

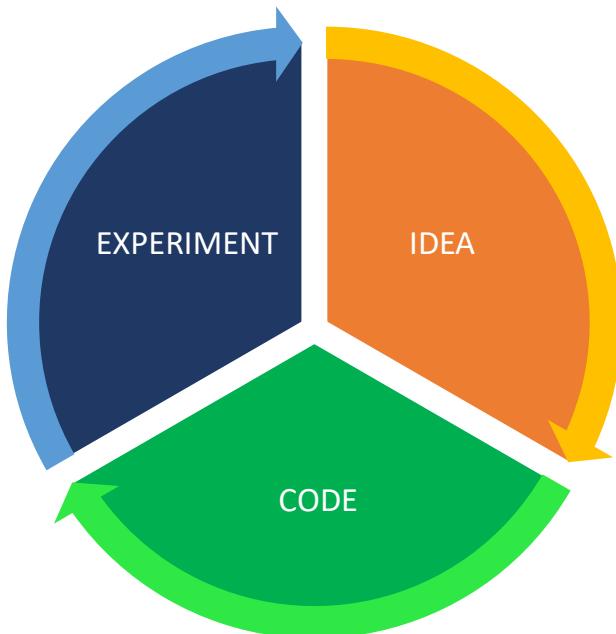
DEEP LEARNING SPECIALIZATION

COURSERA-DEEPMLEARNING.AI-ANDREW NG

These are notes taken from [Deep Learning Specialization](#) by deeplearning.ai in [Coursera](#). All the screenshots shown here are taken from the course videos. I would advise you to first watch the videos of the course in order to completely understand the notes.

COURSE 2

Improving Deep Neural Networks: Hyperparameter tuning, Regularization and Optimization



Hyperparameters

- **Learning rate α** → Decides the rate at which gradients the parameters change.

- **No of iterations** → Decides no of times gradient descent Is done to adjust the parameters.
- **No of hidden layers l** → Decides the number of variety of features.
- **No of hidden units $n^{[l]}$** → Decides the number of activation units.
- **Choice of activation function** → Decides the nature of features produced in the NN units.

Other than these common hyperparameters there are other parameters that we can tune by considering them as hyperparameters, which we would see as we go through the course.

All the hyperparameters are tuned to decide which ones are best to get the best results. It's very rare to have knowledge of the best value for all these parameters in advance.

So, what we do is we experiment with the code that we written using the particular parameter values. If we fail, we get an idea about what changes to make then we code again and we test again.

**In this course we would learn about how to approach these Idea
→Code →Experiment method to get the best model**

About this Course

This course will teach you the "magic" of getting deep learning to work well. Rather than the deep learning process being a black box, you will understand what drives performance, and be able to more systematically get good results. You will also learn TensorFlow.

After 3 weeks, you will:

1. Understand industry best-practices for building deep learning applications.
2. Be able to effectively use the common neural network "tricks", including initialization, L2 and dropout regularization, Batch normalization, gradient checking.
3. Be able to implement and apply a variety of optimization algorithms, such as mini-batch gradient descent, Momentum, RMSprop and Adam, and check for their convergence.
4. Understand new best-practices for the deep learning era of how to set up train/dev/test sets and analyse bias/variance
5. Be able to implement a neural network in TensorFlow.

This is the **second course** of the Deep Learning Specialization.

WEEK 1 - Setting Up Your Machine Learning Application

Train, Dev, Test sets (For small data size)



Usually when we create a machine learning model, we split our dataset into 3 parts, Train set which we use to train our model, Cross validation set which we use to tune our parameters and experiment with different models and choose which is the best and Test set which is used to finally test our model. And before big data came in, we used to have data in medium sizes like in 10,000 or 1000 etc. So, we used to give more percentage of data for testing and cross validation. But when we work with **Big data**, we have data in range of millions or even billions. So, we usually divide our data into portions like **99.5%** train, **4%** cross validation and **1%** test. So, it basically depends on your size of data that determines the ratio at which you divide the dataset into train, cross validation and test data.

Train, Dev, Test sets (For Big Data)

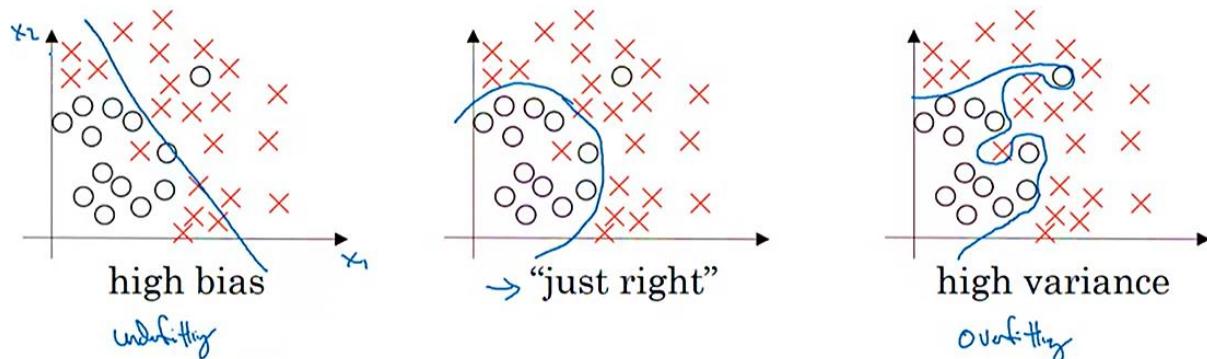


Mismatched Train, Test and Dev set

When people divide train/test and dev set they might tend to choose different sources for different sets in order to collect maximum data for training set. This is not good since images

collected from different sources might have different resolutions and might not give true accuracy when we use data from different source for test set. So, its always good to divide data collected from similar sources into train, test and dev sets.

Bias and Variance



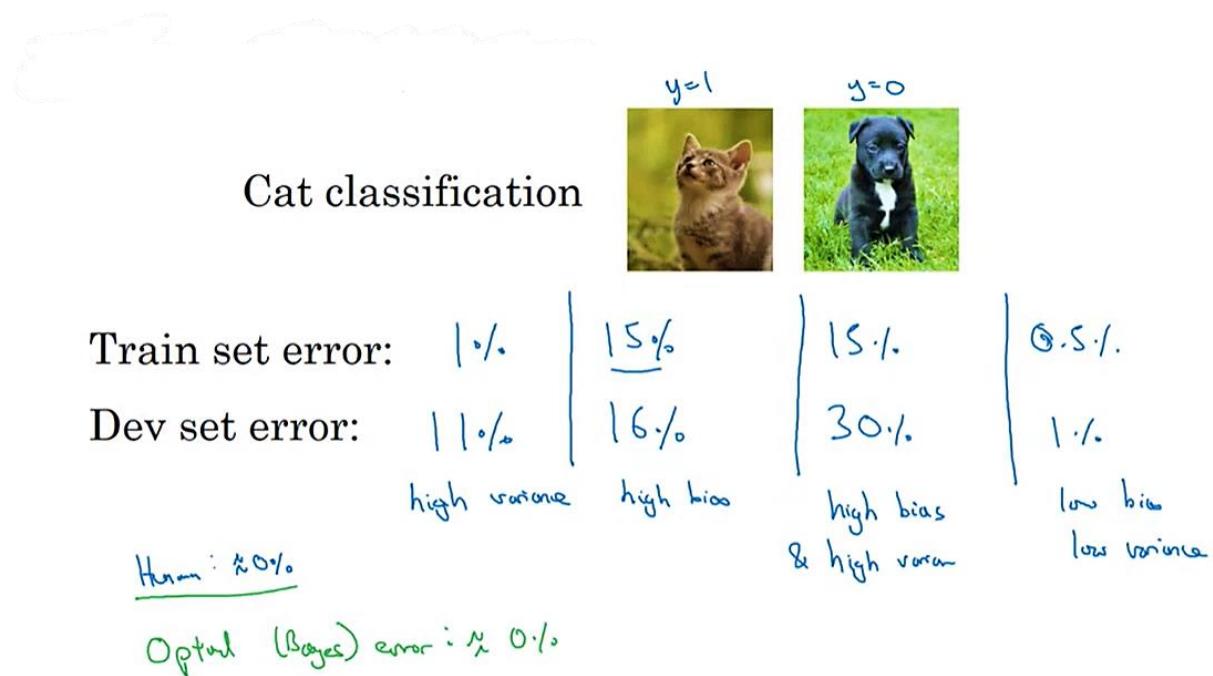
Assumptions

1. The base error is 0
2. The data is from same distribution

In the above data we can see that if we try to fit a logistic regression straight line onto the data it might not fit the data very well and we can say it has high bias as it might predict lot of the observations wrongly. The last figure shows fitting the data with that unusual curve which might fit the training data very well but not the test data. It might be able to predict all the values in the training set correctly but it might perform poorly on the test data. The middle figure shows the better fit to the data as it covers all the points very well and even though there are one or two mismatches it would perform better when we try it on test data.

- **High bias** → Fit the training data poorly and dev data too(**underfit**)

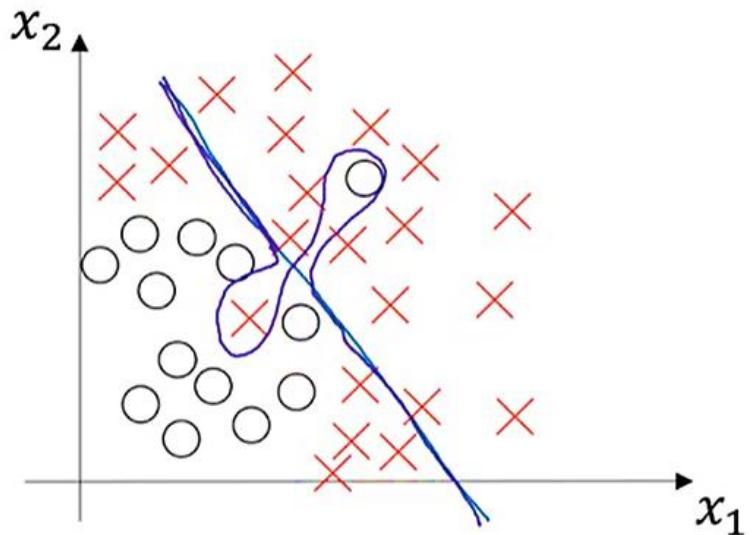
- **High variance** → **Overfit** the training data and not good fit on dev data



In this example we classify the images into cat or not cat. We can see that we got 1% train set error and 11% Dev set error. Which shows high variance as it is fitting the training data perfectly but not the Dev data. The other case where there is a 15% error for the training data and 16% error for the dev data it means that your model is not fitting even the training data well. It shows that the model is underfitting and ha bias.

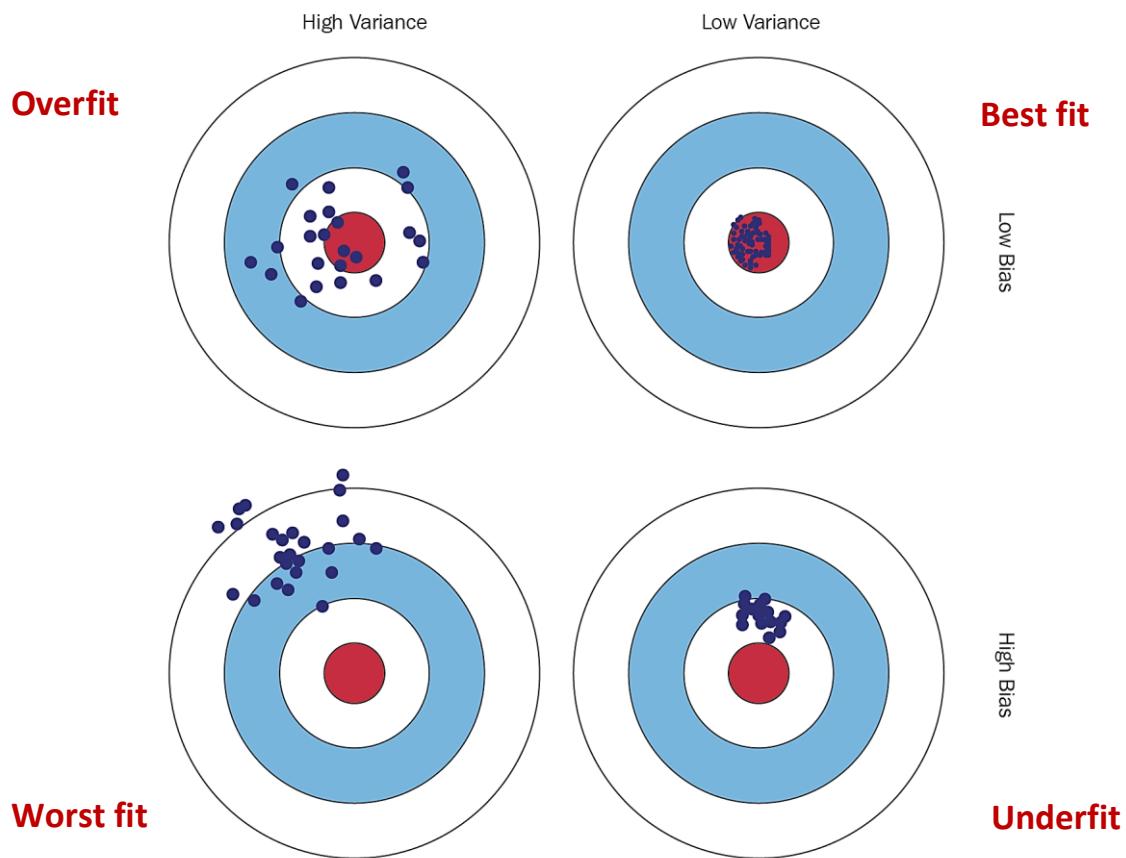
When it has 15%error on training data and 30% error on dev data it is both underfit and overfit. Let's see an example to better understand this.

High bias and high variance



We can see that the line is underfitting the data since it is linear almost in all portions but it must be a curve to fit it perfectly. At the same time, it takes a weird path in the middle and overfit the data in middle.

And finally, if your model has very low bias and variance then it is said to be good fit.



<http://scott.fortmann-roe.com/docs/BiasVariance.html>

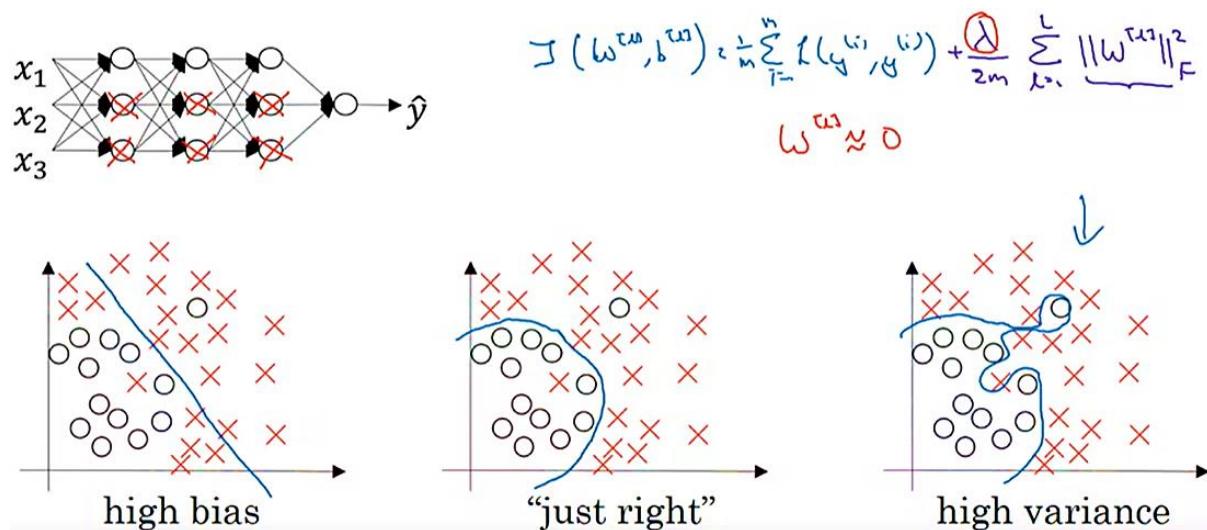
Ways to reduce high bias or variance

High Bias	High Variance
Bigger Neural Network	More Data
Train Longer	Regularization
Search for a better NN architecture	Search for a better NN architecture

If we have a model that is underfitting the training it with more layers of neural network would try to fit it better and might reduce bias. Also, we can try to train it bit longer to see its best fit. If we have high variance in our model, we can try to add more data to it so it won't overfit it, also we can try regularization to reduce overfitting. Using a bigger neural network won't overfit the data if we do regularization.

Bias variance trade-off is a term used to describe the increase in variance/bias as we decrease variance/bias.

Why Regularization reduces Overfitting?



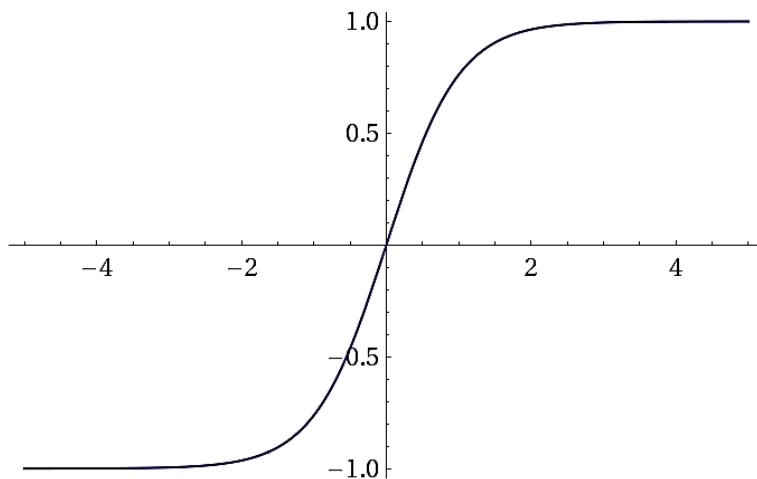
If we are including regularization term in our cost function our cost function looks like this:

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

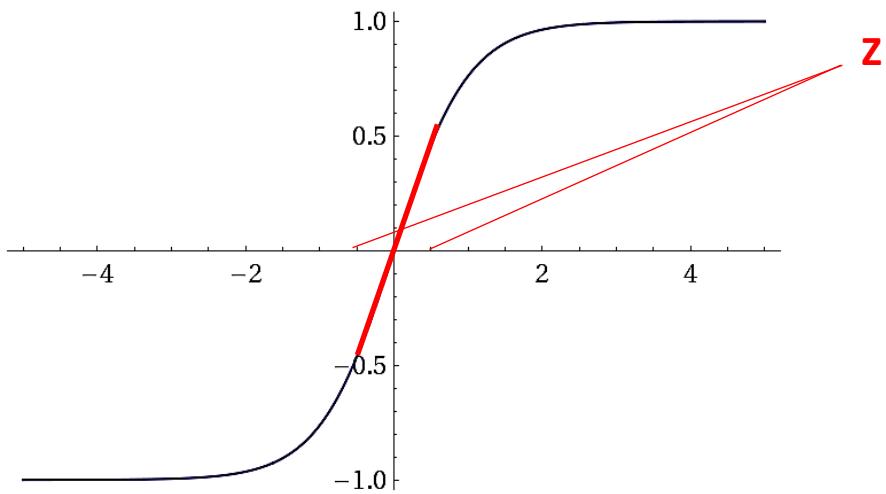
So, the parameter update would be:

$$\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha(d\mathbf{W}^{[l]} + \frac{\lambda}{m} \mathbf{W}^{[l]})$$

Due to the addition of the new regularization term if we increase λ the term $\mathbf{W}^{[l]}$ would be reduced. And some of the weights that contribute less to the model might get closer to zero. So effectively it might reduce the effect of few of the neural network nodes on the overall model and thereby reducing overfitting because most of the times overfitting might be caused when we use deeper neural networks and train our data. But we might need these deeper neural networks as it would give us additional features that we won't get from shallow networks.



If we are using the tanh function as activation function for any layers then we know as Z increases or decreases so much the functions gradient would become lesser to 0. But if we are regularizing the cost function then we would decrease \mathbf{W} thereby decreasing \mathbf{Z} as $\mathbf{Z} = \mathbf{WA} + \mathbf{b}$. Therefor for some neurons the activation function would behave as linear as the value of Z would be lying in a small region close to 0. Thereby reducing the formation of complex curves which overfit the data.



Regularization

Now let's see how to do regularization

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

$$\|\mathbf{w}^{[l]}\|^2 = \sum_{i=1}^{n^l} \sum_{j=1}^{n^{[l-1]}} (w_{i,j}^{[l]})^2$$

The norm $\|\mathbf{W}^{[l]}\|_F^2$ is called Forbenius norm.

$$\text{So } \mathbf{dW}^{[l]} = \frac{1}{m} \mathbf{dZ}^l \mathbf{A}^{[l-1]} + \frac{\lambda}{m} \mathbf{W}^{[l]}$$

$$\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha (\mathbf{dW}^{[l]} + \frac{\lambda}{m} \mathbf{W}^{[l]})$$

This is also called weight decay.

Here we can see that how regularization is controlling the weights further and thereby reducing the chance of overfitting. It doesn't allow the weights to decrease or increase too much.

Explanation

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

$$\|W^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (W_{ij}^{[l]})^2 = \text{np.sum(np.multiply}(W^{[l]}, W^{[l]}))$$

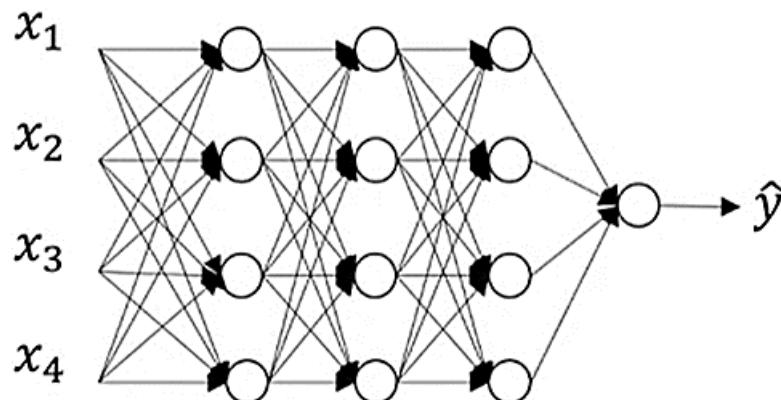
= For $i = 1 \text{ to } L$:

Sum $\left(\begin{bmatrix} W_{11}^{[l]} & W_{12}^{[l]} & \dots & W_{1n^{[l-1]}}^{[l]} \\ W_{21}^{[l]} & W_{22}^{[l]} & \dots & W_{2n^{[l-1]}}^{[l]} \\ W_{31}^{[l]} & W_{32}^{[l]} & \dots & W_{3n^{[l-1]}}^{[l]} \\ \vdots & \vdots & & \vdots \\ W_{n^{[l]}1}^{[l]} & W_{n^{[l]}2}^{[l]} & \dots & W_{n^{[l]}n^{[l-1]}}^{[l]} \end{bmatrix} \right)_{n^{[l]} \times n^{[l-1]}}$

$$\begin{aligned} dW^{[l]} &= \frac{dL}{dW^{[l]}} = \frac{dL}{dz^{[l]}} \times \frac{dz^{[l]}}{dW^{[l]}} + \frac{d}{dW^{[l]}} \left(\frac{\lambda}{2m} \sum_{i=1}^L \|W^{[i]}\|^2 \right) \\ &= dz^{[l]} \times A^{[l-1]} + \frac{\lambda}{m} \|W^{[l]}\| \end{aligned}$$

$$\left\{ \begin{array}{l} dW^{[l]} = \frac{dL}{dz^{[l]}} \times \frac{dz^{[l]}}{dW^{[l]}} \\ \frac{dz^{[l]}}{dW^{[l]}} = \frac{d}{dW^{[l]}} (W^{[l]} A^{[l-1]} + b^{[l]}) \\ = A^{[l-1]} \end{array} \right\} \text{For non regularized}$$

Dropout regularization



In dropout regularization we randomly eliminate some neurons from different layers in the network in order to get simpler network.

How is it done?

Let's say we are applying dropout to layer l .

```
Keep_prob = x
db = np.random.rand(a[l].shape[0], a[l].shape[1]) > Keep_prob
db would be matrix with random True, False values
a[l] = np.multiply(a[l], db)
a[l] /= keep_prob
```

How does it work?

Let's say $x = 0.5$, and $a^{[l]}$ is of size $3, 3$

Let $a^{[l]} =$

```
[[1, 2, 5],
 [4, 5, 8],
 [8, 9, 5]]
```

```
db = np.random.rand(a[l].shape[0], a[l].shape[1]) > Keep_prob
```

This would create matrix db as follows:

```
[[False False False]
 [ True  True  True]
 [ True  True False]]
```

```
a[l] = np.multiply(a[l], db)
```

After this $a^{[l]} =$

```
[[0, 0, 0],  
 [4, 5, 8],  
 [8, 9, 0]]
```

We can see that half the values are turn to zero. Now as we know this would mean that when we multiply $a^{[l]}$ with W and find Z it would be a smaller value than previous $a^{[l]}$ so in order to counter this loss we divide $a^{[l]}$ by the value we use to set probability keep_prob

```
a^{[l]}/= keep_prob
```

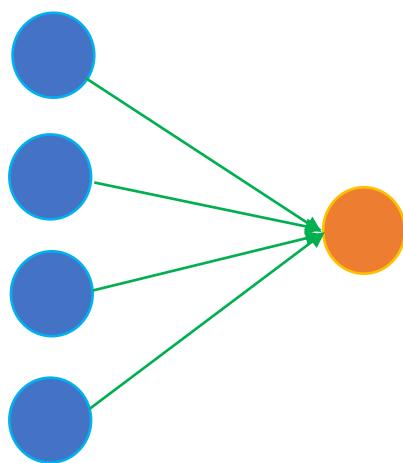
So $a^{[l]}$ would now look like:

```
[[ 0., 0., 0.],  
 [ 8., 10., 16.],  
 [16., 18., 0.]]
```

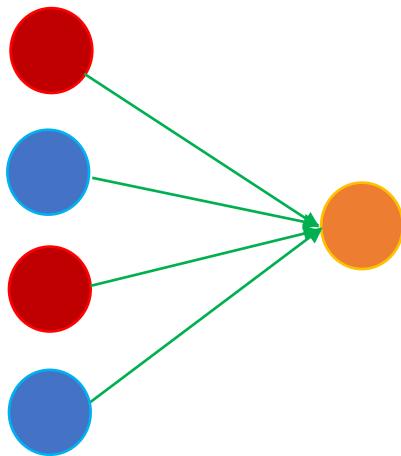
Which is scaled up to counter the loss

Why does it work?

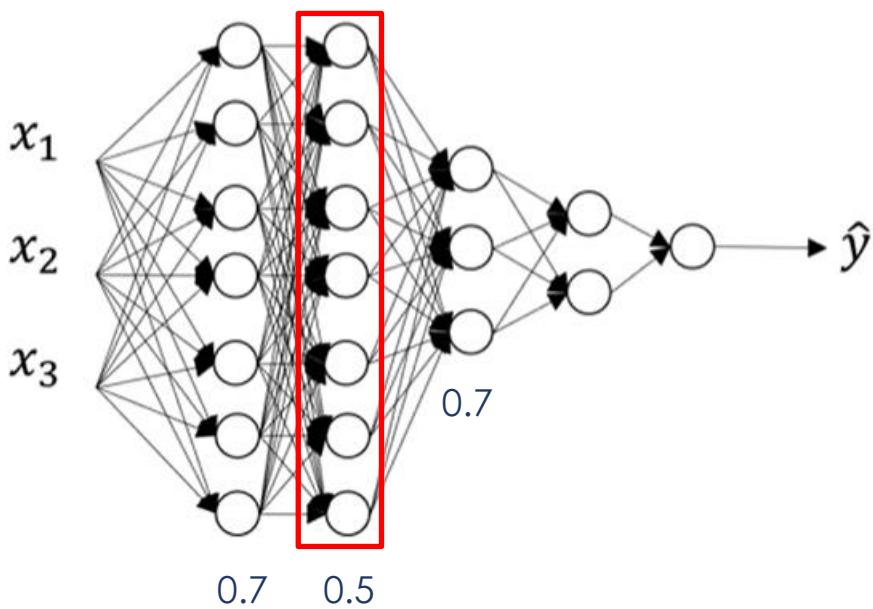
Consider a simple neural network



Let's say we apply dropout to eliminate few of these nodes



Now what happens is if we eliminate few neurons the weight assigned for each node would become smaller as it can't rely on any particular node more as any node can be eliminated. Therefore the overall norm of weights would be less. Therefore it reduces overfitting as it reduces weights distributed little bit.



If we consider this neural network, we can see that 2nd layer is layer with most networks and is more prone to overfitting, therefore we apply less value of keep_prob for that layer and more for others.

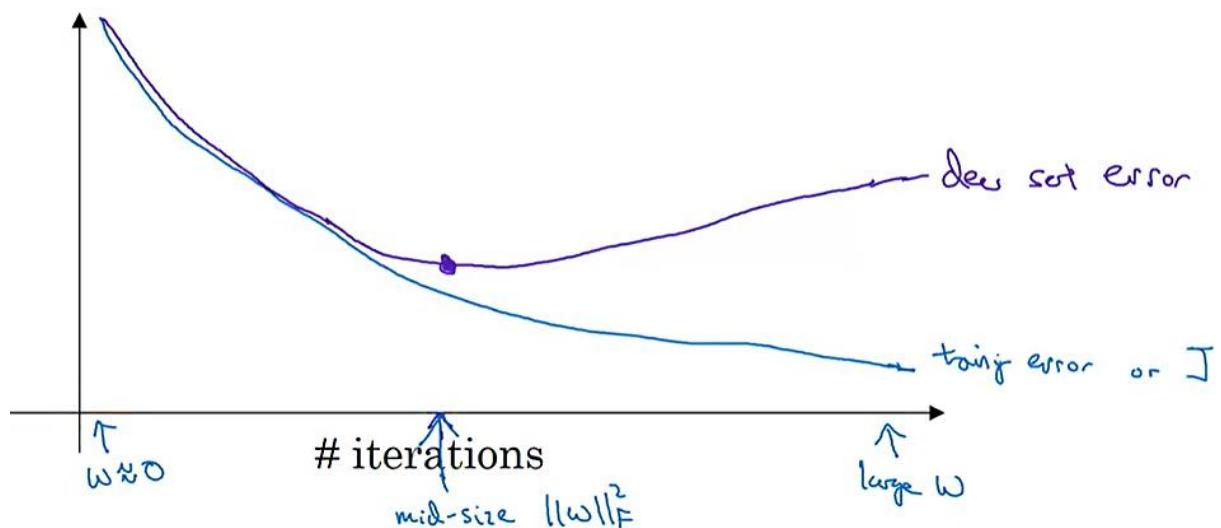
Other techniques to reduce overfitting

Data Augmentation



One way to reduce overfitting is to collect more data. But most of the times it's a costly process. So, one thing we can do is to take our data and make few changes like the orientation or dimension or maybe distort if our data is set of images and use it as additional data.

Early stopping



In early stopping we will be plotting graphs of cost function for dev set and training set. And after few iterations we would observe the dev set error stop reducing and going up. We would stop the training at a point where both training set and dev set are at their best low point. This will effectively reduce the weights as the weights would only keep increasing as the number of iterations increase and by stopping it in middle, we would reduce it from increasing very high and would do similar job as L2 regularization.

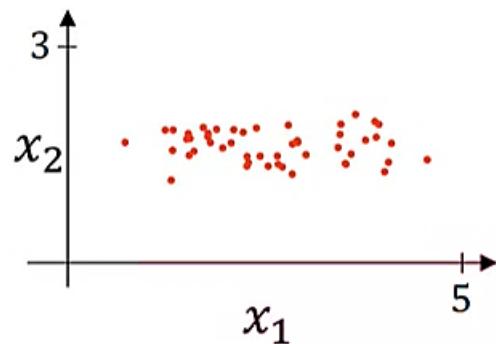
Normalizing inputs

What?

Normalizing is the process of making all the input features to a similar scale. That is to have similar value of mean and standard deviation for all the features.

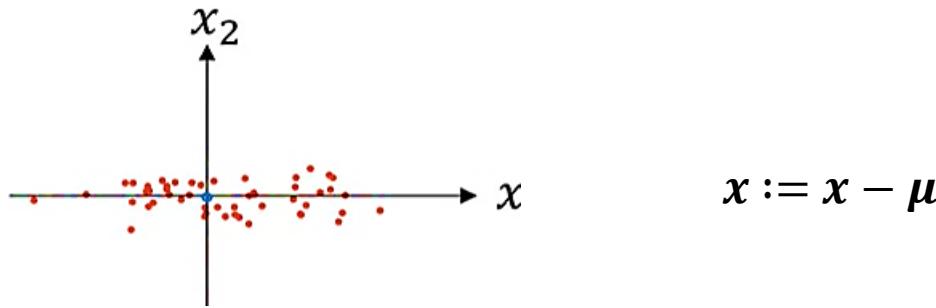
How?

Let's take an example of data with different scales for different features.

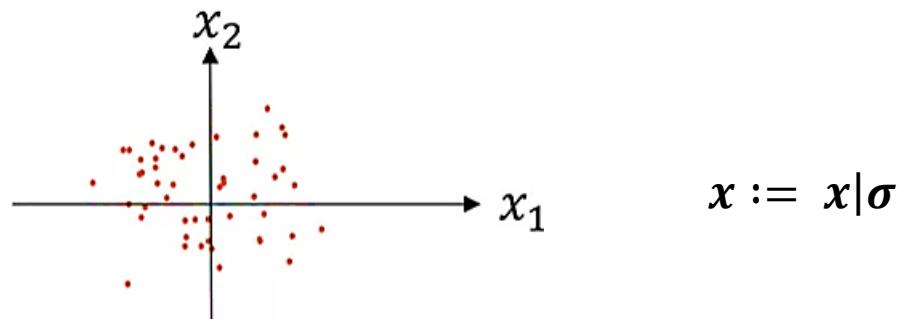


Here we can see that it doesn't have a common mean and the deviation for x_1 is larger than x_2 .

1st we would make the mean of data to 0.



2nd we would make the standard deviation equal by dividing with std.



Now we can see that our both features have same mean and variance.

Steps:

```

u = (1/m) * np.sum(x)
x = x-u
var = (1/m) * np.sum(np.power(x, 2))
std = np.sqrt(var)
x = x/std
    
```

Why?

Usually when we collect data for a particular purpose it might have different features with different scales. For example, weight of a person might have scale from 20-100 and his salary might have scale from 100-1000000. So, both these features are having different scales which might contribute more towards the activation functions final output than the other feature. So, our model might be predicting the results based on effectively fewer features than what's given.

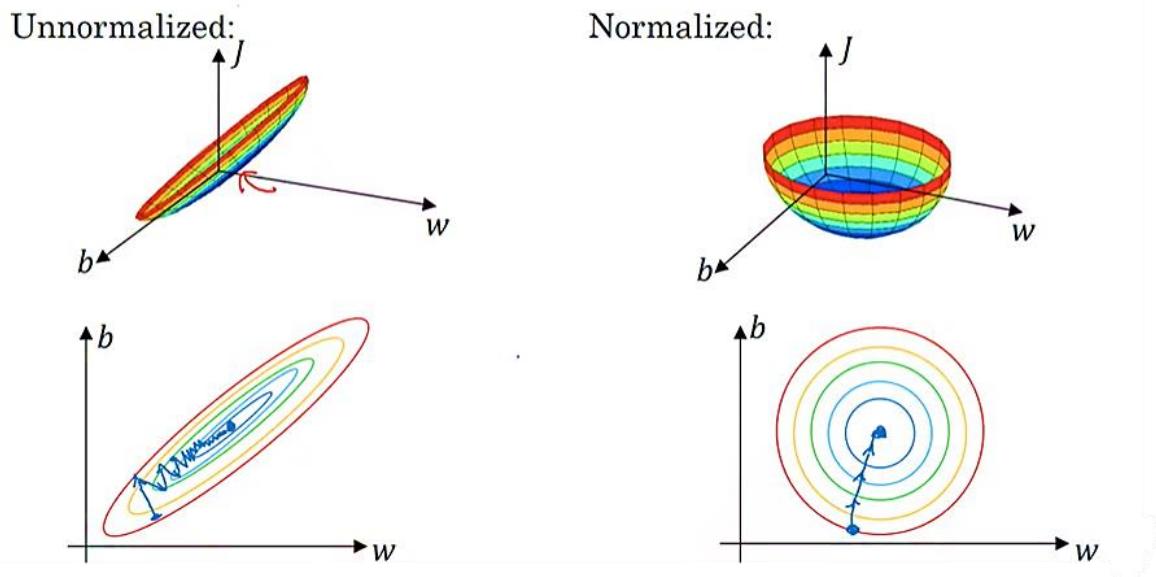
Let's again consider the example above where we want to find out whether a person is going to buy an exercise machine which is

costly based on his weight and salary. So, as we know the prediction should depend on both the weight and the salary of the person. But since the salary has a higher scale of values our model might trend to give that feature more importance than the weight. So, we might end up getting wrong predictions for people with less weight and more money or vice versa.

Also, another factor is that normalisation helps gradient descent easier.

Why normalize inputs?

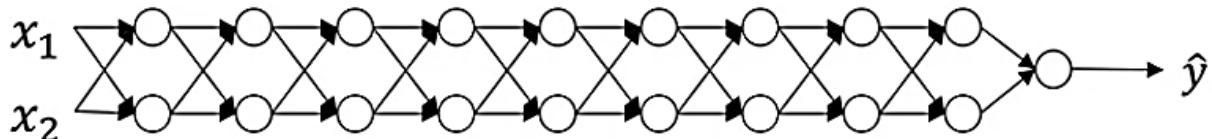
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$



We can see the when we normalize the cost function graph would look more symmetric and therefore gradient decent would be much faster as we can use bigger values of learning rate.

Vanishing/Exploding gradient descent

When we train deep neural networks like the one above, the gradient descent might get exponentially high and we call it exploding gradient descent. When it gets exponentially small, we call it vanishing gradient descent.



Let the weight for each layer be $W^{[1]}, W^{[2]}, W^{[3]}, \dots, W^{[L]}$. Let's take $b = 0$ here. And let the activation function $g(x) = z$ then we have:

$$\begin{aligned}\hat{y} &= W^{[L]} A^{[L-1]} \\ &\downarrow \\ z^{[L-1]} &= W^{[L-1]} A^{[L-2]} \\ &\downarrow \\ z^{[L-2]} &= W^{[L-2]} A^{[L-3]} \\ &\vdots \\ &\downarrow \\ \hat{y} &= W^{[L]} W^{[L-1]} W^{[L-2]} \dots W^{[2]} W^{[1]} x\end{aligned}$$

let $W^{[l]} = \begin{bmatrix} x & \dots & 0 \\ 0 & \dots & x \end{bmatrix}$, then $\hat{y} = W^{[L]} \begin{bmatrix} x & 0 \\ 0 & x \end{bmatrix}^{L-1}$

(1. IF $x > 1$, then $\hat{y} \approx 1 \cdot 5^{L-1} x$ ($x = 1.5$)

(2. FOR $x < 1$, $x = 0.5$ $\hat{y} \approx 0.5^{L-1} x$)

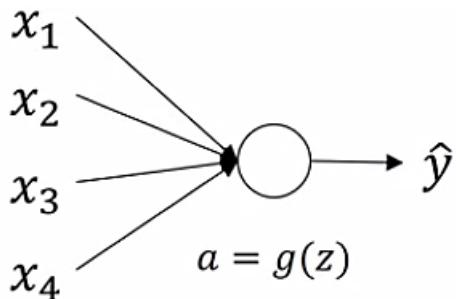
$\left\{ \begin{array}{l} W^{[l]} \text{ has different dimension} \end{array} \right.$

As we can see if the values for weight is greater than 1 the output value increases exponentially or if weights are less than 1 then it decreases exponentially. Similarly, we can show that the gradients increase and decrease exponentially as we change the weights. This can make the learning very difficult as a very less value means the gradient descent would take very small steps and if its too large then gradient descent wouldn't converge to local minima. To reduce this, we have to initialize weights carefully.

Weight Initialization for Deep Learning Networks

Here we would see a way in which we can initialize weights so that we can reduce the vanishing/exploding problem.

Let's consider a single neuron neural network for this.



Here $\mathbf{z} = \mathbf{w}_1 \mathbf{x}_1 + \mathbf{w}_2 \mathbf{x}_2 + \mathbf{w}_3 \mathbf{x}_3 + \dots + \mathbf{w}_n \mathbf{x}_n$

So, in order to reduce the problem, the weights should be smaller for larger number of features since \mathbf{z} is the sum of all these terms, i.e. when n increases the weight values should decrease. One way to achieve this is to set the variance of the randomly initialized values of weights to be neither too large than 1 nor too smaller than one. So, we would set **var(w) = 1/n**. So that when n increases the values of weights would be smaller.

In general, we set:

$$\mathbf{W}^{[l]} = \text{np.random.randn}(n^{[l]}, n^{[l-1]}) * \text{np.sqrt}(1/n^{[l-1]})$$

When $g(z) = \text{ReLU}$ we would use $2/n^{[l-1]}$ as variance.

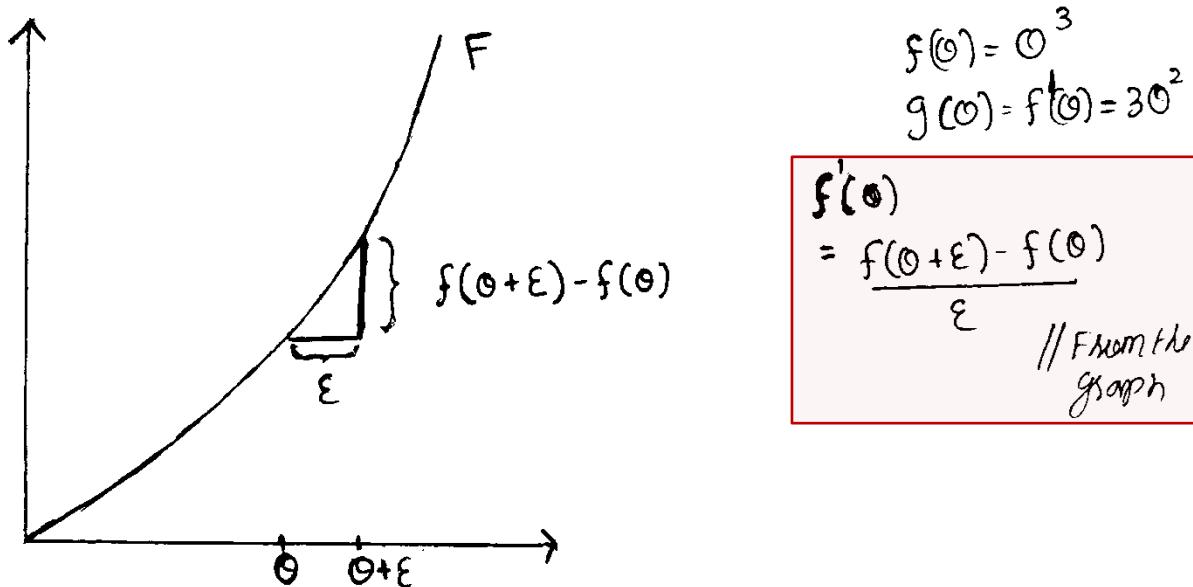
Methods for gradient checking

There are two methods to find gradient of a function:

$$\frac{f(\theta+\varepsilon)-f(\theta)}{\varepsilon} \quad \& \quad \frac{f(\theta+\varepsilon)-f(\theta-\varepsilon)}{2\varepsilon}$$

Let's see which one performs better in finding gradient by taking an example function $f(\theta) = \theta^3$ and take the positive side of the graph to find gradient.

1st finding gradient with one triangle.



$$\text{Let } \theta = 1 \text{ & } \epsilon = 0.01$$

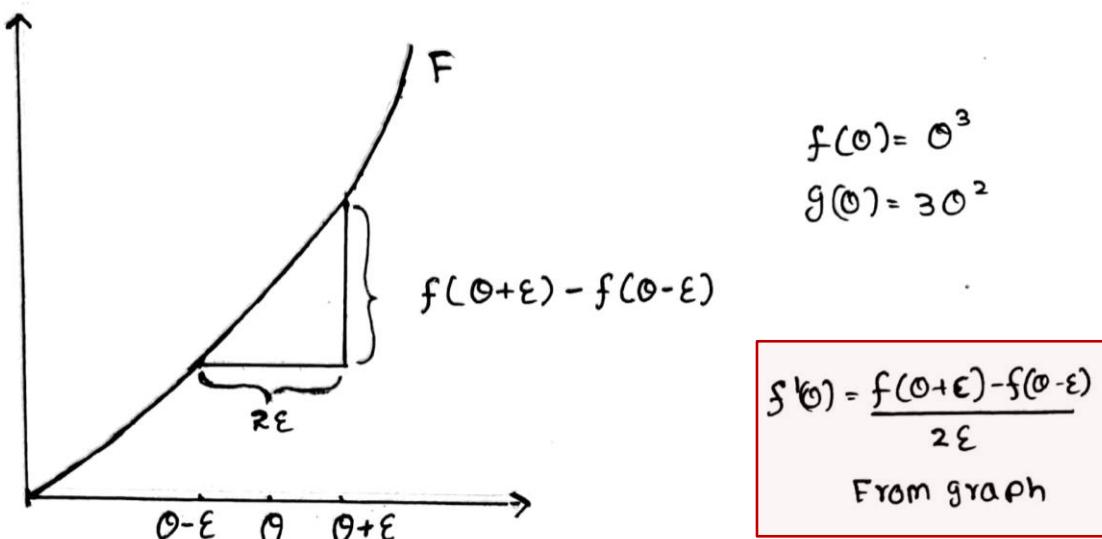
$$\text{then } \theta + \epsilon = 1.01, f(1) = 1, f(1.01) = 1.030301$$

$$\therefore f'(\theta) = \frac{1.030301 - 1}{0.01} = 3.0301 \quad // \theta = 1$$

$$g(\theta) = 3 \times 1^2 = 3 \quad // \theta = 1$$

$$f'(\theta) - g(\theta) = 0.0301 \quad -①$$

Here we can see that the error in finding gradient is 0.0301 when we use triangle with ϵ width. Next, we would find gradient using triangle with width 2ϵ .



Let $\theta = 1$ & $\epsilon = 0.01$

then $\theta + \epsilon = 1.01$, $\theta - \epsilon = 0.99$, $f(1.01) = 1.030301$,
 $f(0.99) = 0.970299$

$$\therefore f'(1) = \frac{1.030301 - 0.970299}{2 \times 0.01} = 3.0001$$

$$g(1) = 3$$

$$f'(1) - g'(1) = 0.0001 \quad \text{--- (2)}$$

From the two pictures we can see that finding gradient using triangle with width 2ϵ is much more accurate than the other.

Gradient checking

How it works?

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

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

Then $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, W^{[L]}, b^{[L]})$ would be equal to $J(\theta_1, \theta_2, \theta_3, \dots, \theta_L) = J(\theta)$

Let $\epsilon = 10^{-7}$

For $i = 1$

$$d\theta_{approx}[i] = \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_i + \epsilon, \dots, \theta_n) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots, \theta_n)}{2\epsilon}$$

The value of $d\theta_{approx}$ should be \approx to $d\theta$ to check this we do:

Check

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

For $\epsilon = 10^{-7}$ if the difference value we got $= 10^{-7}$ then gradients are all correct.

If it is less than or equal to about 10^{-3} then we need to check for errors in calculation of gradients.

Gradient checking implementation notes

- Don't use in training – only to debug
- If algorithm fails grad check, look at components to try to identify bug.
- Remember regularization.
- Doesn't work with dropout.
- Run at random initialization; perhaps again after some training.

WEEK 2 - Optimization Algorithms

Mini Batch Gradient Descent

Assumption: The dataset is large (BIG Data)

In mini batch gradient descent, we split our dataset of $m=5,000,000$ training samples into different batches containing let's say 1000 examples each. We use the term $X^{\{t\}}$ to represent one batch. So, with 5,000,000 training samples we would have 5000 batches each having 1000 training samples.

$$X^{\{t\}}_{(n_x, 1000)}, Y^{\{t\}}_{(1, 1000)}$$

repeat
{

For $t = 1$ to 5000

{ forward propagation
 $Z^{[1]} = W^{[1]} X^{\{t\}} + b^{[1]}$
 $A^{[1]} = g^{[1]}(Z^{[1]})$
 \vdots
 $A^{[L]} = g^{[L]}(Z^{[L]})$

$$J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^{1000} L(y^{(i)}, A^{[L]}) + \frac{\lambda}{2 \times 1000} \sum_{l=1}^L \|W^{[l]}\|_F^2$$

backward propagation

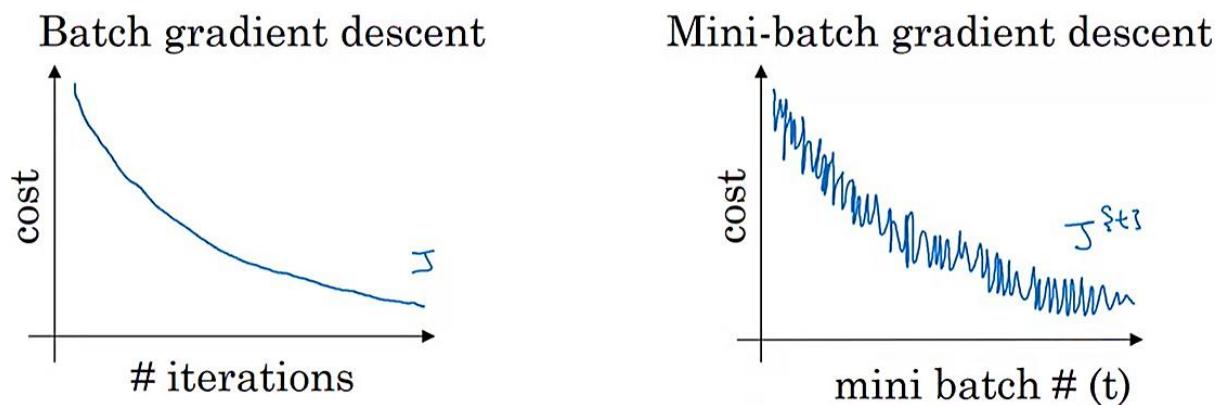
$$W^{[l]} := W^{[l]} - \alpha dW^{[l]}, \quad b^{[l]} := b^{[l]} - \alpha db^{[l]}$$

}

}

We can see that by dividing the training samples into different batches we can now run gradient descent on a batch of training samples multiple times. Means after one epoch (one iteration) we have already trained the weights 5000 times using small batches of the training samples. Whereas in batch gradient descent after one epoch we have only changed the weights once.

Cost function curve



The oscillations in the curve in the case of mini-batch gradient descent is due to the different costs for different batches. But it would also trend down following a similar curve.

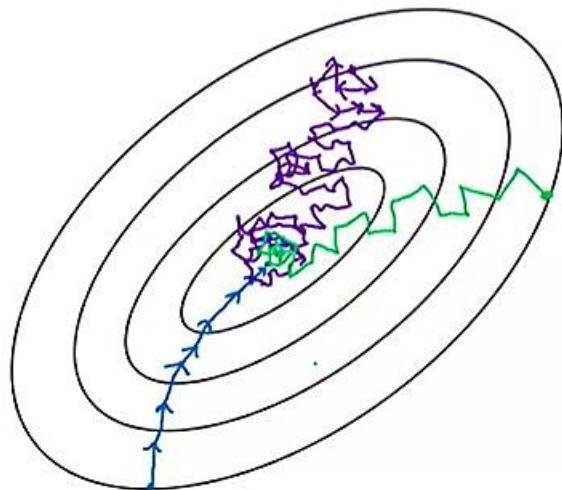
Types of batch gradient descent

If size of batch is equal to m then it is batch gradient descent. The whole training sample is taken as a batch.

If we take each batch as one training sample from the training set the it is called Stochastic gradient descent.

If we take batch size as anywhere between 1 and m then it is mini batch gradient descent.

These three gradient descent methods have difference in the way in which they converge towards the minima.



In batch gradient descent the cost reduces gradually and converge to the minima in time.

In stochastic gradient descent the cost oscillates few times while moving towards the minima and keep on circling the minima while never converging completely

In mini batch gradient descent, the cost oscillates while moving towards the minima and it would come really close to the minima but never completely converge. These oscillations can be reduced by changing the learning rate.

Choosing mini-batch size

- If size of training data is $m \leq 2000$ then do batch gradient descent. If it is greater than that then choose mini-batch.
- When we choose the size for batch choose the size as powers of 2. For example, like **64, 128, 512, 1024** etc.
- Also, when we choose size of the batch, we have to make sure that it fits in the CPU memory.

Exponentially weighted averages

Let's take an example of the temperatures in London for a year.

$$\theta_1 = 40^{\circ}\text{F}$$

$$\theta_2 = 49^{\circ}\text{F}$$

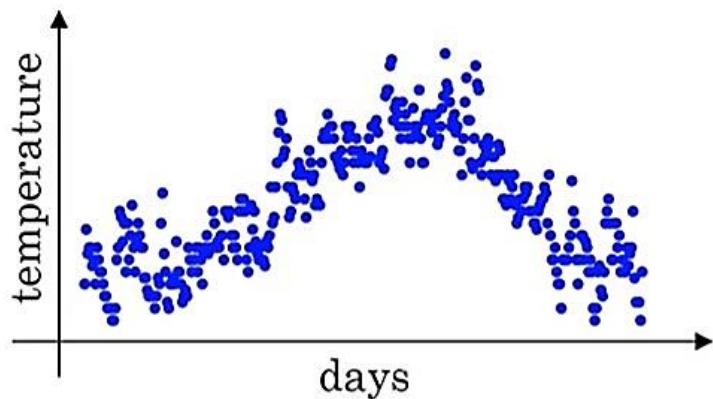
$$\theta_3 = 45^{\circ}\text{F}$$

⋮

$$\theta_{180} = 60^{\circ}\text{F}$$

$$\theta_{181} = 56^{\circ}\text{F}$$

⋮



So, if we need to find the trend of this data by calculating the moving average in a particular window, we would do like this.

$$v_0 = \text{temperature of day 1}$$

$$v_1 = \beta v_0 + (1-\beta) \theta_1,$$

$$v_2 = \beta v_1 + (1-\beta) \theta_2,$$

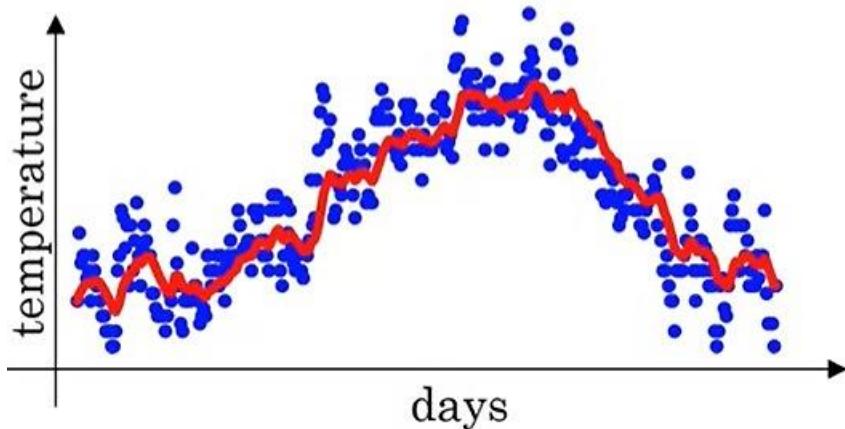
⋮

$$v_t = \beta(v_{t-1}) + (1-\beta)\theta_t$$

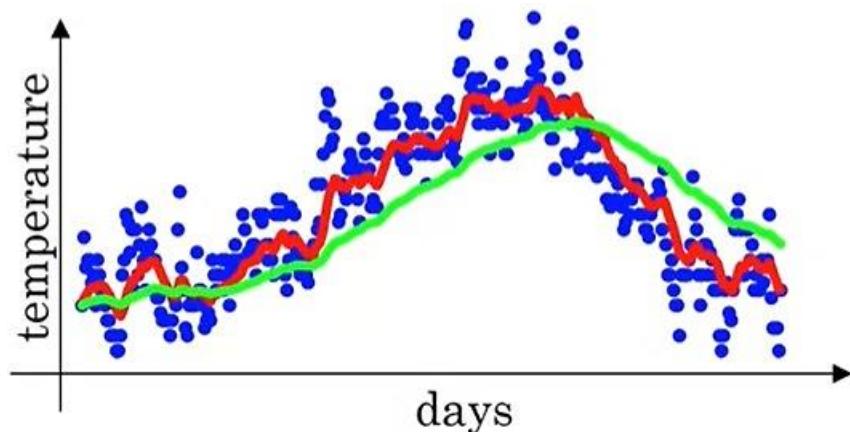
θ_t is the original temperature
of the day t

Where v_t is average over $\frac{1}{1-\beta}$ days of temperature.

When $\beta = 0.9$ we are averaging over 10 days of temperature. The graph would look like this:

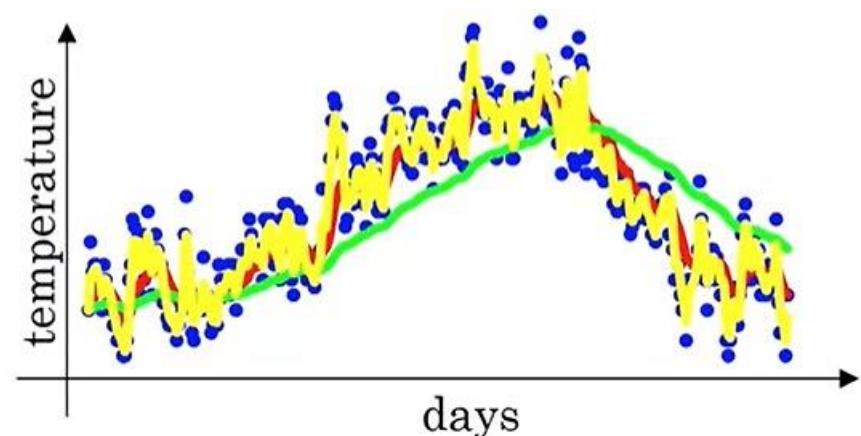


When $\beta = 0.98$ we are averaging over 50 days of temperature. The graph would look like this:



The change of the green line and its squeezed towards right is because we are giving more weight to the past day's temperature than the current day's one. So, the average temperature would be little late to adapt to the new changes.

When $\beta = 0.5$ we are averaging over 2 days of temperature. The graph would look like this:



We can see that this is not a great way to find the trend as it can also include some outliers which wouldn't help in predicting the trend. So, we can say in this case the red line looks more better choice.

Explanation & implementation

Consider:

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

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

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

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

$$\begin{aligned}
 v_{100} &= (0.1\theta_{100} + 0.9v_{99}) \\
 &\quad \downarrow \\
 &= (0.1\theta_{99} + 0.9v_{98}) \\
 &\quad \downarrow \\
 &= (0.1\theta_{98} + 0.9v_{97}) \\
 &\quad \vdots \\
 &\quad \downarrow \\
 &= (0.1\theta_{02} + 0.9v_{01})
 \end{aligned}$$

$$\therefore v_{100} = 0.1\theta_{100} + 0.1 \times 0.9 \cdot 0.99 + 0.1 \times 0.9^2 \cdot 0.98 + \dots + 0.1 \times 0.9^{t-1} \cdot \theta_0$$

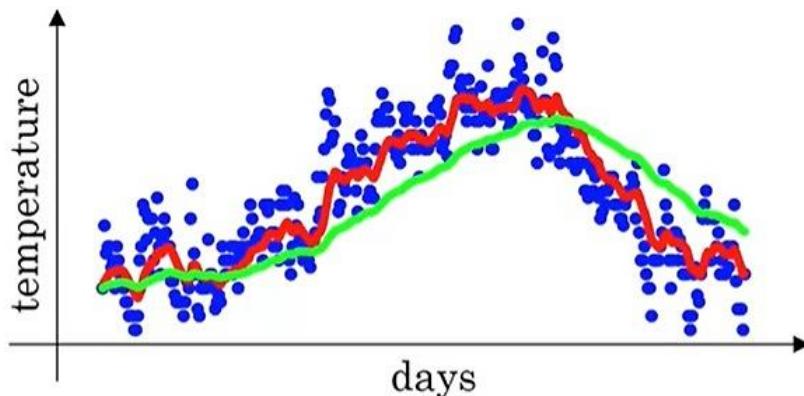
The number of days in the window to find the moving average is calculated as $1/(1-\beta)$. So, for $\beta = 0.9$ no of days = 10. The reason is because when calculating the value of v_{100} above in the example we have to calculate the weighted sum of previous days temperatures. So, when the weight reaches 0.9^{10} it is a small number almost equal to 0.35 which is approximately $1/e$. For a window of n days the value of β is calculated as $(n-1)/n$.

Implementation

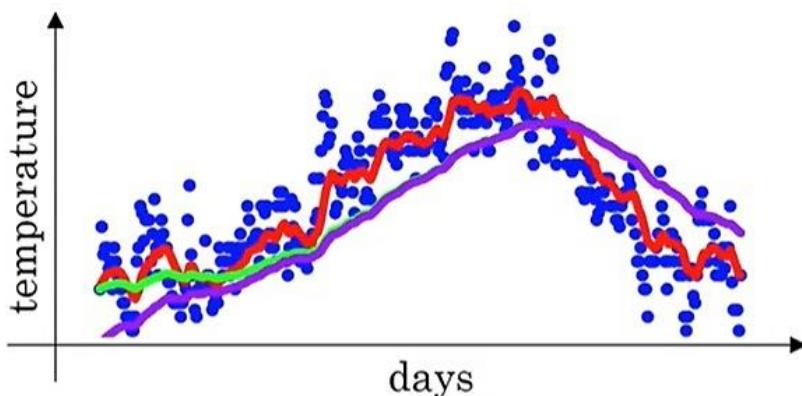
 $v_{\theta} := 0$
 $v_{\theta} := \beta v + (1 - \beta) \theta_0,$
 $v_{\theta} := \beta v + (1 - \beta) \theta_1,$
 \vdots
 Generally:
 $v_{\theta} := 0$
 Repeat {
 Get next θ_t
 $v_{\theta} := \beta v_{\theta} + (1 - \beta) \theta_t$
 }

Bias Correction

We know that for $\beta = 0.9$ and 0.98 we get lines like these



But actually, without bias correction the line we would get for $\beta = 0.98$ is given below as violet coloured line:



We can see that the line starts well below the actual average value is supposed to be. The reason is:

$$v_0 = 0$$

$$\begin{aligned} v_1 &= 0.98 v_0 + 0.02 \theta_1 \\ &= 0 + 0.02 \theta_1 \end{aligned}$$

So if θ_1 = lets say 10° , v_1 would be 0.8°

Similarly

$$\begin{aligned} v_2 &= 0.98 \times 0.02 \theta_1 + 0.02 \theta_2 \\ &= 0.0196 \theta_1 + 0.02 \theta_2 \end{aligned}$$

which would again be a small value

This is the reason why we see that very low starting point for the line in the graph.

To solve this we find $\frac{v_t}{1-\beta^t}$ instead of v_t

$$\text{So for } t=2, 1-\beta^2 = 1 - (0.98)^2 = 0.0396$$

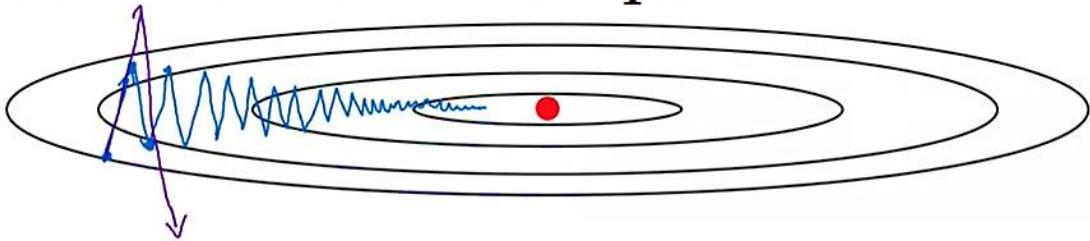
$$\frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

Here we can see that $0.0196 + 0.02 = 0.0396 \therefore$ we are actually taking the weighted average of θ_1 and θ_2 hence removing the bias

Also we can see that as t increases β^t decreases \therefore reducing its effect on the values as t increases

Gradient Descent with Momentum

Gradient descent example



When we do mini-batch or batch gradient descent on our data, the cost function reduces to minima as shown above in blue line. We can see that it **oscillates** up and down while moving towards the minima. This would increase the learning time and is inefficient. And if we increase the learning rate more it can go out of the path and start to **shoot up and down** as shown in purple line.

To reduce this, we can average out the positive and negative values in the vertical direction, thereby reducing the oscillations. It won't affect the learning in horizontal direction as the direction for all the values are the same.

So, we update our gradient descent steps like this:

$$v_{dw=0}$$

$$v_{db=0}$$

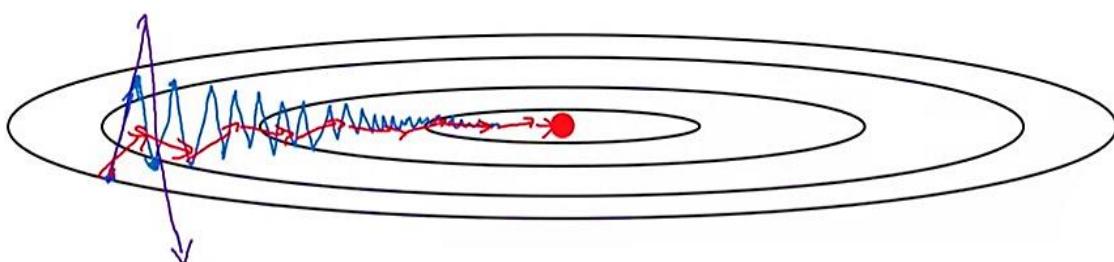
On iteration t :

Compute dW, db on the current mini-batch

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

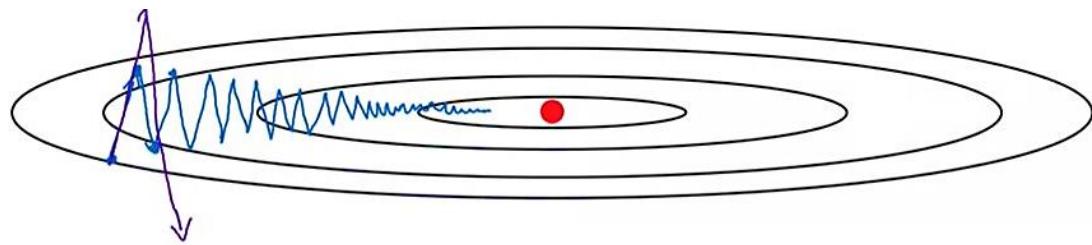
$$v_{db} = \beta v_{db} + (1 - \beta) db$$

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



So, by applying exponentially weighted moving average to our gradient descent steps we have averaged out the oscillations to make our learning faster. Also, we can use β too as a hyperparameter. But usually setting $\beta = 0.9$ works mostly fine.

RMSprop (Root Mean Square Propagation)



This is another way to reduce the oscillations and making the learning faster. Here we update our terms for gradient descent like this:

First let's imagine that the term that makes the oscillations on vertical axis is b and the horizontal is w . In normal cases it might be combinations of lots of features. So, we have to increase w and decrease b .

On iteration t :

(compute dw, db on current mini batch)

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

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

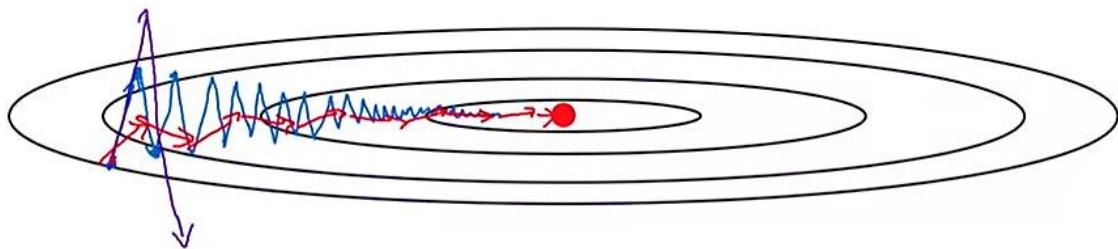
$$w := w - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}, b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$$

where ϵ is just a small term to make sure that the denominator is never zero.

Let's imagine that b is in the horizontal axis and w is in the vertical axis. So, if we need to reduce the oscillations, we need to change

the value of b slowly, for this we have to subtract it with smaller values. Here we can see that we are finding the exponential average using the square of the derivative term. So that if the slope is larger its square would be much larger. If it is smaller the square would be much smaller. In our case the slope for b is large. So db would be large. So db when squared would be larger. Then we are dividing the derivative term which we are going to subtract from b using the weighted average term which is large thereby reducing the value of b slowly.

On the other hand, since w is small the value of weighted average would be small too therefore when we divide the derivative with smaller value, if it is smaller than 1 it would increase the derivative term therefore getting larger changes and faster learning horizontally.



Now we would see an algorithm that combines the use of gradient momentum and RMSprop together for better optimization,

Adam (Adaptive Moment Estimation) Optimization

Adam optimization combines both momentum and RMSprop algorithms to optimize the gradient descent. It is implemented as follows:

How does Adam work?

1. It calculates an exponentially weighted average of past gradients, and stores it in variables v (before bias correction) and $v^{corrected}$ (with bias correction).
2. It calculates an exponentially weighted average of the squares of the past gradients, and stores it in variables s (before bias correction) and $s^{corrected}$ (with bias correction).
3. It updates parameters in a direction based on combining information from "1" and "2".

The update rule is, for $l = 1, \dots, L$:

$$\left\{ \begin{array}{l} v_{dW^{[l]}} = \beta_1 v_{dW^{[l]}} + (1 - \beta_1) \frac{\partial \mathcal{J}}{\partial W^{[l]}} \\ v_{dW^{[l]}}^{corrected} = \frac{v_{dW^{[l]}}}{1 - (\beta_1)^t} \\ s_{dW^{[l]}} = \beta_2 s_{dW^{[l]}} + (1 - \beta_2) \left(\frac{\partial \mathcal{J}}{\partial W^{[l]}} \right)^2 \\ s_{dW^{[l]}}^{corrected} = \frac{s_{dW^{[l]}}}{1 - (\beta_2)^t} \\ W^{[l]} = W^{[l]} - \alpha \frac{v_{dW^{[l]}}^{corrected}}{\sqrt{s_{dW^{[l]}}^{corrected}} + \epsilon} \end{array} \right.$$

where:

- t counts the number of steps taken of Adam
- L is the number of layers
- β_1 and β_2 are hyperparameters that control the two exponentially weighted averages.
- α is the learning rate
- ϵ is a very small number to avoid dividing by zero

Hyperparameter choices:

α – needs to be tuned

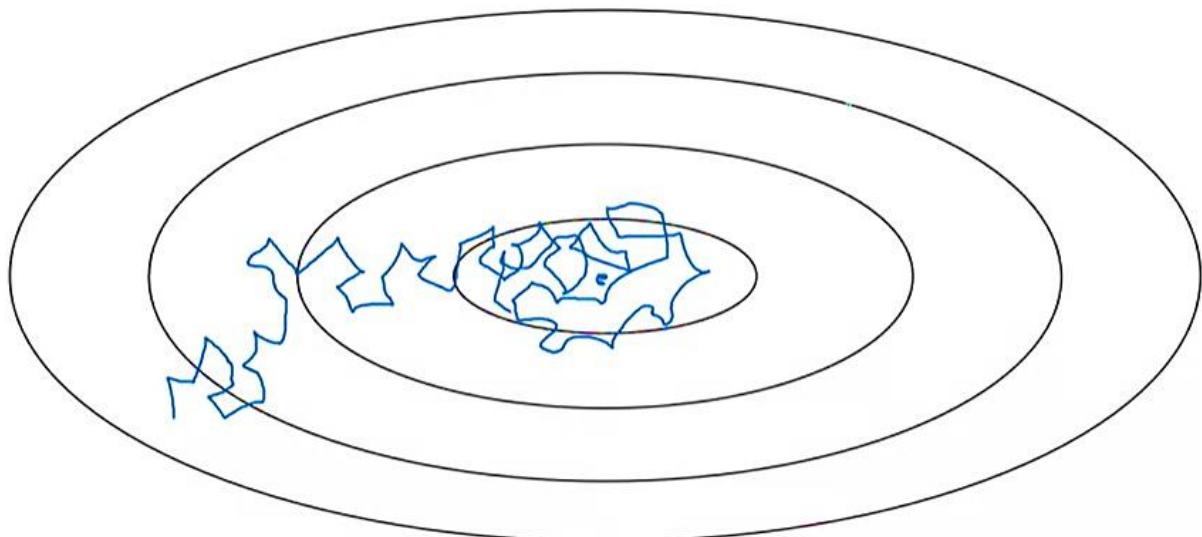
$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

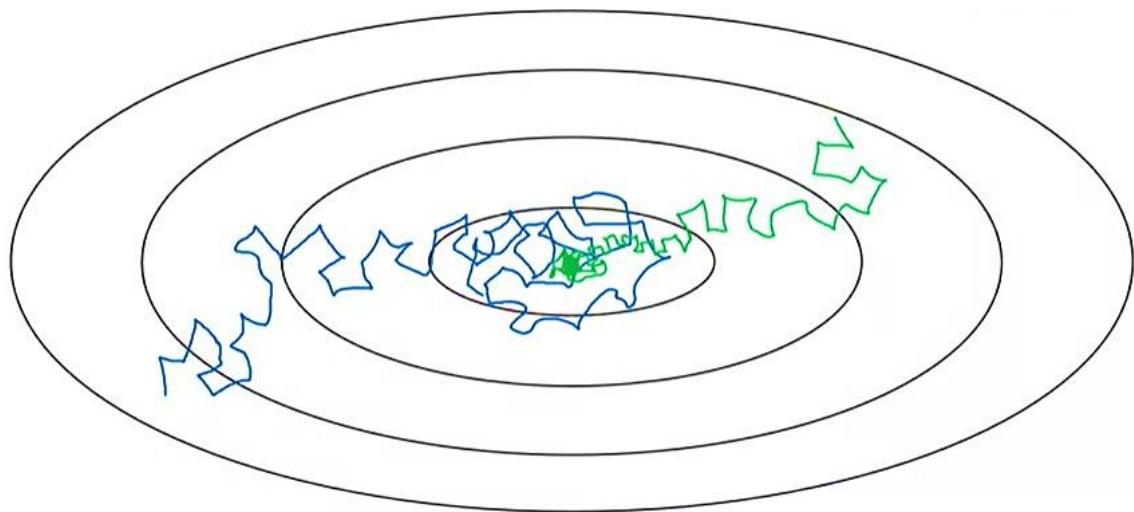
$$\epsilon = 10^{-8}$$

Learning Rate Decay

Suppose we are doing gradient descent on a small mini-batch of training set. Then the steps would tend to be little bit noisy and it would go down towards the minima but it wouldn't converge to it and instead **it would revolve around a larger region around the minima.**



So, one thing we can do to get it closer to the minima is to reduce the learning rate as the number of epochs increases. So, for **smaller learning rate** the steps we take would be slower and smaller and the **steps would revolve around a much smaller region close to the minima.**



By doing this we can afford to have larger learning rate in the beginning of the training and as we get closer to the minima, we can reduce the learning rate to get better convergence.

The different commonly used ways to do this are:

$$(1) \alpha = \frac{1}{1 + \text{decay rate} * \text{epoch-num}} \cdot \alpha_0$$

$$(2) \alpha = \text{num}^{\text{epochnum}} \cdot \alpha_0 ; \text{ num} \leq 1$$

eg: $0.95^{\text{epochnum}} \cdot \alpha_0$

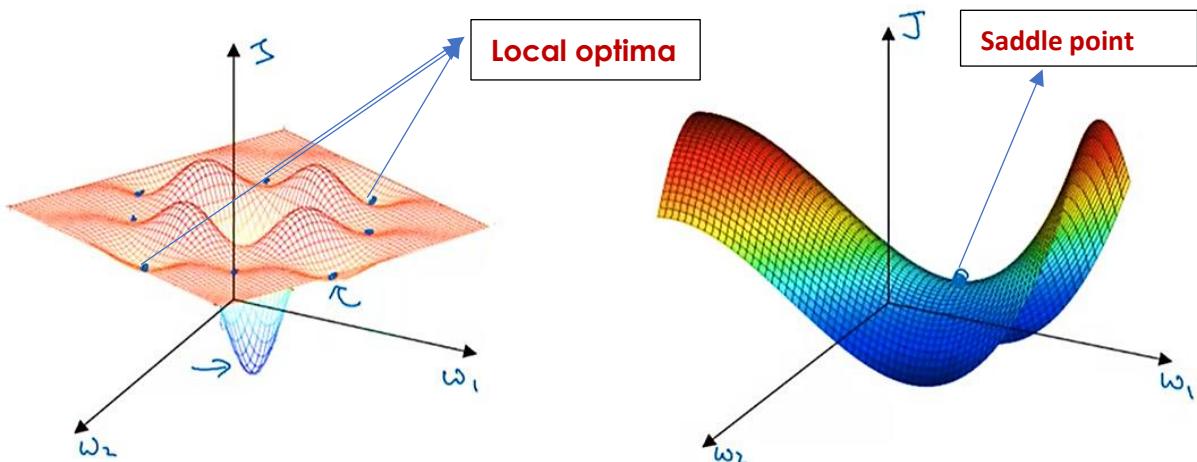
$$(3) \alpha = \frac{k}{\sqrt{\text{epochnum}}} \cdot \alpha_0$$

$$(4) \alpha = \frac{k}{\sqrt{t}} \cdot \alpha_0$$

We can try out different values for α_0 , decay rate, num and k and tune it as hyperparameters.

The problem of local optima

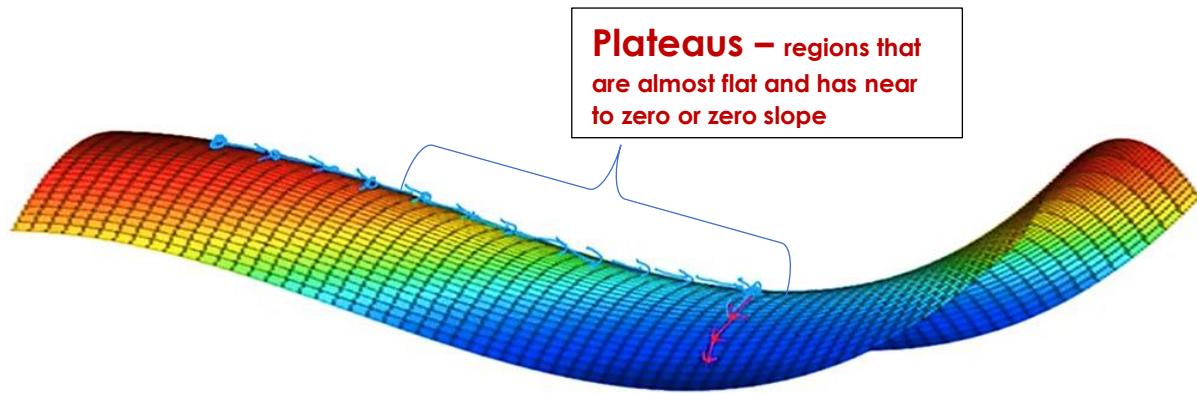
- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow



Getting stuck at local optima is a problem that occurs mostly when we have dimensions like 2 or 3 for the training data features. It

usually won't happen for a training set with lots of dimensions. Instead there can be saddle points where derivative would be zero.

But one problem that can happen is because of plateaus. These are regions which are almost flat and has almost zero slope. So, it would take long time to get from these regions to global optima.



This problem can be reduced by using the optimization algorithms that we learned about.

WEEK 3 – Hyperparameter Tuning

Hyperparameters that are most important to be tuned:

α

β

$\beta_1, \beta_2, \varepsilon$

Layers

Hidden units

Learning rate decay

Mini-batch size

RED → Always tune

ORANGE → 2nd Most important hyperparameter to tune

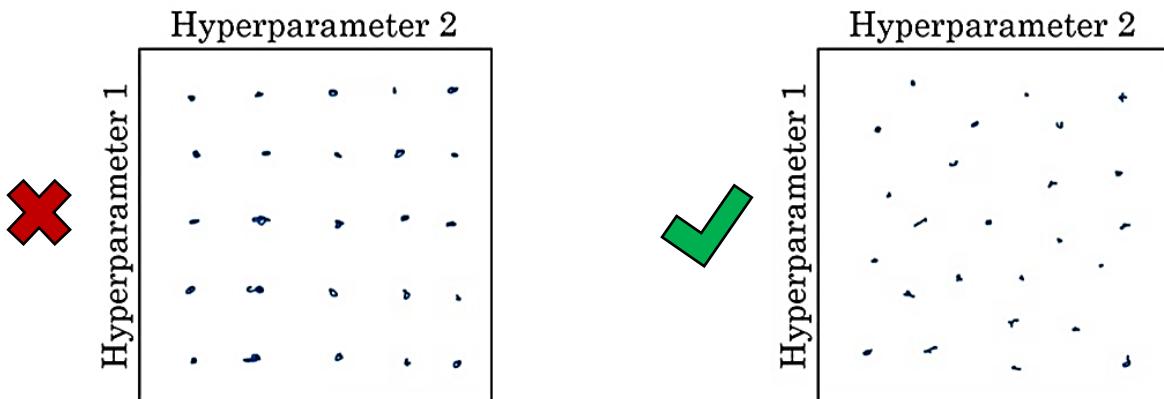
VIOLET → 3rd Most Important

BLACK → Usually not tuned ($0.9, 0.99, 10^{-8}$)

Choosing parameters from different samples

In early days of deep learning people tend to choose parameters by placing them as samples in grid form and then trying out each of the pairs, as shown in the figure left.

Try random values: Don't use a grid

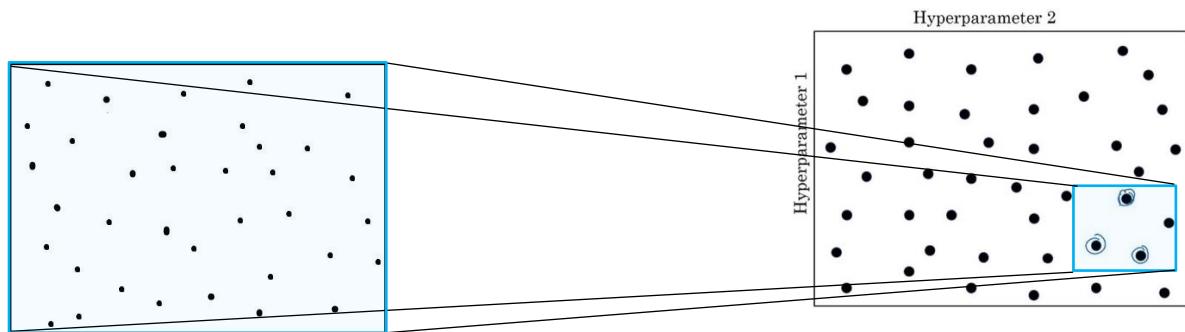


One disadvantage of putting hyperparameters in a grid for is that we would be able to test few different values of a particular hyperparameter even after we tests the entire grid. For example, if the grid contains 25 samples then each hyperparameter would be having 5 different values. So, we would only be seeing 5 different values after looking at 25 samples.

And for example, let's say hyperparameter 1 is learning rate alpha and hyperparameter 2 is epsilon. We know that alpha is more important than epsilon. So even after we are seeing 5 different values of alpha the value of epsilon might not be changing much as it has less importance.

So, when we test hyperparameters in sample we often place them randomly instead of in a grid format to get 25 unique samples with 25 unique values for each hyperparameters. In this case after testing all 25 samples we would be seeing 25 different values.

Another thing that works well in finding hyperparameter is **coarse to fine tuning**. We would focus on the area where more samples are giving good results and we take more samples in that region and test with it further.

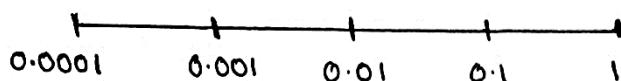


Using appropriate scale for Hyperparameter tuning

While searching for hyperparameters like $n^{[l]}$ or L we can use linear scale and choose values uniformly at random. Like for $n^{[l]}$ it can range from 50 to 100 and for L it can be around 2 to 4 or so. But when we are dealing with hyperparameters like α or β if we choose to tune uniformly at random in a linear scale then we would end up wasting lots of resources tuning a particular range closer to 1 more than the other range of values. So here we use logarithmic scale to tune.

Set $\alpha = 10^r$

Here we will be checking values in log scale



i.e from $10^{-4} \rightarrow 10^0$

Generally we take values from 10^a to 10^b

$$\text{here } a = \log_{10} 0.0001 = -4, b = \log_{10} 1 = 0$$

$$\therefore r = [a, b] = [-4, 0]$$

$$r = -4, \alpha = 10^{-4} = 0.0001$$

$$r = 0, \alpha = 10^0 = 1$$

Set $\beta = 1 - 10^\gamma$

β range from 0.9 to 0.999

$1 - \beta$ range from 0.1 to 0.001

$$= 10^{-1} \text{ to } 10^{-3}$$

$\therefore \gamma$ range from -1 to -3 ($a = \log_{10} 0.1, b = \log_{10} 0.001$)

$$\text{when } \gamma = -1 \quad \beta = 1 - 10^{-1} = 0.9$$

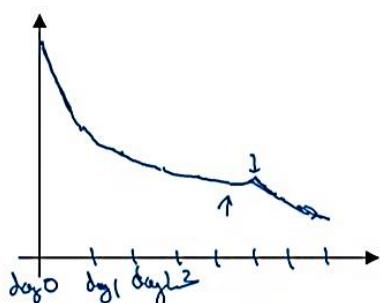
$$\gamma = -3 \quad \beta = 1 - 10^{-3} = 0.999$$

Hyperparameter choices can vary over time due to changes in datasets or changes in sources of data. Therefore, it is always good to retrain your model with the latest data time to time.

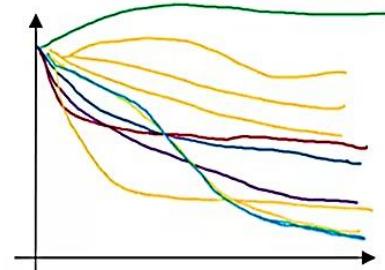
Two ways of training:

1. One model train over long time changing hyperparameters → If you have less computational power and has big dataset.
2. Multiple models with different hyperparameters trained overtime and choosing the best one → When you have lot of computational resources.

Babysitting one model

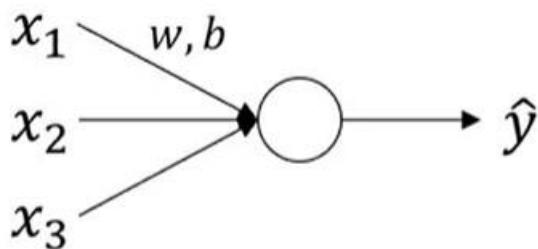


Training many models in parallel



Batch Normalization

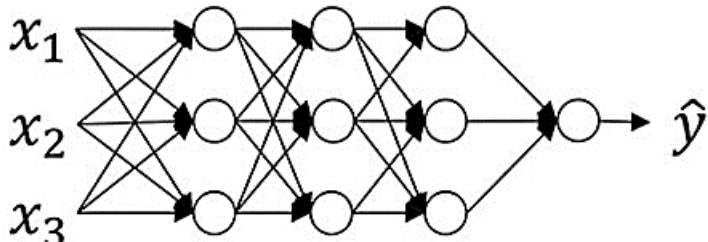
We have seen that for neural network model that resembles logistic regression which doesn't have any hidden layers like the one shown below we normalize only the input features.



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

This would transform the features to have similar mean and variance.

But in the case of neural net frameworks with hidden layers as shown below we have to normalize the hidden layer features too.



Here we need to normalize $a^{[l-1]}$ to improve the learning of $w^{[l]}$ and $b^{[l]}$. We normalize $z^{[l-1]}$ instead of $a^{[l-1]}$.

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

$$x = x - \mu$$

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

$$x = x / \sigma^2$$

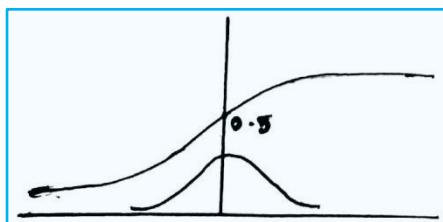
$$\mu = \frac{1}{m} \sum_{i=1}^m z^{i}$$

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

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

ϵ is used to make sure that denominator is never 0

But by doing this would set the mean and variance forcefully to 0 and 1. For example if we are using the activation function sigmoid then we would lose the advantage of the non linear functions capability to learn complex features because our features would be laying in a region closer to 0.



So if we don't want that to happen then we do the following:

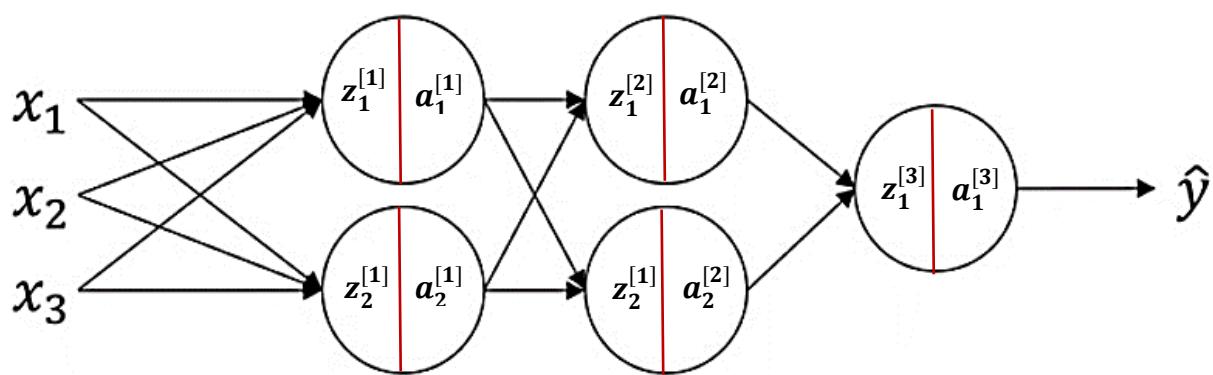
$$\tilde{z}^{[l](i)} = \gamma z_{\text{norm}}^{[l](i)} + \beta$$

Where γ and β are learnable parameters which we would be learning using gradient descent maybe with some optimizers.

If $\gamma = \sqrt{\sigma^2 + \epsilon}$ and $\beta = \mu$ then $\tilde{z}^{[l](i)} = z^{[l](i)}$. Here in this case we are calculating the identity. But if we choose other values then we can normalize the hidden layer features with different mean and variances.

Fitting batch normalization to neural network

While training the model there are mainly two steps involved, forward propagation and backward propagation. We apply normalizations before we do forward propagation and we apply optimization before we do backward propagation. So, we can divide the forward propagation further into two halves, 1st finding the linear function Z and then finding the activation function a. We apply normalization to the z value we calculated. The whole process is depicted below.



$$X \xrightarrow{w^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow[B.N]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \xrightarrow{a^{[1]} = g^{[1]}(\tilde{Z}^{[1]}) w^{[2]}, b^{[2]}} Z^{[2]} \xrightarrow[B.N]{\beta^{[2]}, \gamma^{[2]}} \tilde{Z}^{[2]} \xrightarrow{d^{[2]} = g^{[2]}(\tilde{Z}^{[2]})}$$

Parameters:

$$\begin{aligned} & w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]} \\ & \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]} \quad \left\{ \text{dim} = (n^{[0]}, 1) \right. \end{aligned}$$

In batch gradient descent

$$\begin{aligned} X & \xrightarrow{\{w^{[1]}, b^{[1]}\}} Z^{[1]} \xrightarrow[B.N]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \xrightarrow{g^{[1]}(\tilde{Z}^{[1]}) = a^{[1]} w^{[2]}, b^{[2]}} Z^{[2]} \xrightarrow{\dots} \\ & X^{[2]} \xrightarrow{\dots} \\ & \vdots \\ & X^{[t]} \xrightarrow{\dots} \end{aligned}$$

For $t=1$ to no of batches:

Forward pass on $X^{[t]}$

For each hidden layer use B.N to find \tilde{Z}^t then a^t .

Backprop to compute $d w^{[t]}, d \beta^{[t]}, d \gamma^{[t]}$

Update parameters $w^{[t]}, \beta^{[t]}, \gamma^{[t]}$

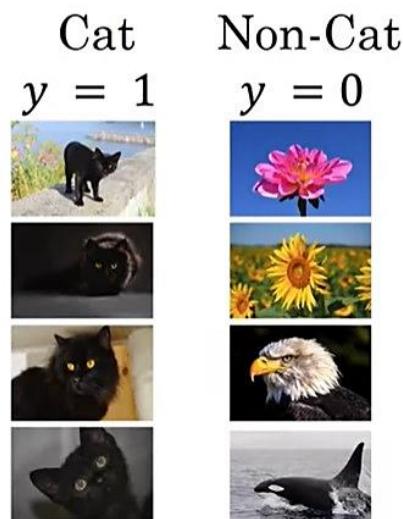
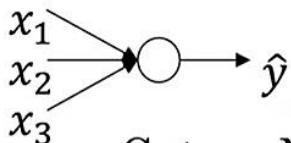
We omit b because since we are using batch normalization the term b would be cancelled out when we normalize the values.

Why batch norm works

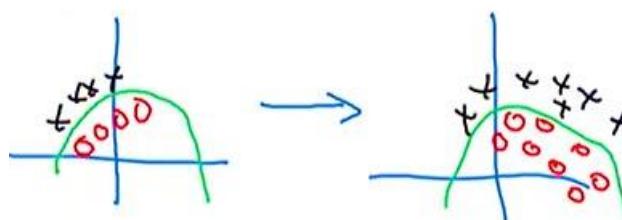
We know that normalizing our inputs would reduce the impact of high valued features on lower valued ones and it would also speed up the learning. But there are more reasons why we should do batch normalization and why it works.

Let's imagine we are training our deep learning model which classifies pictures to cat or non-cat. Our training dataset only contains black cats while we were training. In this case our model will not be performing well on a test set where the cat images are coloured.

Learning on shifting input distribution



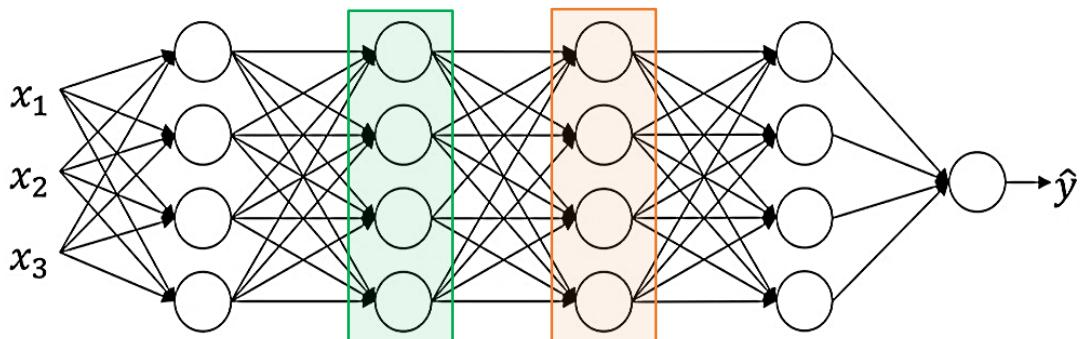
We can see that the distributions are totally different from each other. In this case we can't expect our model to fit a line as shown below which can work well on both the distributions.



This problem of change in the distribution of input values is termed as covariate shift, where changes in distribution would result in

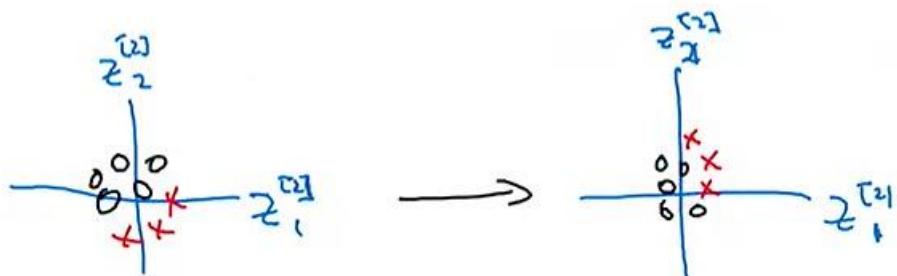
retraining the model. Also, if the ground truth function which maps the input to output changes, we have to retrain the model.

Why this is a problem with neural networks?



Consider **layer 3** of this neural network. It gets input from the former layers as $a_1^{[2]}, a_2^{[2]}, a_3^{[2]} \text{ and } a_4^{[2]}$ and it helps the parameters $W^{[3]}, b^{[3]}, W^{[4]}, b^{[4]}$ and $W^{[5]}, b^{[5]}$ to learn from these features and try to map them to the final output \hat{y} and make it closer to y . But these features $a_1^{[2]}, a_2^{[2]}, a_3^{[2]} \text{ and } a_4^{[2]}$ are again learned by the features from earlier layers and it changes when the parameters $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$ changes. So, if these parameters change as they get updated it would eventually change the features passed to the third layer. So the parameters of that layer and the proceeding layers have to adapt to this new change in distribution of features it gets to learn from.

What normalising the features does is that it wouldn't allow the features to change its distribution by keeping its mean and variance same as before even after the values are changed when we do gradient descent and updating. So even though the values are changing it would still be having the same distribution as shown in the below graph. So this would speed up the learning.



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 some noise to each hidden layer's activations.
- This has a slight regularization effect.

So as the noise is added to hidden layer activations the model cannot rely upon any of the nodes as any of it can have more noise in the future. So, this would make the weights to be spitted among the different nodes more uniformly without giving any particular node more weight.

Also, this noise would get reduced as we use bigger batch sizes. The noise occurs due to the change in z to \tilde{z} and the mean and variance we find are only taking into consideration the values present in the batch and not the entire training set. So as batch size decreases the noise increases and vice versa.

Batch Norm at test time

At test time we won't be having batches to do mini batch. We would probably have single training sample to do normalization. To in order to do that we would find the mean and variance of the distribution as we train the data and we use that to do the mini batch normalization. But instead of finding the mean and variance by running the dataset on the whole model we would keep track of the mean and variance as exponentially weighted averages and would use these as mean and variance to normalize the values at test time.

So similarly, like we found the average of temperature changes over a given window, we would find the mean and variance of the data over a number of iterations.

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

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

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

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

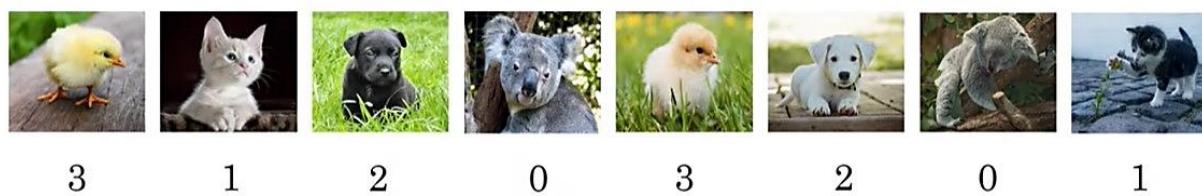
This is the equation to do batch normalization. Now we would find μ and σ^2 by finding exponentially weighted average across mini-batches $X^{\{1\}}$, $X^{\{2\}}$, $X^{\{3\}}$ etc. as $\mu^{\{1\}[l]}$, $\mu^{\{2\}[l]}$...etc. and similarly the variance too and use this to normalize the test inputs.

Softmax Regression

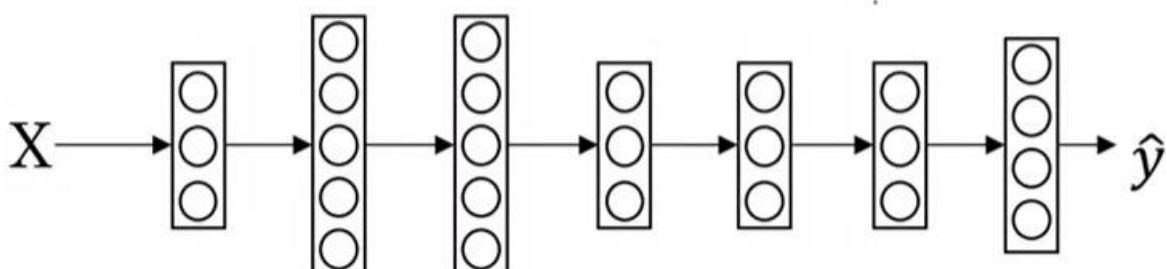
Softmax is a generalization of logistic regression, here we classify our output to one among more than one class or c number of classes.

Before when we used logistic regression to predict cat-vs non cat. Here we can take an example of 3 animals. Class 3 for baby chick's, class 2 for dogs, class 1 for cats and class 0 for other animals. Here we use Softmax as the activation for the whole final layer vector to predict which class our output belongs to. In Softmax regression the number of units in output layer = $n^{[L]} = c$, the number of classes.

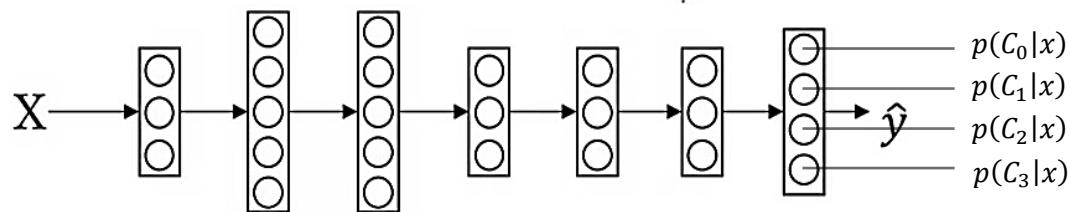
Recognizing cats, dogs, and baby chicks



In this case we would use a neural network with 4 output units.



Here the outputs will be probability of each class being the correct output given x



Here c_1, c_2, c_3, c_4 represents the four classes.

In order to generate these probabilities, we use a layer called Softmax layer in the output layer. Let's see how this works:

First, we would calculate a temporary variable named $t = e^{z^{[L]}}$. The final output $a^{[L]}$ will be equal to $\frac{e^{z^{[L]}}}{\sum_{i=0}^4 t_i}$ so $a_i^{[L]} = \frac{t_i}{\sum_{i=0}^4 t_i}$. Lets take an example to explain this.

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}, \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$\sum_{i=0}^4 t_i = 176.3$$

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

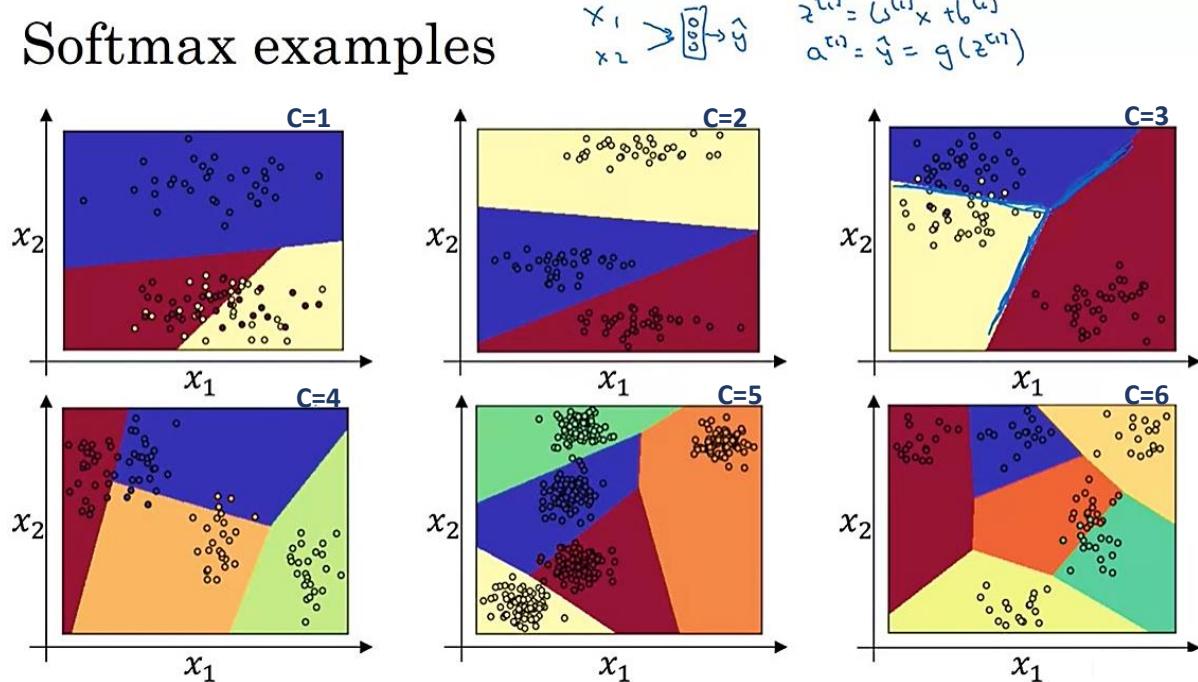
Last Layer after activation

$$\begin{array}{l}
 \begin{array}{ccc}
 \text{Input} & \rightarrow & \frac{148.4}{176.3} = 0.842 \\
 \text{Input} & \rightarrow & \frac{7.4}{176.3} = 0.042 \\
 \text{Input} & \rightarrow & \frac{0.4}{176.3} = 0.002 \\
 \text{Input} & \rightarrow & \frac{20.1}{176.3} = 0.114
 \end{array} \\
 + \\
 + = 1
 \end{array}$$

Here we can see the different probabilities of different cases being the output that is closer to the true value.

Now let's see some examples of Softmax function applied to a neural network with no hidden layer to classify different classes of outputs

Softmax examples



Understanding softmax

$$\begin{aligned}
 & (4, 1) \\
 & z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \\
 & \underbrace{\qquad\qquad\qquad}_{\text{"Soft max"}}, \quad c = 4 \quad g^{[L]}(\cdot) \\
 & a^{[L]} = g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \\
 & \underbrace{\qquad\qquad\qquad}_{\text{"hard max"}}, \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}
 \end{aligned}$$

Softmax regression generalizes logistic regression to C classes.

If $C=2$, softmax reduces to logistic regression. $a^{[L]} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$

The name Softmax is used in contrast to Hardmax where we give 1 to the element with highest probability and 0 to the rest. Also another thing to note is that if the number of classes $c = 2$ then

instead of calculating two outputs we can calculate the probability of just one to understand the probability of the other since both would sum up to one. Therefor Softmax with $c = 2$ is essentially logistic regression.

Now let's take a look at the loss function for Softmax classifier:

$$\text{let } Y = \begin{bmatrix} y^{(1)}, y^{(2)}, \dots, y^{(m)} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix}_{(4, m), c=4}$$

$$\text{let } \hat{Y}^i = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$\therefore \hat{Y} = \begin{bmatrix} \hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(m)} \end{bmatrix} = \begin{bmatrix} 0.3 & \dots & \dots \\ 0.2 & \dots & \dots \\ 0.1 & \dots & \dots \\ 0.4 & \dots & \dots \end{bmatrix}_{(4, m), c=4}$$

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

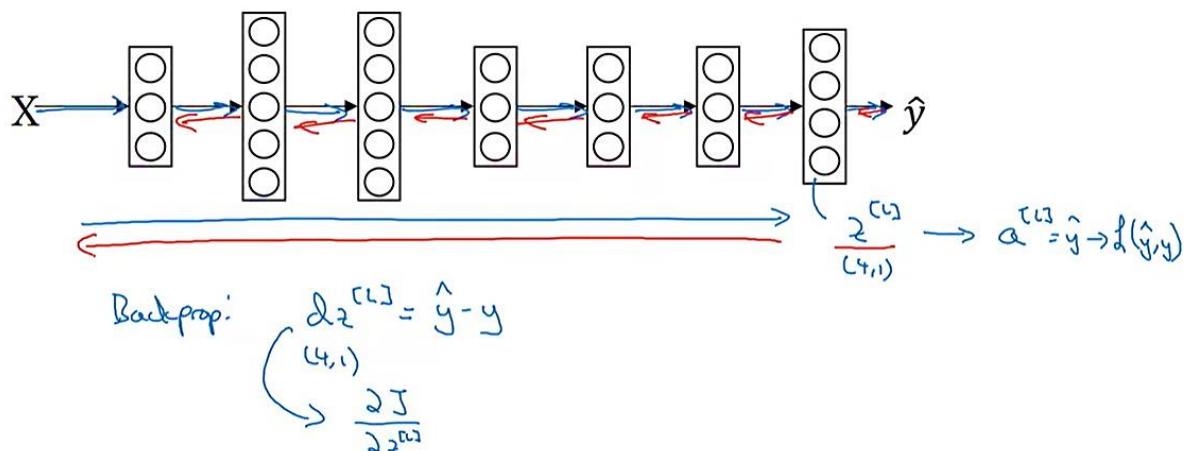
for $j=2$ $y_j = y_2 = 1$ for y^i and all others are 0

$$\therefore L(\hat{y}^i, y^i) = - \log \hat{y}_2$$

so to reduce loss we need to increase \hat{y}_2 ie increase \hat{y}_2 to a much bigger value so that the probability would be close to 1.

$$\text{in general cost: } J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

Gradient descent with softmax



Gradient descent with Softmax is exactly same as that of gradient descent we did before.

*Note

For those who are trying to do TensorFlow calculations as shown by Andrew Ng in the lecture video on their own system, please note that Andrew Ng is using TensorFlow 1 instead of 2 which the latest at the time I write this. So, follow these steps if you encounter any trouble during working out those codes shown in video.

*Assumption: - OS is windows 10

1. Install latest visual C++ distribution from Microsoft
2. If conda install doesn't work try pip install as shown in TensorFlow documentation.
3. No need to install older version instead do this while importing to get version 1 methods:

```
import tensorflow.compat.v1 as tf
tf.disable_v2_behavior()
```

