

# Hyperparameters

$\alpha$

$\beta$  0.9

$\beta_1, \beta_2, \epsilon$   
0.9 0.99

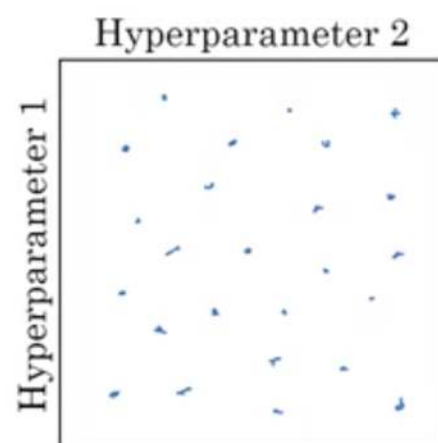
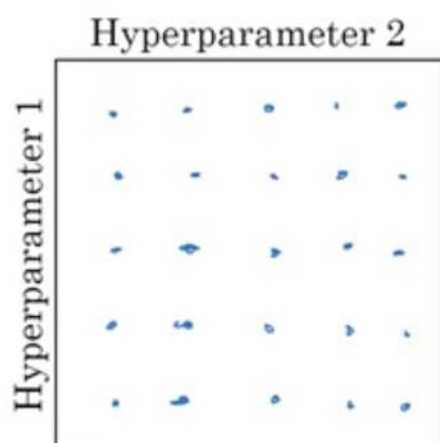
#layers

#hidden units

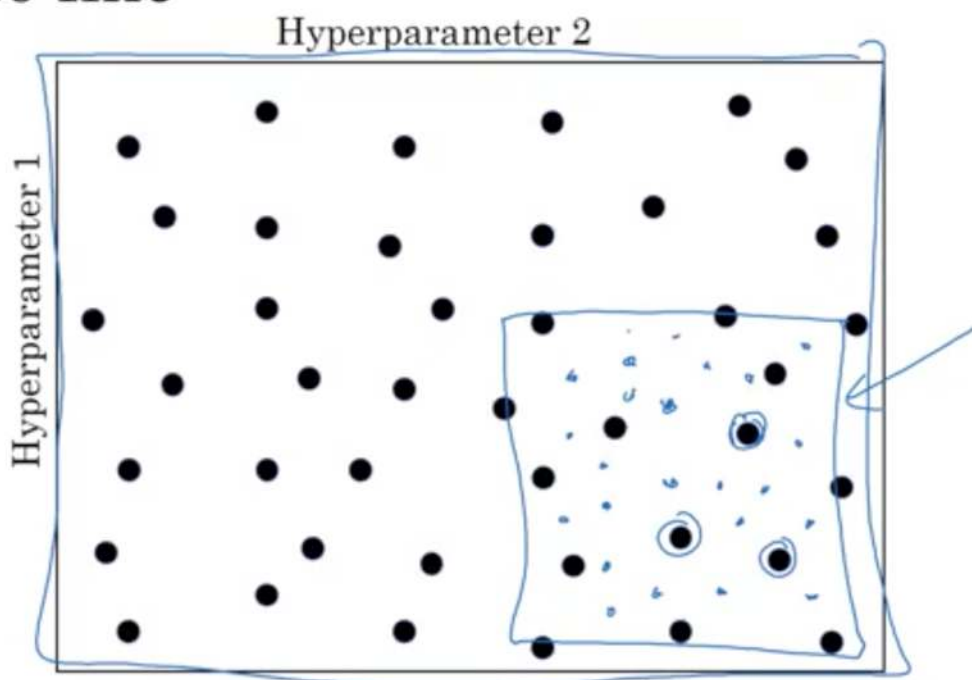
learning rate decay

mini-batch size

# Try random values: Don't use a grid

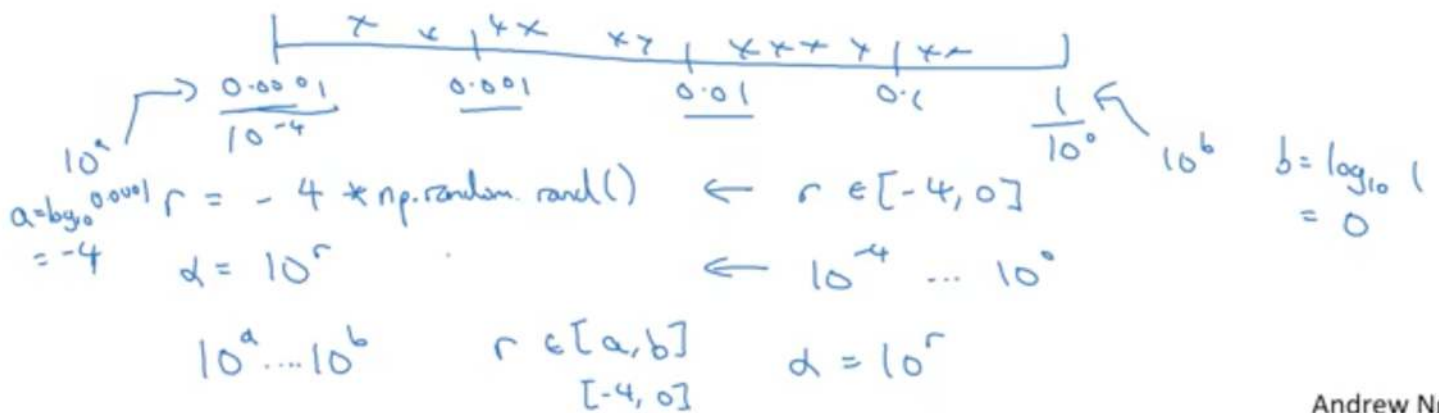
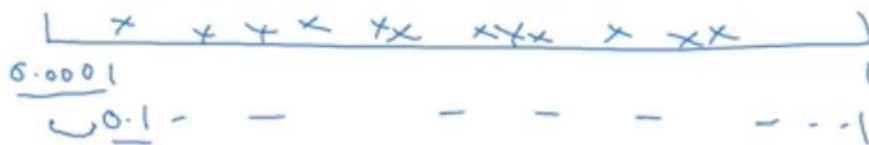


# Coarse to fine



## Appropriate scale for hyperparameters

$$\alpha = 0.0001, \dots, 1$$

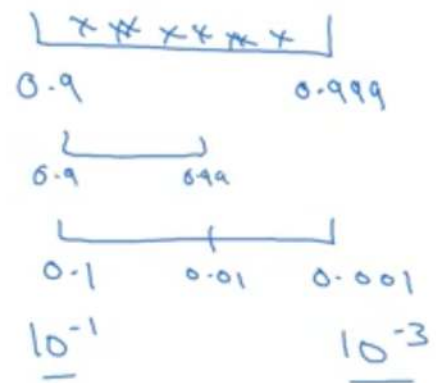


## Hyperparameters for exponentially weighted averages

$$\beta = 0.9 \quad \dots \quad 0.999$$

$\downarrow$                        $\downarrow$   
 10                      1000

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

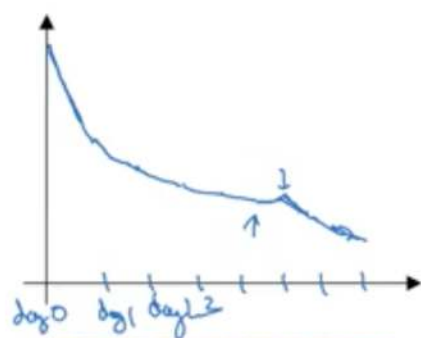


$$r \in [-3, -1]$$

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

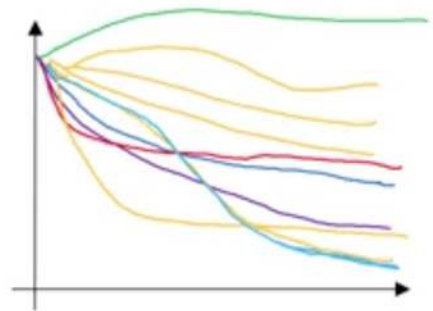
$$\beta = 1 - 10^r$$

## Babysitting one model



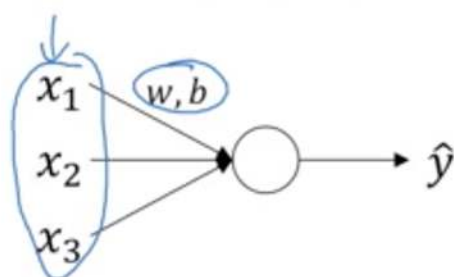
Panda

## Training many models in parallel



Caviar

# <sup>x</sup> Normalizing inputs to speed up learning

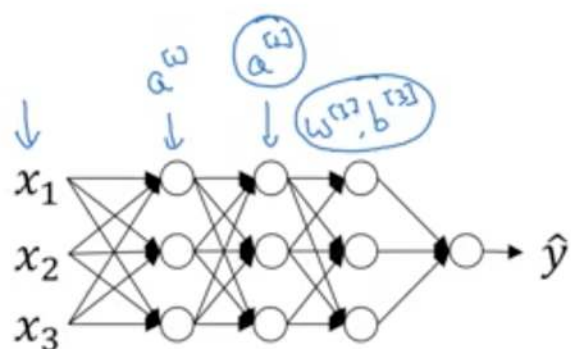
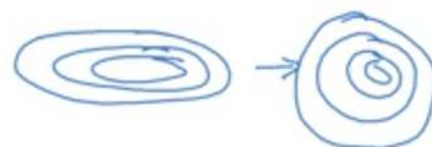


$$\mu = \frac{1}{m} \sum_i x^{(i)}$$

$$X = X - \mu$$

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

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



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

Normalize  $z^{[2]}$



# Implementing Batch Norm

Given some intermediate values in NN

$z^{(1)}, \dots, z^{(n)}$   
 $z^{(i)}$

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

Use  $\hat{z}^{(i)}$  instad of  $z^{(i)}$

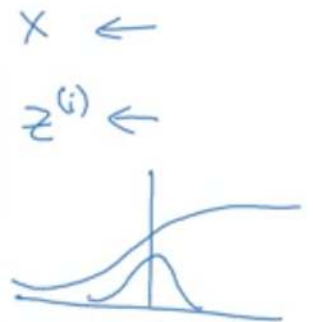
If

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

$$\beta = \mu$$

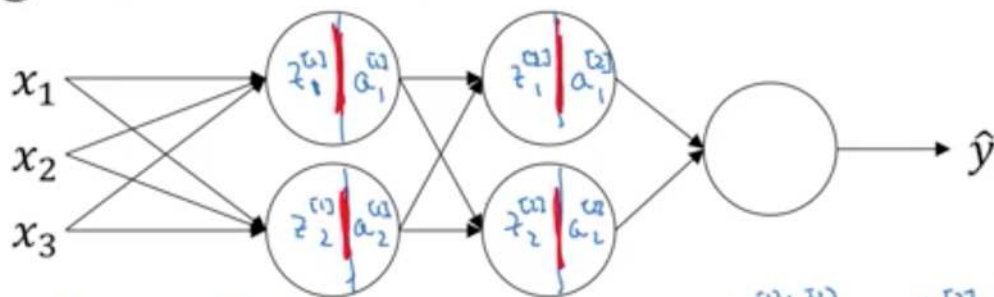
then  $\hat{z}^{(i)} = z^{(i)}$

learnable parameters of model.





# Adding Batch Norm to a network

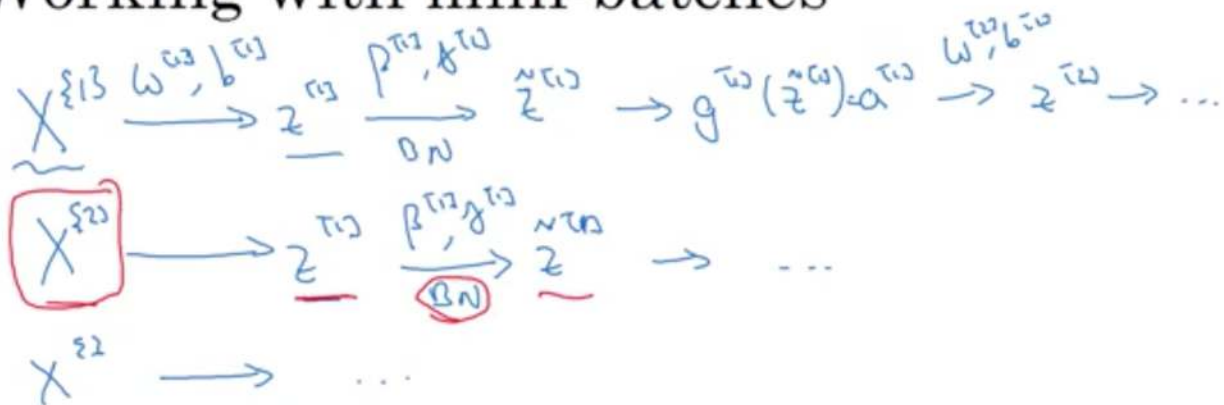


$$x \xrightarrow{W^{(1)}, b^{(1)}} \underline{z^{(1)}} \xrightarrow[\text{Batch Norm (BN)}]{\beta^{(1)}, \gamma^{(1)}} \underline{N^{(1)}(z)} \rightarrow a^{(1)} = g(N^{(1)}(z)) \xrightarrow{W^{(2)}, b^{(2)}} \underline{z^{(2)}} \xrightarrow[\text{BN}]{\beta^{(2)}, \gamma^{(2)}} \underline{N^{(2)}(z)} \rightarrow a^{(2)} \rightarrow \dots$$

Parameters:  $\{W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots, W^{(L)}, b^{(L)}\}$   
 $\rightarrow \beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)}, \dots, \beta^{(L)}, \gamma^{(L)}$   
 $\rightarrow \beta$

$d\beta^{(2)}$   
 $\beta = \beta^{\text{old}} - \alpha d\beta^{(2)}$   
 tf.nn.batch-normalization ←

## Working with mini-batches



Parameters:  $W^{[1]}, \cancel{b^{[1]}}, \beta^{[1]}, \delta^{[1]}$ .

Below the parameters, the dimensions are indicated:

- $z^{[1]} (n^{[1]}, 1)$  is aligned under  $W^{[1]}$ .
- $\hat{z}^{[1]} (n^{[1]}, 1)$  is aligned under  $\beta^{[1]}$ .
- $\delta^{[1]} (n^{[1]}, 1)$  is aligned under  $\delta^{[1]}$ .

The  $\cancel{b^{[1]}}$  parameter is crossed out with a red X.

Red arrows point to the equations below:

$$\rightarrow \underline{z^{[1]}} = W^{[1]} a^{[0]} + \cancel{b^{[1]}}$$

$$z^{[1]} = W^{[1]} a^{[0]}$$

$$\rightarrow \hat{z}^{[1]} = \delta^{[1]} z_{\text{norm}}^{[1]} + \beta^{[1]}$$

The  $\cancel{b^{[1]}}$  term in the first equation and the  $\beta^{[1]}$  term in the third equation are highlighted with red boxes and red arrows.

## Implementing gradient descent

for  $t = 1 \dots \text{num MiniBatches}$   
Compute forward pass on  $X^{(t)}$ .

In each hidden layer, use BN to replace  $\underline{z}^{(l)}$  with  $\tilde{z}^{(l)}$ .

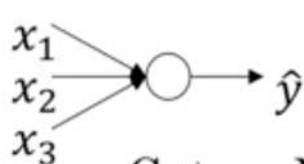
Use backprop to compute  $\underline{dw}^{(l)}$ ,  ~~$d\beta^{(l)}$~~ ,  $d\gamma^{(l)}$ ,  $\underline{df}^{(l)}$

Update params 
$$\left. \begin{aligned} W^{(l)} &:= W^{(l)} - \alpha dw^{(l)} \\ \beta^{(l)} &:= \beta^{(l)} - \alpha d\beta^{(l)} \\ \gamma^{(l)} &:= \dots \end{aligned} \right\} \leftarrow$$

Works w/ momentum, RMSprop, Adam.

the parameters Beta and Gamma that Batch Norm added to algorithm.

# Learning on shifting input distribution



Cat

Non-Cat

$y = 1$  ✓

$y = 0$



$y = 1$  ✓

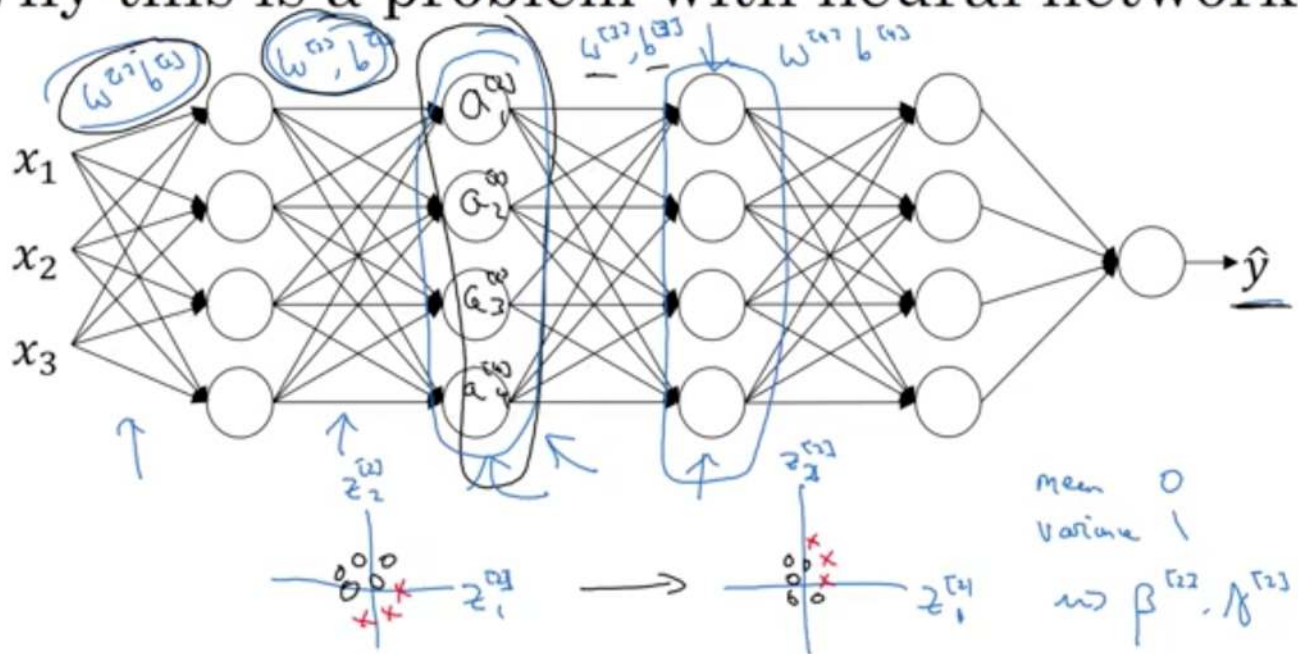
$y = 0$



"Covariate shift"

$\underline{x} \rightarrow y$

Why this is a problem with neural networks?



## Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.  $\times$
- 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.  $\tilde{z}^{[l]}$   $\leftarrow 128$   $\times^{t+1}$   $z^{[l]}$
- This has a slight regularization effect.



## <sup>x</sup>Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.  $X^{(t)}$   
 $z^{(t)}$
- 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.  $\mu, \sigma^2$
- This has a slight regularization effect.

mini-batch: 64  $\rightarrow$  512



## Batch Norm at test time

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

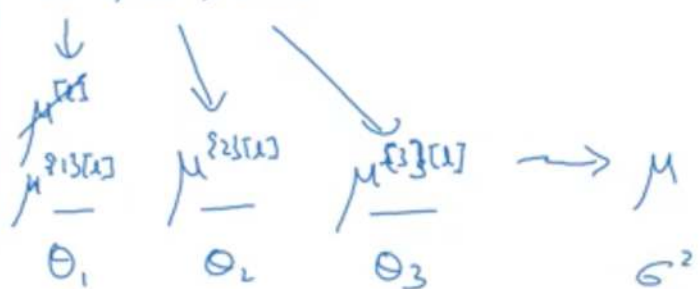
$$\rightarrow \sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

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

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

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

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





## <sup>x</sup> Batch Norm at test time

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

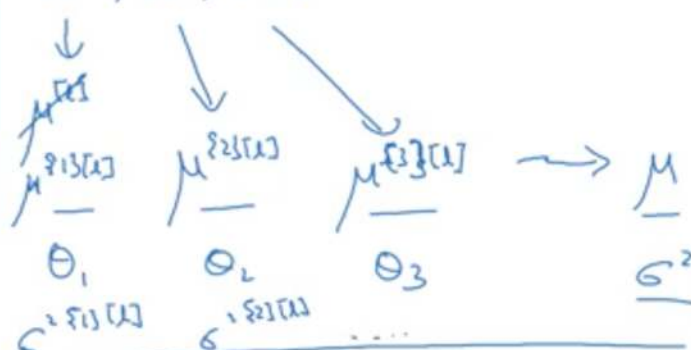
$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

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

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

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

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



$$\tilde{z} = \gamma \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta$$

## Batch Norm at test time

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

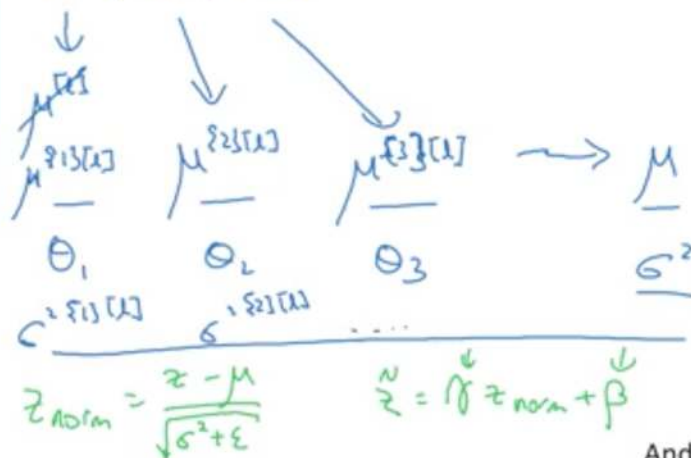
$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

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

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

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

$X^{[1]}, X^{[2]}, X^{[3]}, \dots$



# Recognizing cats, dogs, and baby chicks, <sup>other</sup><sub>0</sub>



3

1

2

0

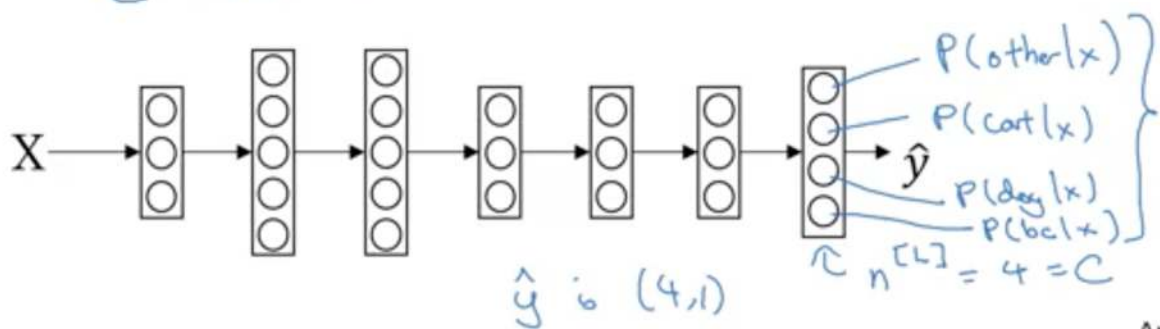
3

2

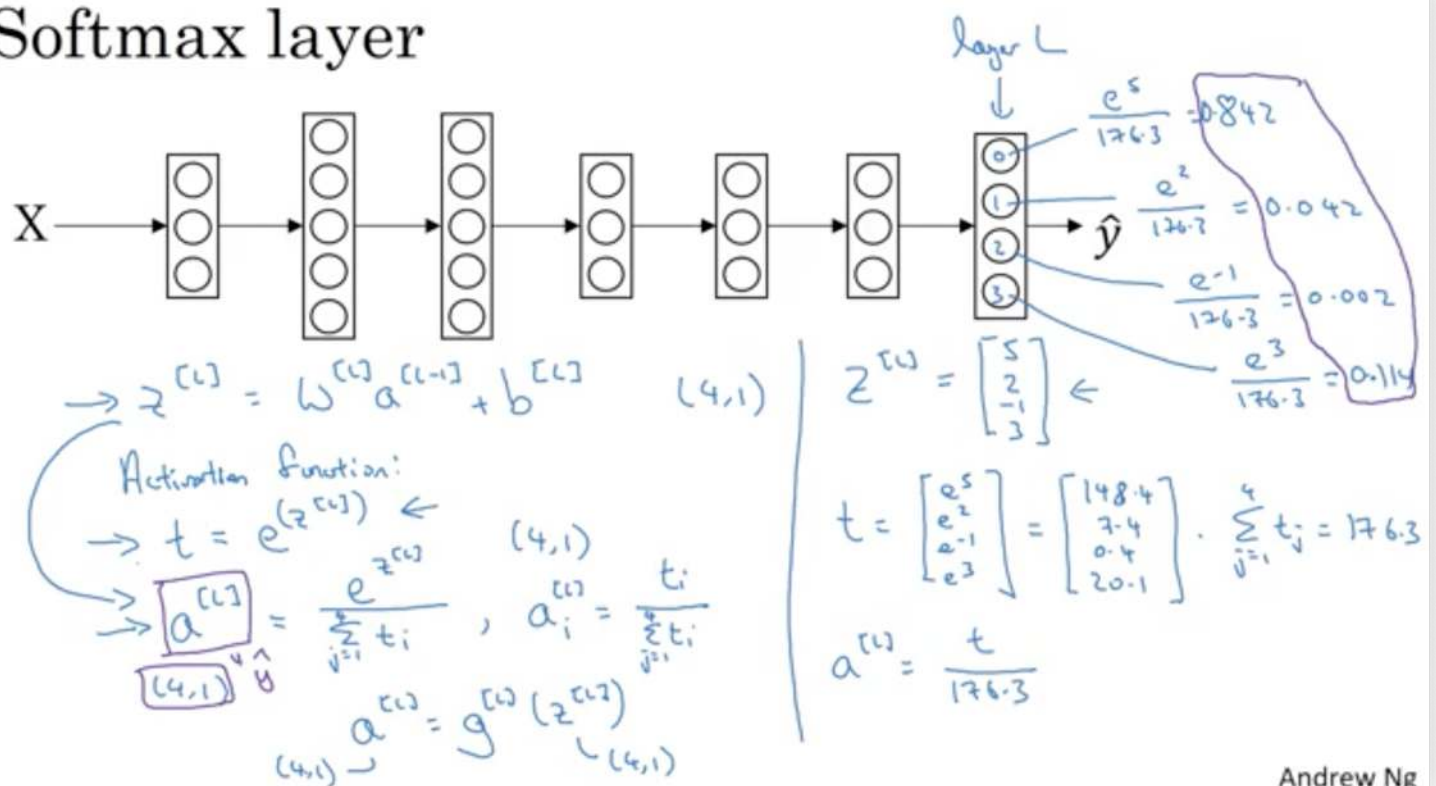
0

1

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



# Softmax layer

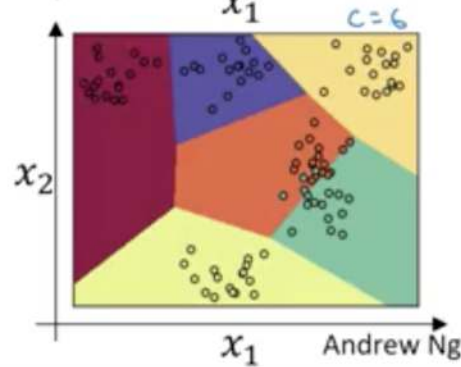
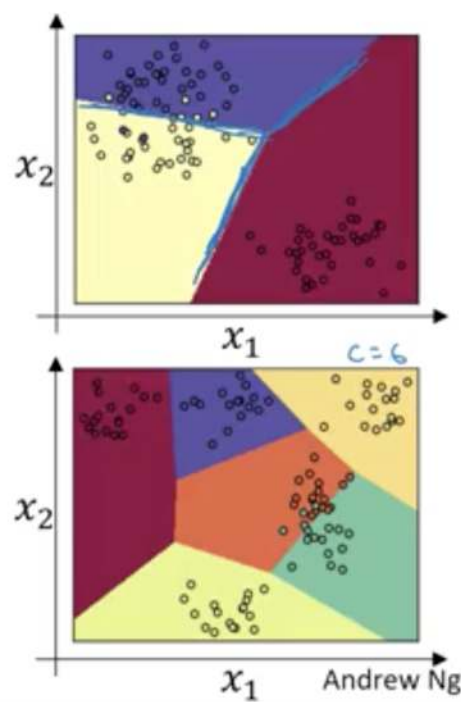
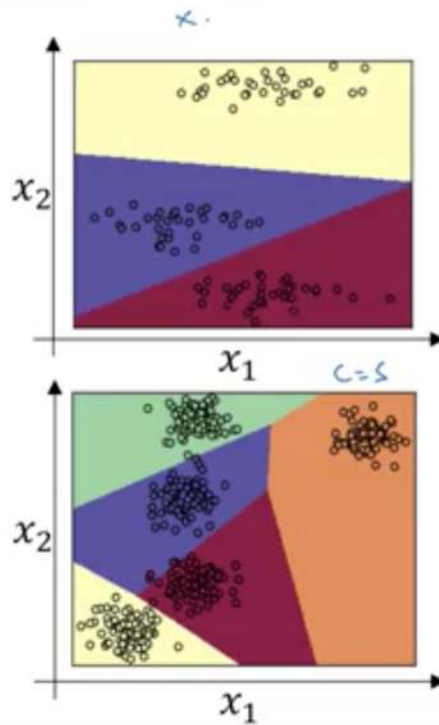
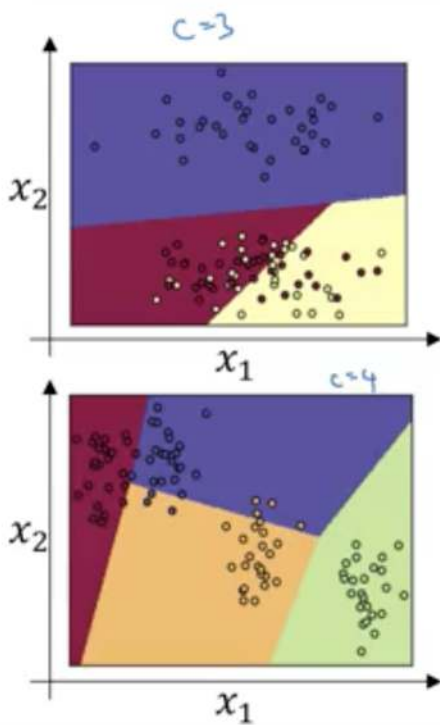


# Softmax examples

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \hat{y}$$

$$z^{(1)} = W^{(1)}x + b^{(1)}$$

$$a^{(1)} = \hat{y} = g(z^{(1)})$$



# Understanding softmax

(4,1)

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

$C=4$        $g^{[L]}(\cdot)$

"Soft max"

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

"hard max"

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

Softmax regression generalizes logistic regression to  $C$  classes.

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



# Loss function

$$y^{(1)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \leftarrow \text{cat } y_2 = 1$$

$$y_1 = y_3 = y_4 = 0$$

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

$\underbrace{\hspace{10em}}_{\text{small}} \quad \quad \quad \uparrow$

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

$\text{Make } \hat{y}_2 \text{ big.}$

$C = 4$

$a^{(1)} \approx \hat{y}^{(1)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$

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

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

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

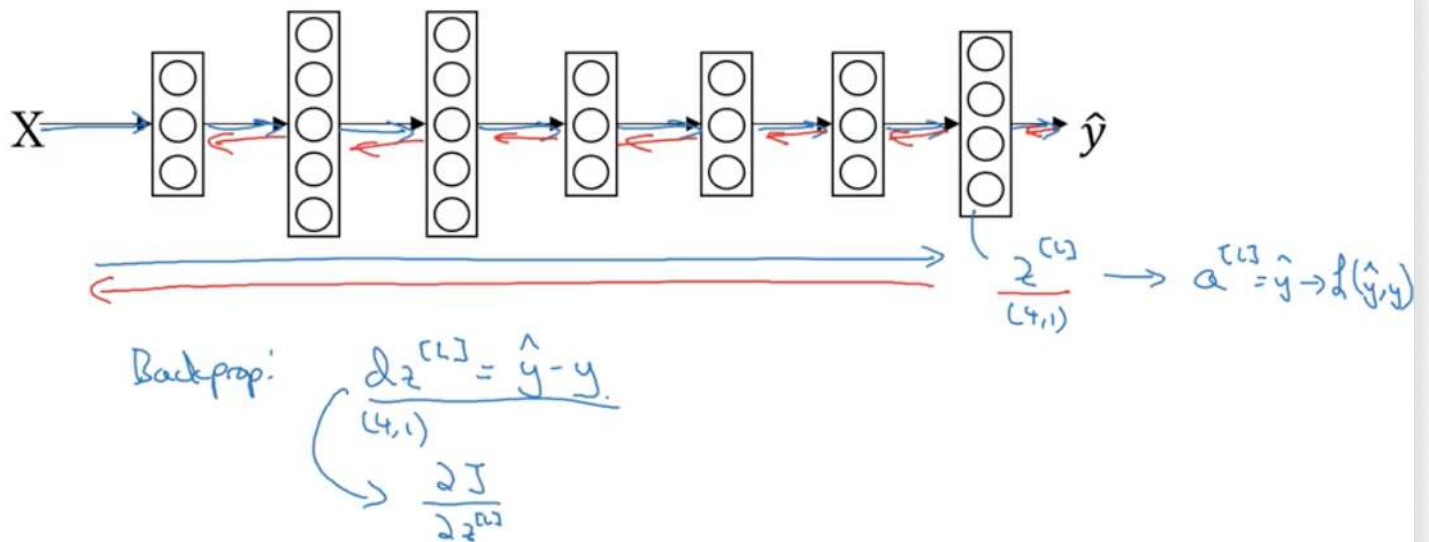
$(4, m)$

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

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

$(4, m)$

## Gradient descent with softmax





# Deep learning frameworks

- 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)