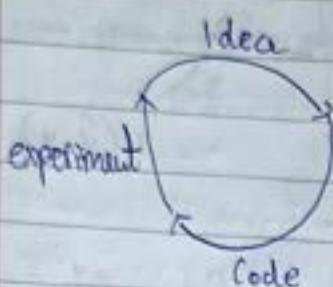


Applied deep learning is a very empirical process.



Applying DL is an iterative process by doing trial and error:

- Set initial α value. Start with guess for hyperparameters $\alpha = 0.01$
- Train and Evaluate: Observe the behavior of cost function J :
 - if J converges slowly, α might be too small
 - if J diverges, α might be large
- Adjust the hyperparameters based on observations

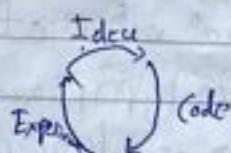
What does this have to do with the brain?

Course 2: Improving Deep NNs: hyperparameter tuning, regularization, optimization
Setting up ML applications

Train / Dev / Test sets

Applied ML is a highly iterative process. When training a neural network, we need to make a lot of decisions:

- no of layers NN will have
- no of hidden units you want each layer to have
- What activation function you want to use for each layer
- What learning rate to use



- Very hard to guess the hyperparameters correctly the very 1st time
- Applied DL is an iterative process where you undergo the cycle many times to have a good choice of network for the application.
- One thing we can determine is how quickly you can make progress around the cycle effectively.

Course - 2 Improving Deep NN - Hyperparameter tuning, Regularization, Optimization

Technique

Data

Train set

- Holdout

- Test

- Validation

- Dev set

Workflow:

- Train the algorithms on train set
- Use the dev set to see which different models perform the best on your dev set.
- When you have the final model that you want to evaluate. Use the best model you found on the test set to get an unbiased estimate of how well the algorithm is doing.
- Previous era we used: → 70% train, 30% test
 → 60% train, 20% dev, 20% test
- This is ok if we have examples 10 - 10000
- In the Big Data era where we got 1 million examples, the trend is that your dev, test set has to be smaller percentage of the total. Remember: Dev set only has to be big enough for you to evaluate different algorithm choices and decide which works better. ex, if we 1 mil ex for dev set we only need 10000 ex.
~~so~~
 → Split for 1 mil ex, 99% train, 1% dev, 1% test
 we can also use 99.5% train, 0.25% dev, 0.25% test
 or 0.5% train, 0.4% dev, 0.1% test

Mismatched train/test distributions:

ex ~~so~~ We build an app to detect photos of cats.

Train set: Cat pictures from webpages
 (This has professional, high res, nicely framed pictures)

Distribution
 of data
 is different

Test / Dev set: Cat pictures from ^{users}
 using the app (users will take photo from their own camera)
 Which maybe blurrier, lower res)

→ ~~so~~ Rule thumb: Make sure the dev and tests come from the same distribution

- ex for rule of thumb: if cat train ~~8x's~~ are from internet and dev/test ex's are user uploaded from the user's phone they will mismatch. Make sure ~~you~~^{that} dev and test set are from the same distribution
- The dev set rule is to try some good models you created.
- If you don't need an unbiased estimate, it's fine not to have a test set. Since we don't have the test set, we train on train set and then try different model architectures. Evaluate them on dev set and then use to iterate and try to get good model. Since we fit dev on dev set it no longer gives an unbiased estimate of performance
- In the above case we call dev set as test set

Bias and Variance

Logistic regression

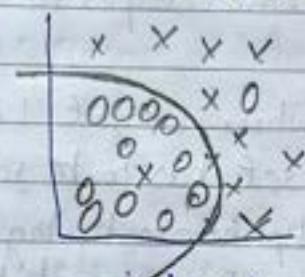


high bias

Underfitting

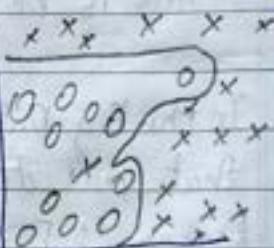
- Doesn't fit very good to data has a high bias
- Underfits

cff with medium complexity



just right

complex classifier
Deep NN (a lot of hidden units)



high variance

Overfitting

- May fit data perfectly
- Not a great fit
- High variance and overfits

ex, Dog / Cat classification

(Case 1: 1 f (cats), 0 f (dog))

- Train set error : 1%

- dev set error : 11%

- We have high variance model has high dev set error and low train set error.

(Case 2:

Train err: 15%

Dev set err: 18%

- Underfitting

- High bias

- Doesn't do well on both train, dev set

(Case 3:

Train err: 0%

Dev set err: 30%

- High variance

- Low variance

(Case 4:

Train err: 0%

Test err: 1%

- High variance

- Low variance

- These assumptions come from the baseline that human has 0% error or more generally they are the optimal error / Bayes error.

Case 5:

Train error: 15% ↗ Pay attention to

Dev error: 16% ↗ the metrics

human error: 15% ↗

so the above case is ↗ low variance because the human err itself is 15%.

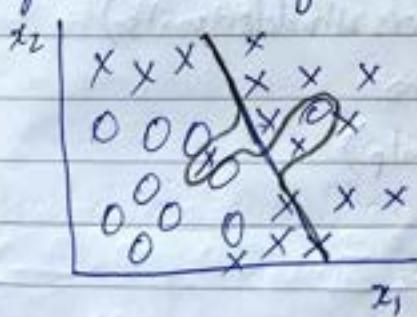
In the case of analysing bias / variance when no clf can do well:

- ex, we have very blur images, even human or the system can't do well ~~then~~
- ~~Bayes~~ Then the bayes error is too high

• Note:

- By looking at J_{train} it gives an idea of how well your fitting to the train set it tells if we will have a bias problem.
- Looking at how high error goes when you go from train set to dev set, that gives you sense of how ~~bad~~ bad the variance problem is.
- All this is under the assumption that ~~Bayes error is small~~
- Bayes error is small
- Train and Dev sets are drawn from the same distribution

High bias and high variance:



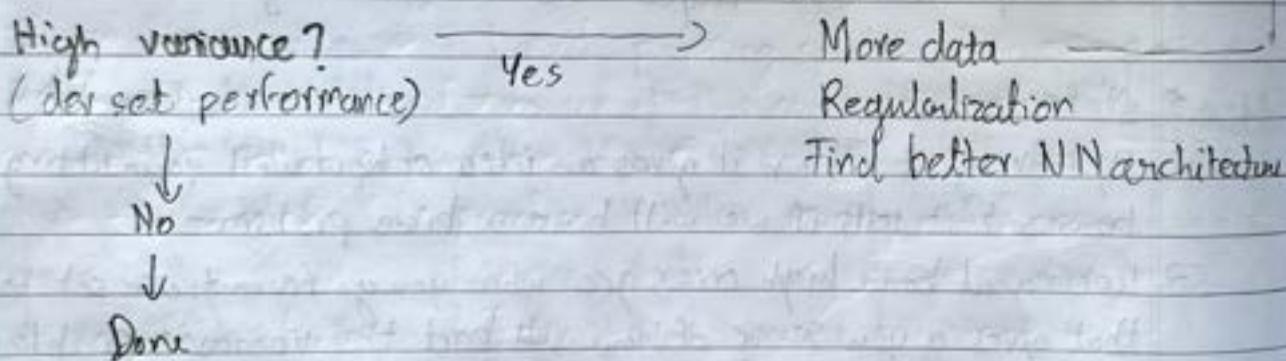
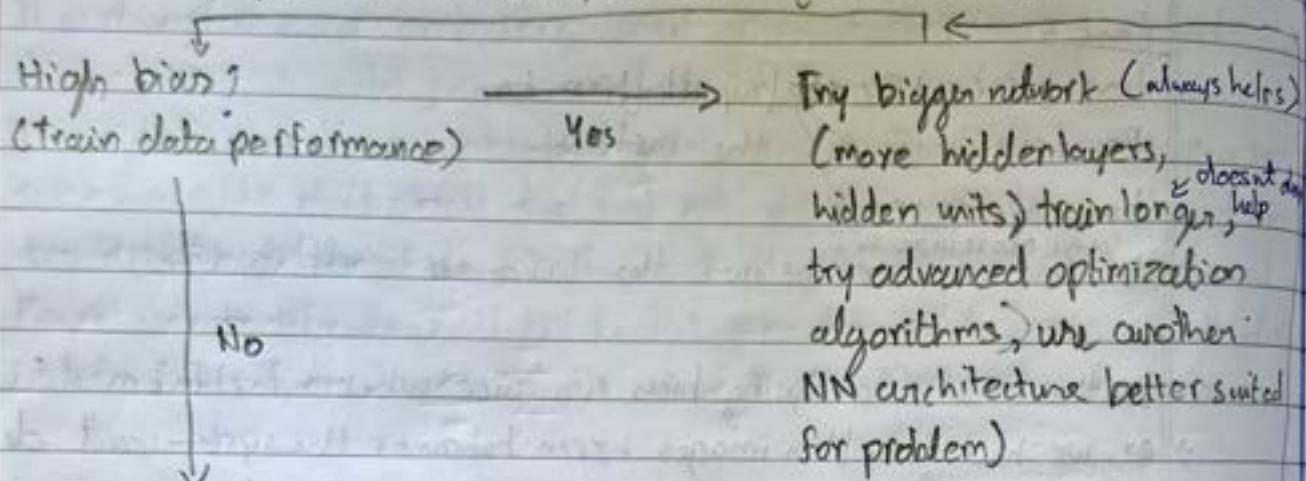
Black line denotes high bias and high variance

High bias, because by being a linear clf, it doesn't fit the quadratic shape well. But by having too much flexibility ^{in middle} it gets ~~those~~ those ex, high variance.

Basic Recipe for ML

After training an initial model we ask if algorithm has:

repeat twice until you can fit the train set well

Key Points

- Use the train ^{and} dev set to determine if you got a bias or variance problem, then you can do the appropriate things.
- ~~Handle~~ High bias problem:
 - Getting more train data would help (not efficient to do either)
 - Try bigger network (more hidden layers & hidden units)
 - Try advanced optimization algorithms
 - NN architecture search (May help)

High variance:

- More data
- Regularization
- NN architecture search

- Bias-Variance tradeoff: → Increase bias and reduce variance
→ Reduce bias and increase variance
- In the modern Big Data era - Pre DL era, not many tools were there that reduce bias or reduce variance without hurting the other one.
- In modern Big Data era, you can train a bigger network and so, as long as you can keep getting more data isn't always the case for either of these, but if it is ~~the~~ the case, then getting the bigger net always reduces bias without necessarily hurting Variance, as long as you regularize it appropriately.
- Getting more data ~~per~~, always reduces variance and doesn't hurt bias
- With bigger net ^{and} more data we now have tools to drive down bias or drive down variance without hurting the other thing that much. This is one of big reasons that DL is useful for supervised learning and we have a lesser tradeoff where you carefully have to balance bias and variance but sometimes you ~~get~~ more options to reduce bias or variance without hurting the other.
- Training a big net never hurts but the cost of training a NN that's too big is just computational time as long as you regularize it.

Regularizing your NN

Regularization

- If NN is overfitting the data, we have high variance.
- We can use To address this we can try:
 - regularization
 - Get more training data, ^{that's reliable} but that's expensive and not always possible

Note: We will use Regularization on Logistic Regression

Logistic Regression : Recall the cost func $\min J(w, b)$, $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

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

To add regularization we add $\frac{\lambda}{2m} \|w\|_2^2$ to the cost fn.

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

[$\|w\|_2^2 = \sum_{i=1}^{n_x} w_i^2 = w^T w$]

- Why don't we regularise b ?
- w is a high dimensional parameter, especially if it's high variance problem.
 - w has a lot of parameters, while b is a single no.
 - If we were to add regularization to b it won't make a difference because b is only 1 parameter

$$L_2 \text{ regularization} : \frac{1}{2m} \|w\|_2^2 \Rightarrow \|w\|_2^2 = \sqrt{\sum_{j=1}^n w_j^2} = \sqrt{w^T w}$$

(most commonly used)

$$L_1 \text{ regularization} : \frac{\lambda}{2m} \sum_{j=1}^m |w_j| = \frac{\lambda}{2m} \|w\|_1 \quad (\text{we can use either } \lambda/m \text{ or } \lambda/2m, m \text{ is just a scaling constant})$$

- If we use L_1 reg, w will end up being sparse. This means w vector will have lots of zeros in it. This can help compress the model because the set of parameters are 0, then you need less memory to store the model. (In practice it helps making model sparse only a little)

λ : regularization parameter

- This is set using holdout/validation set / dev set
 - Try a variety of values for λ , see which gives the best result in terms of trading off between doing well on train set ~~vs~~ ^{to be} vs setting the norm of parameters small to prevent overfitting.
 - Lambda λ is another hyperparameter that has to be tuned
- Note: lambda is a ^{Python} keyword

Regularization for NN

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

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^{n^{(l+1)}} (w_{i,j}^{(l)})^2$$

Frobenius norm \rightarrow

$$w = (w^{(1)} \ n^{(2-1)})$$

no of hidden units in layer l and \rightarrow prev layer $l-1$

rows of the matrix should be no of neurons in current layer $n^{(l)}$

cols of the wt matrix should equal no of neurons in the prev layer $n^{(l-1)}$

6.3 Computing gradient descent:

$$\frac{dw^{(l)}}{w^{(l)}} = \text{(from backprop)} + \frac{\lambda}{m} w^{(l)} \quad \frac{\partial J}{\partial w^{(l)}} = dw^{(l)}$$

$$w^{(l)} := w^{(l)} - \alpha \frac{\partial J}{\partial w^{(l)}}$$

→ with regularization coeff of w slightly < 1 in which case its called weight decay. L2 regularization is called weight decay

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

$$= w^{(l)} - \frac{\alpha \lambda}{m} w^{(l)} - \alpha \text{(from backprop)}$$

This term shows whatever wt mat $w^{(l)}$ is we are making it smaller

$$= w^{(l)} \left(1 - \frac{\alpha \lambda}{m} \right) - \alpha \text{(from backprop)}$$

we are multiplying wt mat w with $\left(1 - \frac{\alpha \lambda}{m} \right)$ which will be < 1 , this is why L2 reg is called weight decay

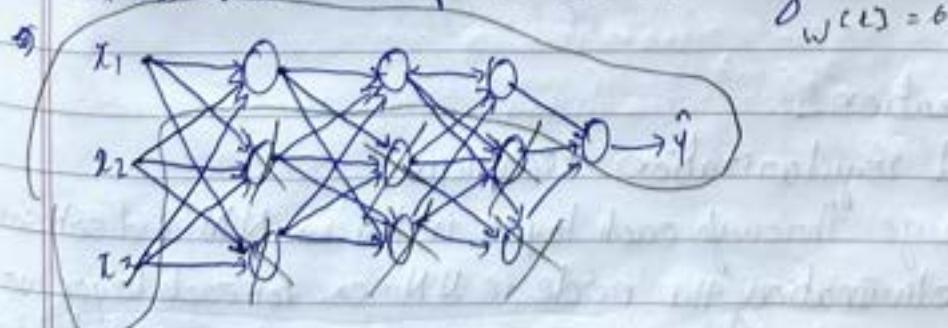
Why does regularization prevent overfitting?

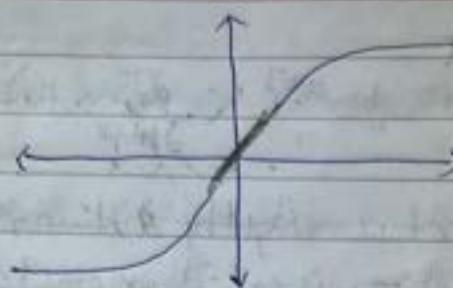
→ Why does shrinking the L2 norm with parameters cause less overfitting?

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

→ if we make the regularization term very small

→ If we make the λ very big, then weight matrices W will be reasonably close to zero, ~~it will set the weights~~ effectively zeroing out the impact of hidden units. The simplified NN becomes a much smaller NN, eventually almost like logistic regression. We will end up with a small NN therefore is less prone to overfitting.





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

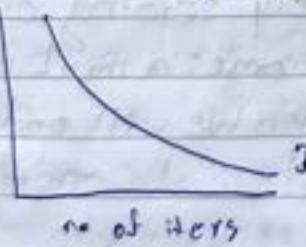
- Taking activation function ~~to~~ $g(z) = \tanh(z)$ as or, if z is large then weights w is small and subsequently z ends up taking relatively small values, where g and z will be roughly linear which ~~is not able~~ to fit those very complicated decision boundary i.e less able to overfit.
- Every layer is roughly linear

Implementation tip:

When implementing regularization,

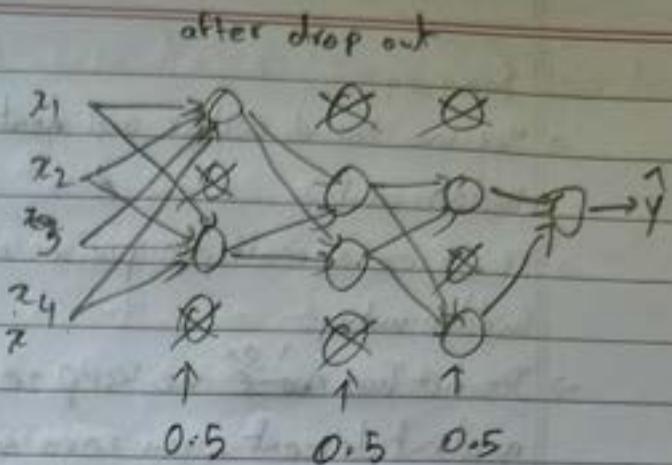
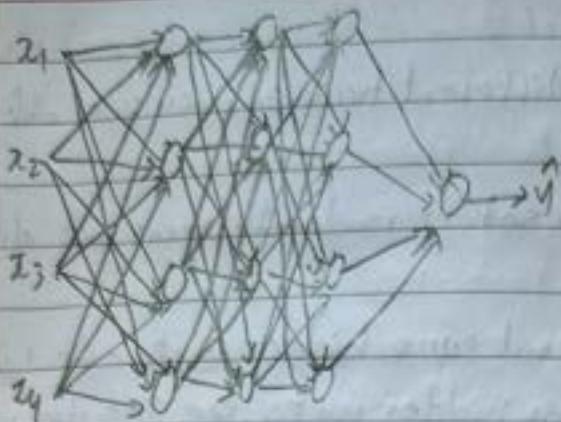
$$\text{cost fn: } J(w^{(1)}, b^{(1)}) = \sum L(y^{(1)}, \hat{y}^{(1)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w_l\|_F^2$$

- We have $\frac{\lambda}{2m} \sum_{l=1}^L \|w_l\|_F^2$ which penalizes the weights being too large
- When we implement gradient descent, one of the steps to debug gradient descent is to plot cost function J as the no. of ~~iterations~~^{iterations} for gradient descent, you want to see cost fn J decrease monotonically after every ~~elevation~~^{iteration} of gradient descent. Make sure you plot J with reg term, without it won't decrease monotonically



Dropout Regularization

- Another powerful regularization technique
- With dropout we go through each layer in the ~~NN~~ NN and set some probability for eliminating the node in NN. ex, for each layer we go through the nodes we set $p = 0.5$, it will have 0.5 chance to keep and drop the node, then we remove all ingoing and outgoing connections from that node
- After that we have a small diminished net ~~and~~ and after that we can do backprop or forward pass on diminished net.



- For train ^{cost} ex you would train it on the neural ~~net~~^{smaller} net which has a regularizing effect

Implementing dropout (Inverted dropout):

- Illustrating with layer $l = 3$ $\text{keep_prob} = 0.8$ $0.2 \rightarrow \text{prob of}$ eliminating
- Set ~~a~~ dropout vector, d_3 , \uparrow probability of keeping node
- $d_3 = \text{np.random.rand}(\alpha_3.\text{shape}[0], \alpha_3.\text{shape}[1]) < \text{keep_prob}$
- d_3 will be a matrix, each ex have a hidden unit there is a ~~80%~~^{80%} chance that will correspond to 1, 20% chance its 0
- ~~We will~~ take activations α_3 in 3rd layer
- $\alpha_3 = \text{np.multiply}(\alpha_3, d_3)$ $\# \alpha_3 * = d_3$
- $\alpha_3 / = \text{keep_prob}$ \rightarrow inverted drop out
- ex, we have 50 units in 3rd hidden layer
 - α_3 would be $(50, 1)$ or $(50, m)$ dim
 - 80% - keep, 20% - drop each unit
 - We ~~drop~~^{keep} 10 units
- $\alpha_3^{(4)} = w^{(4)} \alpha^{(3)} + b^{(4)}$ (In order to not reduce $\alpha_3^{(4)}$, we can take)
 - reduced by 20% $\alpha^{(3)} / = 0.8$ because it roughly bumps it back to $\approx 20\%$. This does not change the expected value of α_3)
- Inverted dropout technique is $\alpha_3 / = \text{keep_prob}$, ensures expected value of α_3 remains the same which makes test time easier because you have less scaling problem.

Summary:

- You use the d vector and ~~mask~~ for different train ex, you zero out different hidden units.
- If we make multiple passes on trainset it randomly zeros out different hidden units.
- It's not for ~~one~~^{1^{ex}} we keep zeroing out same hidden units - On 1st iter of gradient descent, you zero out some hidden units, On 2nd iter we zero at a different pattern of hidden units.
- vector d_3 for 3rd layer decides what to zero out, only showing prob here.

Making predictions at test time:

- At test time given some x we want to make a pred

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

- ~~No~~ dropout in test time

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

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

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

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

⋮

⋮

⋮

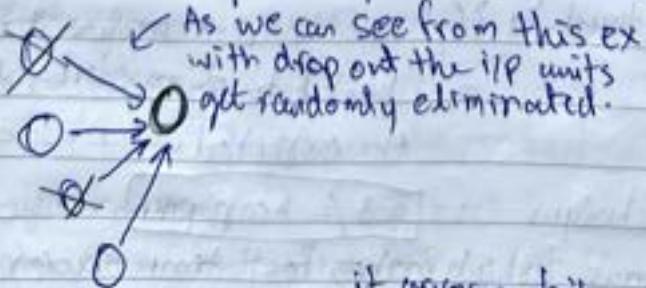
Note: We are not using dropout in test time because while making preds we don't want o/p to be random

- if we were using dropout we would be adding noise to preds. Also if we dropped out hidden units randomly its computationally inefficient and it gives the same result.

Understanding Dropout

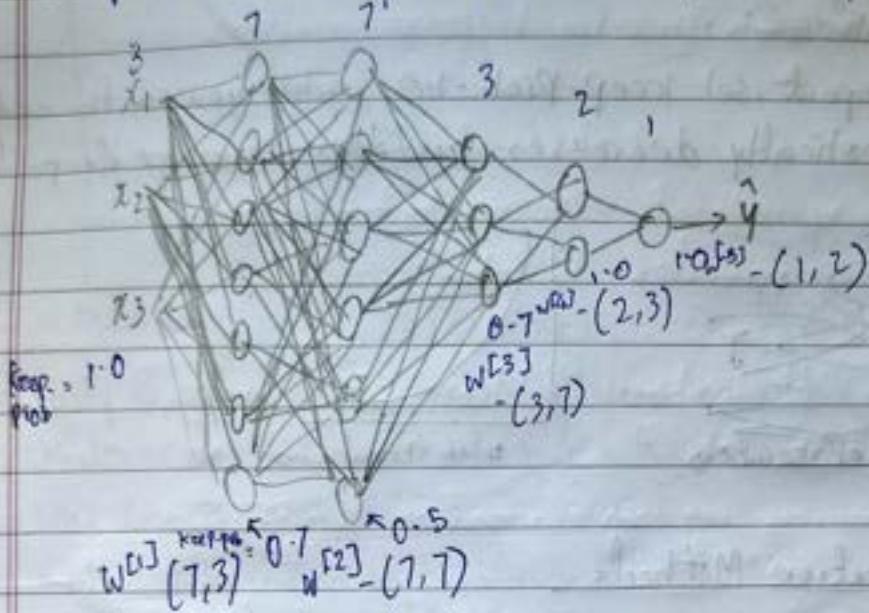
Why does dropout work?

Intuition: Can't rely on any one feature so have to spread out weight



- By spreading the weights ~~the input~~ of weight to each of 4 i/p.
- Spreading the weights would shrink the ^{squared} norm of weight similar to L2 regular helping to prevent overfitting.

- L2 regularization can be applied in different ways can be adaptive to scalable to different ilp.
- Its possible to vary keep-prob by the layer

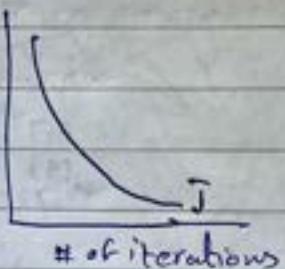
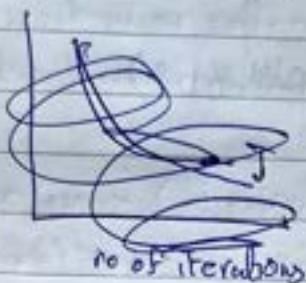


- $W^{[2]}$ has the largest set of parameters and has the biggest weight. To reduce overfitting for $W^{[2]}$ we might have keep-prob relatively low (keep-prob = 0.5). Different layers can have keep-prob = 0.7 (overfitting won't be a concern) for every layer.
- keep-prob is different layer
- layers that have keep-prob = 1 we won't drop them.
- For layers that overfit and have a lot of parameters keep prob is smaller and you can apply powerful form of dropout (similar to regularizer)
- We can also apply dropout to ilp layer but it isn't done often in practice, it usually uses keep-prob = 1.0 or 0.9.
- For layers you are more worried about overfitting, you can set lower keep-prob for some layers, you apply dropout - The downside is that we will more hyperparameters to search for using cross validation
- An alternative can be for some layers you apply dropout and some layers where you don't apply dropout. Then just have 1 hyperparameter to search for
- Many of 1st successful implementations of dropout are in computer vision. In CV the ilp sizes are too big so never have enough data prone to overfitting.

- Downside of dropout: cost function J is no longer well defined, on every iteration we ~~are~~ zeroing out a bunch of nodes, harder to check performance of gradient descent - ~~so~~

Soln for the above: .

- Turn off dropout, set keep-Prob = 1.0 , when you run the code, monitor J monotonically decreases and then turn on dropout.



Other Regularization Methods

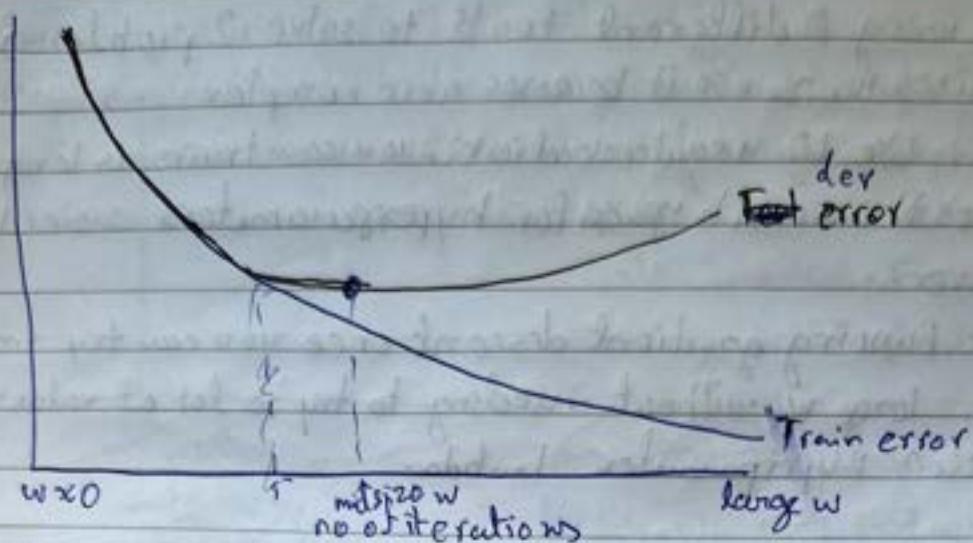
Data Augmentation

Data Augmentation:

- Getting more training data is expensive sometimes its not possible
- We can ~~not~~ augment the train data by flipping it horizontally and we can add to our train set ~~in~~ ^{however} doubling the train set
- ~~Train set is~~, a bit redundant it isn't as good as collecting a new independent example
- We can ~~do~~ do this without ^{any} expense, ~~so~~
 - ex, - random flips, - randomly zooming
 - random crops
 - random ~~distortion~~ distortion and translation of images
- These methods can be used to generate more fake training ex. ^{augmented data}
- ~~These ex~~ This is an inexpensive way to give your algo more data and sort of regularize it to reduce overfitting

Early stopping:

- Ex, we run gradient descent, plot the train error J_{train} on train set it should monotonically decrease
- We can plot the ~~test~~ set error ~~error~~



→ Early stopping stops training your neural network halfway and takes whatever dev set error at that value.

Why does it work?

- We didn't run many iterations for the NN because $w \approx 0$ because we want to avoid overfitting.
- With random initialization w has small random values and as we iterate w gets bigger and bigger parameter values
- Early stopping: stops training midway to get midsize w .

~~Penalties~~

Orthogonalization: ML Process comprises of several steps:

- Optimize cost fn J
 - Gradient descent, RMSprop, Adam
 - Not to overfit
Regularization
- You have 1 set of ~~parameters~~ ^{top} for optimizing ~~the~~ cost fn, we only care finding w, b , so that $J(w, b)$ is as small as possible. We only focus on this.
- Reducing variance is a separate task, you use a separate set of tools for it. This principle is orthogonalization.

Downside: → It couples the tasks of optimizing cost fn J and not to overfit
 → We can't work on those problems independently because by stopping gradient descent early - ~~we're breaking what we're trying to do~~ doesn't do a great job with ~~optimizing~~ optimizing cost fn ~~and we~~ Additionally we are trying not to overfit

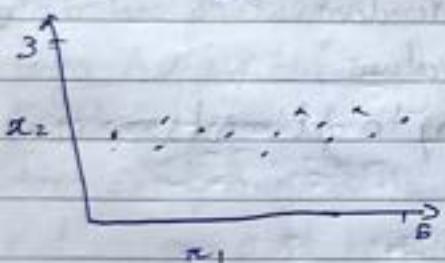
- Instead of using 2 different tools to solve 2 problems you will use one that mixes the 2, so it becomes more complex
- Alternative: Use L2 regularization, you can train as long as possible. It also makes the search space for hyperparameters easier to decompose and search over.
- Advantage: Running gradient descent once you can try small w , medium w , large w without needing to try a lot of values for L2 regularization hyperparameter lambda.

Setting Up your Optimisation Problem

Normalizing Inputs

- To speed up training we can normalize the inputs

Normalizing train sets



After subtracting mean

Normalization corresponds to 2 steps:

- ① Subtract the mean

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x := x - \bar{x}$$

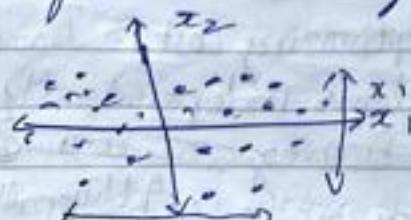
- ② Normalize the variance

x_1 has more variance than x_2 .

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

$$x := \frac{x - \bar{x}}{\sigma}$$

After doing ② we will get this

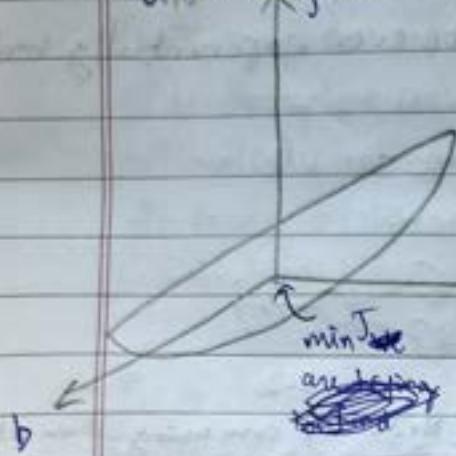


Variance of x_1 and x_2 is 1

Tip: If you use this \bar{x}, σ to scale training data we ~~use~~ ^{some} \bar{x}, σ to normalize the test data. Don't normalize train and test set separately. We want the obs to go through the same transformation on training and test samples

Why normalize inputs?

Unnormalized

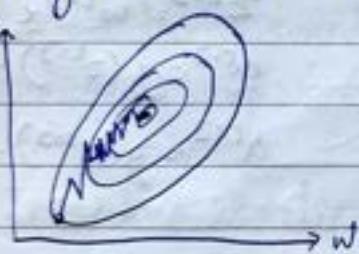


$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

- cost fn looks very elongated
- if features ~~are~~ are on very different scales ~~so fast~~. For ex,
- x_1 ranges from (1, 1000)

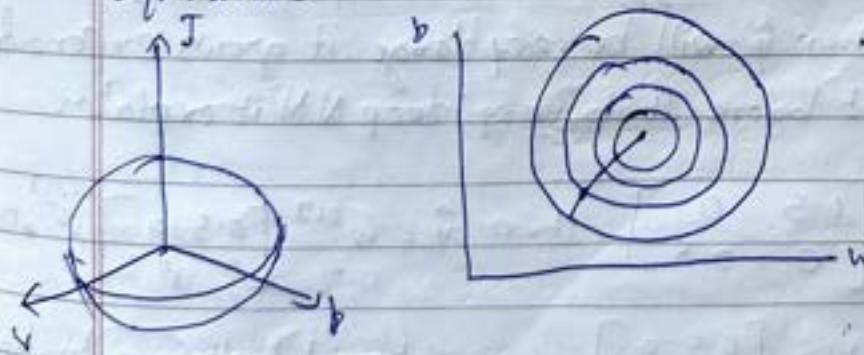
~~-~~ Ratios for range of values of parameters for w_1, w_2 will end up taking different values

When we plot the contours of the function we can have a very elongated ~~function~~ function.



- If we run gradient descent we need to use a small learning rate because the gradient descent might need a lot of steps to oscillate back and forth, ^{before} finally reaches the minimum.

Normalized: if normalize features the cost fn will look more symmetric



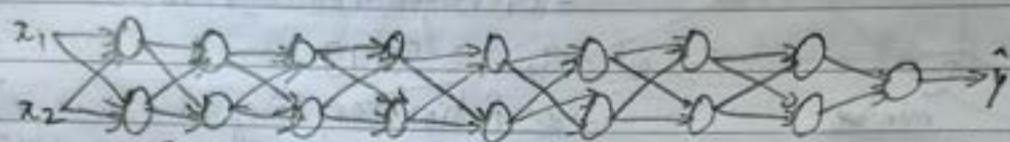
- When we have more spherical contours the gradient descent can go straight to the minimum
- We can take longer steps + for gradient descent rather than needing to oscillate

Intuition: cost function is easier to optimise when features are in similar scale. ex instead of $x_1 \rightarrow [0, 1000], x_2 \rightarrow [0, 1]$ makes it easier and faster to optimise

→ We can set all of them to have ~~to~~ mean and ~~a~~ variance, sets all features to similar scale and normalizes them.

Vanishing / Exploding Gradients

- When you're training a deep NN your derivatives or slopes can be either very very large or small (maybe even exponentially small) and this makes training difficult
- Ex,



$$g(z) = z \quad b^{[L]} = 0$$

$$\hat{y} = w^{[L]} w^{[L-1]} w^{[L-2]} \dots \boxed{w^{[3]} w^{[2]} w^{[1]} x} \rightarrow a^{[3]}$$

(Note: we are using linear activation)

$$w^{[L]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

$$\hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{[L-1]} x = 1.5^{L-1} x$$

$$z^{[L]} = w^{[L]} x$$

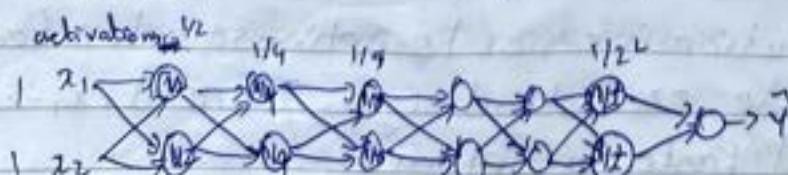
$$a^{[1]} = g(z^{[1]}) = w^{[1]} x$$

$$a^{[2]} = g(z^{[2]}) = g(w^{[2]} a^{[1]})$$

- The weights if they're little bit bigger than 1 or identity matrix then deep units will grow exponentially
- if L is very high then \hat{y} will be very large, it grows exponentially like 1.5 to number of layers. If its very deep NN it explodes

$$\rightarrow \text{if we replace } W = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \text{ then } \hat{y} = w^{[L]} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{[L-1]} x = 0.5^L x$$

Ex, $x_1, x_2 = 1$, $w = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$ (weights < 1) then activations become



the activation values decrease exponentially

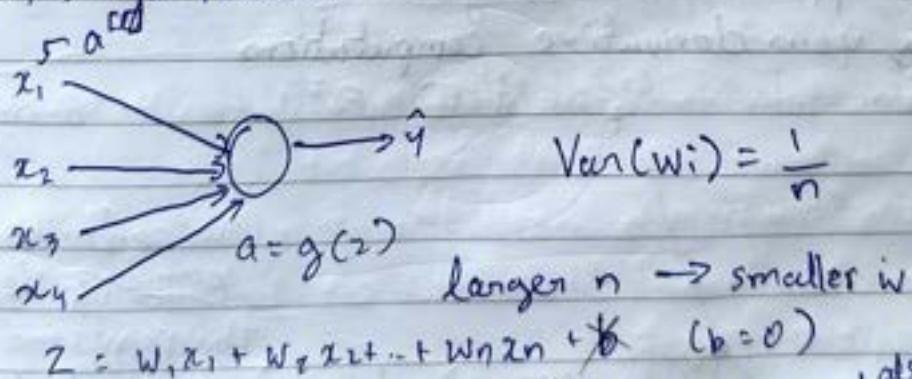
Important note: The weights W , if they are just a little bit bigger than 1 or the identity matrix then with deep NN the activations will explode. If W is little less than identity, then with deep NN, ~~the~~ the activations decrease exponentially.

- Similarly we show that the derivatives or gradients will also decrease or increase exponentially as function of no of layers.
~~= with deep NN if activations increase~~

→ In a deep NN if activations or gradients increase or decrease exponentially as a function of L then these values can get really big or really small. This makes the training difficult especially if gradients are exponentially smaller than 1, ~~gradient~~ gradient descent takes small steps. It will take a long time for gradient descent to learn anything.

Weight Initialization for Deep Networks

- A partial solution to exploding or vanishing gradient is a better or more careful choice of random initialization of weights
- ex, Single Neuron



$$\text{Var}(w_i) = \frac{1}{n}$$

larger $n \rightarrow$ smaller w

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b \quad (b=0)$$

→ In order to prevent z from blowing up and ^{also} not becoming too small, larger the n the smaller the w . We can set variance = $\frac{1}{n}$ and we can set $W^{(l)} = \text{np.random.randn}(\text{shape of mat}) * \text{np.sqrt}\left(\frac{1}{n^{(l-1)}}\right)$

$$\text{np.sqrt}\left(\frac{1}{n^{(l-1)}}\right)$$

* no of features fed in unit at l

If we use ReLU activation we can set, $\text{Var}(w_i) = \frac{2}{n}$

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

- So if all features ~~are completely~~^{have} $w=0$, $\sigma=1$, $\text{var}=1$ it would cause σ to take similar scale and it doesn't solve but it reduces the vanishing/exploding gradient problems.
- It tries to set each of the weight matrices w so that its not too much bigger than 1 or not too less than 1. So it doesn't explode/vanish too quickly.
- This is by ReLU.

Tanh Variant:

~~$w^{[l]}$~~ $w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}(1/n^{[l-1]})$
This is called Xavier Initialization

Yoshua Bengio Variant:

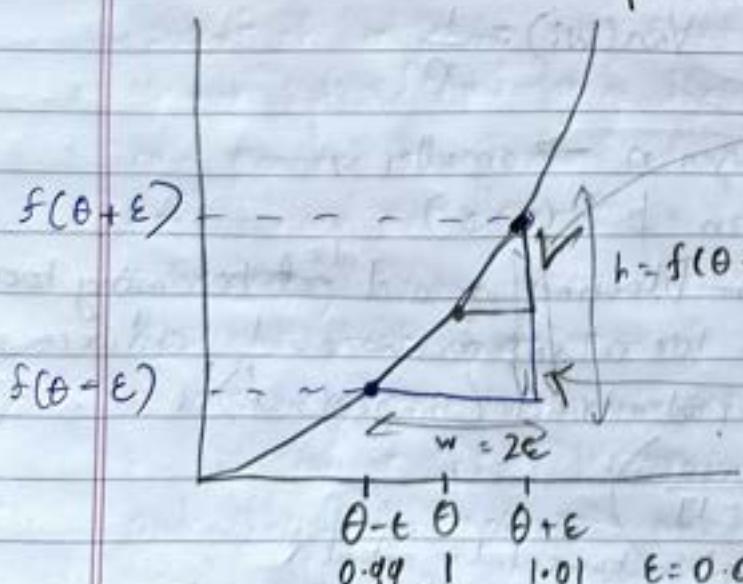
$$w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}(2/(n^{[l-1]} + n^{[l]}))$$

- Variance can be hyperparameter $\text{var}(w) = 1/n$ or $2/n$, but its not as important as other hyperparameters.

Numerical Approximation of gradients

Checking your derivative computations

$$f(\theta) = \theta^3$$



Instead this triangle and computing slope. We can the points marked in blue $b = f(\theta + \epsilon) - f(\theta - \epsilon)$ and compute the slope at

This slope gives a better approx of θ

$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta)$$

$$\frac{1.01 - 0.99}{2(0.01)} = 3.001$$

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

approx error = 0.0001

from prev-slide $g(\theta) = 3 \cdot 0.301$ error = 0.03

- This method for gradient checking and backprop, it's slow. In practice it's worth using the method because it's more accurate.

Two-sided difference formula is more accurate:

- Two-sided case, $f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$ error term $\sim O(\epsilon^2)$
- One-sided case, $f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}$ error $\sim O(\epsilon)$

Goal: Numerically verify implementation of derivative of fn is correct
and to check if there is a bug in backprop implementation

Gradient Checking

Gradient check for a NN:

- Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape into a big vector θ .
 \downarrow $J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$
- Take $d\theta^{[1]}, d\theta^{[2]}, \dots, d\theta^{[L]}$ and reshape into a big vector $d\theta$.
 \downarrow \checkmark \downarrow concatenate

Is $d\theta$ the gradient of $J(\theta)$?

Gradient Checking (Grad Check)

J is a fn of giant parameter θ $J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots)$

To implement this:

for each i : # each component of θ

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

$$\approx d\theta^{[i]} = \frac{\partial J}{\partial \theta_i} \quad | \quad \text{we get 2 vectors} \approx d\theta$$

Check $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta\|_2} \approx \epsilon = 10^{-7} \approx 10^{-3}, 10^{-5} \approx \text{great}$ (check bugs)

Gradient checking implementation notes

- Don't use in training - only to debug $d\theta_{approx}[i]$ vs $d\theta$
- If algorithm fails grad checking, look at the components and try to ~~debug~~ identify the bug.
 $db^{(l)}, dw^{(l)}$
- Remember regularization.

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

$d\theta$ = gradient of J wrt θ

- Doesn't work with dropout. keep-prob = 1.0
- ~~Random~~ Run at random initialization perhaps after some ^{again} training
 $w, b \approx 0$

Week-2 - Optimization Algorithms

Mini Batch Gradient Descent

Batch vs Mini batch Gradient descent:

- Vectorization allows you to efficiently compute on m ex.

$$X = [x^{(1)} \ x^{(2)} \ \dots \ x^{(m)}] \quad Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \\ (n_x, m) \quad (1, m)$$

- However if m is very large then it can be slow.

What if we have $m = 5,000,000$?

- While implementing gradient descent on the whole train set, you can process the ~~train~~ entire train set before you take one little step of gradient descent

- By doing this you can get a faster algorithm if you let gradient descent start to make some progress even before you finish processing the full train set at $m = 5 \times 10^6$

- Let's say we split the training set into smaller 'baby' train sets and these baby train sets are called minibatches.

- Each baby train set has 1000 ex.

$$X = [x^{(1)} | x^{(2)} | \dots | x^{(1000)} | x^{(1001)} | \dots | x^{(2000)} | \dots | \dots | x^{(m)}]$$

$x^{(1)} \in \mathbb{R}^{n_x, 1000}$ $x^{(2)} \in \mathbb{R}^{n_x, 1000}$ $x^{(5000)} \in \mathbb{R}^{n_x, 1000}$

↑ minibatch

$$Y = [y^{(1)} | y^{(2)} | \dots | y^{(1000)} | y^{(1001)} | \dots | y^{(2000)} | \dots | \dots | y^{(m)}]$$

~~We have~~ $y^{(1)} \in \mathbb{R}^{1, 1000}$ $y^{(2)} \in \mathbb{R}^{1, 1000}$ $y^{(5000)} \in \mathbb{R}^{1, 1000}$

~~minibatch~~: $x^{(t)}, y^{(t)}$

We have 5000 mini batches each having 1000 ex.

minibatch \hat{x}^t : $x^{(t)}, y^{(t)}$

$x^{(t)}: 1^{\text{st}} \text{ train ex}$
 $z^{(t)}: 1^{\text{st}} \text{ layer of train}$
 $x^{(t)}, y^{(t)}: \text{mini batches}$

- Batch gradient descent: This is the gradient descent where we process the entire training set all at the same time. We process the entire batch of train ex all the same time.
- Minibatch gradient descent: We process single mini batch X_t, Y_t at the same time rather than processing your entire ~~batch~~ of train set X, Y at the same time.

Minibatch Gradient descent:

repeat

For $t = 1, \dots, 5000$ i

Forward prop on $X^{(t)}$

$$Z^{(t)} = W^{(t)} X^{(t)} + b^{(t)}$$

$$A^{(t)} = g^{(t)}(Z^{(t)})$$

:

$$A^{(t)} = g^{(t)}(Z^{(t)})$$

Compute cost $J^{(t)}$

$$J^{(t)} = \frac{1}{1000} \sum_{i=1}^{1000} L(y^{(ci)}, a^{(ti)}) + \frac{1}{2 \cdot 1000} \sum \|W^{(t)}\|_F^2$$

inside the for loop implement

1 step of gradient descent using $X^{(t)}, Y^{(t)}$. ($m = 1000$)

Vectorized implementation
(1000 ex)

We are doing
the same
work
simpler

Backprop to compute gradients wrt $J^{(t)}$ using $(X^{(t)}, Y^{(t)})$

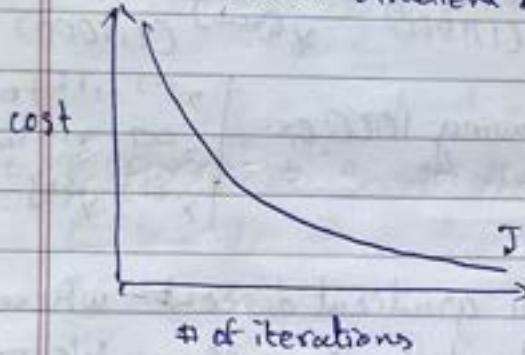
$$W^{(t)} := \frac{\partial J}{\partial W^{(t)}} = \frac{\partial L}{\partial W^{(t)}}, b^{(t)} := \frac{\partial J}{\partial b^{(t)}} = \frac{\partial L}{\partial b^{(t)}}$$

"1 epoch": 1 pass through the training set

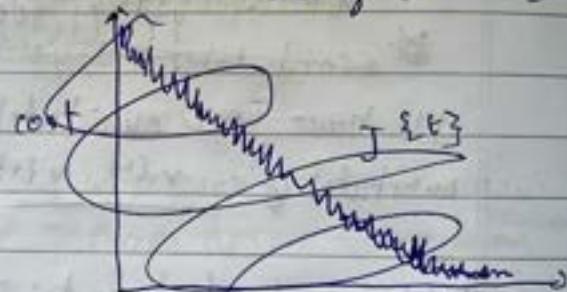
- With batch gradient descent single pass through train set allows you to take only 1 step in gradient descent.
- Mini batch gradient descent a single pass through the train, 1 epoch allows you to take 5000 gradient descent steps

Understanding MiniBatch Gradient Descent

Batch Gradient Descent



mini Batch gradient descent



On Batch gradient descent the cost fn J should decrease on every iteration

On mini batch gradient descent :

- When we plot cost fn $J^{(t)}$, it uses $X^{(t)}, y^{(t)}$, $J^{(t)}$ computed using $X^{(t)}, y^{(t)}$ then on every iteration we are training on a different train set, q. ~~one~~ mini batch.

The trend is downward but noisy, because maybe $X^{(t)}, y^{(t)}$ is relatively easy mini batch. So cost might be a bit lower, and $X^{(t)}, y^{(t)}$ mislabeled ex, cost will be higher. So we get these oscillations while running batch gradient descent.

Choosing your mini batch size :

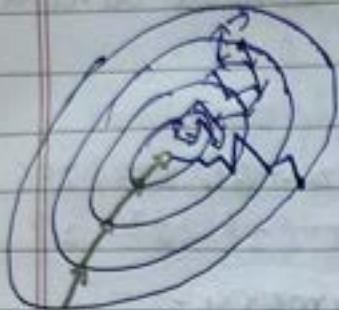
If mini batch size = m : Batch gradient descent $(X^{(t)}, y^{(t)}) = (X, y)$

If mini-batch size = 1 : Stochastic gradient descent every ex has its own minibatch $(X^{(t)}, y^{(t)}) \rightarrow (X^{(1)}, y^{(1)}) \dots (X^{(n)}, y^{(n)})$

In practice : Somewhere between 1 and m

too small too large

on the contours of
ex the cost function we are trying to minimise



- batch gradient descent
- stochastic gradient descent
- somewhere between [1, m]

- Batch gradient descent can start somewhere and can take relatively low noise, relatively long steps. It goes towards the minimum.
- Stochastic gradient descent starts at some point. On every iteration of gradient descent with a single ex, most of the times you'll hit a global ^{minimum}. Sometimes it'll head in wrong direction; if that ex points in ^{bad} direction.
- Stochastic gradient can be extremely noisy and on an avg it can ~~not~~ take you to a good or wrong directions. It doesn't converge, it oscillates and wanders around the region of minimum.

- Why do we use ~~batch~~ batch size $\rightarrow [1, m]$?

- If we use batch gradient descent, minibatch size = m. We are processing a huge training set on every iteration. Disadvantage is it will take too long per iteration, i.e. we have a big train set. If we have a small train set then batch gradient descent is fine.

- If we use stochastic gradient descent, we make progress after 1 ex. The noisiness can be reduced by just using a smaller learning rate. Disadvantage is we lose all the speed up from vectorization because we are processing a single ex at a time which is inefficient.

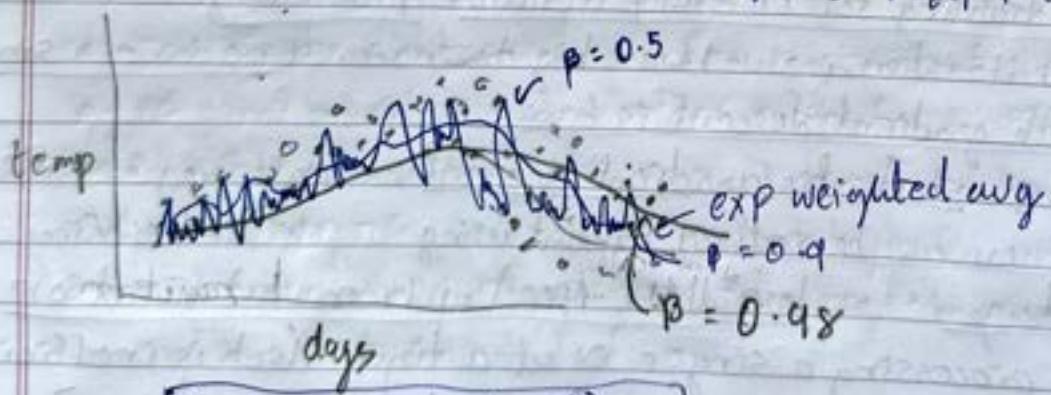
- In practice if we use ~~batch~~ mini batch size $\rightarrow [1, m]$, where minibatch size is not too big or too small it gives you faster learning. We get a lot of vectorization ex. minibatch size = 1000 we can vectorize across 1000 ex which is going to be much faster than processing 1 ex at a time. We can also make progress without needing to wait to process the entire train set. It doesn't guarantee to always head towards the minimum but it tends to head more consistently in direction of minimum then consequent descent. It doesn't always oscillate in a very small region; for this we can reduce learning rate.

Guidelines: to choosing $[1, m]$ for minibatch set:

- If train set is small, use batch gradient descent ($m \leq 2000$)
- Bigger train set (~~Typical~~): typical minibatch size $\rightarrow 64, 128, 256, 512$
 $2^6, 2^7, 2^8, 2^9$
- Make sure minibatch fit in CPU/GPU memory.
 $x^{(i)}, y^{(i)}$

Exponentially weighted Averages or Exponential weighted moving Avg
ex Temp in London

$$\begin{aligned} \theta_1 &= 40^\circ\text{F} \quad 4^\circ\text{C} & V_0 &= 0 \\ \theta_2 &= 49^\circ\text{F} \quad 9^\circ\text{C} & V_1 &= 0.9V_0 + 0.1(\theta_1) \\ \theta_3 &= 45^\circ\text{F} & V_2 &= 0.9V_1 + 0.1(\theta_2) \\ \vdots & \vdots & V_3 &= 0.9V_2 + 0.1(\theta_3) \\ \theta_{180} &= 60^\circ\text{F} \quad 15^\circ\text{C} & \vdots & \\ \theta_{181} &= 56^\circ\text{F} & V_t &= 0.9V_{t-1} + 0.1(\theta_t) \end{aligned}$$



$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$\beta = 0.9 \approx 10$ days temperature

We can think of V_t as averaging over $\frac{1}{1-\beta}$ days

$$\frac{1}{1-\beta} = 10$$

$$\beta = 0.98$$

$$\frac{1}{1-\beta} = 50$$

- As β is increased we get a smoother plot since we are averaging more days of temperature. Curve shifted to right

Latency: time delay bet moment an ilp is provided to a model and model produces output

classmate

Date _____
Page _____

- By averaging over a larger window the exp weighted avg formula adapts more slowly when temperature changes so there is a bit more latency
- When $\beta = 0.98$, it gives a lot of weight to the previous value and smaller weight 0.02 to the current value.
- When $\beta = 0.5$, $1/(1-0.5) = 2$ days
it's like averaging over just 2 days temp. ~~like~~ averaging over a shorter window. It's more noisy and more susceptible to outliers. But it adapts more quickly to temp changes

Understanding exponentially weighted avg

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t \quad \beta = 0.9$$

~~$V_{100} = \theta_{100}$~~

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

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

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

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

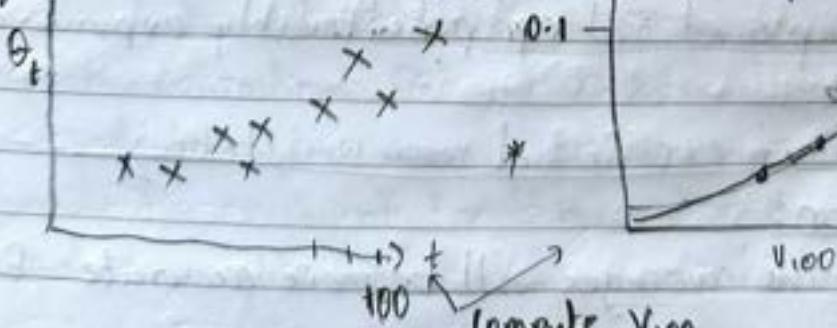
$$\nwarrow V_{99} = 0.9 V_{98} + 0.1 \theta_{99}$$

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

$$\nwarrow 0.1 \theta_{98} + 0.9 V_{97}$$

$$= 0.1 \theta_{100} + 0.1(0.9) \theta_{99} + 0.1(0.9)^2 \theta_{98} + 0.1(0.9)^3 \theta_{97} + 0.1(0.9)^4 \theta_{96} \dots$$

Visualizing this



compute V_{100}
we take elementwise
product these 2 and
sum it

How many days are we going over?

$$0.9^{10} \approx 0.33 \approx 1/e$$

10 days

$$(1 - e^{-\frac{t}{\tau}})^N = \frac{1}{e}$$

Implementing exp weighted avg :

$$V_0 = 0$$

$$V_1 = \beta V_0 + (1-\beta) \theta_1$$

$$V_2 = \beta V_1 + (1-\beta) \theta_2$$

$$V_3 = \beta V_2 + (1-\beta) \theta_3$$

 \vdots

$$V_0 := 0$$

$$V_1 := \beta V_0 + (1-\beta) \theta_1$$

$$V_2 := \beta V_1 + (1-\beta) \theta_2$$

 \vdots

$$V_\theta = 0$$

Repeat {

Get next θ_t

$$V_\theta := \beta V_\theta + (1-\beta) \theta_t$$

}

Advantages :-

→ takes very little memory, you just need to keep just 1 row no in computer to compute expwa memory and override it with the formula - $V_\theta := \beta V_\theta + (1-\beta) \theta_t$ based on the latest values.

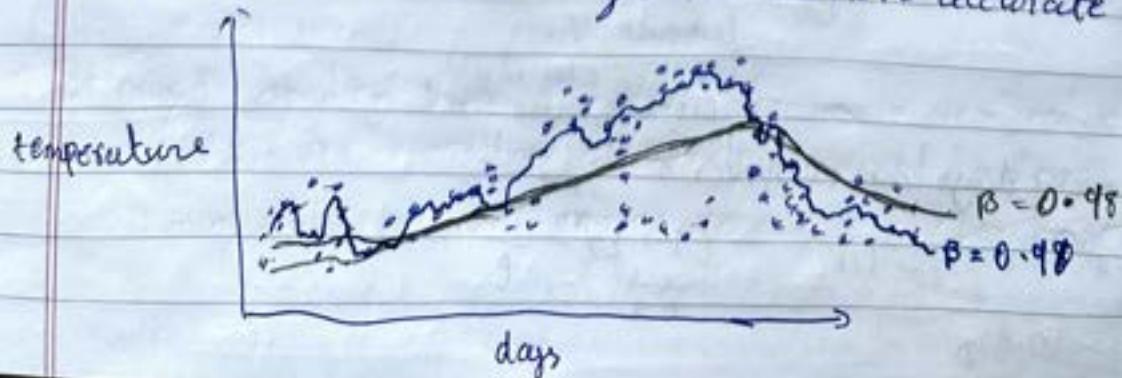
Disadvantage :-

→ Not really the best or most accurate way compute avg. ex. if we compute the moving window for the last 10 days, we divide by 10, that usually gives better estimate. But the disadvantage is of explicitly keeping temp around and sum of the last 10 days it requires more memory, this more complicated and computationally expensive.

Bias correction in exponential weighted avg

Bias Correction:

→ Computations of the averages will be more accurate - Bias Correc



$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

? When we use $\beta = 0.98$ we won't actually get the black curve we get blue color!

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$$V_0 = 0$$

$$V_1 = 0.98V_0 + 0.02\theta_1$$

$$V_2 = 0.98V_1 + 0.02\theta_2$$

$$= 0.98(0.02)\theta_1 + 0.02\theta_2$$

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

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

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

$$\frac{V_2}{0.0396} = \frac{0.0196\theta_1 + 0.02\theta_2}{0.0396}$$

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

assuming θ_1, θ_2

are tve no's, V_2 → Assuming θ_1, θ_2 are tve no's, V_2 is not a good estimate because is not a very good its less than θ_1 and θ_2 . There is a way to modify the estimate which estimate since it'll makes it much better and accurate shown on the right. be less than θ_1, θ_2 , → This becomes the weighted avg and removes the bias.

→ As t becomes large, β^t approaches 0; ~~which makes~~ →

→ When t is large enough bias correction makes no difference. ~~both~~ this line and this line overlap.

→ During the initial phase of learning when you're warming up your estimates, bias correction helps you obtain a better estimate of your temp.

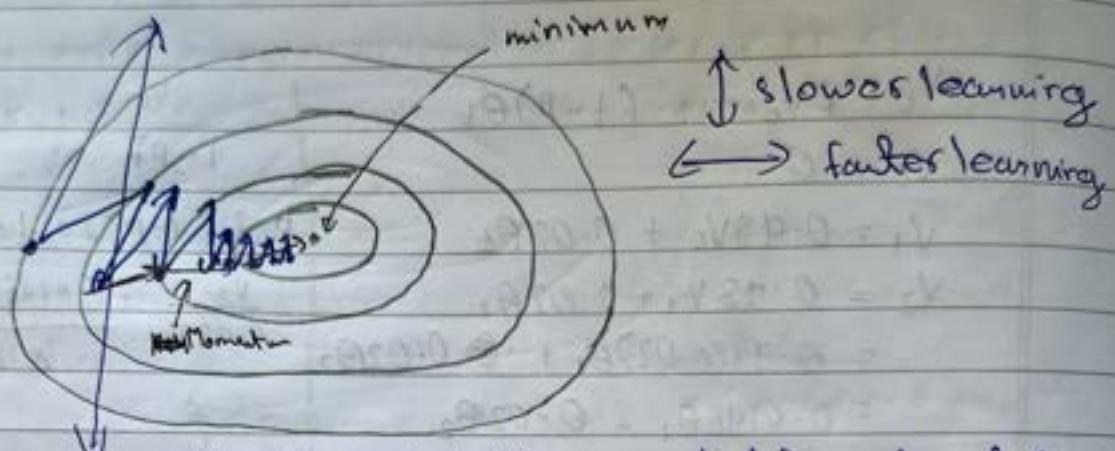
Gradient Descent with Momentum

- Momentum / Gradient Descent with momentum works faster than standard gradient descent -
- Compute the exponential weighted avg of the gradients and use that gradient to update weights instead.

Gradient Descent example:

3:50

ex. we are optimizing a cost fn that has contours.



- if we take an iteration of batch or minibatch gradient descent
- The gradient descent takes a lot of steps and slowly oscillate towards the minimum.
- Up and down oscillations prevent you from using a large learning rate
 - It might overshoot and end up diverging. So to prevent oscillations it forces you to use not ~~too~~ large learning rate.
- Another way to view this problem is that on vertical axis you want a bit small because we don't want oscillations and on horizontal axis we want faster learning

Momentum:

On iteration t:

Compute dW, db on current minibatch

$$Vdw = \beta Vdw + (1-\beta) dw$$

$$Vdb = \beta Vdb + (1-\beta) db$$

$$Vw = \beta Vw + (1-\beta) w$$

$$w := w - \alpha Vdw, b := b - \alpha Vdb$$

- It ^{we} avg the gradient → the oscillations in vertical direction tends to 0. In the vertical direction, you want to slow things down this will avg out +ve and -ve no's, so avg is closer to 0.

- On the horizontal direction, the derivatives are pointing to the right of the horizontal direction, so the avg in horizontal direction will be pretty big.
- In this algo, with a few ~~steps~~^{all} iterations we will find that gradient descent with momentum ends up eventually taking smaller oscillations in vertical direction, but more directed to just moving quickly in horizontal direction. This allows the algorithm to take a more straightforward path or to damp out the oscillations in the path to minimum.

Momentum intuition: (Optional)

if we try to minimise a bowl shaped fn



$$V_{dw} = \beta V_{dw} + (1-\beta) dW$$

friction

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

velocity acceleration

→ The derivative terms dW, db can be thought as providing acceleration to a ball your rolling downhill

$$w := w - \alpha V_{dw}$$

$$b := b - \alpha V_{db}$$



→ Imagine taking a ball and little ball is rolling downhill, it imparts acceleration to the ball

- Rather than gradient descent just taking every single step independently of all previous steps. Little ball can roll downhill and gain momentum it can accelerate down the bowl and gain momentum.

Implementation details

$$V_{dw} = 0, V_{db} = 0$$

On iteration t:

Compute dW, db on the current mini-batch

$$v_{dw} = \cancel{\beta v_{dw}} + \beta v_{dw} + (1-\beta) dW \quad | \quad V_{dw} = \beta v_{dw} + dW$$

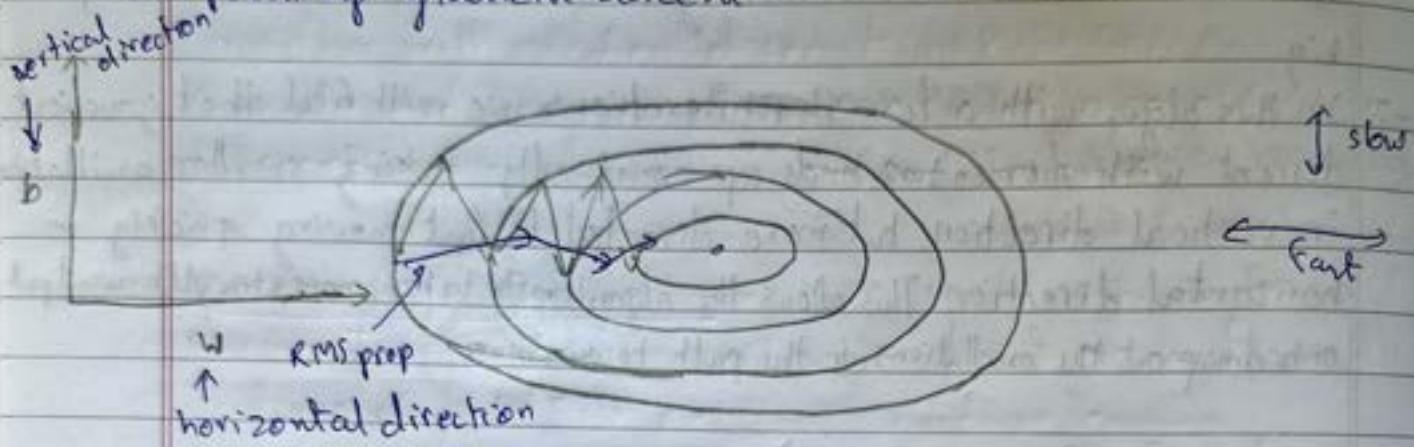
$$v_{db} = \cancel{\beta v_{db}} + \beta v_{db} + (1-\beta) d_b$$

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

Hyperparameters: α, β $\beta = 0.9$

RMSprop (Root Mean Square prop)

→ Speeds up gradient descent



- We want to slow down learning in the b direction or vertical
- ~~We want~~ Speed up learning in the horizontal direction
- RMS prop does this on iteration t

On iteration t :

Compute dW, db on current mini batch

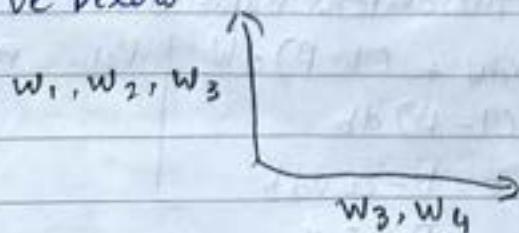
~~$Sdw = \beta_1 Sdw + (1 - \beta_1) dW^2$~~ element wise summing

$$Sdb = \beta_1 Sdb + (1 - \beta_1) db^2$$
 small

$$w := w - \alpha \frac{dW}{\sqrt{Sdw}} \quad b := b - \alpha \frac{db}{\sqrt{Sdb}}$$

- We want Sdw to be small and Sdb to be large
- The net effect is updates in vertical direction (dW) are divided by larger no that helps damp out oscillations. Updates in horizontal direction (db) are divided by a smaller no.

- Updates in horizontal and vertical direction keep going
- Use larger α to get faster learning without diverging in vertical direction
- In practice we are in a very high dimensional space of parameters observe below



- In practice dW and db are high dimensional vectors - But the idea is that in dimensions you're getting these oscillations, you end up computing a larger sum.
- A weighted away for these squares and derivatives we damp out the directions in which there are oscillations
- Hyperparameter: β_2
- What if \sqrt{dW} is too small?
 - ~~This could~~ things can blow up
 - Add a very small ϵ to denominator for ensuring numerical stability
- RMSprop is similar to momentum, has the effects of damping out oscillations in gradient descent, in minibatch gradient descent. Maybe allow you to use larger α and speeds up your learning speed of algo.

Adam Optimization algorithm

- Adam algo is a combination of RMSprop and momentum together implementation

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

on iteration t :

compute dW, db using current minibatch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dW, V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dW^2, S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

momentum β_1

RMSprop β_2

$$V_{dw}^{corrected} = V_{dw} / (1 - \beta_1^t), V_{db}^{corrected} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{corrected} = S_{dw} / (1 - \beta_2^t), S_{db}^{corrected} = S_{db} / (1 - \beta_2^t)$$

$$w := w - \alpha \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}} \quad b := b - \alpha \frac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected} + \epsilon}}$$

Hyperparameters choice

→ α : needs to be tuned $\rightarrow \Sigma: 10^{-3}$

→ $\beta_1: 0.9$ (dW)

→ $\beta_2: 0.999$ (dW^2)

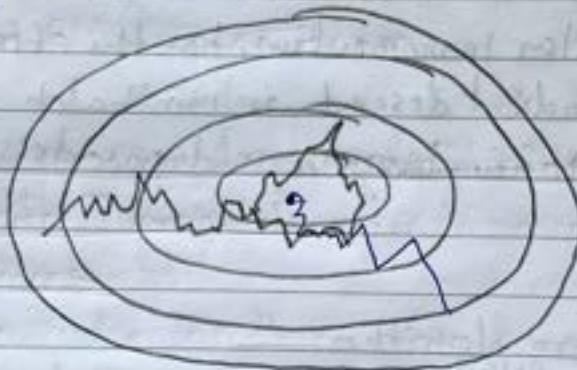
Adam: Adaptive Moment for estimation

- Combines the effect of gradient descent with momentum together with gradient descent with RMSprop.
- Commonly used learning algorithm, proven to be effective many different NN's of wide variety architectures

Learning rate Decay

ex we implement mini batch gradient descent with ^{small} mini-batches of 64, 128 ex -

→ Slowly reduce α



- As we iterate the steps are noisy, tending to the minimum and it wont converge because ~~we~~ are using a fixed α or there is noise in different batches
- We can ^{slowly} reduce α during the initial ~~step~~ of learning, so we can take much bigger steps but then as learning approaches convergence, then having slower learning rate allows you to take smaller steps.

Implementation

1 epoch = 1 pass through data

$$\alpha = \frac{1}{1 + \text{decay Rate} \times \text{epoch number}}$$

Epoch	α
1	0.1
2	0.067
3	0.05
4	0.04

$$\begin{array}{|c|c|c|} \hline & x^{(1)} & x^{(2)} & \dots \\ \hline \end{array} \rightarrow \begin{matrix} \text{epoch 1} \\ \text{epoch 2} \end{matrix}$$

α_0 : initial learning rate

$$\alpha_0 = 0.2$$

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



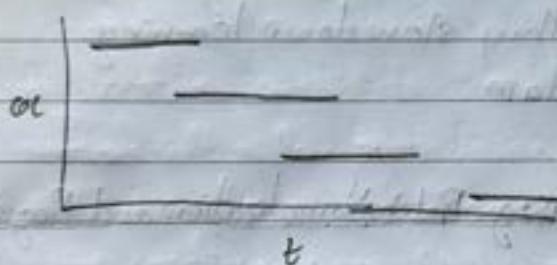
→ ~~Hyperparameters~~ Hyperparameters : α_0 , decay rate

Other learning rate decay methods :

$$\alpha = 0.95^{\text{epoch_num}} * \alpha_0 \leftarrow \text{exponential decay}$$

$$\alpha = \frac{k}{\sqrt{\text{epoch num}}} \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \alpha_0 \leftarrow \text{Polynomial decay}$$

discrete staircase



smooth decay

~~piecewise constant~~
~~discrete staircase~~

The Problem of local Optima

- Algorithms get stuck in local optimum rather than finding its way to global optimum.
- In NNs most points of zero gradients are not local optima.
- Most points of zero gradients in cost function are called saddle point
- A function of high dimensional space if gradient is zero then in each direction it can be a convex ^{function} light function or concave ^{function} light

convex
concave

U
∩
- We are more likely to run into saddle points in high dimensional spaces.
- Derivative is zero at saddle points

Problem of Plateau

Problem of Plateaus :

- Plateaus slowdown learning
- Region where the derivative is very close to 0 for long time

Key Points :

- You won't get stuck at local optima as long as your training a reasonably large ~~network~~ NN; since a lot of parameters and cost function J is defined over a relatively high dimensional space
- Plateaus ~~are~~ an problem and they slowdown learning.
Solv: Use RMSProp or Adam

Week - 3 : Hyperparameter tuning, Batch Normalization and Programming Framework

Hyperparameter Tuning

Tuning Process

Hyperparameters :

1. α (learning rate)
2. $P \approx 0.9$ (momentum term) 0.9 is usually a good estimate
3. B_1, B_2, ϵ (Adam optimization) → we ^{never} fine tune this. $B_1 = 0.9, B_2 = 0.999, \epsilon \approx 10^{-8}$ we these by default.
4. no of layers
5. no of hidden units
6. learning rate decay
7. minibatch size

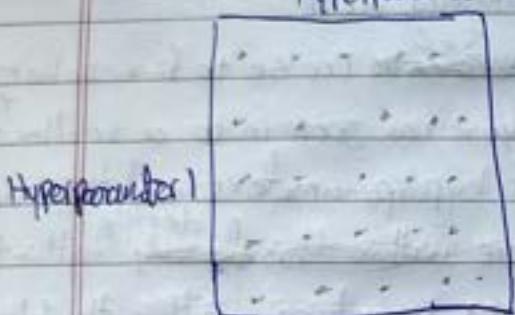
Most important hyperparameters to finetune : (order of importance)

1. α
2. P
3. no of hidden units
4. minibatch size
5. no of layers
6. learning rate decay

How do we select a set of values to explore?

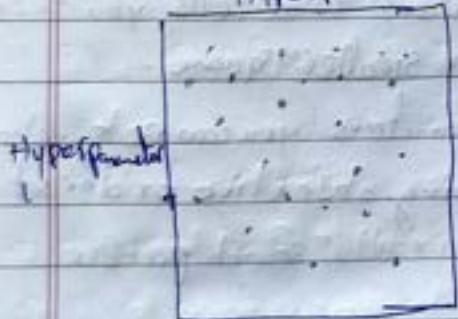
- In earlier ML algorithms, it's a common practice sample points in a grid
- Try random values; Don't use a grid
- In earlier ML algorithms, if we had 2 hyperparameters, hyperparameter-1 and hyperparameter-2. We sample the points on a grid and systematically explore the values.

Hyperparameter-2



→ We are placing in a 5×5 grid in this ex all 25 points and then we pick whichever hyperparam that works best. This works only if we have a small no of hyperparameters

- In DL choose points at random, we choose the hyperparameter-2



→ Choose 25 points

→ Try out hyperparameters on the randomly chosen ^{set of} points.

→ Reason why we do this is because its difficult to know in advance which hyperparameters are going to be the most important for the problem.

→ Some hyperparameters can be more important than others

- ex, hyperparameter-1 is α and hyperparameter-2 is ϵ . The choice of α matters more than ϵ . If we used sampling from grid, tried 5 diff values of α and we will find 5 values of ϵ gives you the same answer
- We ~~had~~ trained 25 models and only got ~~5~~ 5 trial values for the

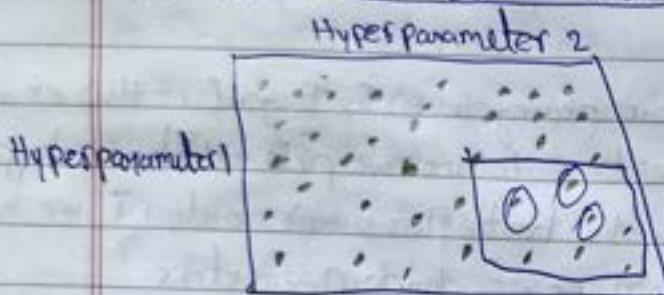
- If we tried sampling at random then we tried 25 distinct values of α and you are more likely able to find a value that really works well

- In practice we search over ^{right} more than 3 hyperparameters, its hard to know in advance which are more important for applicatn. So sampling instead

- Sampling in random rather than grid allows to richly explore the ~~set~~ possible values for the most important hyperparameters

Another common practice to find hyperparameters:

~~from coarse to fine~~



- Let's say we found these points working well

- In ~~Ir~~ Coarse to fine scheme you can zoom in to a smaller set of hyperparameters and then sample more ~~with this~~ densely within this

or maybe at random but focus more resources searching within the blue square if we think it's a ~~bad~~ ^{best} setting; the hyperparameter ~~is~~ maybe in this region

- We coarse the sample of the entire square that tells you to focus smaller square, we sample more densely in smaller square.
- By trying diff values of diff hyperparameters you can pick whatever value allows you to do best on your training objective or loss or any hyperparameter you're trying to optimise in the search process.

Using an appropriate scale to pick hyperparameters

Picking hyperparameters at random:

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

For no of hidden units we can pick a no at random to search for the particular hyperparameter.

- no of layers $l = 2 - 4$

→ Sampling x at random 2, 3, 4 might be reasonable

→ GridSearch to evaluate the values 2, 3, 4

- The above examples is where we can apply sampling at random over the range you're contemplating is reasonable. But this is not true for all hyperparameters

Appropriate scale for hyperparameters:

ex, we search for α :

$\alpha = 0.0001, \dots, 1$
we sample the values ^{uniformly} at random

$0.0001, \dots, 1$

90% of values → This means 90% of the resources are used to search we sample between 0.1 to 1.

- 10% of resources to search between 0.0001 to 0.1
- This isn't right.

- Use log scale to search for hyperparameters instead of linear scale.

$0.0001, 0.001, 0.01, 0.1, 1$

- Using sampling at random on logarithmic scale → We have more resources dedicated to searching between $0.0001 - 0.001, 0.001 - 0.01, 0.01 - 0.1, 0.1 - 1$

implementing in Python:

$$r = 4 * np.random.rand() \leftarrow r \in [-4, 0]$$

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

ex if we are sampling between 10^a and 10^b on log scale

$$10^a = 0.0001$$

~~$10^b = 1.0$~~

$$a = \log_{10} 0.0001$$

$$b = \log_{10} 1.0 = 0$$

$$a = -4$$

After this sample uniformly at random $r = [-4, 0]$

$$r = 10^r$$

Hyperparameters for exponentially weighted avg:

$$\alpha, \beta = 0.9, \dots, 0.999$$

We use \log_{10} values we avg over 1000 values

0.9

0.9

0.999

0.999

$$1 - \beta = 1 - 0.9 = 0.1$$

$$0.1 \dots 0.001$$

0.1

0.01

0.001

$$r \in [-3, -1]$$

10^{-1}

10^{-2}

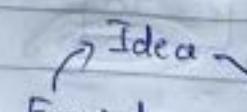
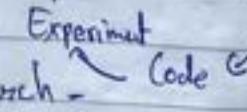
$$1 - \beta = 10^r$$

$$\beta = 1 - 10^r$$

- We spend an ~~some~~ amount of resources exploring from 0.9 - 0.99 and 0.99 - 0.999
- Why is it a bad idea to use linear scale?
- When $\beta \approx 1.0$, the sensitivity of your results changes even with very small changes to β . If β goes from 0.9 to 0.9005, doesn't change the results. If β goes from 0.999 to 0.9995 it will have a huge impact on what the algorithm is doing.
- $\frac{1}{1-\beta}$ formula is sensitive to small changes in β , when $\beta \approx 1$, what the whole sampling process does is it causes you to sample more densely in the region where $\beta \approx 1$ or ~~near~~ $1 - \beta \approx 0$.
- This will be more efficient on how distribute the samples to explore the space of possible outcomes more effectively.

Hyperparameters Tuning in Practice : Pandas vs Lawler:

Retest hyperparameters occasionally:

- There is a lot of cross fertilization over diff domains. ex ConvNet, ResNet in CV 
- People from diff application domains read research papers from other application domains to look for inspiration for cross fertilization 
- Intuitions do get stale. Re-evaluate occasionally

ex, we have a logistics problem;

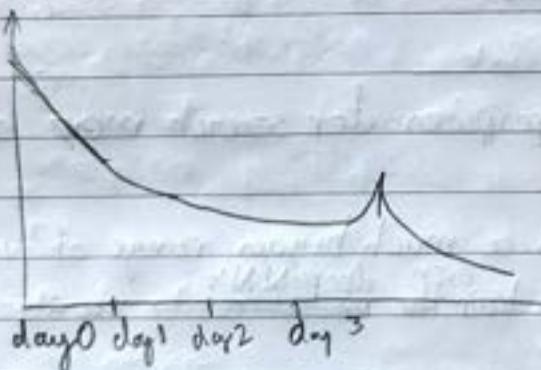
- we have a good set of hyperparameters, kept developing the algorithm or maybe ~~update~~ the data gradually changes over time because of that hyperparameters get stale

- Retest hyperparameters every several months.

How do we search for hyperparameters?

There are 2 approaches:

1. Babysitting one model : (Panda)
- Usually done when we have huge dataset but not a lot of computational resources to train a lot of models at the same time.
- Not a lot of CPU and GPU.
- You can afford training only one model or a very small no of models at a time
- We will babysit the model as its training.



Day 0: Random parameter initialization and start training. Gradually see the learning curve, cost fn J decrease gradually

Day 1: Learning is ~~good~~. We try increasing α a little bit and see how it goes

Day 2: goes well, may ~~increase~~ decrease the momentum term. or decrease α .

- Every day we will observe and nudge up or down the parameters
- Maybe one day, α is too big, we go back to the previous days model.
- We Babysit the model one day at a time as it trains over a few days or weeks.

2. Training many models parallelly (Cavion)
 - Ideal if we have enough computing resources
 - Simultaneously train models with varied hyperparameters
 - Each model produces a ~~diff~~ ^{distinct} learning curve
 - At the end of training, compare the performance and select the best model.

Choosing Between approaches

- Best approach depends on computational resources
- Cavion approach → is preferred when abundant computing resources are available to test multiple ~~models~~ hyperparameters at once.
- In scenarios with massive datasets and resource intensive models the puzzle approach is preferred where one invests significantly in fine-tuning a single model might be more feasible.

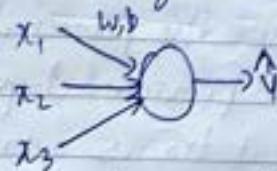
Batch Normalization

Normalizing Activations in a network

- Batch Normalization makes hyperparameter search more easier
- Makes NN more robust
- The choice of hyperparameters is a much bigger range of hyperparameters that work well and allows you to train ^{even} ~~deep~~ MN's easily quickly.

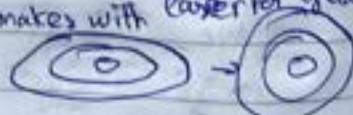
Normalizing inputs to speed up learning:

ex, In logistic regression we normalize input features to speed up learning

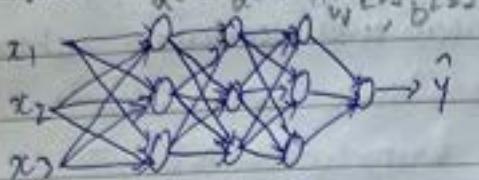


$$\left. \begin{array}{l} \mu = \frac{1}{m} \sum x^{(i)} \\ x = x - \mu \\ \sigma^2 = \frac{1}{m} \sum (x^{(i)})^2 \end{array} \right\} \rightarrow \text{This makes contours of gradient descent from elongated to more round shapes with easier gradient optimization}$$

This works in terms of normalizing $x = x / \sigma$
the i/p features values into NN after regression



what if we have a Deep NN?



- we have i/p features and activations
~~so it we want to normalize o~~

Can we normalize $a^{(2)}$ so as to train $w^{(3)}, b^{(3)}$ faster?

$\rightarrow a^{(2)}$ is i/p to next layer & it affects training of $w^{(3)}, b^{(3)}$

\rightarrow Batch Norm will normalize the values of ~~$a^{(2)}$~~ $z^{(2)}$ not $a^{(2)}$.

\rightarrow In practice normalizing $z^{(2)}$ is done more often.

Implementing Batch Norm :

Given some intermediate values in NN $z^{(1)} \dots z^{(m)}$

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

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

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

$\leftarrow z \text{ with mean } 0 \text{ and std dev } 1$

IF

$$z = \sqrt{\sigma^2 + \epsilon}$$

$$\beta = \mu$$

$$\text{then } \tilde{z}^{(i)} = z^{(i)}$$

tilde

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

learnable parameters of model

\rightarrow we use this for hidden units
 (we don't want hidden units to have

0 mean and 1 std dev because hidden
 units should have different dist)

Update param: γ, β

We use $\tilde{z}^{(l+1)}$ instead of $z^{(l+1)}$

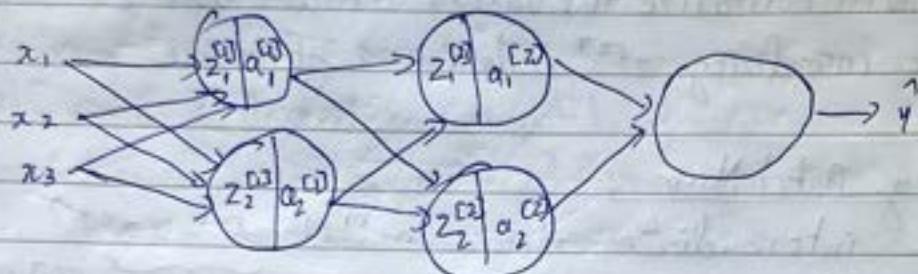
process not only to i/p layer

- BatchNorm applies this type of normalization to normalize mean but even some values deep in some hidden layer in NN. It applies and variance of some hidden units values $z^{(i)}$. This normalization to normalize mean and variance of some hidden units
- Difference bet training i/p ~~and~~ hidden unit $z^{(i)}$ is we don't want our values to be forced to have mean 0, var 1. ex if we have sigmoid we don't want the values to be clustered there. You might want them to have larger variance or mean that's different from 0 in order to take advantage of non-linearity of sigmoid rather than having all values in linear regime

Parameters
 γ and β make sure your $z^{(l)}$ values have the range you want. If your hidden units have standardized mean and variance, where they are controlled by γ and β .

Fitting Batch Norm into a NN

Adding BatchNorm to a network:



$$X \xrightarrow{w^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow{\text{BatchNorm(BN)}} \tilde{Z}^{[1]} \xrightarrow{\gamma^{[1]}} \alpha^{[1]} = g^{[1]}(\tilde{Z}^{[1]}) \xrightarrow{w^{[2]}, b^{[2]}} Z^{[2]} \xrightarrow{\text{BN}} \tilde{Z}^{[2]}$$

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$ } Update parameters using optimization algo Adam, RMSprop etc.
 $\gamma^{[1]}, \beta^{[1]}, \gamma^{[2]}, \beta^{[2]}, \dots, \gamma^{[L]}, \beta^{[L]}$

We can implement in tensorflow like,
`tf.nn.batch_normalization`

Working with mini batches:

$$X^{[1]} \xrightarrow{w^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow{\text{BN}} \tilde{Z}^{[1]} \xrightarrow{g^{[1]}(\tilde{Z}^{[1]}) = \alpha^{[1]} w^{[2]}, b^{[2]}} Z^{[2]} \xrightarrow{\dots}$$

$$X^{[2]} \xrightarrow{\dots} Z^{[1]} \xrightarrow{\text{BN}} \tilde{Z}^{[1]} \xrightarrow{\dots}$$

Parameters: $w^{[1]}, \gamma^{[1]}, \beta^{[1]}, \gamma^{[2]}$
 $(n^{[1]}, 1) \quad (n^{[1]}, 1) \quad (n^{[1]}, 1)$
 $z^{[1]}$
 $(n^{[1]}, 1)$

$$\begin{aligned} Z^{[1]} &= w^{[1]} a^{[1-1]} + b^{[1]} x^{[1]} \\ Z^{[2]} &= w^{[2]} a^{[1-1]} \\ \tilde{Z}^{[1]} &= \gamma^{[1]} Z^{[1]} + \beta^{[1]} \end{aligned}$$

Implementing gradient descent :

for $t = 1 \dots \text{num MiniBatches}$

Compute forward prop on $X^{(t)}$

In each hidden layer we BN to replace $\tilde{z}^{(l)}$ with $\hat{z}^{(l)}$

Use backprop to compute $dw^{(l)}, db^{(l)}, d\beta^{(l)}, dy^{(l)}$

Update parameters $\left. \begin{array}{l} w^{(l)} := w^{(l)} - \alpha dw^{(l)} \\ \beta^{(l)} := \beta^{(l)} - \alpha d\beta^{(l)} \\ y^{(l)} := y^{(l)} - \alpha dy^{(l)} \end{array} \right\}$

Works with RMSProp, Momentum, Adam

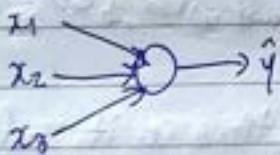
Why does BatchNorm work?

→ Batch Norm normalizes all ^{input} features x to take a similar range of values to speed up learning. It does this for hidden units and input features.

Learning on shifting input distribution :

→ It makes weights ^{that are} deeper in NN more robust to changes in weights to earlier layers

ex. we train a NN on cat detection (Cat $y=1$, Non-Cat $y=0$)



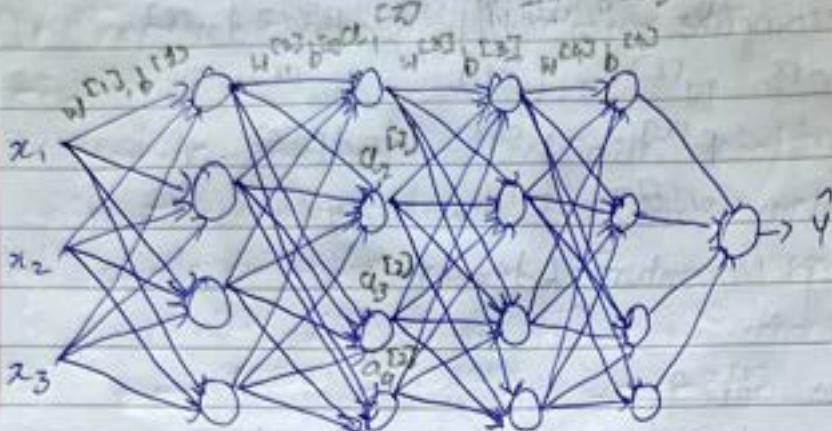
We train the CLF on images of black cat and if test on colored cat the CLF might not do very well. ex it our train set was like this

~~xx -ve ex~~ and if we tried to generalize this to dataset where ~~xx +ve ex~~ trained on data on the left then the model, won't do well on the data on the right even though they are same function that works well

- The idea of data distribution changing is called covariate shift. If you learn $X \rightarrow Y$ mapping, if the dist of X changes then we need to retrain the learning algo.
- This is true when the mapping of $X \rightarrow Y$ doesn't change

- ~~Return the function~~ The need to retain your function becomes worse if mapping $X \rightarrow Y$ changes

How does covariate shift problem apply to NN's?

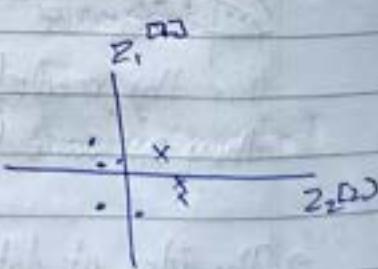
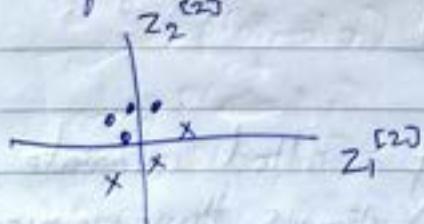


From the perspective of 3rd layer,

- The features are $a_1^{(2)}, a_2^{(2)}, a_3^{(2)}, a_4^{(2)}$
- The job of the 3rd hidden layer is to take the ip features and find a way to map them to \hat{y}
- Gradient descent learns parameters $w^{(3)}, b^{(3)}, w^{(4)}, b^{(4)}$ so it does a good job mapping from $X \rightarrow \hat{y}$
- The net also adapts these values $w^{(1)}, b^{(1)}$ and $w^{(2)}, b^{(2)}$ so the values at $a_1^{(2)}, \dots$ change
- The hidden unit values keep changing all the time so it suffers from the problem of covariate shift.

- BatchNorm reduces the amount of distribution these hidden unit values shift around

Plotting dist values



- The values for $z_1^{(2)}, z_2^{(2)}$ will change when parameters get updated in the earlier layers. Batch Norms ensures mean and variance of $z_1^{(2)}, z_2^{(2)}$ remains the same

- mean = 0, variance = 1 or whatever value governed by μ, σ^2
- This limits the amount that updating parameters in the earlier layers can affect the distribution of values that 3rd layer sees.
- Batch Norm reduces the problem i/p values changing, it makes the values become more stable, so later layers will have a more firm ground to stand on
- If i/p dist changes a little bit ~~is fixed~~, even as earlier layers keep learning the amounts that forces this later layers to adapt to ~~changes~~ earlier layers changes is reduced or if we will, it weakens the coupling between what the early layers parameters has to do and what later layer parameters has to do.
- It allows each layer of the net to learn by itself a little bit more independently of other layers and this speeds up learning in the whole net
- Batch Norm from the perspective of later layers in NN, the earlier layers don't shift as much because they are constrained to the same mean and variance and this makes the job of learning on later layers become easier

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 that minibatch. So similar to dropout it adds noise to each hidden layers activations
- This has a slight regularization effect.

Batch Norm at test time

Recall:

$$\begin{aligned} \bar{m} &= \frac{1}{m} \sum_i z^{(l)} \\ \sigma^2 &= \frac{1}{m} \sum_i (z^{(l)} - \bar{m})^2 \\ z_{\text{norm}}^{(l)} &= \frac{z^{(l)} - \bar{m}}{\sqrt{\sigma^2 + \epsilon}} \end{aligned}$$

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

At test time we may not have minibatches to process at the same time

\bar{m}, σ^2 : estimate using exp weighted avg

$X^{1|3}, X^{2|3}, X^{3|3}$

$\downarrow_{\bar{m}^{1|3|3}}, \downarrow_{\bar{m}^{2|3|3}}, \downarrow_{\bar{m}^{3|3|3}}$

$\sigma^{(l)}$

$\theta_1, \theta_2, \theta_3$

$\sigma^{2(1|3|3)}, \sigma^{2(2|3|3)}$

$\tilde{z} = \gamma z_{\text{norm}} + \beta$

$Z_{\text{norm}} = \frac{z - \bar{m}}{\sqrt{\sigma^2 + \epsilon}}$

- During training time μ and σ^2 are computed on the entire batch of 64 ex.
- In test time we process one ex at a time, we can do that by esting μ, σ^2 from train set.
- In theory we can run whole train set through the final net to get μ, σ^2 .
- In practice, you can implement exp weighted avg where you keep track of μ, σ^2 and ~~exp weight~~, sometimes we call this running avg to get rough estimate of μ, σ^2 and then use these to scale at testing time.

$$\mu_{\text{last}}, \mu_{\text{new}}, \mu_{\text{new}} \rightarrow \mu_{\text{new}}$$

$$\sigma^2_{\text{last}}, \sigma^2_{\text{new}}, \dots \rightarrow \sigma^2_{\text{new}}$$

Multi-class Classification

Softmax Regression

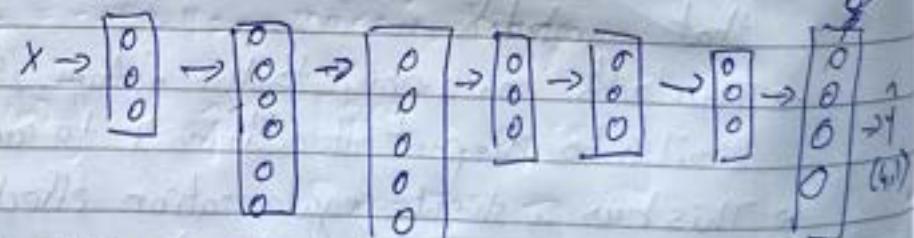
- Softmax Regression is a generalization of logistic Regression
- We try to recognise one of multiple classes instead of just 2 classes ex, Recognising cats, dogs, baby chicks labeled

0 others $C = \text{no of classes} = 4 (0 \dots 3)$

1 cat

2 dog

3 chick



→ No of units in o/p layer is 4 or equal to C .

→ Output layer tells to the probability of each of the classes o/p layer

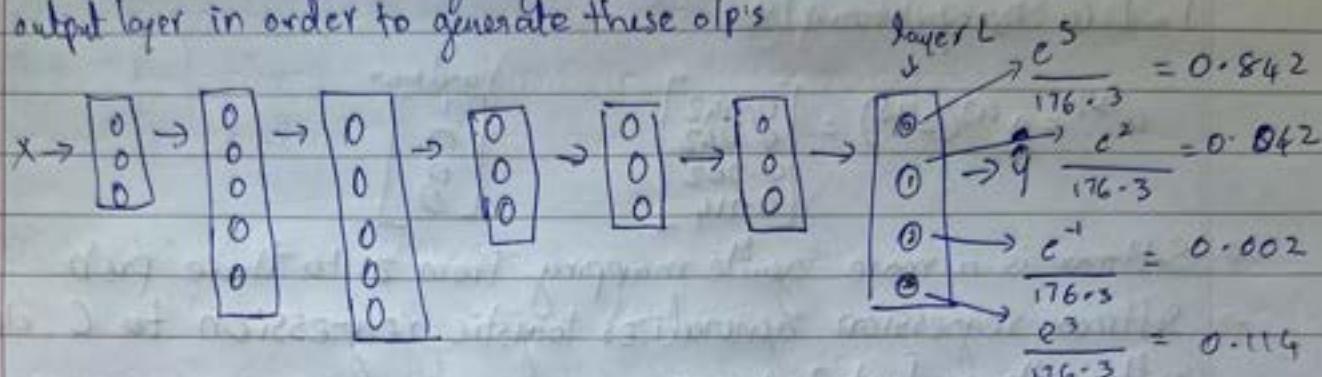
0	$P(\text{others} x)$
0	$P(\text{cat} x)$
0	$P(\text{dog} x)$
0	$P(\text{chick} x)$

The probabilities (y_i) should sum to 1.

~~2 The end~~

Softmax Layer :

The standard model for getting your net to do this is called Softmax layer and output layer in order to generate these o/p's



$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \quad (4, 1)$$

Activation Function

$$t = e^{z^{[L]}} \quad (4, 1)$$

$$\hat{y} = a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^4 t_i}, \quad a_i^{[L]} = \frac{t_i}{\sum_{j=1}^4 t_j}$$

$$a^{[L]} = g^{[L]}(z^{[L]}) \quad \hookrightarrow (4, 1)$$

(4, 1)

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \quad \sum t_i = 176.3$$

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

→ Unusual behavior of softmax activation layer is it needs to be normalized across different possible outputs and needs to take a vector put in output vector.

Training a Softmax Classifier

Understanding softmax :

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}, \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$a^{[L]} = g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

- The biggest element in the z vector is 5, it has the highest prob
- Softmax comes from Hardmax where it looks at the biggest elements in z and 1 and 0's everywhere else.

$$a^{(1)} = g^{(1)}(z^{(1)}) = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

hardmax
 $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

- Softmax is a more gentle mapping from z to these prob
- Softmax regression generalizes logistic regression to C classes rather than just 2 classes.
- If $C=2$, softmax reduces to logistic regression, $a^{(1)} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$. These probs need to sum to 1, we don't need to compute 2 of them, just one
- $a^{(1)} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$

Loss function :

Notice $y_1, y_3, y_4 = 0$ & $y_2 = 1$

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{cat } a^{(1)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$L(\hat{y}, y) = -\sum_{j=1}^4 y_j \log \hat{y}_j \quad J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_m L(\hat{y}^{(m)}, y^{(m)})$$

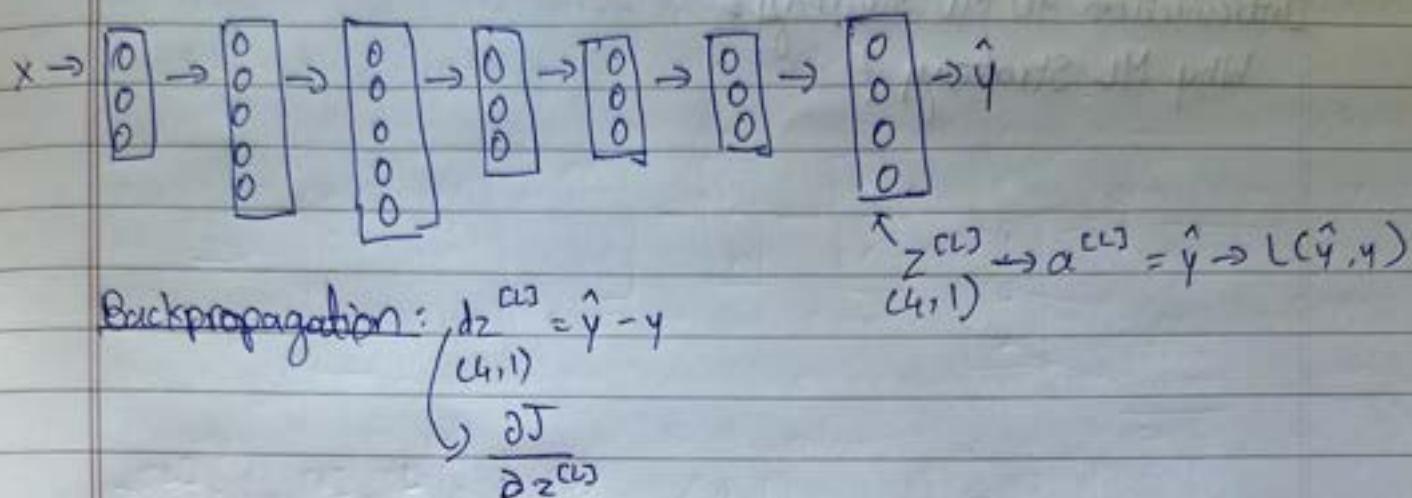
$$= -y_2 \log \hat{y}_2 = -\log \hat{y}_2 \quad \text{make } \hat{y}_2 \text{ bigger}$$

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

$$= \begin{bmatrix} 0 & 0 & 1 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & 0 \dots \\ 0 & 0 & 0 & \dots \end{bmatrix} \quad = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$(4, m)$

Gradient Descent with softmax:



Introduction to Programming Frameworks

Deep Learning Frameworks ↗

- Caffe / Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- Tensorflow
- Theano
- Torch

Choosing DL Frameworks:

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