*Coursera Notes*

**Improving Deep neural Networks: Hyperparameter Tuning, Regularization and Optimization Methods**
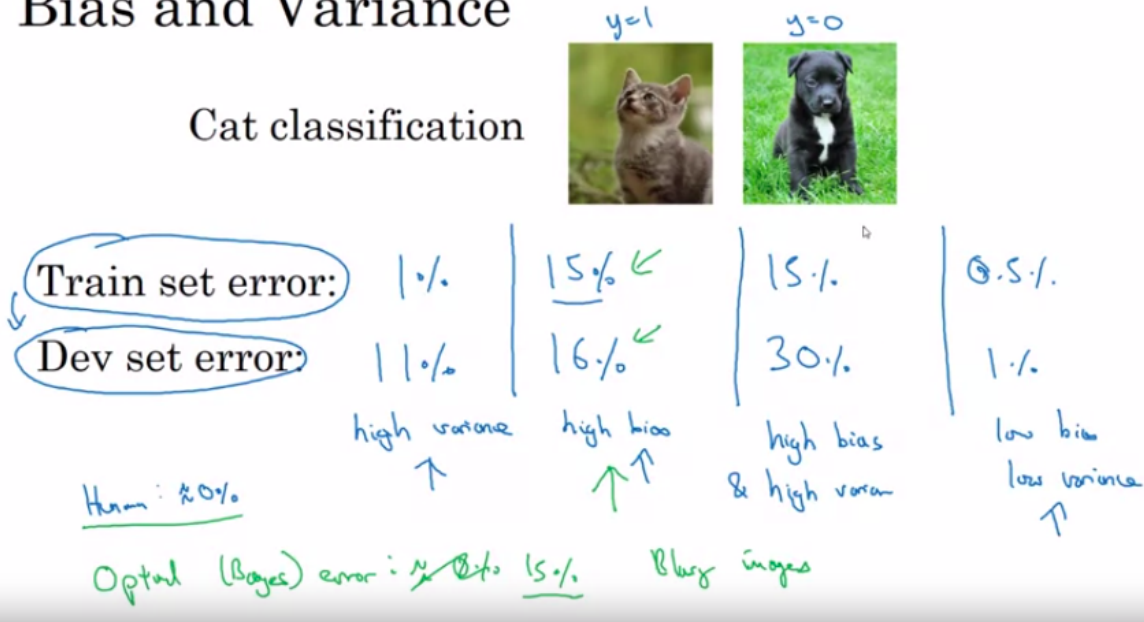
Setting up an ML application:

W1L1: train/dev/test sets
- Smaller datasets: 60/20/20 or 70/15/15
- Big datasets: 98/1/1
- Make sure that your dev and test sets come from the same data distributions or source of inputs
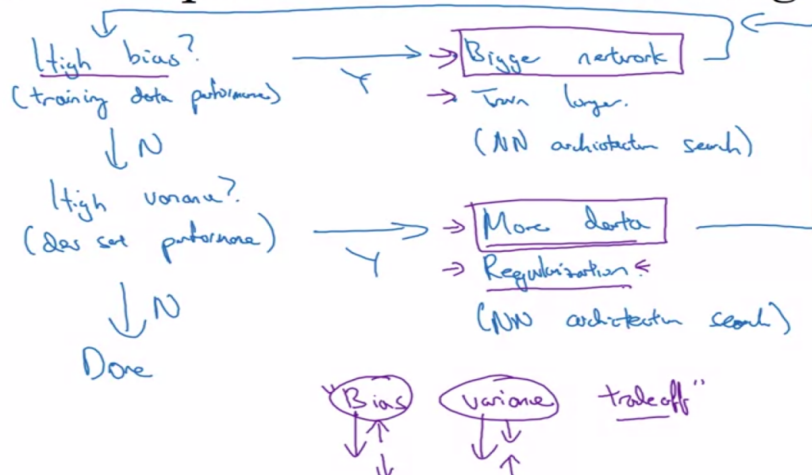
W1L2: bias/variance



W1L3:



Andrew Ng

The vector norm $|\mathbf{x}|_p$ for $p = 1$, 2, ... is defined as

$$|\mathbf{x}|_p \equiv \left( \sum_i |x_i|^p \right)^{1/p}.$$

W1L4: Regularization:
- To prevent high variance problem & overfitting: use regularization
- L2 - generally used - also called as weight decay
- L1 - weights will be sparse

# Neural network

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

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

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

"Frobenius norm"

$$\|\cdot\|_2^2 \qquad \|\cdot\|_F^2$$

$$dw^{[l]} = \boxed{(\text{from backprop}) + \frac{\lambda}{m} w^{[l]}}$$

$$\to w^{[l]} := w^{[l]} - \alpha \, dw^{[l]}$$

$$\frac{\partial J}{\partial w^{[l]}} = dw^{[l]}$$

"Weight decay"

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

W1L5: Why regularization reduces overfitting?
**Regularization** in machine learning is the process of regularizing the parameters that constrain, regularizes, or shrinks the coefficient estimates towards zero. In other words, this technique discourages learning a more complex or flexible model, avoiding the risk of **Overfitting**.

W1L6/7: Dropout regularization
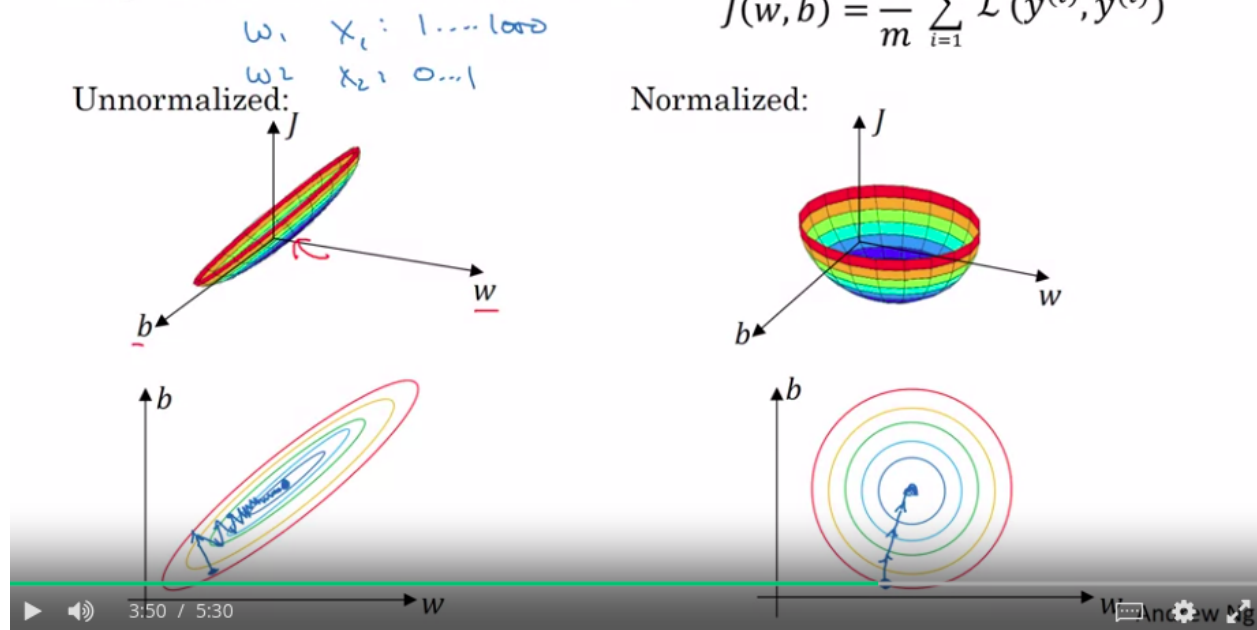- It helps prevent overfitting as the weights have to shrink when some features are dropped.
- In computer vision applications: dropout is used most often
- Downside: cost function J is not well defined, harder to double-check the performance of the gradient descent model.
- To avoid this: Ng runs the code with dropouts off, keep prob = 1, and then turns the dropout on to avoid introducing bugs.

W1L8: few more regularization techniques:
- Data Augmentation:
- Early stopping:

W1L9: Normalizing the inputs



## Why normalize inputs?

$w_1 \quad x_i : 1 \cdots 1000$
$w_2 \quad x_2 : 0 \cdots 1$

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

Unnormalized:

Normalized:

3:50 / 5:30

Andrew Ng

W1L10: vanishing/exploding gradients:
for a deep neural network with L layers (L is large integer):
- If W < Identity matrix : vanishing gradients
- If W > identity matrix : exploding gradients as
  W raised to L-1 gives the hypothesis

W1L11: Weight Initialization: To avoid vanishing and exploding gradients
- Relu : W = np.random.rand(shape) * np.sqrt(1/(n^(L-1))) where n: number of hidden units
- Tanh: Xavier Initialiation

W2L1/2: Mini Batch Gradient Descent
- Mini batch is faster than batch GD.
- Dividing the whole dataset into small batches and update weights for every batch.
- Choosing the size of mini-batch:

# Choosing your mini-batch size

If small tray set : Use batch gradient descent.
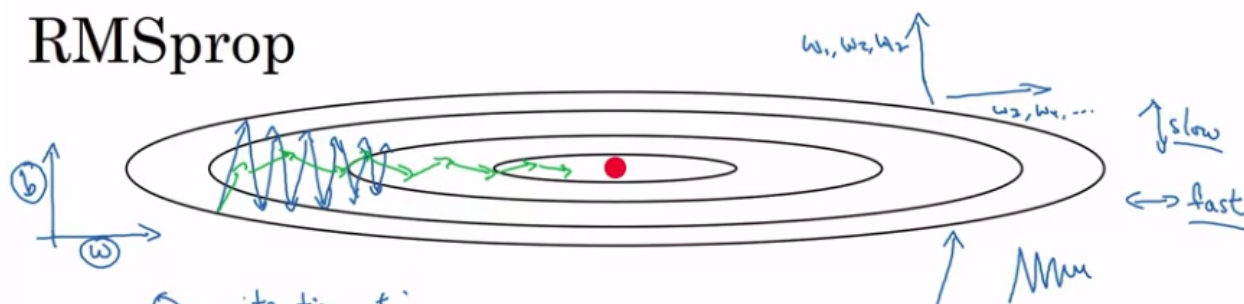$(m \leq 2000)$

Typical mini-batch sizes:

$$64 \, , \, 128, \, 256, \, 512 \qquad \frac{1024}{2^{10}}$$
$$2^6 \qquad 2^7 \qquad 2^8 \qquad 2^9$$

Make sure mini-batch fit in CPU/GPU memory.
$X^{\{t\}}, Y^{\{t\}}$.

W2L4: RMSprop

# RMSprop



On iteration $t$:

Compute $dW, db$ on current Mini-batch

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dW^2 \leftarrow \text{small}$$ element-wise

$$\rightarrow S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \leftarrow \text{large}$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dw}}+\varepsilon} \qquad b := b - \alpha \frac{db}{\sqrt{S_{db}}+\varepsilon}$$

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

W2L5: Adam
Combination of exponential weighted avg and RMSprop

# Adam optimization algorithm

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

On iteratn $t$:

  Compute $dw, db$ using current mini-batch

  $V_{dw} = \beta_1 V_{dw} + (1-\beta_1)dw$ , $V_{db} = \beta_1 V_{db} + (1-\beta_1)db$ ← "momentum" $\beta_1$

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

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

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

  $W := W - \alpha \dfrac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected}}+\varepsilon}$     $b := b - \alpha \dfrac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected}}+\varepsilon}$

Andrew Ng

W2L6: learning rate decay

method1:

# Learning rate decay

1 epoch = 1 pass through data.

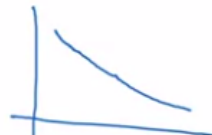$\alpha = \dfrac{1}{1 + decay\text{-}rate * epoch\text{-}num} \alpha_0$ ←

$\boxed{x^{\{1\}} \quad x^{\{2\}} \quad - - - -}$ → epoch 1
→ epoch 2

$\alpha_0 = 0.2$
decay-rate = 1

| Epoch | $\alpha$ |
|-------|------|
| 1 | 0.1 |
| 2 | 0.67 |
| 3 | 0.5 |
| 4 | 0.4 |

# Other learning rate decay methods

formula $\begin{cases} \alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 & \text{— exponentially decay.} \\[2em] \alpha = \dfrac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 & \text{or} \quad \dfrac{k}{\sqrt{t}} \cdot \alpha_0 \\[2em] \end{cases}$

discrete staircase
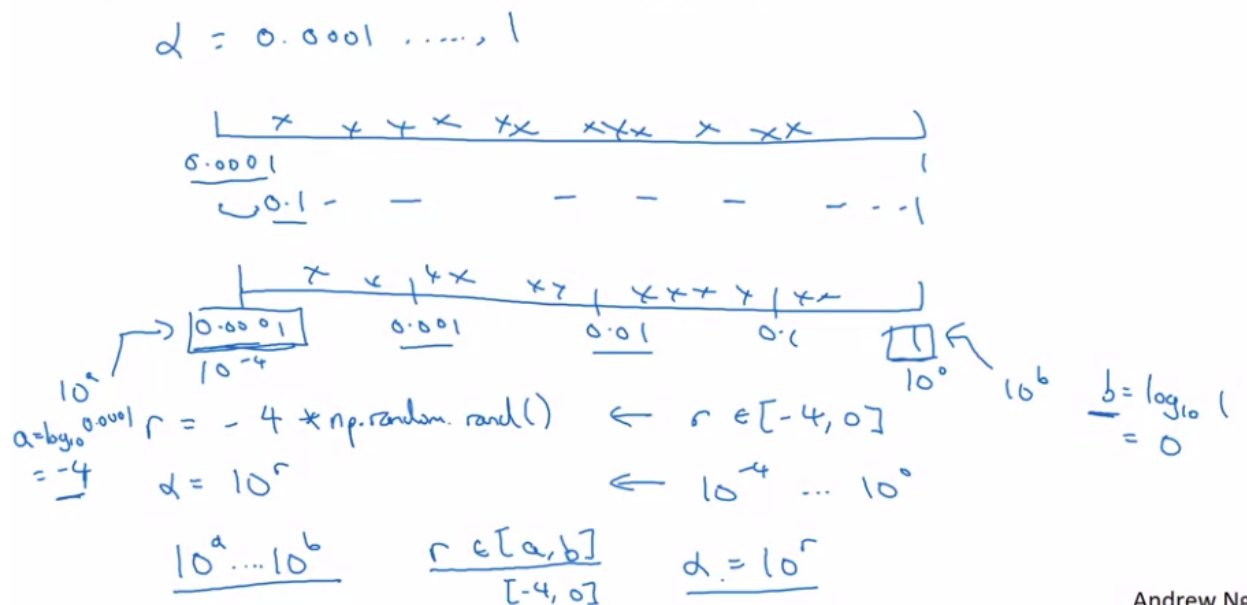
$\alpha$ ⊥___

Manual decay.

W3L1: tuning process
- Order of importance of tuning the hyperparameters:
  learning rate
  momentum
  # hidden units
  mini-batch size
  #layers
  learning rate decay
  adam optimizer parameters: beta1, beta2, epsilon
- Do not use a grid, try random samples

W3L2: picking the right scale for tuning: use logarithmic scale
- Learning rate:

## Appropriate scale for hyperparameters

$\alpha = 0.0001 \ldots, 1$

$0.0001$

$0.1 - - - - - - - - \cdot \cdot \cdot 1$

$0.0001$   $0.001$   $0.01$   $0.1$

$10^a$   $10^{-4}$   $10^0$   $10^6$   $b = \log_{10} 1$
$= 0$

$a = \log_{10} 0.0001$   $r = -4 * np.random.rand()$   $\leftarrow r \in [-4, 0]$
$= -4$   $\alpha = 10^r$   $\leftarrow 10^{-4} \ldots 10^0$

$10^a \ldots 10^b$   $\underline{r \in [a,b]}$   $\alpha = 10^r$
$[-4, 0]$

Andrew Ng

- Beta for exponentially weighted averages

# Hyperparameters for exponentially weighted averages

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

$\downarrow \qquad\qquad \downarrow$

$10 \qquad\qquad 1000$

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

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

$\beta: 0.999 \rightarrow 0.9995$

$\sim 1000 \qquad \sim 2000$

$\frac{1}{1-\beta_R}$

$| \times \# \times \not{\times} \times \not{\times} \times |\Leftarrow$

$0.9 \qquad\qquad 0.999$

$0.9 \qquad 0.99 \qquad 0.999$

$0.1 \qquad 0.01 \qquad 0.001$

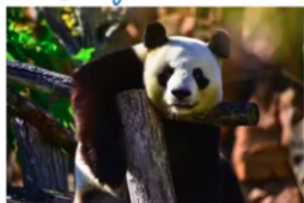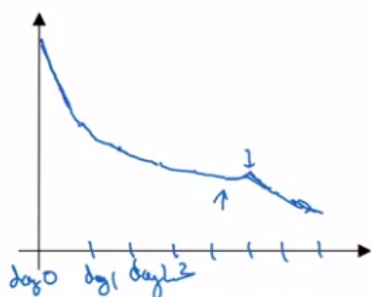$10^{-1} \qquad\qquad 10^{-3}$

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

$1-\beta = 10^r$

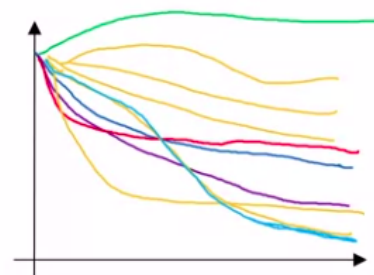$\beta = 1 - 10^r$

An

---

**W3L3: tuning in practice: Application dependent**

## Babysitting one model
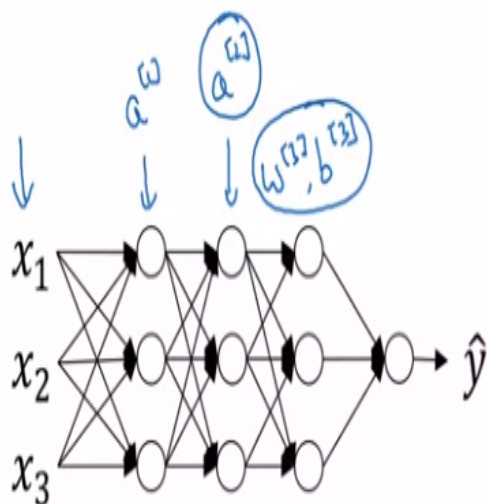


Panda

## Training many models in parallel



Caviar

Andrew

**W3L4: Batch Normalization**
- Normalizing the values of the parameters before applying the activation function.

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

Normalize $z^{[2]}$

Andrew Ng

- Implementing Batch Norm:

# Implementing Batch Norm

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

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

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

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

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

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

If
$$\gamma = \sqrt{\sigma^2 + \varepsilon}$$
$$\beta = \mu$$
then $\tilde{z}^{(i)} = z^{(i)}$

learnable parameters of model.

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

$X \leftarrow$

$z^{(i)} \leftarrow$

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

Andrew Ng

# Working with mini-batches

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

$$\boxed{X^{\{2\}}} \xrightarrow{\quad} Z^{[1]} \xrightarrow[\boxed{BN}]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z} \rightarrow \cdots$$

$$X^{\{3\}} \xrightarrow{\quad} \cdots$$

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

$\quad\quad\quad\quad\quad | \quad\quad\quad | \quad\quad\quad |$

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

$Z^{[l]}$

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

$\rightarrow Z^{[l]} = W^{[l]} a^{[l-1]} + \cancel{\boxed{b^{[l]}}}$

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

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

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

# Implementing gradient descent

for $t = 1 \ldots$ num MiniBatches

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

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

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

$\quad$ Update params $\left. \begin{array}{l} W^{[l]} := W^{[l]} - \alpha \, dW^{[l]} \\ \beta^{[l]} := \beta^{[l]} - \alpha \, d\beta^{[l]} \\ \gamma^{[l]} := \cdots \end{array} \right\} \Leftarrow$

Works w/ momentum, RMSprop, Adam.

W3L5: BN at the test time

# Batch Norm at test time



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

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

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

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

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

$X^{\{1\}}, X^{\{1\}}, X^{\{3\}}, \dots$

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

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

$\sigma^{2\{1\}[\ell]} \qquad \sigma^{\{2\}[\ell]}$

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

Andrew Ng

W3L6: Multiclass classification:
- Softmax Layer: takes a vector as an input and outputs a vector with the same dimensions. It computes the chance of occuring a particular class in the multiclass classification.
- In the example below: there is 84% chance of the predicted class to be class 0.

# Softmax layer



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

Activation function:

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

$a^{[\ell]} = \frac{e^{z^{[\ell]}}}{\sum_{j=1}^{4} t_i}, \qquad a_i^{[\ell]} = \frac{t_i}{\sum_{j=1}^{4} t_i}$

$(4,1)$

$a^{[\ell]} = g^{[\ell]}(z^{[\ell]})$
$(4,1) \qquad \qquad (4,1)$

layer L

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

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

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

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

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

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

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

Andrew Ng

(Notes by Samruddhi Patil) 11

W3L7: training a softmax classifier

# Loss function

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

$$y_2 = 1$$

$$y_1 = y_3 = y_4 = 0$$

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

$$C = 4$$

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

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

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

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

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

$$(4, m)$$

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

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

$$(4, m)$$

Andrew Ng