

# DEEP LEARNING AND NEURAL NETWORK SPECIALIZATION (COURSERA)

COURSE 2

Improving Deep Neural Networks:  
Hyperparameter Tuning, Regularization and  
Optimization

## WEEK 1

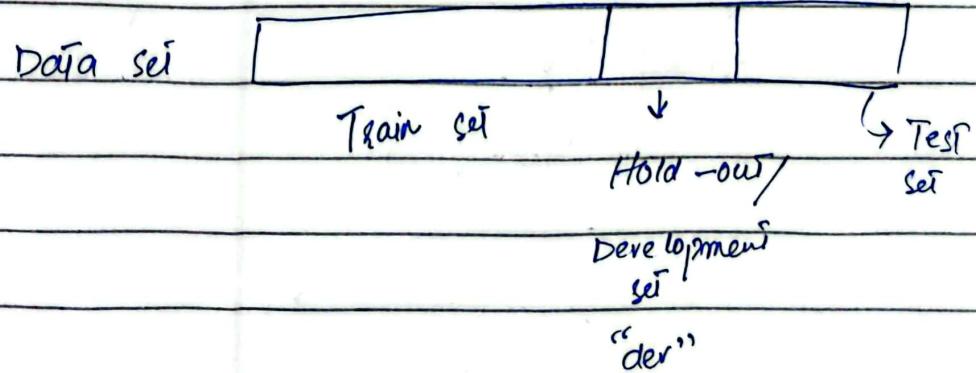
- Setting up your Machine Learning Application:

1) Train / Dev / Test sets

Applied ML is a highly iterative process. We have to consider following

# layers? # hidden units?

# learning rates? # activation functions?



Dev set: It is usually for knowing which algorithm is best to use for these sets. Gradually these sets are decreasing.

Traditional ratio is 60/20/20  
(for small dataset)

- Mismatched train/test distribution

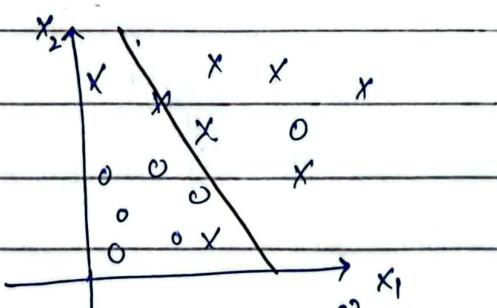
Picture from web pages for train  
and camera for test set

So, we have to have same distribution.

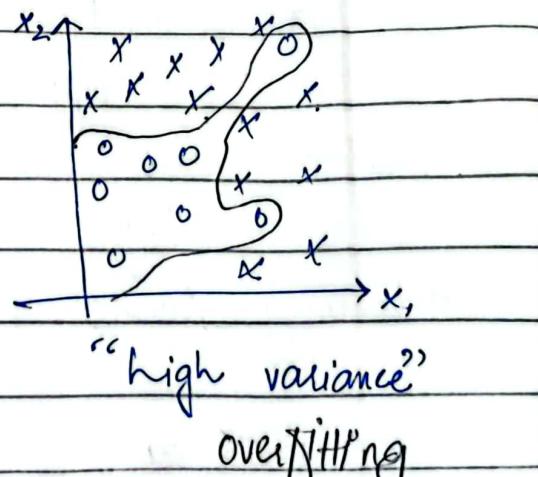
Note:

"Not having a test set might be okay."

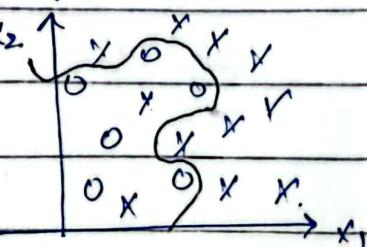
## 2) Bias/Variance:



"high bias"  
underfitting



"high variance"  
overfitting



"just right"

Cat Classification:

$y=1$   
Cat

$y=0$   
Dog



Train set error: 1% }

Dev set error: 11%

You are doing very well on the training set and relatively poorly on development set.  
This looks like you might have overfit the training set. This has high variance.

Train set error: 15% } 15% } 0.5% }

Dev set error: 16% } 30% } 1%

↓ high bias & low bias ?

Because training error is so much then it means it is underfitting the data. So, it has high bias.

This all when bayes error or optimal error is  $\approx 0\%$

### 3) Basic Recipe for Machine Learning:

high bias?

(Training data

performance)



high variance?

(Dev set

performance)

Bigger Network:

more hidden layers

• Train longer.

• NAT architecture search

→ • More data

• Regularization.

## Bias Variance Tradeoff:

In the pre deep learning era, when bias increases variance goes down and when variance increases bias goes down.

But now we can focus on one without having other much. Regularization is a very useful technique for reducing variance. It might increase bias a little bit.

### — Regularizing your Neural Network

#### 1) Regularization:

Frobenius norm formula

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

Note:

The rows "i" of the matrix should be the number of neurons in the current layer  $n^{(l)}$ ; whereas the columns "j" of the weight matrix should equal the number of neurons in the previous layer  $n^{(l-1)}$ .

### — To prevent overfitting.

## Logistic Regression

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

Regularization parameter

$w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \rightarrow L_2 \text{ Regularization}$$

Because you are using Euclidean norm

Why do we just  $w$  not  $b$ ?

Because  $w$  has a lot of parameters as  $n_x$  while  $b$  is only a single parameter

$$\frac{\lambda}{2m} \sum_{i=1}^{n_x} |w_i| = \frac{\lambda}{2m} \|w\|_1, \quad w \text{ will be sparse}$$

which means there will be a lot of zeros.

Note:  $L_1$  Regularization can not be used for suppressing the model.

$\lambda \rightarrow$  is the hyperparameter which we may have to tune

Above will be  $L_2$  regularization for Logistic Regression and now for

Neural Network:

$$J(w^0, \dots, w^{(l)}, b^{(l)}) = \frac{1}{m} \sum_{i=1}^m d(\hat{y}^i, y^i)$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|^2$$

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

$$w: (n^{(l)} \ n^{(l-1)})$$

$$dw^{(l)} = (\text{from back prop}) + \frac{\lambda}{m} w^{(l)}$$

$$w^{(l)} = w^{(l)} - \alpha dw^{(l)}$$

Note:

Sometime L2 regularization is called as weight decay.

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

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

2) Why Regularization Reduces Overfitting?

3) Dropout Regularization:

Other technique from L2 regularization is dropout Regularization.

- We do drop out the nodes.

Implementing dropout ("Inversed dropout")

Illustrate with layer  $l=3$

$d3 = np.random.rand(a3.shape[0],$

$a3.shape[1]) < \text{keep\_prob}$

$\text{keep\_prob}=0.8$

It means there is 20% chance

of eliminating hidden units.

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

We will multiply both so that there will zeros so that matrix will be reduced.

$$a_3 /= 0.8 \text{ or keep-prob}$$

What is happening here.

Let's say we have 50 units then that means 20% will be zeroed out which is approximately that means.

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

↳ reduced by 20%.

So we will divide it by 0.8 so that value of  $z$  can be remain same as by division, value will increase.

Making predictions at test time

No drop out.

#### ④ Understanding Dropout:

Dropout prevents overfitting.

One big drawback that there will be no exact cost function.

#### ⑤ Other Regularization Methods.

- o Data augmentation.

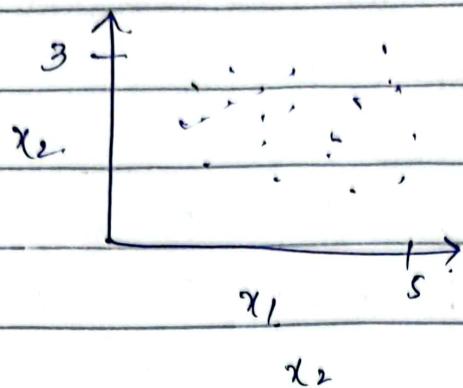
- o (make additional fake data)

- o Early Stopping

- o It is beneficial because we don't have to use Lambda values more often.

## - Setting Up your Optimization Problem:

### 1) Normalizing Inputs



Subtract means:

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

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

Normalize variance

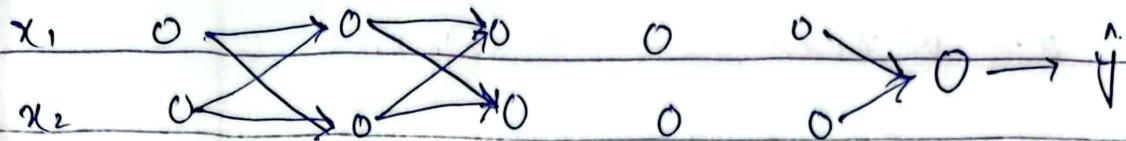
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^i - \bar{x})^2$$

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

Why normalize inputs?

### ② Exploding/Vanishing Gradients

Sometimes slope can be very high or very very small.



$$w^{(1)} \quad w^{(2)} \quad w^{(3)} \quad \dots \quad w^{(l)}$$

$$g(z) = z \quad b^{(l)} = 0$$

$$y = w^{(l)} w^{(l-1)} \dots w^{(1)} x$$

$$z' = w^{(1)} x$$

$$a' = g(z') = z'$$

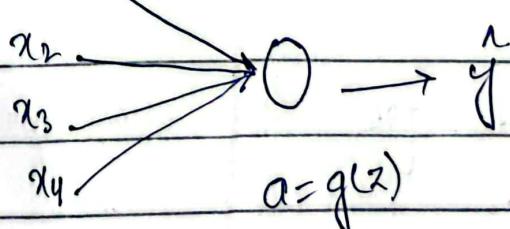
$$a^2 = g(z^2) = g(w^2 a')$$

$$w^l = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \quad y = w^l \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{l-1}$$

### ③ Weight Initialization for Deep Networks

The problem from Exploding/Vanishing gradient is that when there is large neural network then the weights can have very large slope or very small. Now the solution for that is by initializing weights.

Single neuron example



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

large  $n \rightarrow$  small  $w_i$

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

This is for ReLU activation function  
For tanh.

$$\tanh \sqrt{\frac{1}{n^{l-1}}}$$

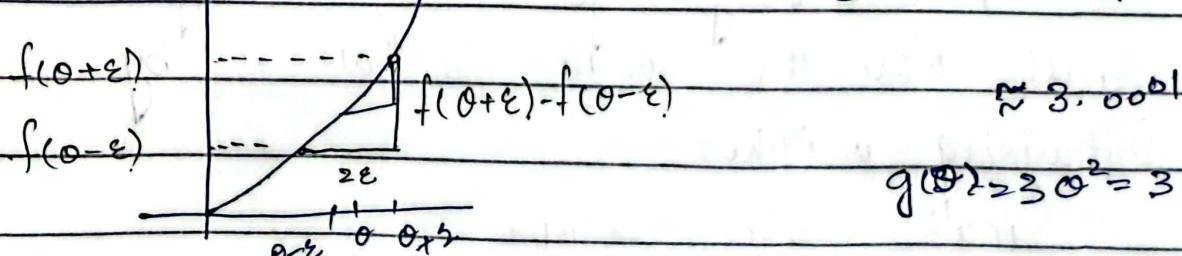
For Xavier activation

$$\sqrt{\frac{2}{n^{l-1} + n^l}}$$

#### (4) Numerical Approximation of Gradients

$$f'(0) = 0^3$$

$$\frac{f(0+\epsilon) - f(0-\epsilon)}{2\epsilon} \approx g(0)$$



$$0.99 \quad 1 \quad 1.01$$

$$\epsilon = 0.01$$

We do use this for less error production

## ⑤ Gradient Checking:

help find bugs in the back prop.

Take  $w^l, b^l, \dots w^p b^p$  and reshape into a big vector  $\theta$

Same for  $dw^l, db^l, \dots db^p$  and reshape into  $d\theta$

Grad check

$$J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots)$$

for each  $i$

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

$$\approx d\theta^{(i)} = \frac{\delta J}{\delta \theta_i}$$

$$\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta\|_2} \approx 10^{-7}$$

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

Then that means

the difference is correct or

u can say there is small value and implementation is correct

## ⑥ 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 drop out.  
(For this we use keep-prob is 1.0).

## WEEK 2

### - Optimization Algorithms

#### ① Minibatch gradient Descent

$$X = [x^1 \ x^2 \ x^3 \ \dots \ x^m] \quad (n \times m)$$

$$Y = [y^1 \ y^2 \ \dots \ y^m] \quad (1 \times m)$$

let's say we have  $m = 500,00,000$

then we will do because it will take a lot of times.

To do well against this, we will divide the data in the mini-section of 1000.

$$X = \underbrace{\left[ x^1 \ x^2 \ \dots \ x^{1000} \right]}_{X^{[1]}} \ \underbrace{\left[ x^{1000} \ \dots \ x^{2000} \right]}_{X^{[2]}} \ \dots \ \underbrace{\left[ x^m \right]}_{X^{(1000)}}$$

same for  $Y$

for  $t = 1 \dots 5000$  1 step of gradient descent using  $x^t, y^t$

Forward prop on  $x^{(t)}$  for 1000

$$z^t = w^{(1)} x^{(t)} + b^{(t)} \quad \left. \right\} \text{vectorize}$$

$$A^t = g^{(1)} z^{(1)} \quad \left. \right\} \text{implementation}$$

$$A^t = g^t z^t$$

Compute cost

$$\hat{J}^{(t)} = \frac{1}{1000} \sum_{i=1}^l d(\hat{y}^i, y^i) + \frac{\lambda}{2 \cdot 1000} \|w^0\|_F^2$$

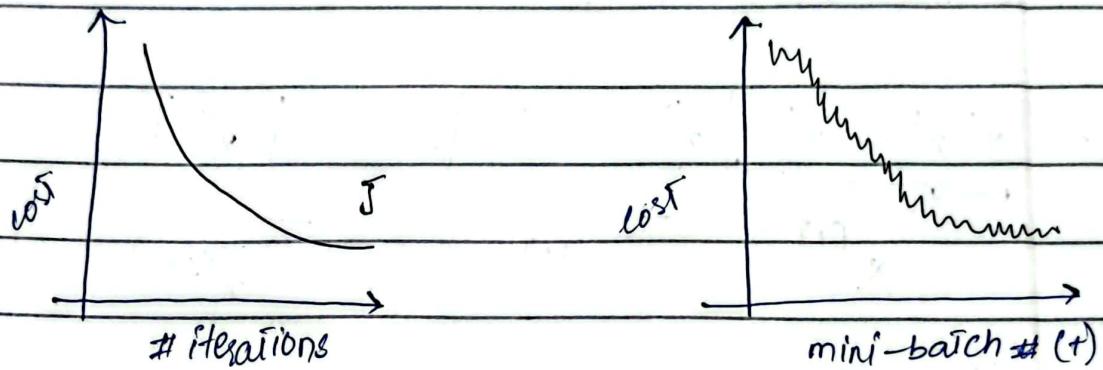
Backprop to compute gradient w.r.t  
 $\tilde{z}^{(t)}$  (using  $x^{(t)}, y^{(t)}$ )

$$w^{(l)} = w^{(l)} - \alpha d w^{(l)}, b^l = b^l - \alpha d b^l$$

② Understanding Mini-batch Gradient Descent:

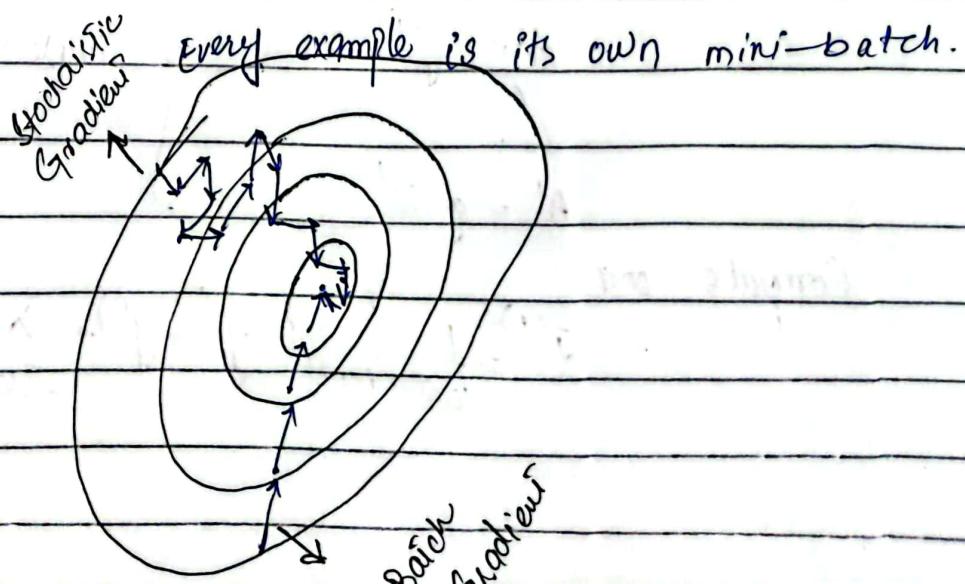
Batch Gradient  
Descent

Mini-batch gradient  
Descent



If mini-batch size =  $m$ : Batch gradient descent  
 $(x^l, y^l) = (x, y)$

" " " size = 1 : Stochastic gradient  
descent



## In-between

Too long per  
iteration.

(Batch Gradient  
Descent)

Stochastic Gradient Descent

(Loose speedup from  
vectorization)

Faster learning

- Vectorization

- make progress  
without needing to wait.

Choosing your mini-batch size:

- If small training set: use batch  
gradient descent

- Typical mini-batch sizes:

64, 128, 256, 512

$2^6, 2^7, 2^8, 2^9$

- Make sure mini-batch fit in CPU/GPU  
memory.

③ Exponentially Weighted Averages:

Because data is a bit noisy

when plotted. So make it look cooler. we will  
do

$$V_0 = 0$$

$$V_1 = 0.9V_0 + 0.1\theta_1$$

$$V_2 = 0.9V_1 + 0.1\theta_2$$

:

$$V_t = 0.9V_{t-1} + 0.1\theta_t$$

So, this will be called as exponentially weighted average.

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

As the values of  $\beta$  increases, then averages also takes place which has slow effect.

#### ④ Understanding Exponentially weighted averages:

- Key component for several optimized algorithms that you used to train your neural networks.

- It takes very little memory.

- This is known as bias correction.

#### ⑤ Bias Correction in exponentially Weighted.

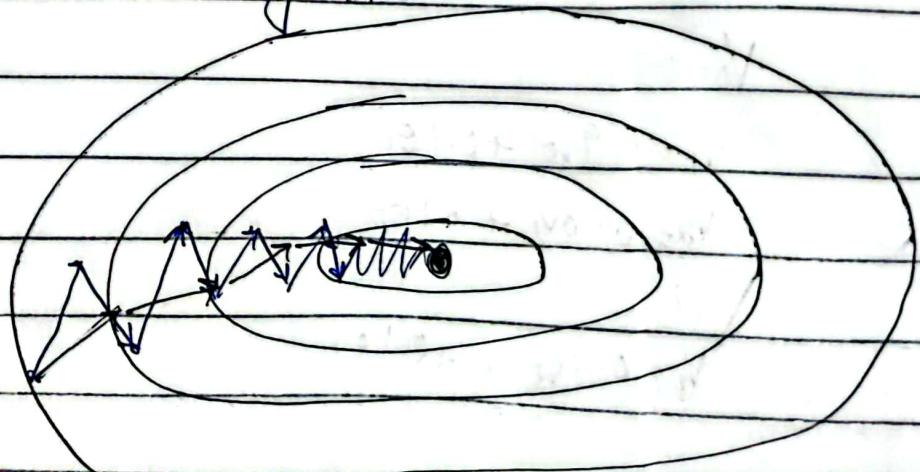
Average :

When we initialize with zero then it starts with very low. So, to resolve that

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

#### ⑥ Gradient Descent with Momentum:

- Always work faster than gradient descent algorithm



This will become up and down oscillation with ~~smaller~~ smaller learning rate.  
We want faster learning, so with Momentum:

On iteration  $t$ :

compute  $d\omega, db$  on current mini-batch

$$V_{d\omega} = \beta V_{d\omega} + (1-\beta) d\omega$$

$$\theta_t = \theta_0 + \beta V_0 + (1-\beta) \theta_t$$

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

$$\omega = \omega - \alpha V_{d\omega}, b = b - \alpha V_{db}$$

Hyperparameters:  $\alpha, \beta$        $\beta = 0.9$

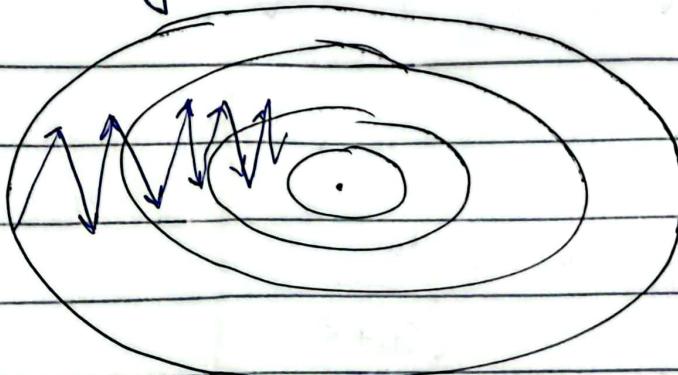
Sometimes, you will see

$$V_{d\omega} = \beta V_{d\omega} + d\omega$$

Because, it does not matter for small values —?

## ⑦ RMSprop

algorithm



We want to speed up in horizontal with smaller vertical.

compute  $d\omega, db$  on current mini-batch

$$S_{d\omega} = \beta S_{d\omega} + (1-\beta) d\omega^2 \xrightarrow{\text{element wise}} \text{small}$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2 \xrightarrow{\text{large}}$$

$$\omega := \alpha \frac{d\omega}{\sqrt{S_{d\omega}}} \quad b := \alpha \frac{db}{\sqrt{S_{db}}}$$

Because  $S_{d\omega}$  is smaller so.

by dividing it will be smaller then  $\omega$  will be larger.

and same for  $b$  but its value will be smaller so give us free. Now you can use larger  $\alpha$  rate.

RMSprop  $\rightarrow$  Root Mean Square prop

We are going to use  $\beta_2$  for RMSprop

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

$$\omega := \omega - \alpha \frac{d\omega}{\sqrt{S_{d\omega} + \epsilon}}$$

(A) Adam Optimization Algorithm

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

On iteration  $t$ :

compute  $dw, db$  using current batch

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

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

After Bias correction

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

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

$$w_i^t = w_i - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corr}} + \epsilon}}$$

$$b^t = b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corr}} + \epsilon}}$$

$\alpha$  needs to tune

$$\beta_1: 0.9. \quad (\text{dw})$$

$$\beta_2: 0.9999 \quad (\text{dw}^2)$$

$$\epsilon: 10^{-8} \quad (\text{By Adam's Paper})$$

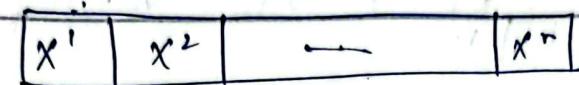
Adam: Adaptive momentum estimation

## ⑨ Learning Rate Decay:

$$\alpha_t = \frac{\alpha_0}{1 + \text{decay Rate} \times \text{epoch Number}}$$

We need to converge because if  $\alpha$  is larger then it will be wandering, but if it is smaller then it will converge.

1 epoch = 1 pass through data



first Epoch  $\hookrightarrow$  second Epoch

Epoch	$\alpha$	$\alpha_0 = 0.2$ decay rate = 1
1	0.1	
2	0.067	
3	0.05	
4	0.04	

Other learning rate decay methods

$\alpha = \alpha_0 \cdot \text{epoch}^{-\frac{1}{n}}$   $\alpha_0 \rightarrow$  exponential decay.

## ⑩ The Problem of Local Optima:

# WEEK 3

## HYPERPARAMETER TUNING, BATCH NORMALIZATION AND PROGRAMMING FRAMEWORKS

### Hyperparameter Tuning

#### ① Tuning Process

By default

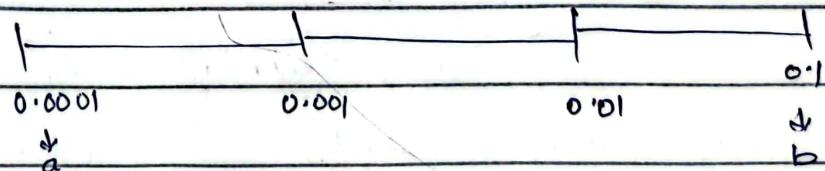
$$\beta_2 = 0.9$$

$$\beta_1, \gamma, \epsilon$$

$$0.4 \quad 0.999 \quad 10^8$$

#### ② Using Appropriate Scale to pick Hyperparameters.

We use random value from sampled for hyper parameter. So for uniform sampling



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

$$\alpha = 10^{-n}$$

$$\lambda \in [a, b]$$

$$\beta = 0.9 \dots 0.999$$



10 values

1000 values

$$1 - \beta = 0.1 \rightarrow 0.001$$

$$\downarrow \quad \downarrow$$
$$10^{-1} \quad 10^{-3}$$

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

$$\beta = 1 - 10^{-x}$$

### ③ Hyperparameters Tuning in Practice: Pandas vs Caviar

- Retest your hyper-parameters once in a several months.
- Babysitting one model (One approach) → (Pandas strategy)
- Training many models in parallel (2nd approach) → (Caviar strategy)

These 2 approaches depends on computation power, we do have

### ④ Normalizing Activations in a Network:

This technique can make very easier search for hyper-parameter but not for all neural network

$$X = X/\sigma$$

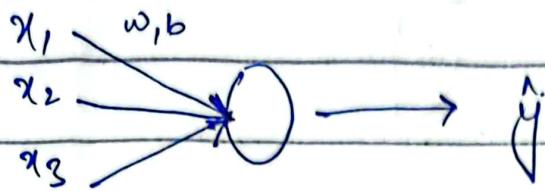
→ Normalizing inputs to speed up learning.

$$\bar{y} = \frac{1}{m} \sum x^i$$

$$X = X - \bar{y}$$

$$S^2 = \frac{1}{m} \sum_{i=1}^m x^i$$

$$X = X/S.$$



When we have a deeper network  
Batch norm will normalize  $z^{(2)}$  for  
 $w^{(3)}, b^{(3)}$  or  $a^{(3)}$ .

Implementing Batch Norm

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

$$S^2 = \frac{1}{m} \sum_i (z^{(i)} - \bar{y})^2$$

$$z_{\text{norm}} = \frac{z^{(i)} - \bar{y}}{\sqrt{S^2 + \epsilon}}$$

- Batch Normalization

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

$\gamma, \beta$  ↴ learnable parameters of model

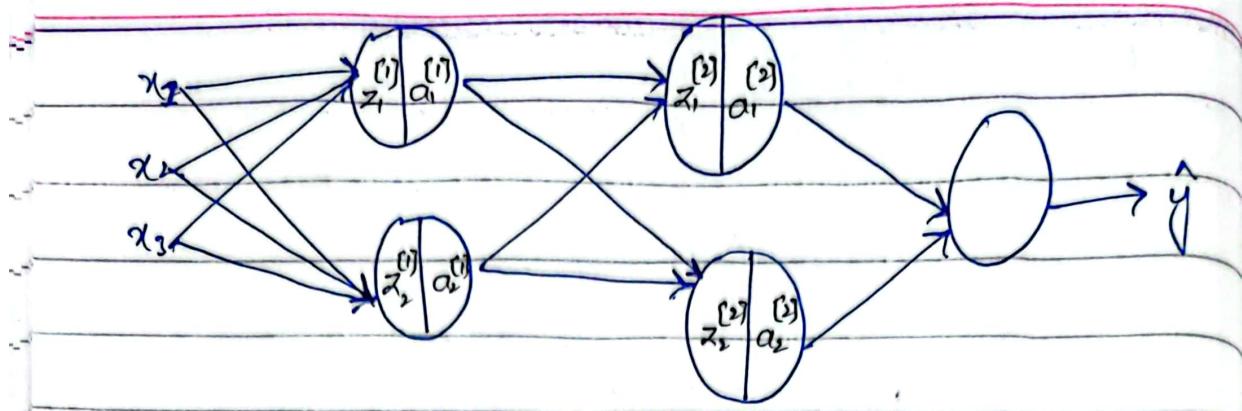
$$\gamma = \sqrt{S^2 + \epsilon}, \beta = 0$$

then

$$\hat{z}^{(i)} = z^{(i)}$$

How does it work from  
scratches?

Fitting Batch Norm into a Neural  
Network:



$$\begin{aligned}
 & \text{Input } x \xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow{\substack{\beta^{(1)}, \gamma^{(1)} \\ \text{Batch Norm}}} \hat{z}^{(1)} \rightarrow a^{(1)} = g^{(1)}(\hat{z}^{(1)}) \\
 & w^{(2)}, b^{(2)} \xrightarrow{} z^{(2)} \xrightarrow{\substack{\beta^{(2)}, \gamma^{(2)} \\ BN}} \hat{z}^{(2)} \rightarrow a^{(2)} = g^{(2)}(\hat{z}^{(2)})
 \end{aligned}$$

Parameters:  $w^{(1)}, b^{(1)}, \dots, w^{(l)}, b^{(l)}$  }  $d\beta^{(l)}$   
 $\beta^{(1)}, \gamma^{(1)}, \dots, \beta^{(l)}, \gamma^{(l)}$  }  $\beta := \beta^{(l)} - \alpha d\beta^{(l)}$

These are different from

Adam's hyper parameters.

For programming we can use

tf.nn.batch\_normalization.

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

Batch Norm. will remove  $b^{(l)}$  by subtracting mean and dividing variance.

$$\hat{z}^{(l)} = w^{(l)} a^{(l-1)}$$

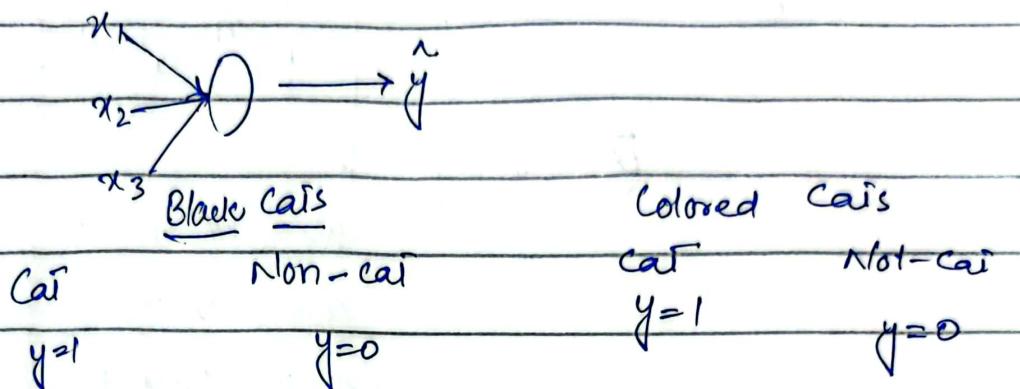
$$\hat{z}^{(l)} = \gamma^{(l)} \hat{z}_{\text{norm}} + \beta^{(l)}$$



$\vec{z}^{(e)}_{(n^{(e)}, 1)} \rightarrow n^{(e)}$  no. of hidden units.

$\vec{s}^{(p)} = (n^{(e)}, 1) \rightarrow$  Dimension.

③ Why does Batch Norm work?



of the distribution of  
if changes then you have  
to retain the model.  
Suppose complex value for  
3rd hidden layer  
 $w^{[3]}, b^{[3]}$  want to go to  $j$

Batch norm will reduce the variance  
distribution for the specific layer. It ensures the  
mean and variance for the layers remains  
same.

④ Batch Norm at Test time:

Batch Norm process mini-batches  
and one mini-batch at a time while  
at test time it has to predict / process  
one example / complete dataset at a time.

This is the equation at  
Training time.

$$y = \frac{1}{m} \sum_i z^i$$

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

$$z_{norm}^i = \frac{z^i - y}{\sqrt{\delta^2 + \epsilon}}$$

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

At test time, we expect  
is one

$y, \delta^2$ : estimate using

exponentially weighted  
average (across mini-  
batch)

$x^{(1)}, x^{(2)}, x^{(3)}, \dots$

$$y^{(1)est}, y^{(2)est}, \dots \rightarrow y$$

$\delta^2$

We will use the average  
across this.

So at test time

$$z_{norm} = \frac{z - y}{\sqrt{\delta^2 + \epsilon}} \quad z = \gamma z_{norm} + \beta$$

But average  $\delta^2, y$

- Multi-class Classification

① Softmax Regression:

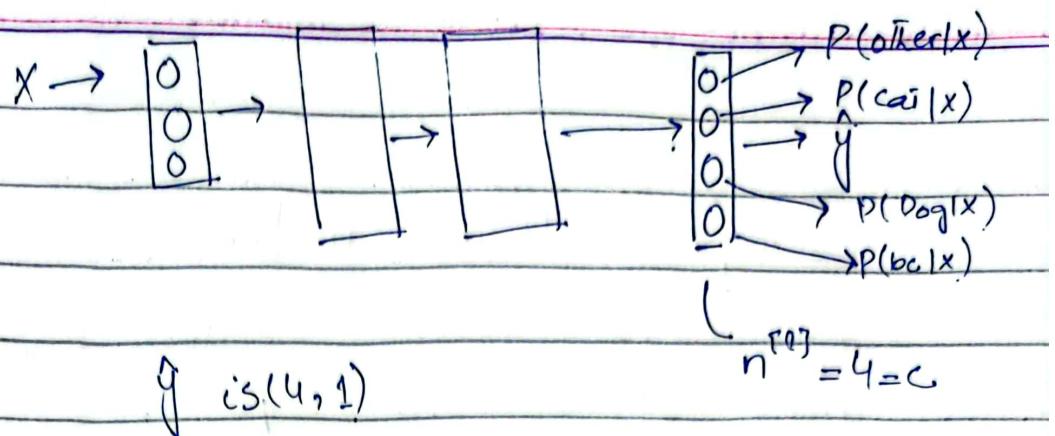
$$a^{(c)} = \frac{e^{z^{(c)}}}{\sum_{i=1}^C t^i}$$

• Generalization of Logistic Regression

When there is more than two classes

then

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



$y$  will give us the probabilities whose sum will be equal to 1.

We use softmax layer for the final layer to output these probabilities.

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

Activation function:

$$t = e^{z^{(l)}} \quad (4, 1)$$

Element wise exponentiation

$$a^{(l)} = \frac{e^{z^{(l)}}}{\sum_{i=1}^4 t_i} \quad (4, 1)$$

$$a_i^{(l)} = \frac{t_i}{\sum_{i=1}^4 t_i}$$

e.g.

$$z^{(l)} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}, \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.1 \\ 20.1 \end{bmatrix}$$

$$\sum_{i=1}^4 t_i = 176.3$$

$$\frac{e^5}{176.3} = P(\text{other}|x) = 0.842$$

$$\frac{e^2}{176.3} = P(\text{cat}|x) = 0.042$$

$$\frac{e^{-1}}{176.3} = P(\text{Dog}|x) = 0.002$$

$$\frac{e^3}{176.3} = P(\text{bird}|x) = 0.114$$

$$a^{(r)} = g^{(r)}(z^{(r)})$$

## ② Training a Softmax Classifier:

$$z^{(1)} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \quad c=4$$

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

Softmax

Hardmax

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

Softmax regression generalizes logistic regression to  $c$  classes.

If  $c=2$ , softmax becomes logistic regression.

Loss function

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow y_2=1 \quad a^{(2)} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

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

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

$$-y_2 \log \hat{y}_2 = -\log \hat{y}_2$$

make  $\hat{y}_2$  large to make loss small

$$L(w_b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^i, y^i)$$

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

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

(4, m)

Gradient descent with softmax.

Backprop:

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

$$\frac{\partial J}{\partial z^{(l)}}$$

## — Introduction to Programming Frameworks

### ① Deep Learning Frameworks :

→ Caffe / Caffe2

→ CNTK

→ DL4J

→ Keras

→ Lasagne

→ mxnet

→ PaddlePaddle

→ TensorFlow

→ Theano

→ Torch

### ② TensorFlow :

$$J(w) = w^2 - 10w + 25$$

cost

we have this function and how we will define it Tensorflow to minimize

To startup Tensorflow

import numpy as np

import tensorflow as tf

① -  $w = \text{tf.Variable}(0, \text{dtype}=\text{tf.float32})$

② - ~~define train step()~~:  
define ~~train step()~~:  $\text{optimizer} = \text{tf.keras.optimizers.Adam}(0.1)$

④ ← with  $\text{tf.GradientTape}()$  as tape:

To compute  $\text{grads}$  (prop) ③ -  $\text{cost} = w^{**2} - 10w + 25$

⑥ - trainable-variables = [w]

⑦ -  $\text{grads} = \text{Tape.gradient}(\text{cost}, \text{trainable-variables})$

⑧ -  $\text{optimizer.apply_gradients(zip(grads, trainable-variables))}$

⑨ -  $\text{print}(w)$

• w is a parameter we want to know so we declare it as a variable.

When there is a dataset then cost function

① ,  $x = \text{np.array}([1.0, -10.0, 25.0]) \text{dtype=np.float32}$

② , def training(x, w, optimizer):  
def cost\_fn():

return  $x[0]^*w^{**2} + x[1]^*w + x[2]$

for i in range(1000):

optimizer.minimize(cost\_fn, Pw)

return w

w=training(x, w, optimizer)

↓  
↓