# Hyperparameter tuning

---

# Tuning process

# Hyperparameters

first importance one circled in red
then second the one in orange
third in importance is the one in blue
He never tunes the parameters for adam algo optim, just use those as
defoult.

this is the most
important hyperparameter
to tune. Next most imp
is the momentum term.

$\alpha$

$\beta$  NO.9  momentum term beta if you are using momentum

$\beta_1$ , $\beta_2$ , $\varepsilon$    for adam optimization algo
0.9      0.999   $10^{-8}$

# layers    maybe you have to pick the number of layer

#hidden units    maybe you have to pick the number of units
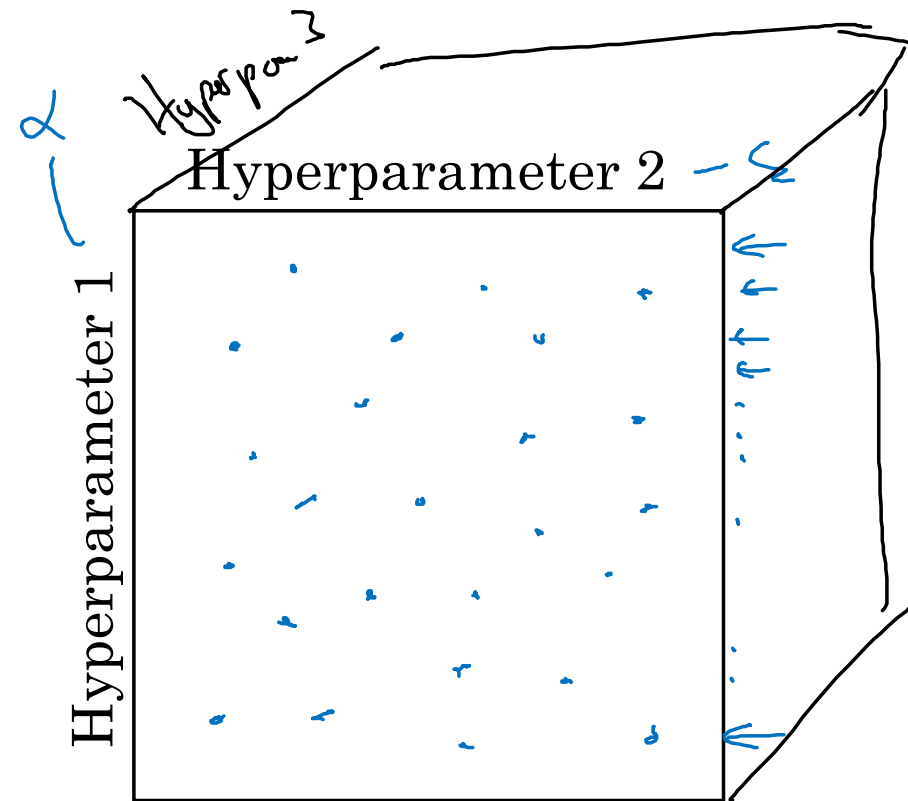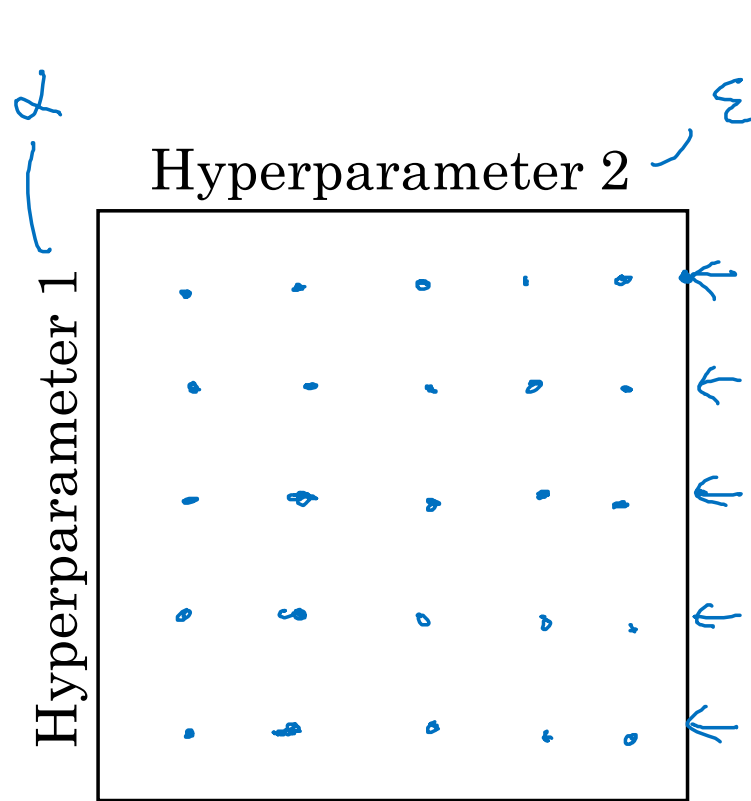
learning rate decay

mini-batch size   you might need to choose the mini-batch size

Andrew Ng

HOW DO YOU SELECT A SET OF VALUES TO EXPLORE??????
IN EARLY GENERATIONS IF YOU HAD TWO HYPERPARAMETERS IT WAS COMMON PRACTICE TO SAMPLE THE POINTS IN A GRID
LIKE BELOW AND SYSTEMATICALLY EXPLORE THESE VALUES, YOU TRY OUT ALL 25 VALUES(CAN BE MORE) AND PICK
WHICHEVER HYPERPARAMETER WORKS BEST. THIS WORKED WHEN NR OF HYPERPARAMETERS WAS SMALL.

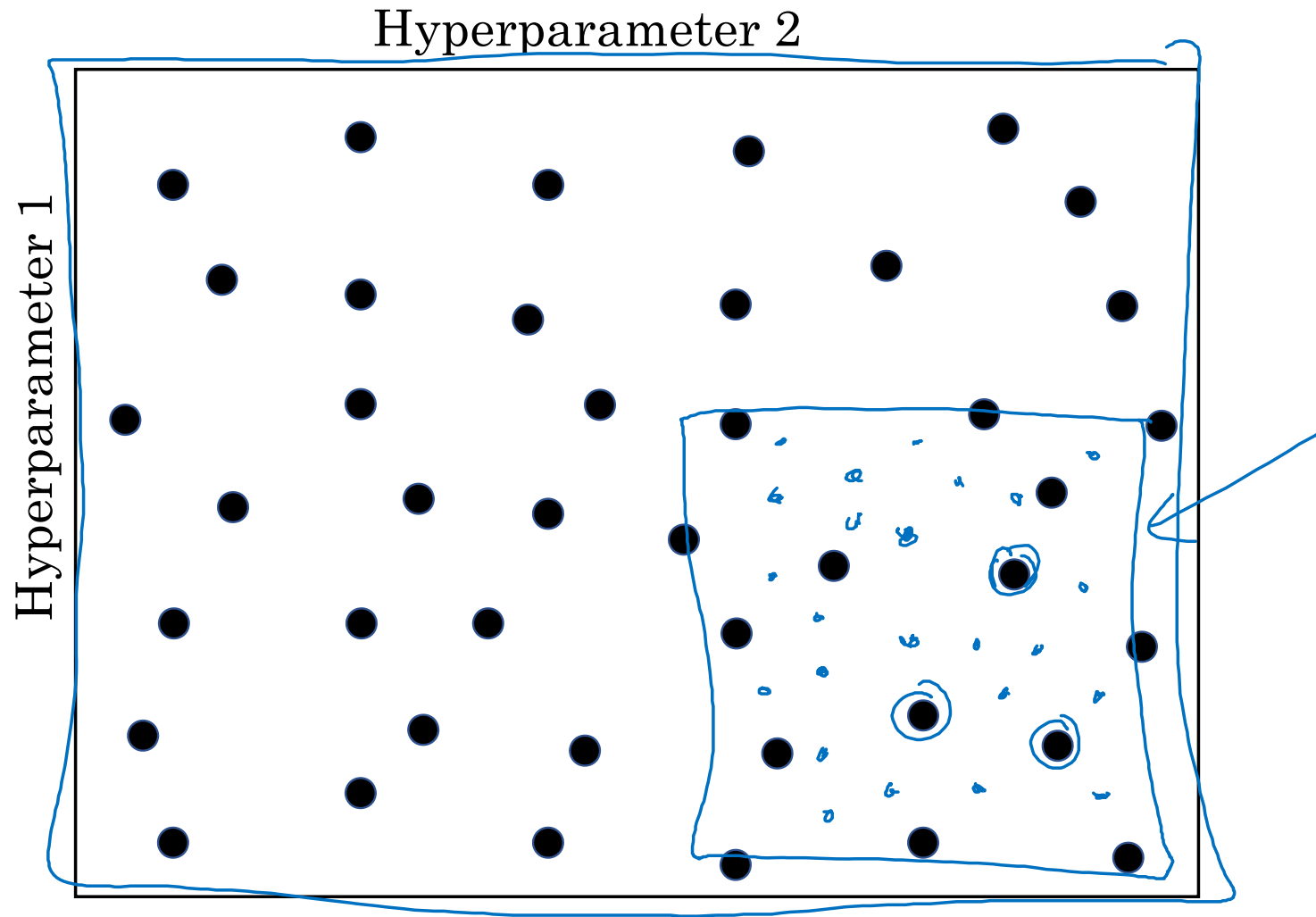# Try random values: Don't use a grid



IN DEEP LEARNING WHAT WE TEND TO DO IS TO CHOOSE POINTS AT RANDOM. SO U CHOOSE 25 POITS AT RANDOM AND THEN TRY OUT HYPERPARAMETER ON
THIS RANDOMLY CHOSEN SET OF POINTS. ITS DIFFICULT TO KNOW IN ADVANCE WHICH HYPERPARAMETER ARE GOING TO BE THE MOST IMPORTANT FOR YOUR
PROBLEM, AND AS WE SAW PREVIEWSLY SOME HYPERPARAMETERS ARE MORE IMPORTNT THEN OTHERS
THATS SAY THE FIRST HYPERPARAMETER IS ALFA AND THE SECOND IS EPSILON THAT IS IN DENOMINATOR OF ADAM ALGO
SO YOUR CHOISE OF ALFA METTERS A LOT WHLE CHOISE OF EPSILON HARDLY. SO WHEN YOU SAMPLE IN THE GRID YOU TRY OUT FIVE VALUES OF ALFA AND YOU
MIGHT FIND THAT ALL OF THE DIFFERENT VALUES OF EPSILON GIVE YOU ESSENTIALLLY THE SAME ANSWER.
SO NOW YOU HAVE TRAINED 25 MODELS AND ONLY GOT INTO TRIAL FIVE VALUES FOR THE LEARNING RATE ALFA, WHICH I THINK IS REALLY IMPORTANT.
WHEREAS IN CONTRAST, IF YOU WERE TO SAMPLE AT RANDOM YOU WILL HAVE TRIED OUT 25 DISTINCT VALUES AND THEREFORE YOU WILL BE MORE LIKELY TO
FIND A VALUE THAT WORKS REALLY WELL. WE EXPLAINED WITH JUST TWO HYPERPARAMETERS WHILE YOU CAN FACE THIS WITH MANY PARAMETERS.
ITS HARD TO KNOW IN ADVANCE WHICH HYPERPARAMETER ARE REALLY IMPORTANT FOR OUR APPLICATION.

Andrew Ng

WHEN YOU SAMPLE HYPERPARAMETERS ANOTHER COMMON PRACTICE IS TO USE A COARSE TO FINE SAMPLING SCHEME.
SO THATS SAY IN THIS TWO DIMENTIONAL EXAMPLE THAT YOU SAMPLE THESE POINTS AND YOU FIND THAT THIS POINT
WORKS BEST AND A FEW AROUND ALSO WORK OK.
IN COARSE OF THE FINAL SCHEME WHAT YOU MIGHT DO IS ZOOM IN TO A SMALLER REGION OF HYPERPARAMETERS AND
SAMPLE MORE DENSLY.

# Coarse to fine

Hyperparameter 2

Hyperparameter 1

Andrew Ng

WE SAW HOW SAMPLING OVER A RANGE OF HYPERPARAMETERS CAN ALLOW YOU TO SEARCH OVER THE SPACE OF HYPERPARAMETERS MORE EFFICIENTLY. IT TURNS OUT THAT SAMPLING AT RANDOM DOES NOT MEAN SAMPLING UNIFORMLY AT RANDOM OVER THE RANGE OF VALUES, INSTEAD IS IMPORTATNT TO PICK THE APPROPRIATE SCALE ON WHICH TO EXPLORE THE HYPERPARAMETERS, IN THIS VIDEO I WANT TO SHOW YOU HOW TO DO THAT.

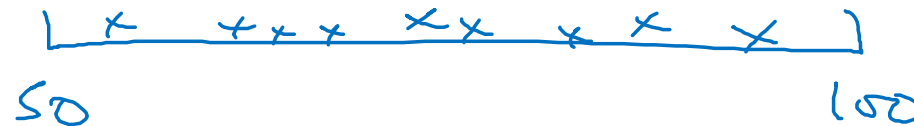deeplearning.ai

# Hyperparameter tuning

---

# Using an appropriate scale to pick hyperparameters

# Picking hyperparameters at random

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

SAY YOU THINK A GOOD RANGE OF VALUES IS 50 - 100

PICKING HERE UNIF RAND MIGHT BE OK

50            100

YOU THINK NUMBER OF LAYER IS BETWEEN 2 TO 4

$\rightarrow$ #layers        $L:$   2 — 4

2 , 3 , 4     ALSO HERE PICKING AT RANDOM
ITS OK, ALSO BUILDING A GRID SEARCH
WHERE EXPLICITLY EVALUATE THE VALUES
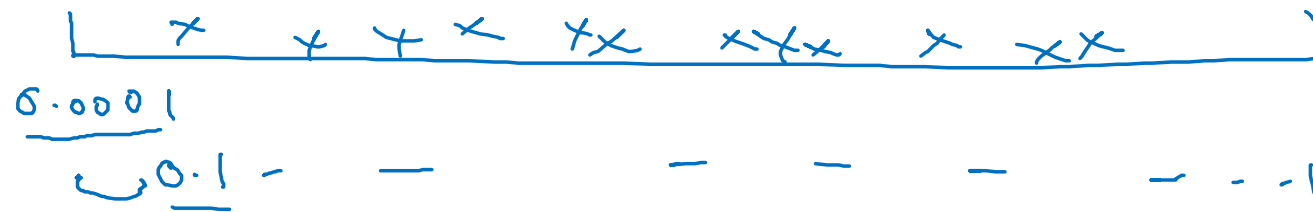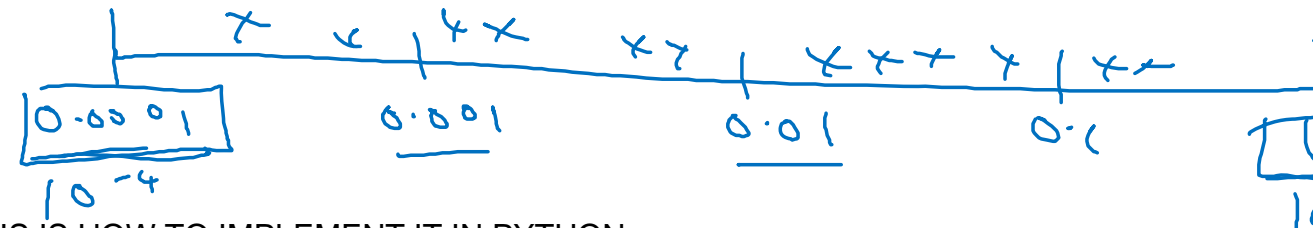2 3 4 ARE MIGHT BE REASONABLE
THIS IS NOT TRUE FOR ALL HYPERPARAMETERS

Andrew Ng

# Appropriate scale for hyperparameters

$$\alpha = 0.0001 , ....., 1$$

THIS IS RANGE YOU SUSPECT



SAMPLING UNIFORMLY AT RANDOM

$$0.0001$$

$$0.1$$

ITS BETTER TO SEARCH ON A LOG SCALE AND YOU HAVE MORE RESOURCES SEARCHING BETWEEN 0.001 ND 0.001 AND '.001

$$0.00001 \quad 0.001 \quad 0.01 \quad 0.1$$

$$10^{-4}$$

$$10^0$$

$$10^6$$

THIS IS HOW TO IMPLEMENT IT IN PYTHON

$$r = -4 * np.random.rand()$$

$$\leftarrow r \in [-4, 0]$$

r will be rand nr between -4 and 0

$$\alpha = 10^r$$

$$\leftarrow 10^{-4} ... 10^0$$

alfa will ne rand nr between above values

$$10^a$$

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

$$= -4$$

$$b = \log_{10} 1$$

$$= 0$$

$$10^a .... 10^b \qquad r \in [a, b] \qquad \alpha = 10^r$$

$$[-4, 0]$$

Andrew Ng

# Hyperparameters for exponentially weighted averages

$\beta = 0.9 \quad \cdots \quad 0.999$

you suspect beta has this range



0.9 is like averaging over 10 last days

while 0.999 is like averaging over 1000 days.

$1-\beta = 0.1 \quad \cdots \quad 0.001$

$\beta: 0.9000 \rightarrow 0.9005 \} \sim 10$

$\beta: 0.999 \rightarrow 0.9995$

$\sim 1000 \qquad \sim 2000$

we need to give more weight to changes near 0.999 since it has higher implications in the number of averaging days.

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

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

$1-\beta = 10^r$

$\beta = 1 - 10^r$

Andrew Ng

deeplearning.ai

# Hyperparameters tuning

---

# Hyperparameters tuning in practice: Pandas vs. Caviar

# Re-test hyperparameters occasionally

Idea



Experiment                    Code

- NLP, Vision, Speech, Ads, logistics, ….

- Intuitions do get stale. Re-evaluate occasionally.

Andrew Ng

# Babysitting one model



Panda ←

# Training many models in parallel



Caviar ←

Andrew Ng

Pand approach if u dont have much CPU otherwise go for caviar approach.

This makes the search for the hyperparameter much easy and makes the NN much more rubust. The choise of hyperparameters is a much bigger rnge of hyperparameters that work and also allows to train very deeo NN.
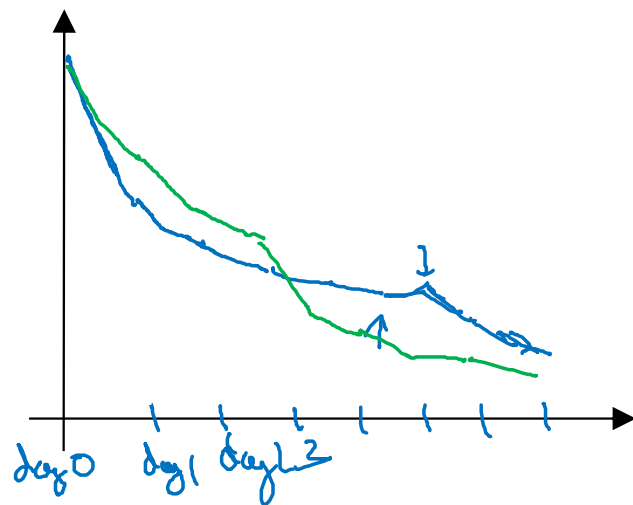Does not work allways for all NN but when it works it can make training very fast.

# Batch Normalization

---

# Normalizing activations in a network

# Normalizing inputs to speed up learning



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

$$X = X - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_i x^{(i)2} \quad \leftarrow \text{element-wise}$$

$$X = X / \sigma^2$$

Normalizing imput features can speed up learning

For any hidden layer can we normalize the values of any hidden layer to train faster

Can we normalize $\dfrac{a^{[2]}}{}$ so as to train $W^{[3]}, b^{[3]}$ faster

Normalize $\dfrac{z^{[2]}}{}$ ↑

should we normalize the value before activation function, so z2 or whether after applying activation function a2.
In practice normalizing z2 is done much often, so thats the version I'll present.

Andrew Ng

So here is how to implement batch norm to normalize the values in the hidden layer.

# Implementing Batch Norm

given some intermediate values in NN say those z1 z2 ...

Given some intermediate values in NN

$$z^{(i)}, \ldots, z^{(m)}$$

we gone ommit this l for semplification

$z^{[l](i)}$

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

compute mean var

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

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

← normalize

epsilon for numerical stability
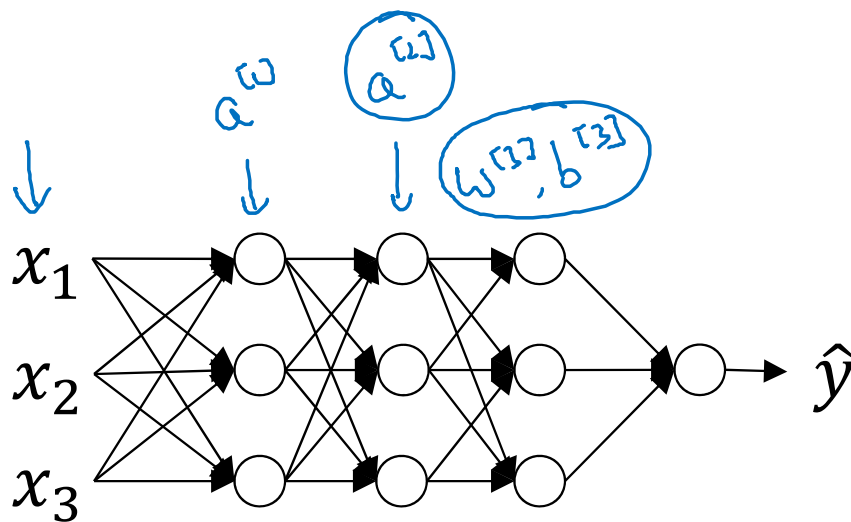
If

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

$$\beta = \mu$$

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

we saw how normalizing features x helps learning algo

$X$

$z^{(i)}$

What BatchNorm does it applys that to values in hidden layer

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

learnable parameters of model.

we dont want hidden layers to have mean 0 and var 1 allways, so we compute z tilda. gamma and beta are learnable parameters of your model.

Use $\tilde{z}^{[l](i)}$ instd of $z^{(i)}$

The effect of gamma and beta lets you set the mean and va fo z i whatever you want.

Here gammma and beta are learnable parameters of your model so if you are using grad descent or rms prop or adam you would update the weights of your nn just like for the weights of the NN.
The effect of gamma end beta is that it allows you to set of z tilde to be whatever you want it to be.
What it does its ensures that the hidden units have standardized mean and variance where mean and var are controlled by gamma and beta parameters that can make mean 0 and var 1 or other values as well, so it normalize hidden values to fixed mean and var, can be 0 and 1 or some other value

Andrew Ng

We see how to fit Batch norm ito a deep NN and how to make it work for the many different layers of NN
We will also give some intuition why Batch norm to train the NN.

We have seen how to fit batch norm for a single hidden layer thats see how it fits in the training of a deep network

# Batch Normalization

## Fitting Batch Norm into a neural network

deeplearning.ai

# Adding Batch Norm to a network
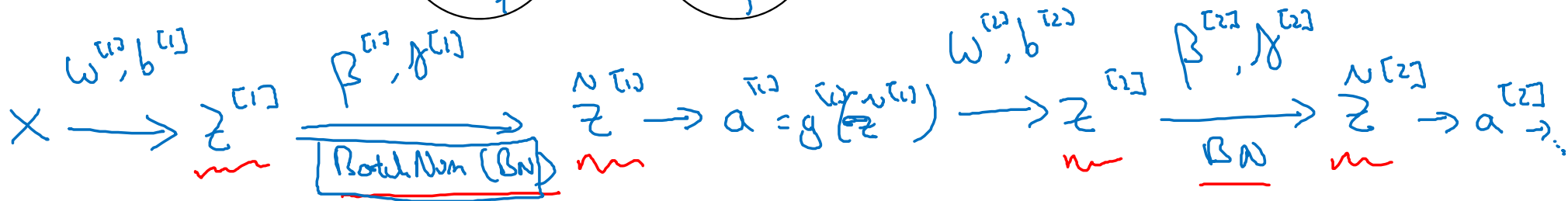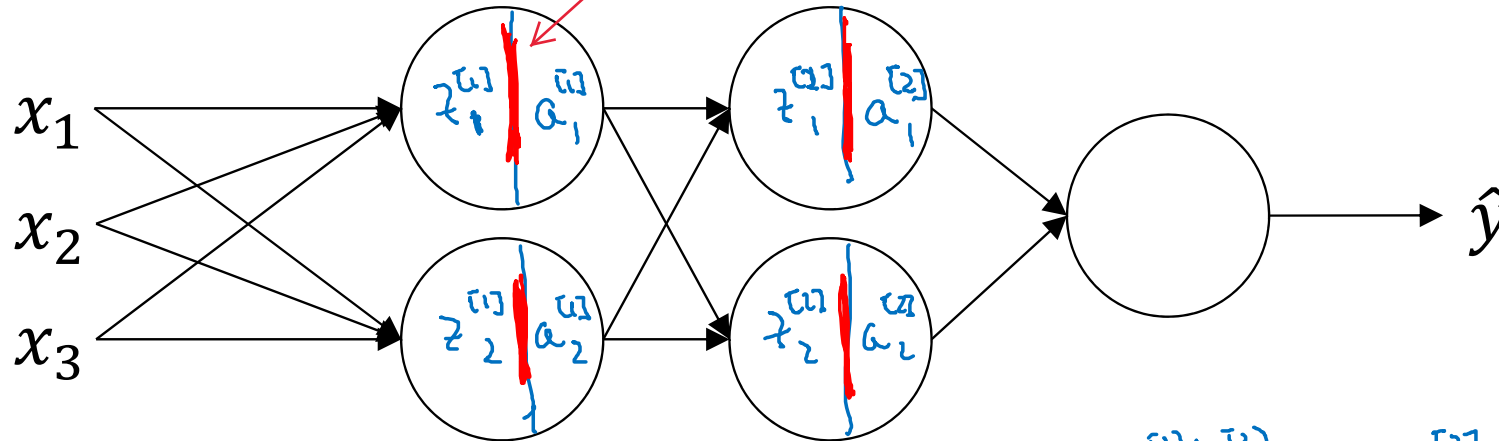
So first calcs the z and then applies activation function to compute a



Neural network diagram with inputs $x_1$, $x_2$, $x_3$ and hidden units showing $z_1^{[1]}, a_1^{[1]}$, $z_2^{[1]}, a_2^{[1]}$, $z_1^{[2]}, a_1^{[2]}$, $z_2^{[2]}, a_2^{[2]}$ leading to output $\hat{y}$

$$X \xrightarrow{\omega^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\boxed{\text{Batch Norm (BN)}}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \rightarrow a^{[1]} = g(\tilde{z}^{[1]}) \xrightarrow{\omega^{[2]}, b^{[2]}} z^{[2]} \xrightarrow[\underline{BN}]{\beta^{[2]}, \gamma^{[2]}} \tilde{z}^{[2]} \rightarrow a^{[2]} \rightarrow \ldots$$

$$\text{Parameters: } \left. \omega^{[1]}, b^{[1]}, \omega^{[2]}, b^{[2]}, \ldots, \omega^{[L]}, b^{[L]}, \atop \rightarrow \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \ldots, \beta^{[L]}, \gamma^{[L]} \right\}$$

$$\rightarrow \beta$$

$$d\beta^{[l]} \qquad \beta^{[l]} = \beta^{[l]} - \alpha \, d\beta^{[l]}$$

using tensorflow to implement batch norm:

$$\boxed{tf.nn.batch\text{-}normalization \Leftarrow}$$

these betas have nothing to do with the betas that we computed for momentum or expon weighted avg.
You can use grad descent or rms prop or adam to update the betas.

Andrew Ng

If you are using a Deep Learning programming framework usually you wont have to implement the Batch Norm step on Batch Norm layer yourself.
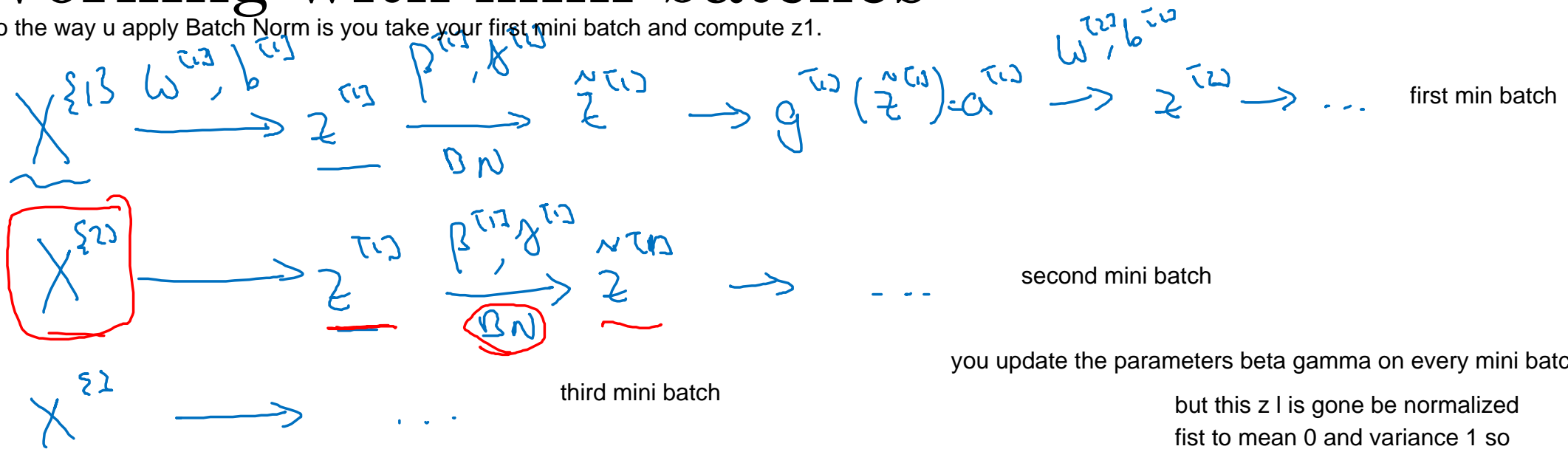In tensorflow you can implement Batch Normalization with the above function.
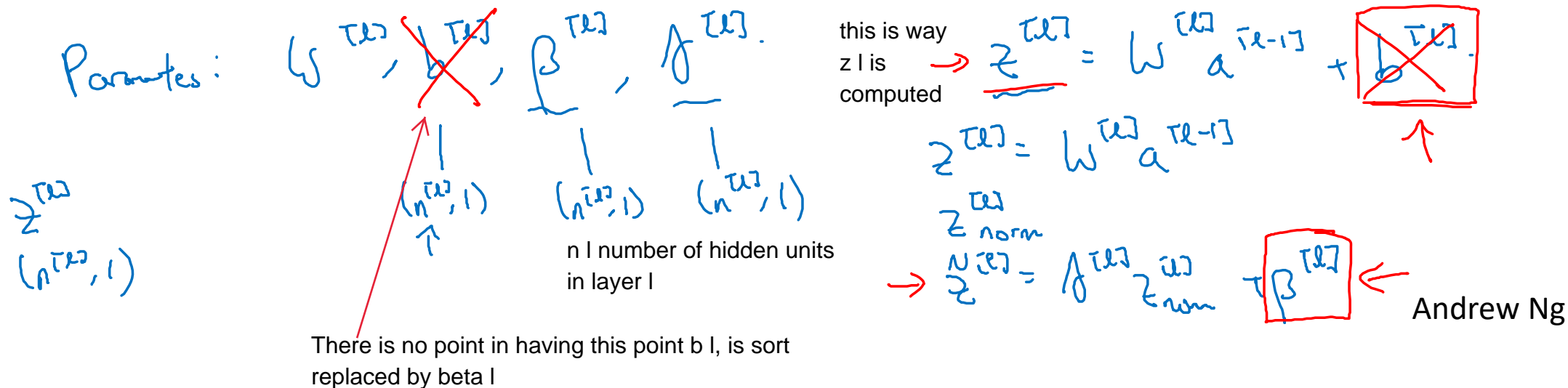
Now so far, we've talked about Batch Norm as if you were training on your entire training site at the time as if you are using Batch gradient descent.
In practice Batch Norm is usually applied with mini-batches of your training set.

# Working with mini-batches

So the way u apply Batch Norm is you take your first mini batch and compute z1.

$$X^{\{1\}} \xrightarrow{W^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow[BN]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow g^{[1]}(\tilde{Z}^{[1]}) = a^{[1]} \xrightarrow{W^{[2]}, b^{[2]}} Z^{[2]} \rightarrow \dots$$

first min batch

$$X^{\{2\}} \longrightarrow Z^{[1]} \xrightarrow[BN]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow \dots$$

second mini batch

$$X^{\{3\}} \longrightarrow \dots$$

third mini batch

you update the parameters beta gamma on every mini batch

but this z l is gone be normalized fist to mean 0 and variance 1 so adding that constant b is irrelevant

Parmeters: $W^{[l]}, \cancel{b^{[l]}}, \beta^{[l]}, \gamma^{[l]}$.

$(n^{[l]}, 1)$  $(n^{[l]}, 1)$  $(n^{[l]}, 1)$

n l number of hidden units in layer l

$Z^{[l]}$
$(n^{[l]}, 1)$

this is way z l is computed $\rightarrow$ $Z^{[l]} = W^{[l]} a^{[l-1]} + \cancel{b^{[l]}}$

$Z^{[l]} = W^{[l]} a^{[l-1]}$

$Z^{[l]}_{norm}$

$\rightarrow \tilde{Z}^{[l]} = \gamma^{[l]} Z^{[l]}_{norm} + \beta^{[l]} \leftarrow$

There is no point in having this point b l, is sort replaced by beta l

Andrew Ng

So thats put all together and describe how to implement grad descent usikng Batch Norm.
Assume we are using mini- batch grad descent.

# Implementing gradient descent

for $t = 1$ .... num Mini Batches  (number of mini batches)

Compute forward prop on $X^{\{t\}}$.

In each hidden layer, use BN to repar $z^{[l]}$ with $\tilde{z}^{[l]}$.

Use backprop to compute $dW^{[l]}$, ~~$db^{[l]}$~~, $d\beta^{[l]}$, $d\gamma^{[l]}$

Update params $W^{[l]} := W^{[l]} - \alpha \, dW^{[l]}$
$\beta^{[l]} := \beta^{[l]} - \alpha \, d\beta^{[l]}$
$\gamma^{[l]} := \cdots$
$\Bigg\}$ ←

Works w/ momentum, RMSprop, Adam.

Works with momentum, RMSprop, Adam where instead of taking this gradient descent update you can use the updates given by these other algorithms as we discussed in the previews week.
So these other optimization algo can be used to update the parameters beta and gamma that Batch Norm added to algorithm.

Andrew Ng

So we have seen how normalizing imputs can speed up learning.
The intuition is that this is doing a similr think but for values in your hidden units and not just in the inputs.
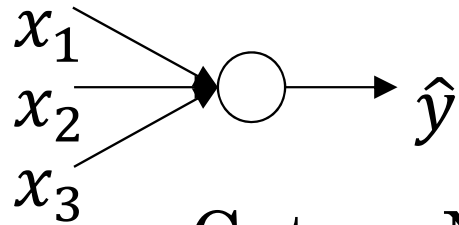
# Batch Normalization

Why does Batch Norm work?

deeplearning.ai

A second reason why batch norm works is it makes weights later or deeper in your network more rubust to changes to weights in earlier layers of the NN, say, in layer one.
To motivate the above that look at this exmple.

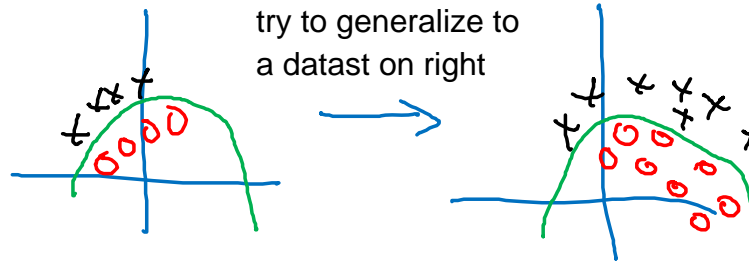# Learning on shifting input distribution

$x_1$

$x_2$

$x_3$

$\hat{y}$

Cat      Non-Cat

$y = 1$   $y = 0$

try to generalize to a datset on right

You might not expect a model trained on data on left to do well on data on right. They might be the same function that actually works, but you wouldnt expect your learning algo to discover that green decision boundary just looking to data on left.

$y = 1$   $y = 0$

So this idea of data distribution changing goes by the name ^covariance shift^

"Covriote shift"

$x \longrightarrow y$

If you have learned some X->Y mapping if the distribution of X changes then you might need to retrain your learning algo. This is true even if the function mapping from x to y reamins unchanged, which it is in this example as the function is telling is cat or not.

you have trained NN on black cats

apply it to colored cats, then classifier might not do well
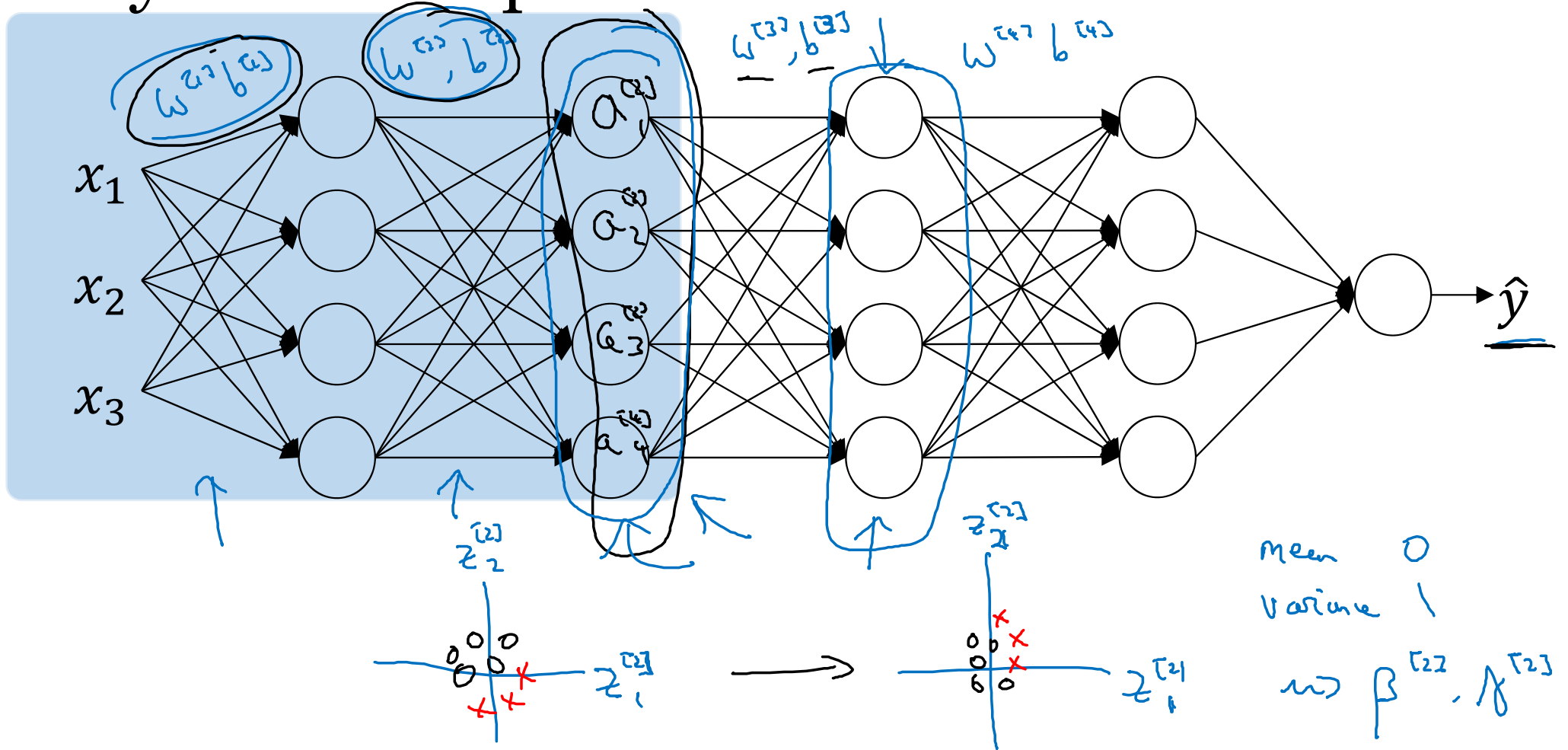
Andrew Ng

How does this problem of covariance shif applies to a NN.
Consider a deep network like this, and thats look at the procss from the perspective of this certain layer, the third hidden layer.
So this network has learned the parameters W3 and b3. From the perspective of the third hidden layer it gets some set of values from the earlier layers, and then has to do some stuff to hopefully make the output y close to groud true value.

# Why this is a problem with neural networks?



if we conver up part on left of third layer and see from perspective of this third layer.

Andrew Ng

# 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.

$\rightarrow \tilde{z}^{[l]}$

$64, 128$

$z^{[l]}$

$X^{\{t\}}$

$\mu, \sigma^2$

Mini-batch : $\underline{64} \longrightarrow \underline{512}$

Batch norm handles data one mini batch at a time, it computes mean and variance on mini batches. So at test time you try to make predictions, try and evaluate the NN, you might not have a mini batch example, you might be processing one single example at a time, so you might need to do something different to make sure your predictions make sense.

# Batch Normalization

---

# Batch Norm at test time

deeplearning.ai

# Batch Norm at test time

m number of example in the mini batch
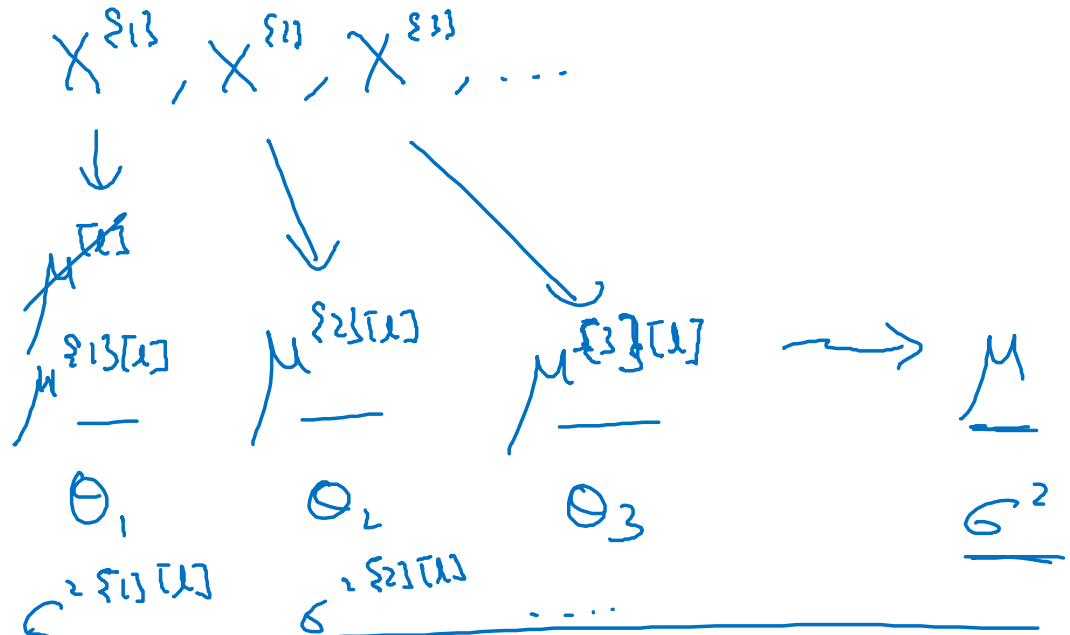
At test time you have to compute mean and variance
and if you have just one example computing the mean
and variance ot that example doesnt make sence.
So we come up with some separate estimate of my and sigma

What u do to estimate those is to comp expon weighted avg the mini batches

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

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

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

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

$\mu, \sigma^2$ : estimate using exponentially weighted average ( across mini-batches ).

$X^{\{1\}}, X^{\{2\}}, X^{\{3\}}, \ldots$

$\mu^{\{1\}}$

$\mu^{\{1\}[l]}$     $\mu^{\{2\}[l]}$     $\mu^{\{3\}[l]}$     $\to \mu$

$\theta_1$     $\theta_2$     $\theta_3$     $\sigma^2$

$\sigma^{2\{1\}[l]}$     $\sigma^{2\{2\}[l]}$     $\ldots$

$z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \varepsilon}}$     $\tilde{z} = \gamma z_{\text{norm}} + \beta$

So here at test time
we use the mean and var
of the exponentially weighted

If we use a framework to compute NN we dont need to worry much about this as its easily done there.
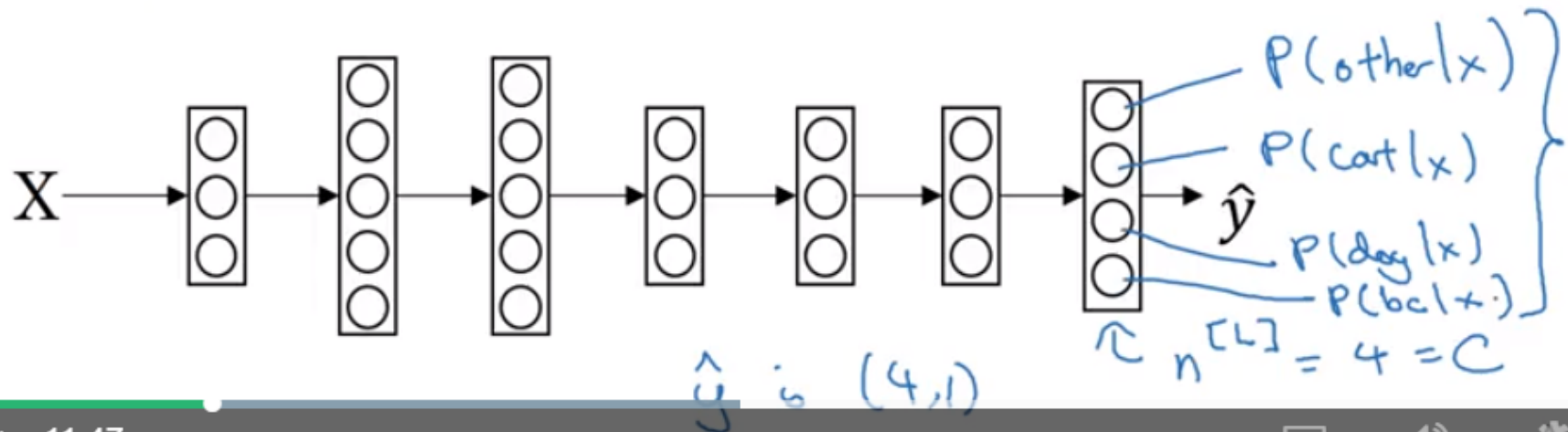
Andrew Ng

Multi-class
classification

---

Softmax regression

deeplearning.ai

# Recognizing cats, dogs, and baby chicks, *other*

*0*

*1*     *2*     *3*



3    1    2    0    3    2    0    1

$C = \#classes = 4$     $(0, \dots, 3)$



$$X \rightarrow \hat{y}$$
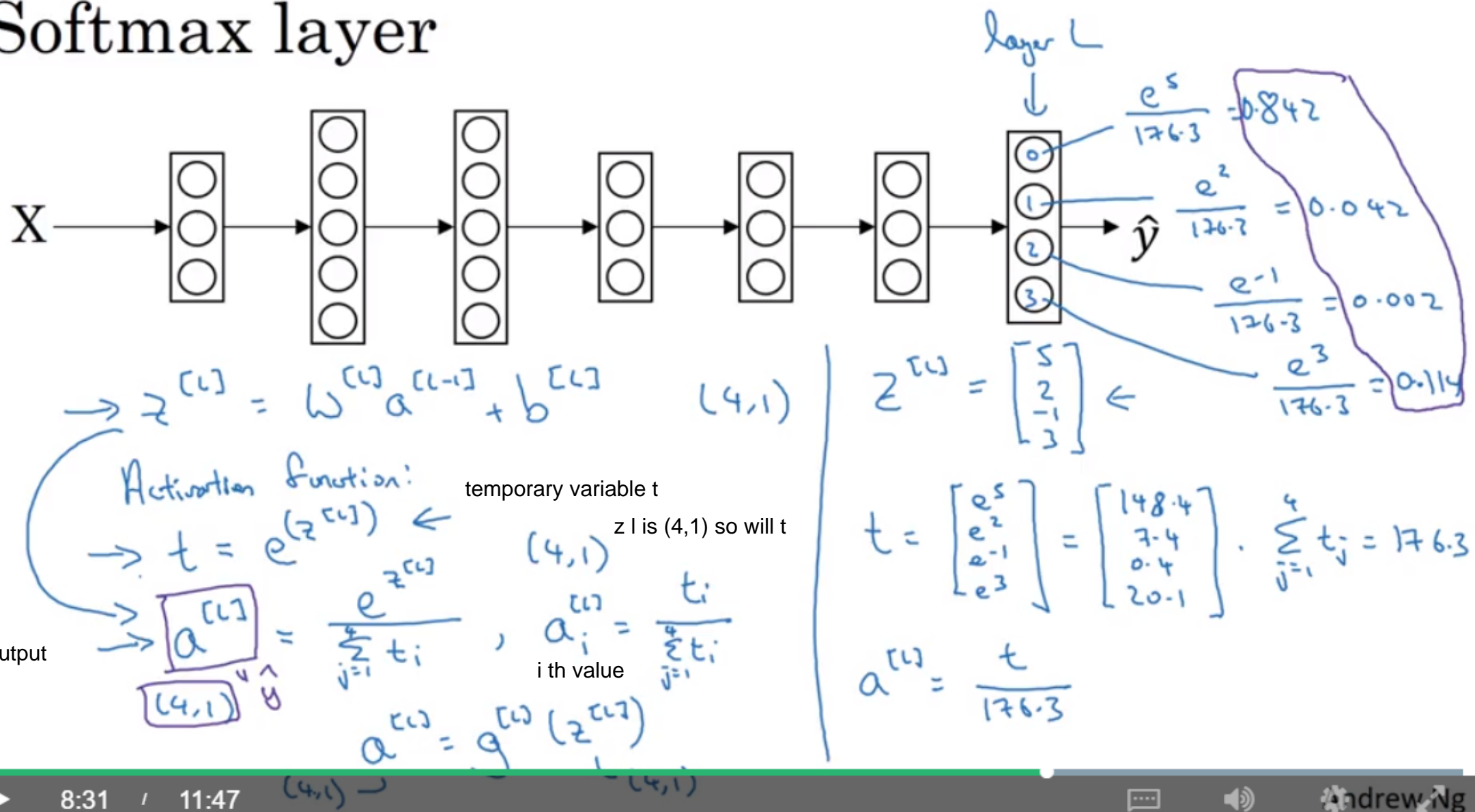
$P(other|x)$
$P(cat|x)$
$P(dog|x)$
$P(bc|x)$

$n^{[L]} = 4 = C$

$\hat{y}$ is $(4,1)$

2:41 / 11:47

Andrew Ng

# Softmax layer

layer $L$



$X \longrightarrow \hat{y}$

$\dfrac{e^5}{176.3} = 0.842$

$\dfrac{e^2}{176.3} = 0.042$

$\dfrac{e^{-1}}{176.3} = 0.002$

$\dfrac{e^3}{176.3} = 0.114$

$$\longrightarrow z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]} \qquad (4,1)$$

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

Activation function:

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

temporary variable t

z l is (4,1) so will t

$$\longrightarrow \boxed{a^{[L]}} = \dfrac{e^{z^{[L]}}}{\sum\limits_{j=1}^{4} t_i} \quad , \quad a_i^{[L]} = \dfrac{t_i}{\sum\limits_{j=1}^{4} t_i}$$

output

$(4,1)$ $\hat{y}$

i th value

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

$(4,1)$ $(4,1)$

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \cdot \sum\limits_{j=1}^{4} t_j = 176.3$$

$$a^{[L]} = \dfrac{t}{176.3}$$
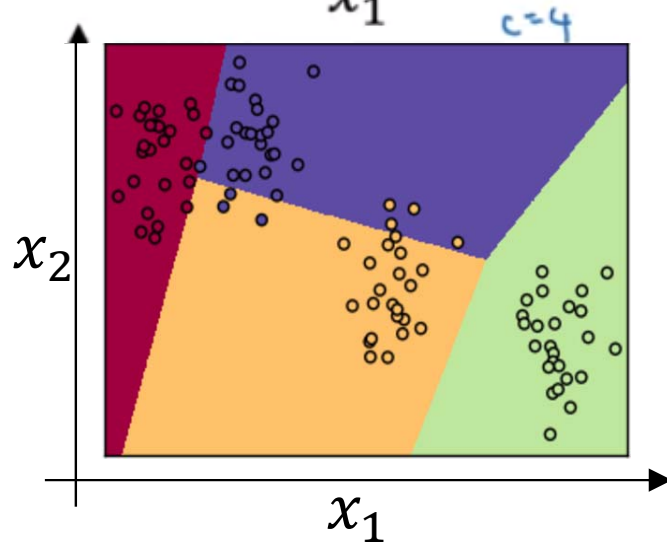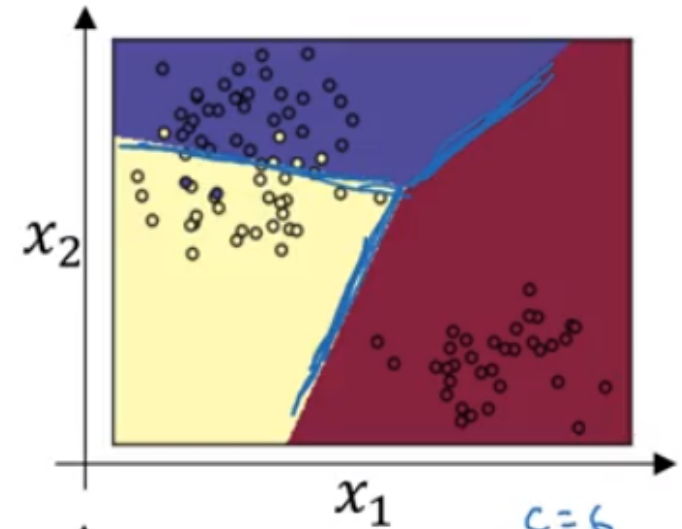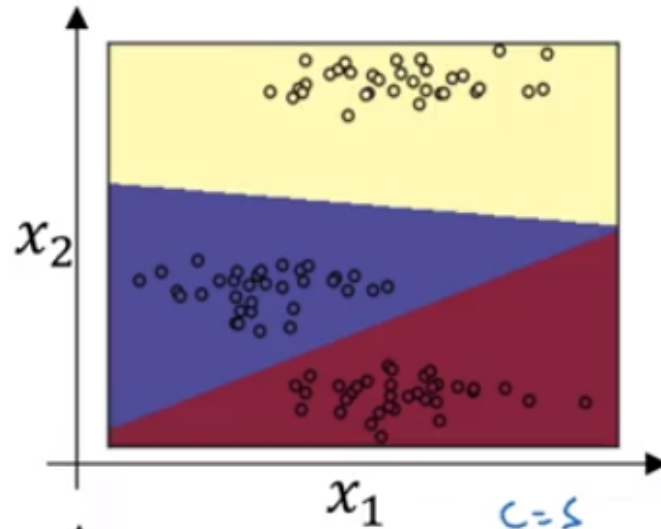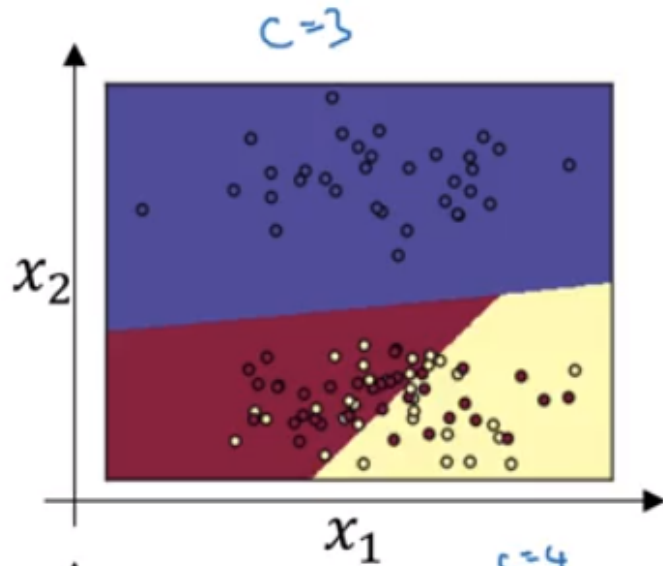
8:31 / 11:47

Andrew Ng

Andrew Ng

So it just described the activation function.
The only strange think is takes a vector and outputs a vector

# Softmax examples

$x_1$
$x_2$ $\rightarrow$ $\hat{y}$

$z^{[1]} = \omega^{[1]} x + b^{[1]}$
$a^{[1]} = \hat{y} = g(z^{[1]})$



c=3

c=4

C=5

c=6

To be noticed that the decision boudary is linear.

We saw ssome example of soft max and the activation function of the softmax classifier.

We will see how to train a NN that uses a Softmax layer

# Multi-class classification

---

# Trying a softmax classifier

deeplearning.ai

# Understanding softmax

since we have four classes z l is a vect of (4,1)

$(4,1)$

we compute t which is this exponentiation function

this would be the activation function

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

$C = 4$

$g^{[L]}(\cdot)$

"soft max"

hard max take one with high prob and set it to 1

"hard max"

we apply act fct to z l and get a

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

biggest element in z was 5 and now the biggest probability is that one .842

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

soft max comes from contrast to hard max which would have mapped z to this vector of 0 and 1, would put one in biggest elem of z.

Softmax regression generalizes logistic regression to $C$ classes.

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

Andrew Ng

Thats see how you would train NN with softmax.
Thats define the loss function u use to train the nn

# Loss function

SAY OUR NN IS OUTPUTTIN THIS Y HAT WITH PROBABILITIES ITS A CAT

$(4,1)$

TAKE EX
WITH THIS
TARGET

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \begin{matrix} - cat \\ - y_2 = 1 \end{matrix}$$

$$y_1 = y_3 = y_4 = 0$$

$(4,1)$

$$a^{[L](i)} = \overset{\wedge}{y}^{(i)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$C = 4$

THIS IS COST ON ALL TRAINING SET COST FCT. U USE
GRAD DESCENT TO MIN THIS

THIS IS
THE LOSS
WE USE:

$$\mathcal{L}(\overset{\wedge}{y}, y) = - \sum_{j=1}^{4} y_j \log \overset{\wedge}{y_j}$$

small

$$J(w^{[i]}, b^{[i]}, \ldots) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\overset{\wedge}{y}^{(i)}, y^{(i)})$$

THE ONLY WAY TO MAKE THE COST SMALL ITS TO MAKE -LOG
SMALL. WHICH MEANS TO MAKE Y HAT BIG

U LEFT
ONLY WITH THIS

$$- y_2 \log \overset{\wedge}{y_2} = - \log \overset{\wedge}{y_2}.$$

Make $\overset{\wedge}{y_2}$ big.

THIS IS
THE VECTORIZED
IMPLEMENTATIO

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \ldots & y^{(m)} \end{bmatrix}$$

$$\overset{\wedge}{Y} = \begin{bmatrix} \overset{\wedge}{y}^{(1)}, & \ldots & , y^{(m)} \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \ldots$$

$$= \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \ldots$$

$(4, m)$

7:29 / 10:07

Andrew Ng

Andrew Ng

# Gradient descent with softmax

We have talked how to implement the forward pro, what about the backpropagation.

$$\frac{z^{[L]}}{(4,1)} \longrightarrow a^{[L]} = \hat{y} \rightarrow \mathcal{L}(\hat{y}, y)$$

Backprop:

$$\underset{(4,1)}{dz^{[L]}} = \hat{y} - y.$$

$$\frac{\partial J}{\partial z^{[L]}}$$

that dz is our
prtial deriv of the cost fct wrt z

The key eq you need to initialize backprop is this expression.
That the derivative wrt z at the loss layer, this turns out, you can compute as that on left.

In this week prog exercise we will start using a deep learning prog framework and in this contest you just need to focus on getting the forward pro right. So we dont need bother about deriv calc as the framework takes care of that.
Thats it for soft max classifications where u characterize inputs into not just one of two classes but one of C different classes          Andrew Ng

As you implement more complex models such as CNN and RNN, or as you start to implement very large models that is not practical to implement everythink yourself.
Its good to know how to do matrix multiplicatio but as u go into more complex model you probably dont want to do that, instead u call a numerical linear algebra library that could do it more efficiently for you.

Programming
Frameworks

_____

Deep Learning
frameworks

deeplearning.ai

# Deep learning frameworks

There are many deep learning frameworks, here are some of
the leading ones

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

Choosing deep learning frameworks
- Ease of programming (development and deployment)
- Running speed
→ - Truly open (open source with good governance)

each of these frameworks has a community of users and developers
They are evolving and you can see the pros and cons of each of these ones

Andrew Ng

deeplearning.ai

Programming
Frameworks

TensorFlow

# Motivating problem

say we have this simple cost fct to min

$$J(w) = \boxed{w^2 - 10w + 25}$$

(cost)

$$\nearrow \quad (w-5)^2$$

$$w = 5$$

this is the solution but thats
say we dont know that
and we try solve with tesnsorflow

$$J(w, b)$$
$$\uparrow \quad \uparrow$$

# Code example

```python
import numpy as np
import tensorflow as tf

coefficients = np.array([[1], [-20], [25]])

w = tf.Variable([0],dtype=tf.float32)
x = tf.placeholder(tf.float32, [3,1])
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0]    # (w-5)**2
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
init = tf.global_variables_initializer()
session = tf.Session()
session.run(init)
print(session.run(w))

for i in range(1000):
    session.run(train, feed_dict={x:coefficients})
print(session.run(w))
```
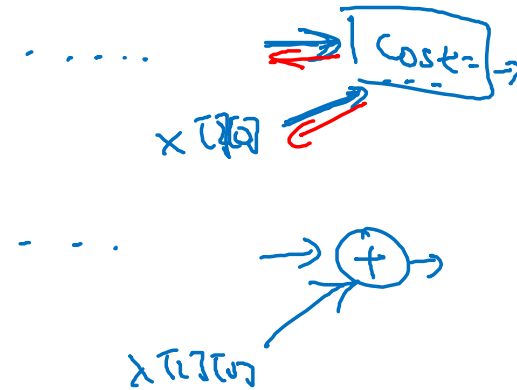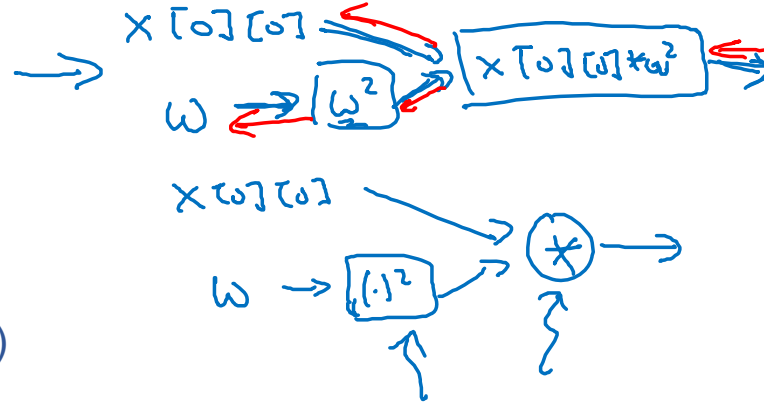
```python
with tf.Session() as session:
    session.run(init)
    print(session.run(w))
```

this is an alternative
form of those three lines
which are quite idiomatic.
does same think.

Andrew Ng