# Hyperparameters

$\rightarrow$ □ $\alpha$ $\longrightarrow$ Priority 1

□ $\beta$ ≈0.9 $\longrightarrow$ Priority 2 ←

**Almost never tune with ADAM** →

$\beta_1, \beta_2, \varepsilon$
0.9  0.999  $10^{-8}$

Priority 3

#layers

#hidden units

learning rate decay

Mini-batch size

---

## Try random values: Don't use a grid



$\alpha$ | Hyperparameter 2 | $\varepsilon$

Hyperparameter 1

$\alpha$ | Hyperparameter 2 | Hyperpara 3

Hyperparameter 1

---

## Coarse to fine



Hyperparameter 2

Hyperparameter 1

# Picking hyperparameters at random

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



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

$2, 3, 4$

OR

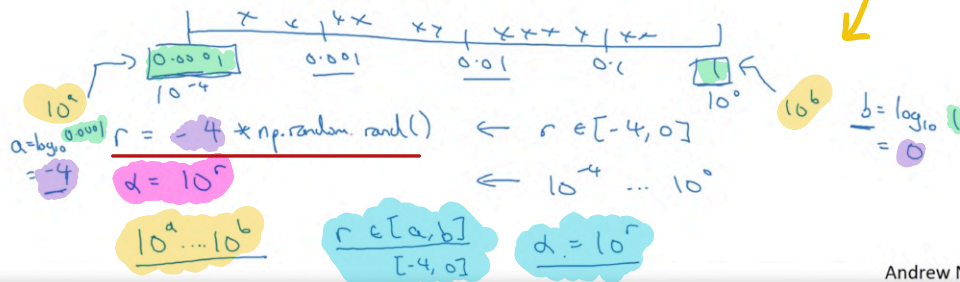# Appropriate scale for hyperparameters

$\alpha = 0.0001, \ldots, 1$



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

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

$a = \log_{10} 0.0001 = -4$

$10^a \ldots 10^b$   $\quad r \in [a, b]$   $\quad \alpha = 10^r$
$\qquad\qquad\qquad [-4, 0]$

$b = \log_{10} 1 = 0$

So $\alpha = 10^r$ where $r \in (-4, 0]$

# Hyperparameters for exponentially weighted averages

$\beta = 0.9 \ldots 0.999$
$\quad\downarrow \qquad\quad \downarrow$
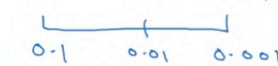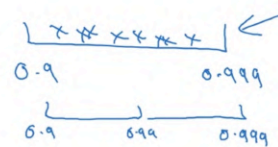$\quad 10 \qquad\quad 1000$

$1 - \beta = 0.1 \ldots 0.001$

$\beta: 0.9000 \rightarrow 0.9005 \quad\} \sim 10$
$\beta: 0.999 \rightarrow 0.9995$
$\quad \sim 1000 \qquad \sim 2000$

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



$0.9 \qquad\qquad 0.999$

$0.9 \quad 0.99 \quad 0.999$

$0.1 \quad 0.01 \quad 0.001$

$10^{-1} \qquad\qquad 10^{-3}$

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

$1 - \beta = 10^r$

$\beta = 1 - 10^r$

# Re-test hyperparameters occasionally

Idea



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

- Intuitions do get stale. Re-evaluate occasionally.

Experiment      Code

---

# Babysitting one model



day 0   day 1   day 2

Panda ←

# Training many models in parallel



Caviar ←

---

# Normalizing inputs to speed up learning



$x_1$   $w, b$

$x_2$

$x_3$ → $\hat{y}$

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

$X = X - \mu$

$\sigma^2 = \frac{1}{m} \sum x^{(i)^2}$   ← elent-wise

$X = X / \sigma^2$

$a^{[1]}$   $a^{[2]}$   $(w^{[3]}, b^{[3]})$

$x_1$

$x_2$

$x_3$ → $\hat{y}$

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

Normalize $z^{[2]}$ ↑

# Implementing Batch Norm

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

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

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

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

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

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

learnable parameters of model.

If $\gamma = \sqrt{\sigma^2 + \varepsilon} \leftarrow$

$\beta = \mu \leftarrow$

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

$x \leftarrow$

$z^{(i)} \leftarrow$

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

learnable

---

# Adding Batch Norm to a network



$x_1$
$x_2$
$x_3$

$z_1^{[1]} \, a_1^{[1]}$ $z_1^{[2]} \, a_1^{[2]}$
$z_2^{[1]} \, a_2^{[1]}$ $z_2^{[2]} \, a_2^{[2]}$

$\hat{y}$

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

Parameters: $\omega^{[1]}, b^{[1]}, \omega^{[2]}, b^{[2]}, \ldots, \omega^{[L]}, b^{[L]},$
$\to \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \ldots, \beta^{[L]}, \gamma^{[L]}$

$\to \beta$

$d\beta^{[l]}$

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

tf.nn.batch-normalization $\leftarrow$

---

# Working with mini-batches

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

$$X^{\{2\}} \xrightarrow{} z^{[1]} \xrightarrow[\text{BN}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \to \ldots$$

$$X^{\{3\}} \xrightarrow{} \ldots$$

Per minibatch

Parameters: $\omega^{[l]}, \cancel{b^{[l]}}, \beta^{[l]}, \gamma^{[l]}$

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

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

$\to z^{[l]} = \omega^{[l]} a^{[l-1]} + \cancel{b^{[l]}}$

$z^{[l]} = \omega^{[l]} a^{[l-1]}$

$z_{norm}^{[l]}$

$\to \tilde{z}^{[l]} = \gamma^{[l]} z_{norm}^{[l]} + \beta^{[l]} \leftarrow$

# Implementing gradient descent

for t = 1 .... num MiniBatches

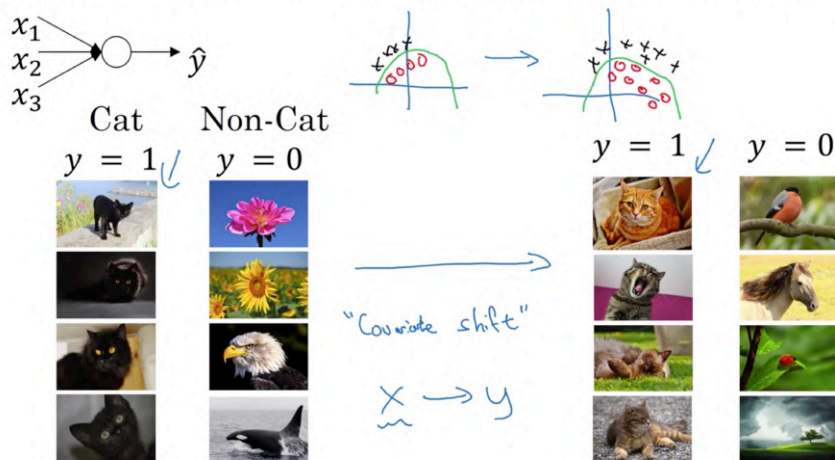Compute forward prop on $X^{\{t\}}$.

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

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

Update params

$$w^{[l]} := w^{[l]} - \alpha \, dw^{[l]}$$
$$\beta^{[l]} := \beta^{[l]} - \alpha \, d\beta^{[l]}$$
$$\gamma^{[l]} := \ldots$$

Works w/ momentum, RMSprop, Adam.

# Learning on shifting input distribution



$x_1$
$x_2$
$x_3$
$\hat{y}$

Cat          Non-Cat
$y = 1$      $y = 0$

$y = 1$      $y = 0$

"Covariate shift"

$x \longrightarrow y$

# Why this is a problem with neural networks?



As changes,
values in $a^{[2]}$ also change

$x_1$
$x_2$
$x_3$
$\hat{y}$

$w^{[2]}, b^{[2]}$   $w^{[3]}, b^{[3]}$   $w^{[3]}, b^{[3]}$   $w^{[4]}, b^{[4]}$

$a^{[2]}_1$
$a^{[2]}_2$
$a^{[2]}_3$
$a^{[2]}_4$

Each layer learns by itself a bit more independently of other layers

$z^{[2]}_2$
$z^{[2]}_1$

$z^{[3]}_3$
$z^{[3]}_1$

Mean 0
Variance 1

doesn't change

$\beta^{[2]}, \gamma^{[2]}$

# Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.

$\to \tilde{z}^{[l]}$   $\{4, 128\}$   $z^{[l]}$

$X^{\{t\}}$
$X^{\{t+1\}}$

- 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 : $64 \longrightarrow 512$

batch norm with dropout for greater
regularization effect of batch norm

I wouldn't really use batch norm as a regularizer, that's really not the intent of batch norm, but sometimes it has this extra intended or unintended effect on your learning algorithm. But, really, don't turn to batch norm as a regularization. Use it as a way to normalize your hidden units activations and therefore speed up learning. And I think the regularization is an almost unintended side effect

Andrew Ng

# Batch Norm at test time

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

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

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

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

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

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

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

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

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

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

Andrew Ng

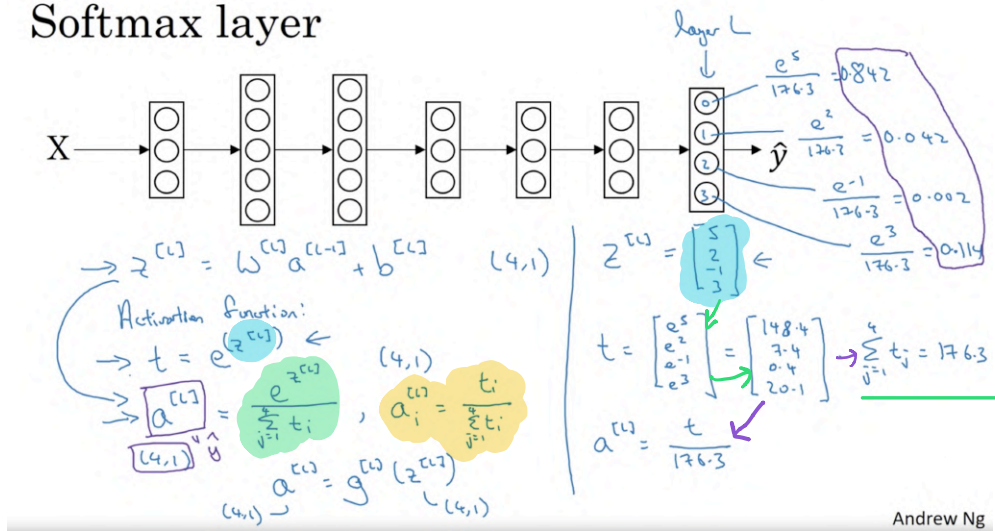# Recognizing cats, dogs, and baby chicks , other

1    2    3    0



3    1    2    0    3    2    0    1

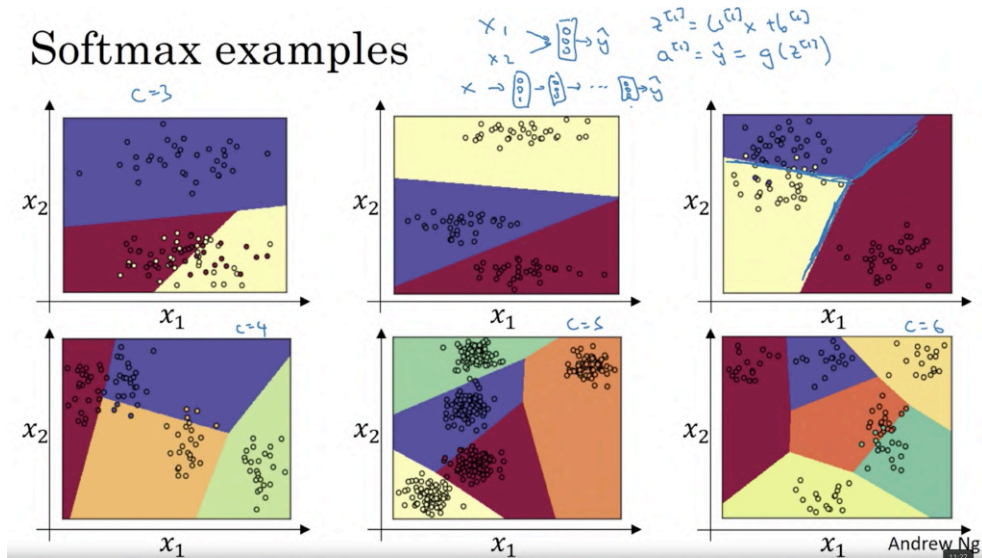$C = \#classes = 4$   $(0, \ldots, 3)$



$P(other|x)$
$P(cat|x)$
$P(dog|x)$
$P(bc|x)$

$\hat{y}$

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

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

Andrew Ng

# Softmax layer



layer $L$

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

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

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

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

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

Activation function:

$\rightarrow t = e^{(z^{[L]})} \leftarrow$

$a^{[L]} = \dfrac{e^{z^{[L]}}}{\sum_{j=1}^{4} t_i}$ , $a_i^{[L]} = \dfrac{t_i}{\sum_{j=1}^{4} t_i}$

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

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

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

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

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

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

# Softmax examples

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

$z^{[1]} = W^{[1]} x + b^{[1]}$

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

$x \rightarrow \square \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow \hat{y}$



$C=3$  $c=4$  $C=5$  $C=6$

# Understanding softmax

$(4,1)$

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

$C = 4$  $g^{[L]}(\cdot)$

"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.

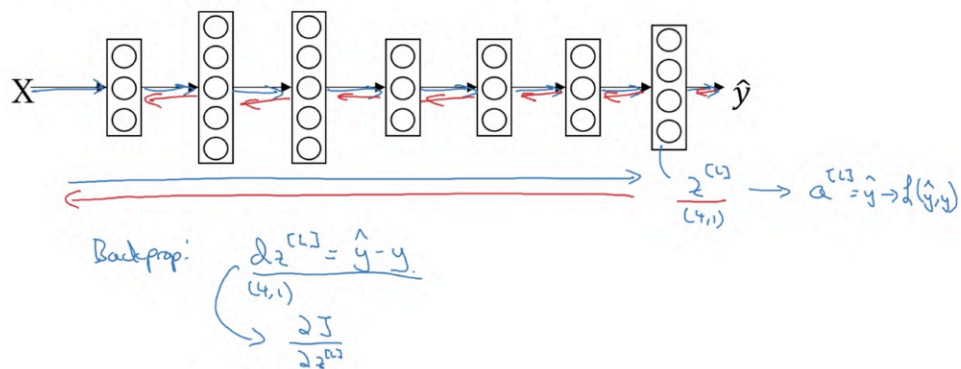If $C=2$, softmax reduces to logistic regression. $a^{[L]} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$

Andrew Ng

# Loss function

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \text{(4,1)} \quad - \text{cat} \quad y_2 = 1$$

$$y_1 = y_3 = y_4 = 0$$

$$a^{[L](i)} \approx \hat{y}^{(i)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \text{(4,1)} \quad C = 4$$

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

small

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

$$- y_2 \log \hat{y}_2 = -\log \hat{y}_2 \quad \rightarrow \quad \text{Make } \hat{y}_2 \text{ big.}$$

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

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

(4, m)

$$\hat{Y} = \begin{bmatrix} \hat{y}^{(1)}, & \dots, & y^{(m)} \end{bmatrix}$$

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

(4, m)

# Gradient descent with softmax



$$X \rightarrow \dots \rightarrow \hat{y}$$

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

Backprop:

$$dz^{[L]} = \hat{y} - y$$

(4,1)

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