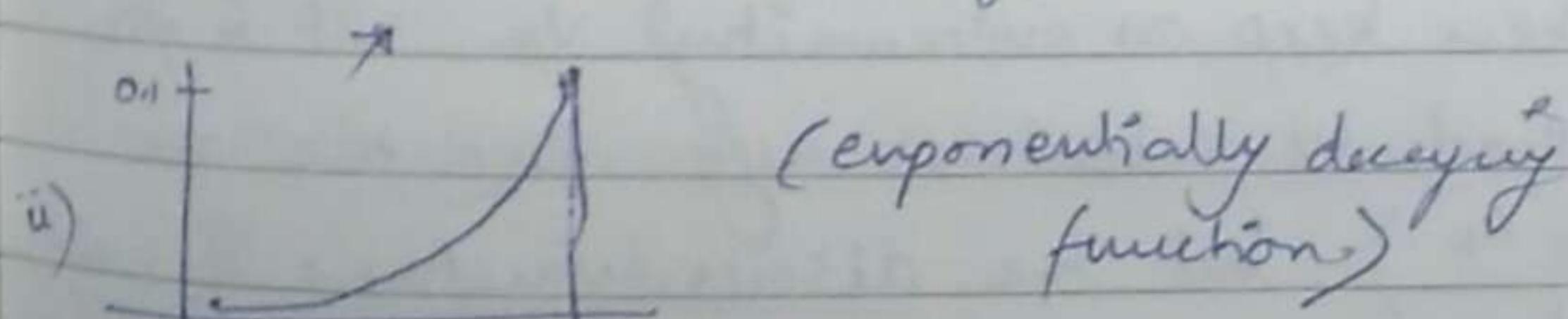
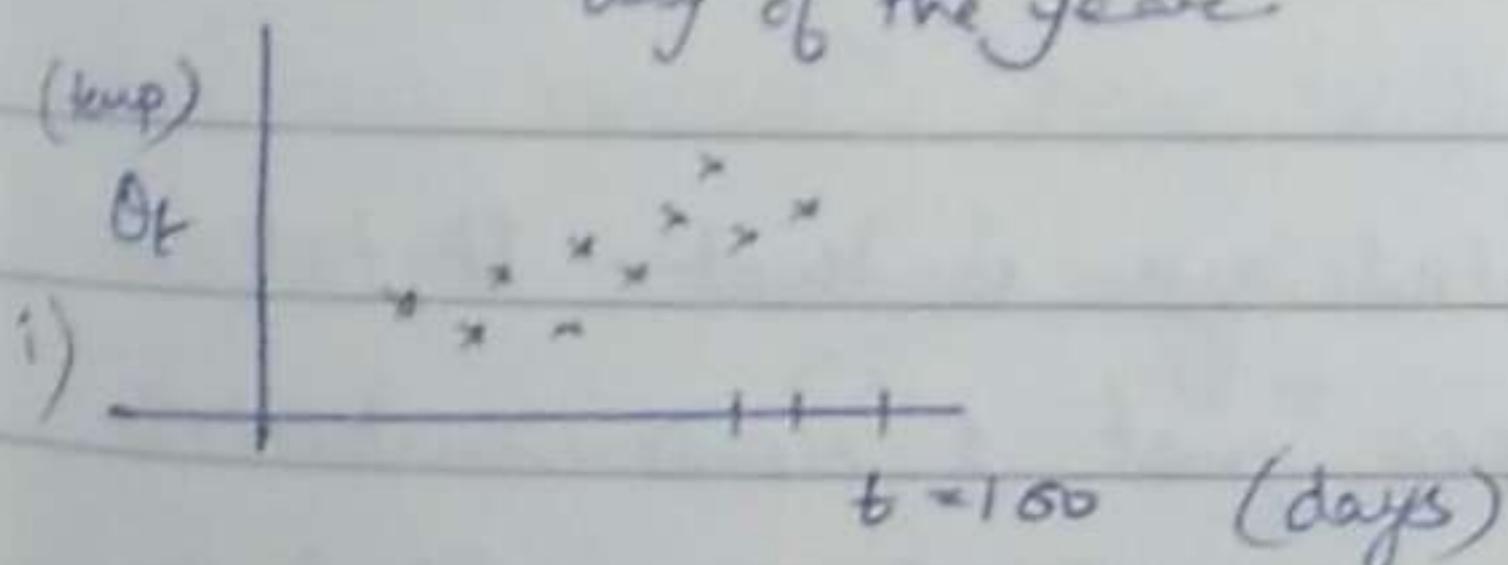
and so on we get,

$$v_{100} = 0.1 \theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1 (0.9)^2 \theta_{98} + \dots$$

calculated on the 100th day of the year

weighted average or sum of θ_{100}

current day's temp.



(exponentially decaying function)

v_{100} = element wise product b/w functions i) and ii). and sum it up.

i.e daily temp (θ_t) multiplied with the exponentially decaying fn (going from 0.1 to $0.1(0.9)$ to $0.1(0.9)^2$ to approx 0)

Note All the coeff. in v_{100} sum up to approximately 1.

Now, we need to know, how many days are being averaged over.

$$\text{well, } (0.9)^{10} \approx 0.35 \approx \frac{1}{e}$$

$$(1-\varepsilon) \gamma_{\varepsilon} = \frac{1}{e} \approx 0.35 \quad (\varepsilon \text{ in this case } = 0.1)$$

This means it takes about 10 days for the weight of exponential fn to decay to around $\frac{1}{3}$ of the peak. Hence, when $\beta = 0.9$ we say that it is as if we are computing an exponentially weighted average that focuses on just the last 10 days temperature, bcoz its after 10 days that the weight decays to less than about a third of the weight of the current day. While when $\beta = 0.98$, $(0.98)^{50} \approx 0.35$, so we say 50 day avg. in case of $\beta = 0.98$. From here, we say that we are averaging over $\frac{1}{1-\beta}$ days $\frac{1}{1-0.98} = 50$.

- Implementing exponentially weighted averages:

$$\begin{aligned} v_0 &= 0 \\ v_1 &= \beta v_0 + (1-\beta) \theta_1 \\ v_2 &= \beta v_1 + (1-\beta) \theta_2 \\ v_3 &= \beta v_2 + (1-\beta) \theta_3 \\ &\dots \end{aligned}$$

$$\begin{aligned} v_0 &= 0 && \left\{ \text{on computer} \right\} \\ \text{Repeat } & \text{ } \\ & \text{Get next } \theta_t \\ v_0 &= \beta v_0 + (1-\beta) \theta_t && \text{|| (update)} \end{aligned}$$

Advantage of this exponentially weighted avg. is that it takes very little memory. We need to keep just 1 word number in computer memory, and then keep on overwriting v_0 . It is for sure not the best method to calculate avg. like the one conventionally used ($\frac{\theta_1 + \dots + \theta_t}{t}$) but the disadvantage of using the conventional is that it requires more memory as we need to keep and store all today's or 50 days temperature and is computationally more expensive.

15/8/18

- Bias correction in exponentially weighted averages:

In the initial days, estimate of the temp. is not very good as

$$v_0 = 0 \\ v_1 = 0.98 v_0 + 0.02 \theta_1 = 0.02 \theta_1$$

$$v_2 = 0.98 \theta_1 + 0.02 \theta_2$$

$$v_2 = 0.0196 \theta_1 + 0.02 \theta_2$$

|| v_2 comes out to be much less than θ_1 or θ_2 , hence not a very good estimate for the first two days of the year.

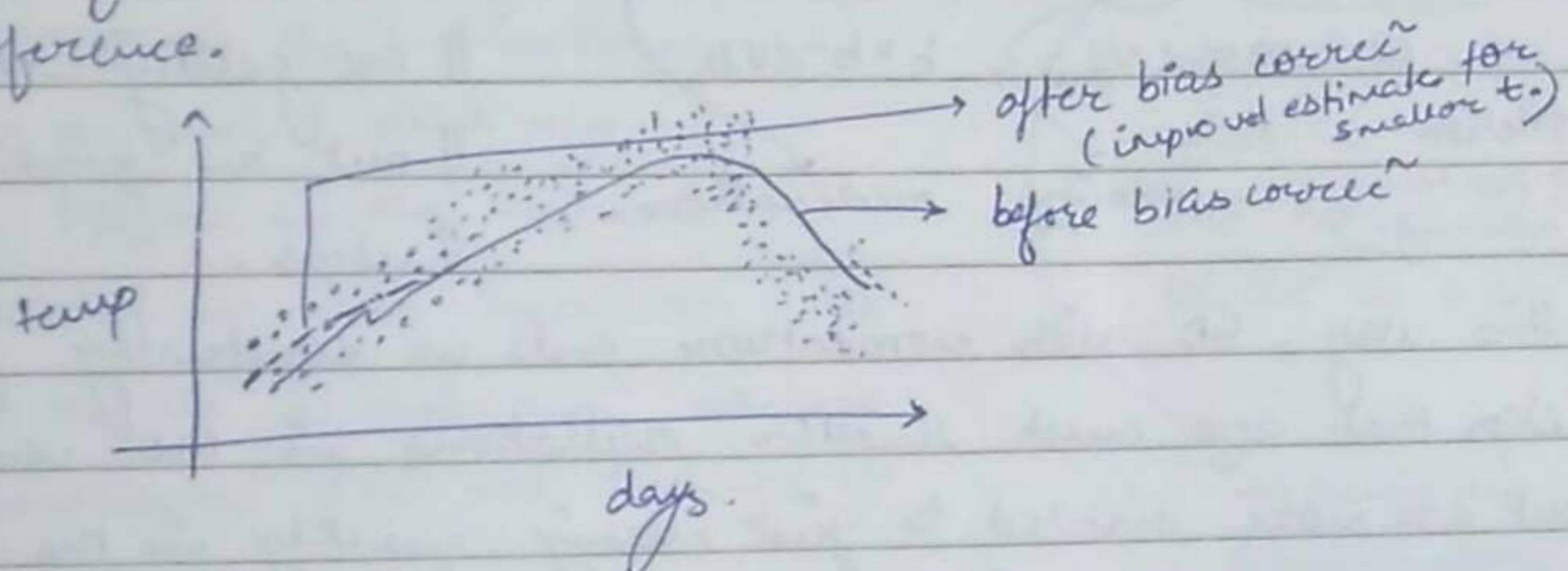
So, an ~~error~~ is made for this bias to remove.

$$\frac{v_t}{1-\beta^t}$$

$$t=2: 1-\beta^t = 1-(0.98)^2 = 0.0396$$

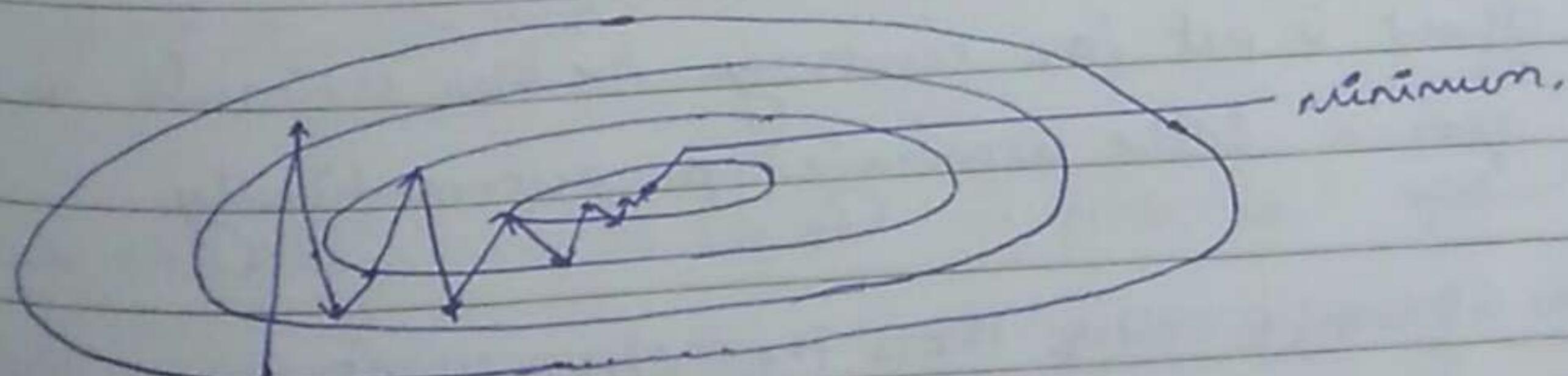
$$\frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396} \quad || \text{bias is removed.}$$

when t is large, $1-\beta^t \approx 1$ and hence the bias correcⁿ makes no big difference.



this bias correcⁿ can help us get a better estimate early on when the moving weighted avg. is warming up and is in its initial phases.

① Gradient descent with momentum: The idea is to calculate the exponentially weighted average of your gradients, and then use that gradient to update the weights instead.



A very large learning rate might end up overshooting or

diverging and so we need to prevent the oscillations from getting too big forces to use a learning rate that's itself not too large.

We want the learning to be slower on vertical axis to avoid those oscillations but on the horizontal axis faster learning is desired.

✓ Momentum:

(mini batch number)

On iteration t :

Compute d_w , db on current minibatch

$$v_{dw} = \beta v_{dw} + (1-\beta) d_w$$

$$v_{db} = \beta v_{db} + (1-\beta) d_b$$

$$w = w - \alpha v_{dw}, \quad b = b - \alpha v_{db}$$

friction!

(to stop from speeding.)

velocity

acceleration

initialize $v_{dw} = 0$

$v_{db} = 0$

{same dimensions as
 w & b respectively}

|| computing moving average of
the derivatives for w you
are getting thus smoothing
out the gradient descent
steps.

This way, GD with momentum ends up eventually just taking steps that are much smaller oscillations in the vertical direction but are more directed to just moving quickly in the horizontal direction, and so this allows our algorithm to take a more straightforward path.

Moving (vertical avg = 0
horiz avg = towards minimum)

2 hyperparameters = α (learning rate), β (to control the exponentially weighted avg.)

$\beta = 0.9$ has so far been the best value, but keep on trying others as well if in case some other values proves to be more robust.

Not much effort is put for removing the bias estimate as the GD learns after a little warming up automatically.

Momentum GD works better than the straight GD.

② RMS prop (Root mean square prop)

We saw that using GD might end up in huge oscillation in the vertical direction even while it's trying to make progress in horizontal direction.

On iteration t:

compute d_w , d_b on current minibatch:

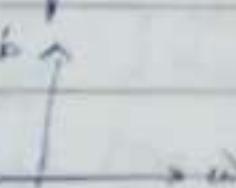
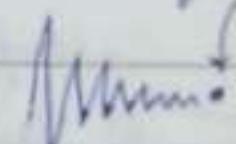
$$S_{dw} = \beta S_{dw} + (1-\beta) d_w^2$$

↳ (element wise square)

$$S_{db} = \beta S_{db} + (1-\beta) d_b^2$$

$$\boxed{w = w - \alpha \frac{d_w}{\sqrt{S_{dw} + \epsilon}}, \quad b = b - \alpha \frac{d_b}{\sqrt{S_{db} + \epsilon}}}$$

optimal
goal

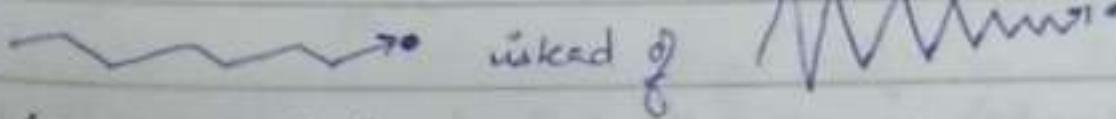


Slope is very large in the b direction, relatively large db and a relatively small d_w . Updates in the vertical dirⁿ are divided by much larger no. $\sqrt{S_{db}}$ to damp the oscillation in that dirⁿ whereas updates in horizontal dirⁿ are divided by a smaller number.

We want learning to be fast (in w dirⁿ) and damp all the oscillation in b dirⁿ.

→ to avoid the denominator being zero.

So, net impact of using RMS prop is that updates will end up as:



(less oscillation in vertical dirⁿ).

and so a larger α can be used without diverging in vertical dirⁿ, and helps in learning fast.

Here we are taking this as 2D, but in practical applicⁿ, the dirⁿ in which we want to damp the oscillation may be a set of many w 's like w_1, w_2, w_3 etc, and the horizontal

$d\vec{w}$ may be w_s, w_b, w , etc and so on bcz in practice it is multidimensional. $d\vec{w}$ and $d\vec{b}$ are both high dimensional.

③ Adam Optimizer Algorithm : (=Momentum + RMS prop)

$$v_{d\vec{w}} = 0, s_{d\vec{w}} = 0, v_{d\vec{b}} = s_{d\vec{b}} = 0 \quad // \text{Initialize}$$

On iterⁿ t :

compute $d\vec{w}, d\vec{b}$ using current mini-batch

$$\left. \begin{array}{l} \text{momentum} \\ \text{like update} \end{array} \right\} \quad \begin{array}{l} v_{d\vec{w}} = \beta_1 v_{d\vec{w}} + (1 - \beta_1) d\vec{w} \\ v_{d\vec{b}} = \beta_1 v_{d\vec{b}} + (1 - \beta_1) d\vec{b} \end{array} \quad \begin{array}{l} // \text{using } \beta_1 \text{ instead of } \beta \text{ to} \\ // \text{differentiate from earlier.} \end{array}$$

$$\left. \begin{array}{l} \text{rmsprop} \\ \text{like update} \end{array} \right\} \quad \begin{array}{l} s_{d\vec{w}} = \beta_2 s_{d\vec{w}} + (1 - \beta_2) d\vec{w}^2 \\ s_{d\vec{b}} = \beta_2 s_{d\vec{b}} + (1 - \beta_2) d\vec{b}^2 \end{array} \quad \begin{array}{l} // \text{usage.} \end{array}$$

$$v_{d\vec{w}}^{\text{corrected}} = \frac{v_{d\vec{w}}}{(1 - \beta_1^t)}, v_{d\vec{b}}^{\text{corrected}} = \frac{v_{d\vec{b}}}{(1 - \beta_1^t)}$$

(after bias correc)

$$s_{d\vec{w}}^{\text{corrected}} = \frac{s_{d\vec{w}}}{(1 - \beta_2^t)}, s_{d\vec{b}}^{\text{corrected}} = \frac{s_{d\vec{b}}}{(1 - \beta_2^t)}$$

$$w = w - \alpha \frac{v_{d\vec{w}}^{\text{corrected}}}{\sqrt{s_{d\vec{w}}^{\text{corrected}}} + \epsilon}, b = b - \alpha \frac{v_{d\vec{b}}^{\text{corrected}}}{\sqrt{s_{d\vec{b}}^{\text{corrected}}} + \epsilon}$$

Adam optimizaⁿ (a combo) is very effective for many different NN of a very wide variety of architecture.

- Hyperparameter choice : α : needs to be tuned

$$\beta_1 : 0.9 \quad (\text{d}\vec{w})$$

$$\beta_2 : 0.99 \quad (\text{d}\vec{w}^2)$$

$$\epsilon = 10^{-8} \quad (\text{doesn't affect much})$$

β_1 and β_2 can also be tuned but are generally not and are used with these robust values only like 0.9 and 0.99 for β_1 and β_2 respectively.

Adam optimizaⁿ can be used in batch and mini-batch GD both.

Adam : Adaptive moment estimation. beoz ↴

β_1 is computing the mean of the derivatives. This is called the first moment. β_2 is used to compute the exponentially weighted avg. of the squares and that's called second moment.

• learning rate decay : To slowly reduce the learning rate over time can help speed up the learning algorithm.

As we iterate, the steps are a little bit noisy and tend towards the minimum but won't exactly converge. It will go on being around the minimum point but would not stop at it bcoz of using some fixed value of α and there is some noise in the mini batches. So, reducing α at later stages like when it reaches close to minimum helps oscillations in a tighter region around this minimum rather than wandering far away, so

α needs to be reduced bcoz in initial stages of learning, one might afford bigger steps but as learning approaches convergence, having a slower learning rate allows to take smaller steps.

1 epoch = 1 pass through the data

$$\alpha = \frac{1}{1 + (\text{decay rate} \times \text{epoch number})} \rightarrow \alpha_0$$

formula for learning rate decay

Epoch	α
1	0.1
2	0.07
3	0.05
4	0.04

$$\text{if } \alpha_0 = 0.2$$

$$\text{decay rate} = 1$$

any constant

less than 1.

$$\alpha = 0.95 \text{ epoch number} \cdot \alpha_0$$

(exponential decay)
of learning rate.

or

$$\alpha = \frac{k}{\sqrt{\text{epoch number}}} \cdot \alpha_0$$

$$\text{or } \alpha = \frac{k}{\sqrt{t}} \cdot \alpha_0$$

epoch number

or

$$\alpha = \begin{cases} \dots & \\ \dots & \\ \dots & \\ t & \end{cases}$$

learning rate decaying in
discrete steps.

(discrete staircase)

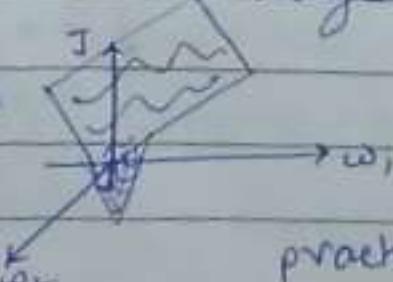
or

Manual decay : If we are training just one model at a time and if your model takes many hours or even days to train, then just watch the model getting trained over a large no. of days and then manually say, it looks like α has slowed down, I am going to $\downarrow \alpha$ a little bit. So, its tuning α by hand hour by hour or day by day. As long as we are training a small no. of models, its feasible.

Note It is sometimes preferred to have just 1 tuned α for learning. Though LRD (decay) helps but is less often used.

* Problem of local optima: There is a problem with optimization algorithms getting stuck in a local optima but with advancement in understanding of deep learning, concept of local minima has also changed.

Earlier we visualized



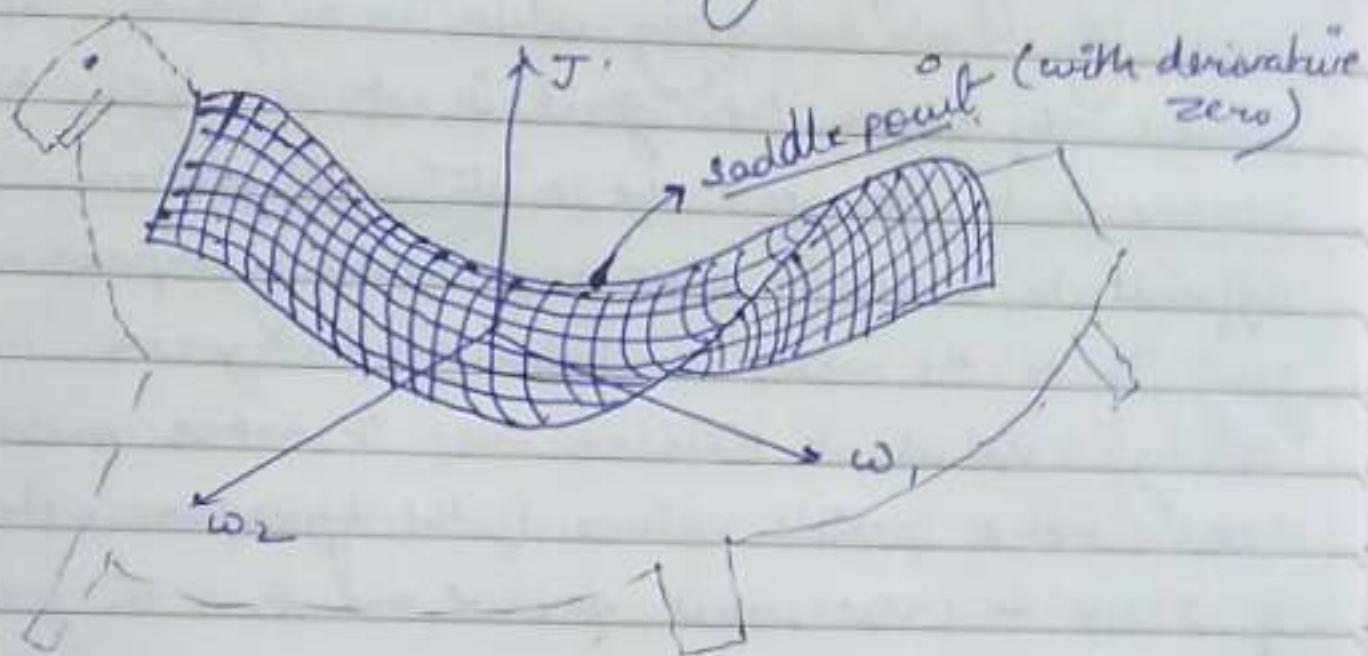
i.e a 2-D plot to see local

optima but in

practice, such low dimensional

problems don't exist and in high dimensional space, it's much more likely to run into a saddle point rather than a global or local optima (where gradient is zero in all the dimensions) which is very hard to occur. Gradient might be zero for some dimensions at a point but not for all and such a point is called saddle point, where the curve function is bending.

called a saddle point bcoz it can be part on the back of the horse.

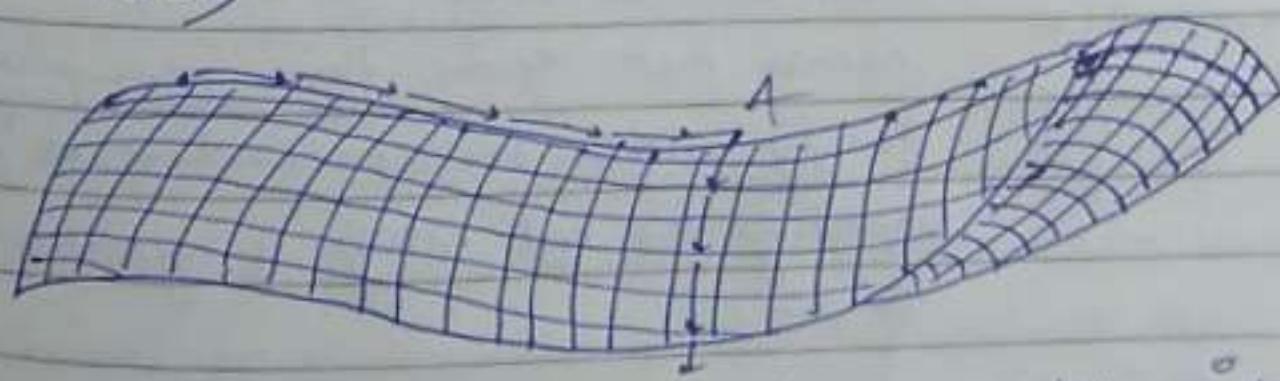


So, our assumptions on low Dimensional space do not always generalize well in high Dimensional space in deep learning.

Q: If local optima is not an issue, then where does the problem lie?

Ans It turns out that plateaus can really slow down learning.

(a region where the derivative is close to zero for a long time.)



So, the gradient might move it long way till point A and then fall off the plateau. So, at these points the optimiser algorithm actually help speed up the process and are of help.

① Week 3 : HYPERPARAMETER tuning and much more

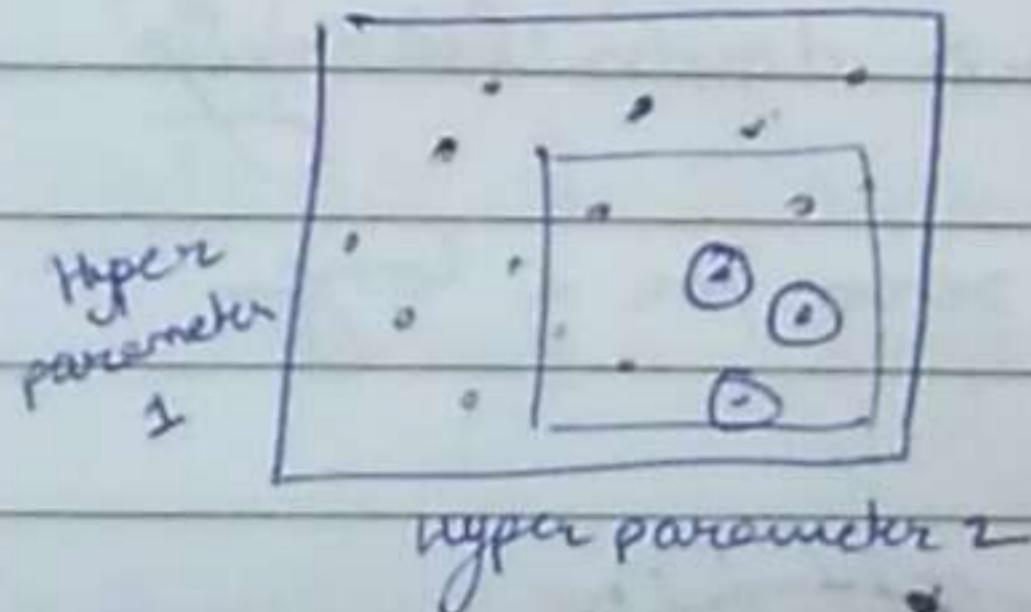
$\alpha > \beta$, #hidden units, minibatch size $> \beta_1, \beta_2, \epsilon$

in top order of the hyperparameter importance
(not a hard and fast order though)

Traditional way of finding the right value of hyperparameters was to try each and every value in a grid but that is feasible only when the grid is not too big but now it's recommended to choose the points at random, becoz it is difficult to know in advance which hyperparameters are going to be the most important for your problem. So, selecting values at random make it more probable to choose more possible values of the hyperparameters - e.g.

if I have to chose value for α and ϵ , in grid method if grid is 5×5 , I will go across 25 values with only 5 different values of α and wasting too much energy on choosing ϵ as it is not affecting the results much, while in random select this won't happen.

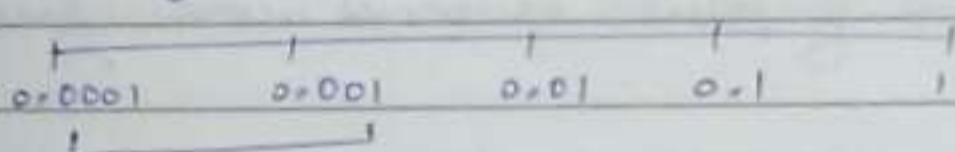
Another common practise is to use a coarse to fine sampling scheme.



If I find that points in a particular region are doing well for the model, then we will zoom in and select the other points from nearby that region thus sampling more density within the space where we are suspecting the best settings of hyperparameters to lie.

Note Sampling at random doesn't mean sampling uniformly at random over the range of valid values instead it's important to pick the appropriate scale on which to explore the hyperparameters

eg suppose we need to choose the no. of hidden units in layer 2.
And we find $n^{(2)} = 50$ to 100 to be the possible range,
then random selec' in b/w 50 to 100 can be reasonable.
If I was to chose no. of layers which I think to be
okay b/w 2 to 4 then uniform sampling b/w 2 to 4
would have been reasonable. In these cases, sampling
at random works BUT
if I am trying to tune α and suspect 0.0001 to 1 to be
the possible range, at that time uniform random
selec' doesn't work, instead it seems more reasonable to
search for hyperparameters on a log scale. Bcoz on a
linear scale b/w 0.0001 to 0.1, 90% of the uniformly
randomly chosen values will lie b/w 0.1 to 1. Thus
diminishing the scope for values b/w 0.0001 to 0.1.
But on a log scale this problem is solved as:



More resources are now available for searching b/w
0-0001 and 0-1

Implementa^o: $\mathcal{E} = -4 \times np.random.rand()$ $x \in [-4, 0]$

30 | 6 | 18

$$\alpha = 10^{\frac{1}{\mu}}$$

[a,b]

$$x \in [-4, 0]$$

Now r can be set uniformly at random b/w a and b .

Hence, appropriate scale (log in above case) is necessary for hyperparameter tuning.

Another difficult thing is selecting hyperparameters for exponentially weighted averages! So, let's say β is suspected to be b/w 0.9 to 0.999 where 0.9 means averaging over the last 10 days while 0.999 is like averaging over the last 1000 days/values.

So, now instead of sampling β , we will sample $1-\beta$ from 0.001 to 0.1, which has same intervals as range of β .

Q: Why is it a bad idea to sample on linear scale?

Ans

When β is close to 1, the sensitivity of the results we get changes even with very small changes to β . If β goes from 0.9 to 0.9005, it's no big deal, this is hardly any change in the results but if β goes from 0.999 to 0.9995, this will have a huge impact on exactly what the algo is doing. In both of the cases, it's averaging roughly over 10 values but in 1st it has gone from an exponentially weighted average over about 1000 examples to now the last 2000 example.

Also, we have $\frac{1}{1-\beta}$ which is very sensitive to small changes in β when β is close to 1.

So, the sampling process helps to sample more densely in the region when β is close to 1 or when $1-\beta$ is close to zero, thus helping in exploring the spaces of possible outcomes more efficiently.

Q How to organize hyperparameter search process?

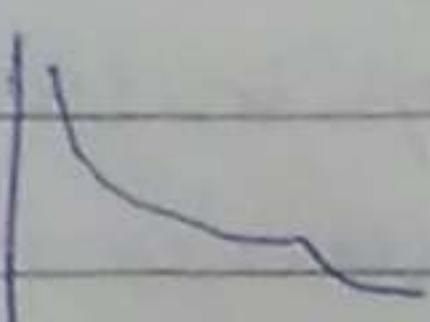
Ans Retest and reevaluate your hyperparameters at least once every several months.

two ways to organize
hyperparameter search
process

Babysit one model (when there is lot of data
but less of computational
resources)

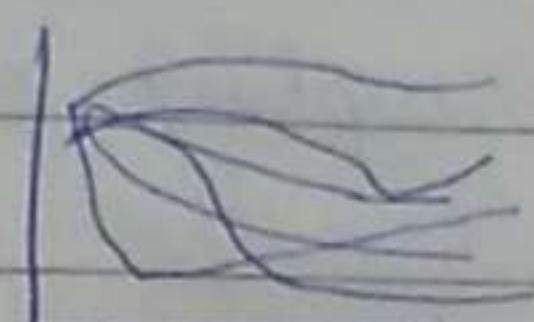
training many
models in parallel.

In babysitting approach, the hyperparameters are tuned every day even when the model is getting trained, while in the other approach, the models are left to run on their own with some setting of the hyperparameters for days or weeks. Many models are run in II and the best one is chosen at the end.



Babysitting
method

(Parde approach)



parallelly
models running
approach.

(Cartier approach)

Panda approach coz they also give one or 2 baby at a time and never pay too much attention to ensure their survival.

Lizard approach coz reptiles give millions of eggs and don't pay much attention afterwards hoping that atleast one of them in a bunch will survive.

So, its the computational resources that decide the kind of approach to be followed.

* Batch Normalization: (makes hyperparameter search problem much easier) (also makes our NN more robust.)

(also helps train very deep NN)



(as seen before)
in logistic regression

Q: for any hidden layer can we normalize the values of $a^{(l)}$ (activations) so as to train $w^{(l)}, b^{(l)}$ faster? It is bcoz $a^{(l)}$ is the input to the next layer and will affect w and b .

Ans: technically we normalize $z^{(l)}$ and not $a^{(l)}$ by default.

Implementation of Batch Norm:

Given some intermediate values in the NN like $z^{(0)(i)}$, so given these values, we compute the mean as follows:

$$\mu = \frac{1}{m} \sum_i z^{(0)(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

$$z_{\text{norm}}^{(0)(i)} = \frac{z^{(0)(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

Normalizing the I/P features in can help learning in a NN. Sates prop extends this to even hidden layers.

for numerical stability
just in case $\sigma^2 = 0$

Now, these z s have mean 0 and standard unit variance.

→ Silde

this makes $\tilde{z}^{(i)}$'s have some fixed mean & variance.

$$\tilde{z}^{(i)(l)} = \gamma z_{\text{norm}}^{(i)(l)} + \beta$$

(γ and β are learnable parameters of the model)

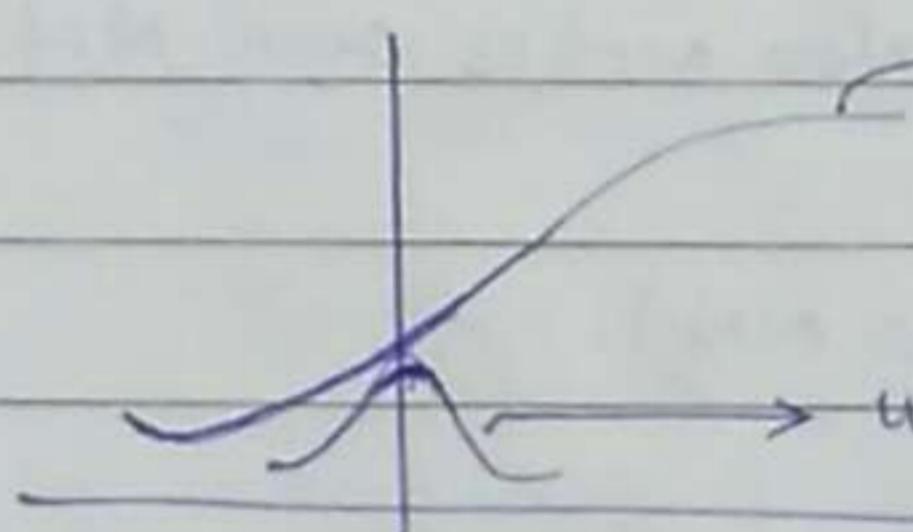
→ This is done bcoz we don't want the hidden units to always have mean 0 and variance 1, maybe it makes sense for hidden units to have a different distribution unlike the i/p features.

Note

$$\text{If } \gamma = \sqrt{\sigma^2 + \epsilon} \text{ and } \beta = \mu, \quad \tilde{z}^{(i)} = z^{(i)}$$

Now new values $\tilde{z}^{(i)}$ will be used instead of $z^{(i)}$ for training purposes and later comput' in NN.

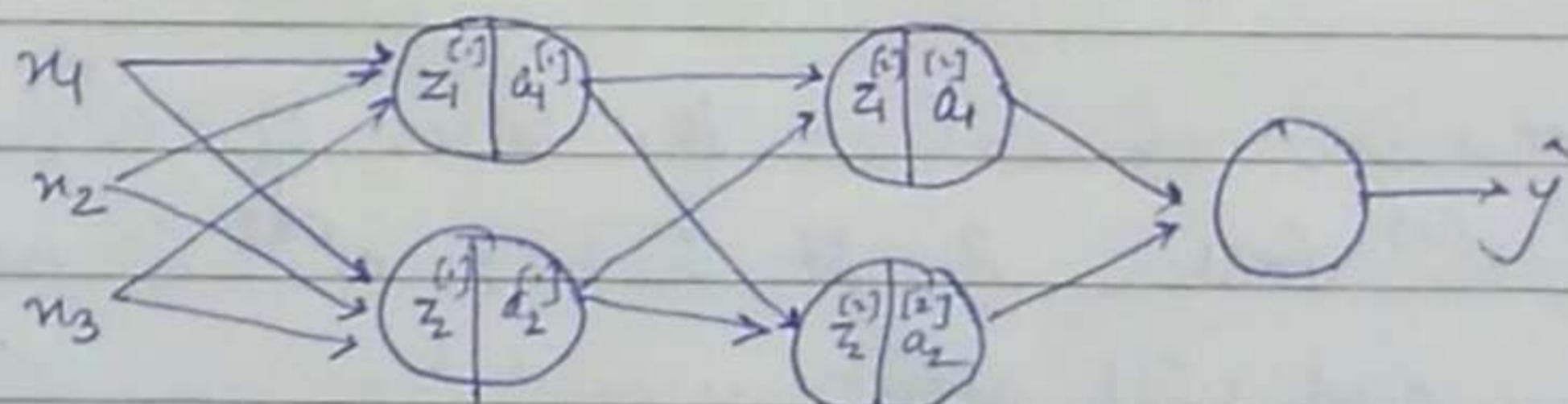
→ when μ and σ^2 have fixed values, activation (sigmoid) has greater range/variance.



→ when $\mu = 0$ and $\sigma^2 = 1$, the range reduces and is concised to such limits which is undesirable.

21/8/18

→ Fitting batch norm into the training of a deep network:



$$X \xrightarrow{w^{[1]}, b^{[1]}} Z \xrightarrow[\text{Batch Norm}]{\mu^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow a^{[1]} - g^{[1]}(\tilde{Z}^{[1]}) \xrightarrow[w^{[2]}, b^{[2]}]{\tilde{Z}^{[2]}} \tilde{Z}^{[2]} \xrightarrow[\text{BN}]{\mu^{[2]}, \gamma^{[2]}} a^{[2]}$$

So, we are using the mean & normalized value \tilde{z} instead of the regular z . Here β is different from one used in momentum.

So, the parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}, \{ \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]} \}$

This batch norm is provided in programming frameworks and can be applied with just 1 line of code as provided in tensorflow or other frameworks.

$$\beta^{[l]} = \beta^{[l]} - \lambda d\beta^{[l]} \quad \text{After using some optimization algorithm.}$$

In case of minibatches, the batch normalization step is done over batch-wise only, considering that much data only for computation.

→ Implementing GD using Batch Norm:

for $i = 1$ to no. of minibatches

 Compute fwd prop on $X^{(i)}$

 In each hidden layer use Batch Norm to replace $z^{(i)}$ with $\tilde{z}^{(i)}$

 // This ensures that the value Z end up with some normalized mean and variance.

 Use back prop to compute $d\omega^{(i)}$, $db^{(i)}$, $d\beta^{(i)}$, $dy^{(i)}$

 Update the parameters. $\omega^{(i+1)} = \omega^{(i)} - \alpha d\omega^{(i)}$

$\beta^{(i+1)} = \beta^{(i)} - \alpha d\beta^{(i)}$

$y^{(i+1)} = y^{(i)} - \alpha dy^{(i)}$

 // This also works with GD with momentum or RMSprop

 // or adam. These algos can also be used to update the

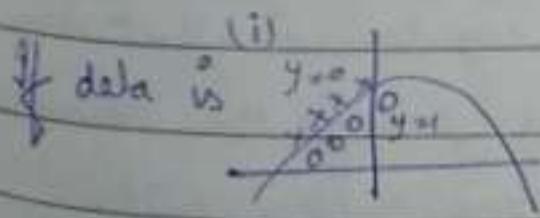
 // parameters that are added by the Batch Norm.

Q Why does Batch Norm work?

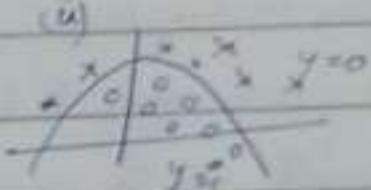
A By normalizing all the input features so to take on a

similar range of values can speed up learning. Batch Norm is a similar thing which extends it to hidden layers and not for the input features.

Also, batch norm makes the weights in the deeper layers of the networks more robust than the ones in earlier layer.

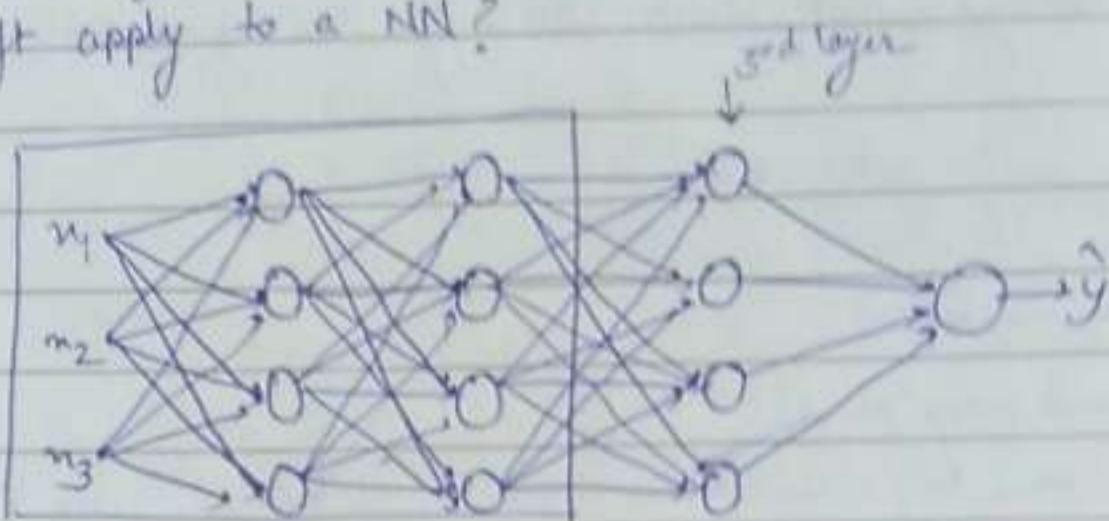


and other data is



Though the fN is same for both the data but model trained on (i) will not be able to identify the fN for (ii). This shift of data is called 'covariate shift'. So, there is a need to retrain the model if data distribution changes and the situation worsens.

even the function shifts. So, how does this problem of covariate shift apply to a NN?



Let's look from the point of view of 3rd hidden layer. This layer is getting some values $a^{(3)} = [a_1^{(3)}, a_2^{(3)}, a_3^{(3)}, a_4^{(3)}]$. Now the job of 3rd layer is to map these as to \hat{y} . The network is adopting the values of $w^{(3)}, b^{(3)}, w^{(2)}, b^{(2)}$ continuously and so as these parameters change, $a^{(3)}$ will also change, so from 3rd hidden layer's perspective, these values are changing all the time and is :- suffering from covariant shift. So, Batch Norm reduces the amount that the distribution of these hidden unit values shifts around. Even if the $z^{(3)}$ changes, the later layers stand firm and don't get affected too much.

Also, batch norm has a slight regularize? effect.

- Each minibatch is scaled by the mean/variance computed on just that mini-batch.
- This adds some noise to the values $z^{(l)}$ within that minibatch. So, similar to dropout, it adds some noise to each hidden layer's activation, becoz the estimates of the mean and standard deviation are noisy.
- So, similar to dropout, batch norm has some regularize? effect becoz by adding some noise to the hidden units, its forcing the downstream hidden units not to rely too much on any one hidden unit.
- Batch norm can be used together with dropouts becoz the regularize? effect is not too impactful.
- Using larger minibatch size, reduces noise and hence reduces regularization.

Note Don't use batch norm for regularization but as a way to normalize the hidden units and \therefore speed up learning, and regularize? is almost an unintended side effect.

Q. Batch norm works on minibatches, but at test time there are single examples at a time to process. There might not be minibatches everytime. So, what should be done to make sure your predict' make sense?

Ans

$$\begin{aligned} \mu &= \frac{1}{m} \sum_i z^{(i)} \\ \sigma^2 &= \frac{1}{m} \sum_i (z^{(i)} - \mu)^2 \end{aligned} \quad \left. \begin{array}{l} \text{averaged over} \\ \text{the minibatch size.} \end{array} \right\}$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

But at test time, taking mean and variance of one example makes no sense.

So, during training time, μ and σ^2 are computed on an entire mini batch of say 64 or some other size. Then to calculate it during test time ∞ one example at a time, just estimate μ and σ^2 from the training set, through exponential weighted avg., also called running average (across mini batches).

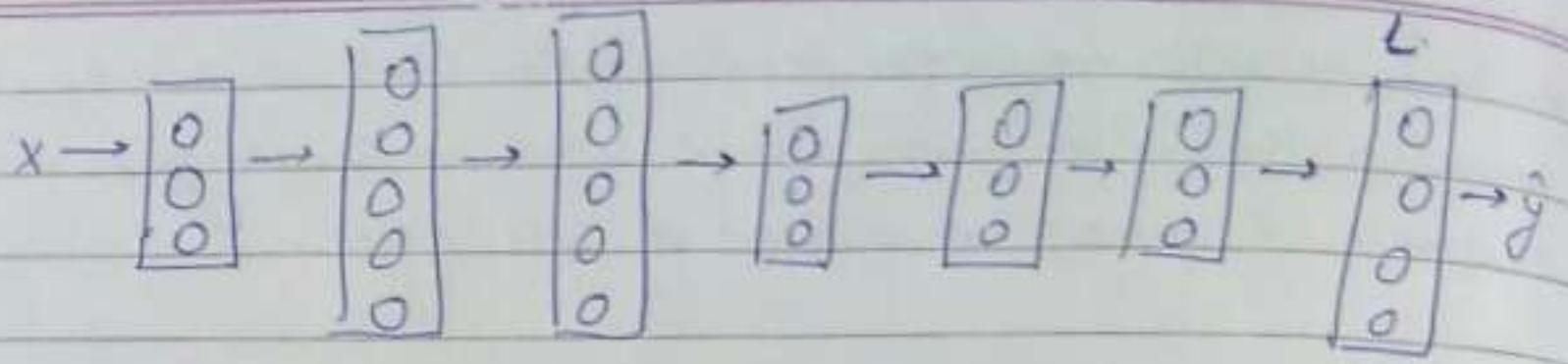
Any way to estimate μ and σ^2 for the hidden unit values z should work fine at test.

Multiclass Classification:

Softmax Regression - (generalization of logistic regression)
to C classes.

$$C = \# \text{ of classes}$$

Suppose $C = 4$ for now, so the NN would be:



The layer L is there to tell the probability

of one of the 4 classes. Also the 4 nos.

in output j should sum to 1 as they

are probabilities. To do this, softmax layer
is used.

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

$$\sum z^{[L]} = (4, 1) \cdot j$$

Activation function:

$$t = e^{z^{[L]}} \quad (\text{element wise exponentiation})$$

dimension of
 $z^{[L]}$ vector and t
vector.

Temporary variable.

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t_j}, \quad a_1^{[L]} = \frac{t_1}{\sum_{j=1}^4 t_j}$$

$$\text{e.g. if } z^{[L]} = \begin{bmatrix} 5 \\ -2 \\ 3 \\ 1 \end{bmatrix}, \quad t = \begin{bmatrix} e^5 \\ e^{-2} \\ e^3 \\ e^1 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 20.1 \\ 2.7 \end{bmatrix}, \quad \sum_{j=1}^4 t_j = 176.3$$

$$\hat{y} = a^{[L]} = \frac{t}{176.3}, \quad \text{so, } a_1^{[L]} = \frac{e^5}{176.3} = 0.842$$

$$a_3^{[L]} = \frac{e^3}{176.3} = 0.002 \quad a_2^{[L]} = \frac{e^{-2}}{176.3} = 0.042$$

$$a_4^{[L]} = \frac{e^1}{176.3} = 0.114$$

So, if X image was input and for which $z^{[L]}$ came to be $\begin{bmatrix} 5 \\ -2 \\ 3 \\ 1 \end{bmatrix}$,
the probability of belonging to each of the classes is 0.842, 0.042,
0.002 and 0.114 respectively. $\hat{y} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$

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

↓
softmax activa²

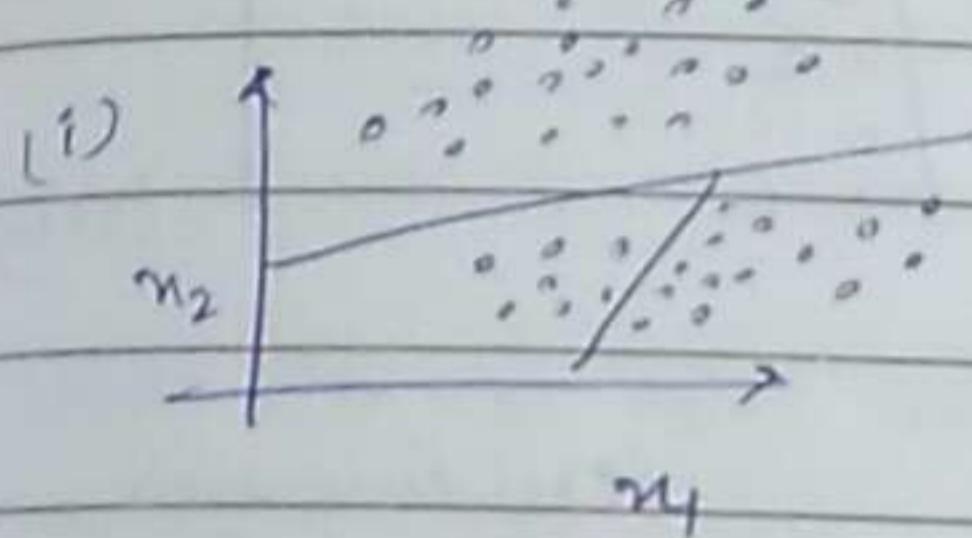
The unusual thing about this particular activation is that it takes $z^{[L]}$ as 4×1 vector

and o/p a 4×1 vector. Previously our activation function used to take in a single row vector q/p like the sigmoid and the ReLU which if a real number and o/p a real number. Softmax is unusual because it needs to normalize across the different possible outputs and it both inputs and outputs a vector.

$$\text{eg } x_1 \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \hat{y} \quad z^{[1]} = w^{[1]} x + b^{[1]}$$

$$x_2 \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad a^{[1]} = \hat{y} = g(z^{[1]})$$

↳ softmax activation



= generalization of logistic regression
with linear decision boundaries &
having more than two classes ($C > 2$).

Softmax classifier can classify data into multiple classes with linear decision boundaries as in (i) when there is no hidden layer. If there is a much deeper network used, then the softmax can learn even more complex non-linear decision boundaries to separate multiple different classes.

- Training a softmax classifier:

22/8/18

Note Softmax is called so because its soft unlike the 'hardmax' which computes t to be $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ instead of $\begin{bmatrix} 1.84 \\ 7.4 \\ 0.4 \\ 90.1 \end{bmatrix}$ from seeing the z vector.

the highest value in z vector's element is mapped to 1 and the rest are mapped to 0.

When $C = 2$, softmax reduces to logistic regression bcoz then it computes $\begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix}$ for the two classes which is effectively equivalent to calculating just 0.8 (single value that LR calculates) bcoz the other value has to be 1-0.8 without any further calculation.

e.g. Suppose $y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ and NN computes \hat{y} as $\begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$ which is not a good computation

$$L(\hat{y}, y) = - \sum_{j=1}^C y_j \log \hat{y}_j = -y_2 \log \hat{y}_2 = -\log \hat{y}_2 \quad (\text{for above case})$$

Now, for reducing the loss J^N , \hat{y}_2 should be big and less than 1
 So, L is the loss on a single training e.g., now the cost J on the entire training set

$$J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) . \quad \text{GD is then used to minimize this } J.$$

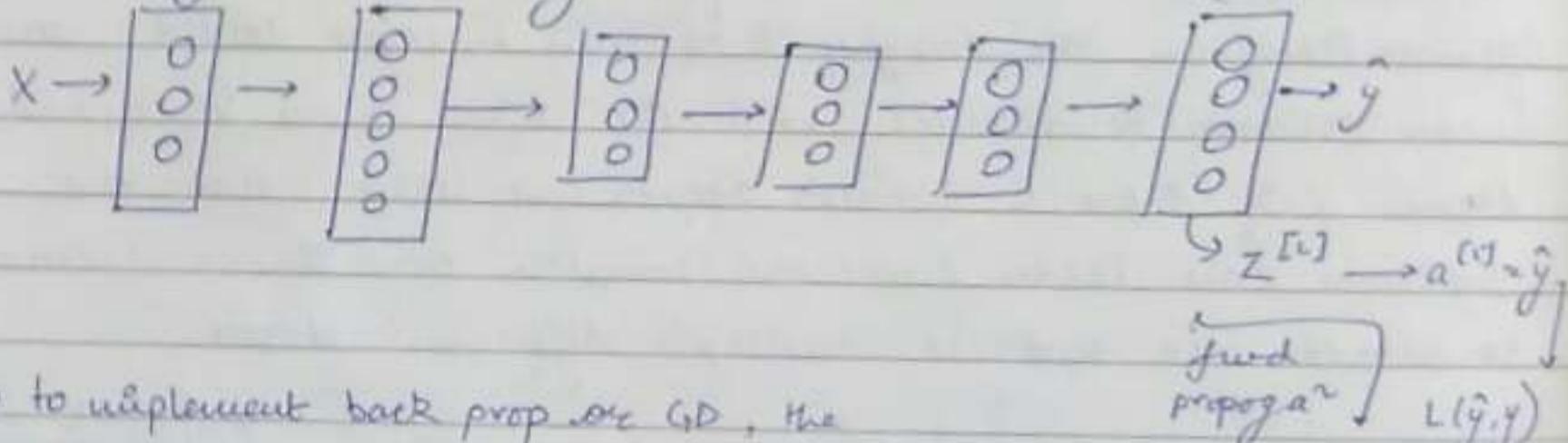
$$\begin{aligned} Y &= [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \\ &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & \dots \\ 0 & 0 & 1 \end{bmatrix} \quad (4 \times m) \end{aligned}$$

$$\begin{aligned} \hat{Y} &= [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \dots \ \hat{y}^{(m)}] \\ &= \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad (4 \times m) \end{aligned}$$

[if $C=4$ as in previous example]

→ GD on sigmoid output layer:

(ex) dimension



Now to implement back prop or GD, the key eqⁿ needed for initialization:

$$\text{backprop: } dz^{(l)} = \hat{y} - y$$

(e, r)
dimension

$$\left[dz^{(l)} = -\frac{\partial J}{\partial z^{(l)}} \right] \quad \begin{array}{l} \text{(using this back prop can be implemented)} \\ \text{to compute all the derivatives} \end{array}$$

Deep Learning Programming frameworks:

These frameworks help us implement complex deep learning models more efficiently. Writing everything and everything from scratch is not practical and very efficient as using frameworks is.
 Some frameworks that make it easy to implement ANN are:-

- ① Caffe / caffe2
- ② CNTK
- ③ DL4J
- ④ Keras
- ⑤ Lasagne
- ⑥ nnunet
- ⑦ TensorFlow
- ⑧ Theano
- ⑨ PyTorch

- choosing deep learning frameworks:
- ease of programming (development and deployment)
- learning speed (some frameworks allow efficient run on large dataset)
- fully open (open source with good governance)

→ Tensorflow:

Basic structure of tensorflow program -

Suppose we want to minimize $J(w) = w^2 - 10w + 25$

(cost)



$(w-5)^2$ and $w=5$ gives minimum values.

But let's see how we can implement something in Tensorflow to minimize this.

Code:

1. import numpy as np
2. import tensorflow as tf
3. $w = \text{tf.Variable}(0, \text{dtype}=\text{tf.float32})$ // w is the parameter to be optimized.
4. cost = $\text{tf.add}(\text{tf.add}(w^2, \text{tf.multiply}(-10, w)), 25)$
5. train = $\text{tf.train.GradientDescentOptimizer}(0.01).minimize(\text{cost})$
6. init = $\text{tf.global_variables_initializer}()$ // learning rate
7. session = tf.Session // to start tensorflow session
8. session.run(init) // to initialize global variables.
9. sess.run(w) // tensorflow is evaluative variable
// to run the learning algorithm.

train is our learning algo with GD Optimizer to minimize J .

Value for session.run(w) = 0 bcoz we haven't initialized w yet.

I run anything yet and so the initialized value is returned for w .

In line 5, some other optimizer can also be used.

10. session.run(train) // will run 1 step of GD.
11. print(session.run(w)) // value for $w=0.1$ now.
12. for i in range(1000): // running 1000 iterations of GD now.
13. session.run(train)
14. print(session.run(w)) // value for $w=4.999$ as we would expect bcoz J was $(w-5)^2$.

|| So tensorflow automatically knows how to take derivatives of add
|| and multiply as well as other f^N which is why we had to
|| implement only feed prop and tf can itself figure out how to do
|| the back prop or GD computation.

Also, $\text{cost} = w^*w - 10^*w + 25$ can also be written in tensorflow
because once w is declared a tensorflow variable, $+, *, -$
operations are overloaded.

Note The previous example had tensorflow minimize a fixed function
of w but one of the f^N we want to minimize is the f^N of our
training set which changes during training. So, to get training
data into a tensorflow program:

$n = \text{tf.placeholder}(\text{tf.float32}, [3, 1])$ || n having 3 values

$\text{cost} = n[0][0]^*w^*w + n[1][0]^*w + n[2][0]$

So now n has become sort of like data that controls the coeff.
of this quadratic fn.

If placeholder tells tf that n is something for which values will be
provided later.

coefficients = np.array([[-1.], [-10.], [25.]])

session.run (train, feed_dict={n: coefficients})

The in between steps remain the same as before code.

So the placeholder allows us to get the training data into the J
conveniently. feed_dict sets the coeff of the eqⁿ as n . this is
it can optimize the given cost for by getting trained
with changing data in just a few lines of code.

Note The cost fn is the heart of the tf program and tf has built
in functions to compute the backward f^N for backprop and
compute derivatives for us. That's why backprop need not be
explicitly implemented.