

# Improving Deep Neural Networks: Hyperparameter Tuning, Regularization and Optimization

- Andrew Ng

Applied NN is a highly iterative process.

# layers, # hidden units, learning rates, activation funcs

Data	Training Set	Dev Set	Test Set
------	--------------	---------	----------

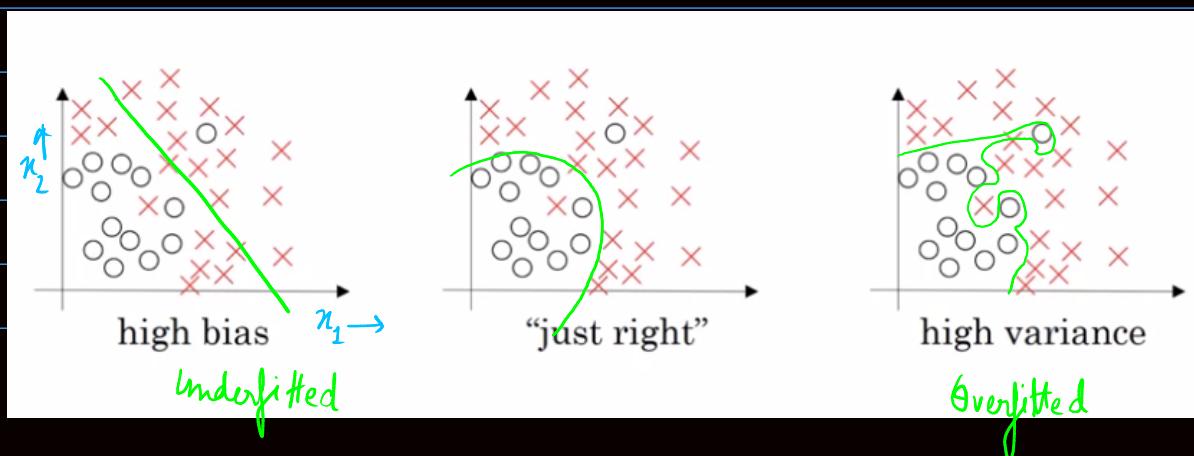
Ratios  $\rightarrow$  60% | 20% | 20% or 80% | 10% | 10%

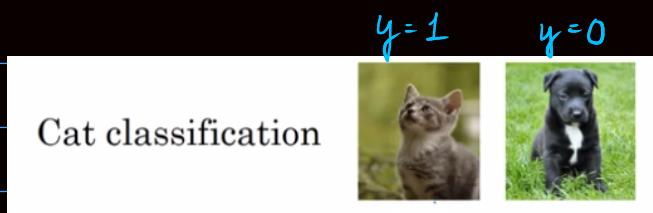
Tip. \* For very big dataset make them split smaller. 10,000 images are enough for dev & test set.

## Mismatched Train / Test Distribution

- \* Make sure the dev and test sets come from the same distribution
- \* Not having a test set might be okay. (only Dev set)

## Bias / Variance





- |  |  |
|--|--|
| ① Train Set Error = 1%.<br>② = 15%.<br>③ = 15%.<br>④ = 0.5%. | Dev Set Error = 11%. → High Variance.<br>or Overfitted<br>= 16%. → High Bias<br>or Underfitted<br>= 30%. → L.B & H.V |
|--|--|
- ↑
- Optimal or Bayes's Error = 0%. {Human Error}

If Optimal Error is 15%. → ②  $\Rightarrow$  L.B & H.V

### BASIC RECIPE FOR ML { Bias - Variance Tradeoff }

- ① High Bias ?  $\xrightarrow{Y}$  ① Bigger n/w   ② Train Longer.  
 (training data performance)
- ↓ N

- High variance ?  $\xrightarrow{Y}$  ① More data   ② Regularization  
 (data set performance)
- ↓
- ③ NN architecture search

DONE

Regularization → Prevents Overfitting

$$\min_{w,b} J(w,b)$$

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

In practice can be omitted

Regularization Parameter  $\lambda$

$$L_2 \text{ Reg } \|w\|^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \leftarrow \text{Used much more often}$$

$$L_1 \text{ Reg } \frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1 \quad \left\{ \begin{array}{l} w \text{ will be sparse} \\ i.e. w \text{ vector will have a lot of zeros which will compress the model} \end{array} \right.$$

$$J(w^{[0]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots) = \frac{1}{m} \sum_{i=1}^n L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2 \quad w : (n^{[0]}, n^{[L-1]})$$

↑

FORBENIUS NORM (L2 Regularization or Weight Decay)

Gradient Descent

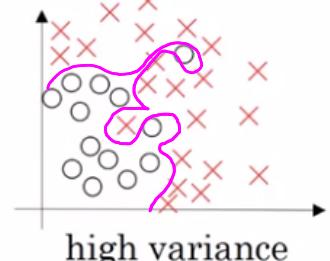
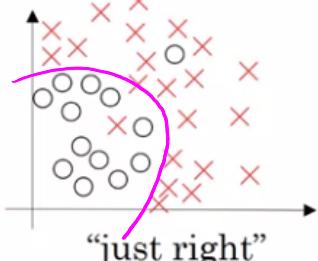
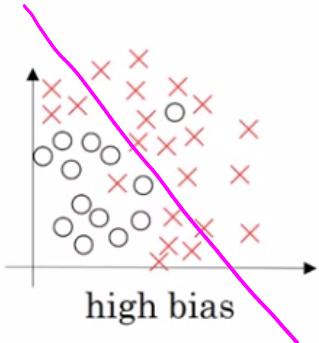
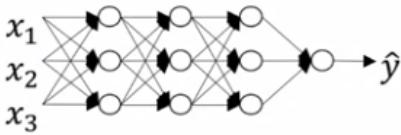
$$dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{2m} w^{[l]}$$

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

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

$$= \left(1 - \frac{\alpha \lambda}{m}\right) w^{[l]} - \alpha (\text{from backprop})$$

# How does regularization prevent overfitting?

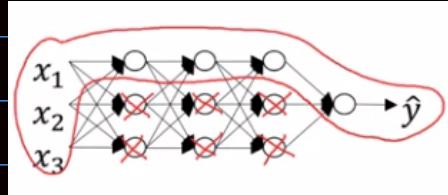


$$J(\mathbf{w}^{[l]}, \mathbf{b}^{[l]}) = \frac{1}{m} \sum_{i=1}^n L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\mathbf{w}^{[l]}\|_F^2$$

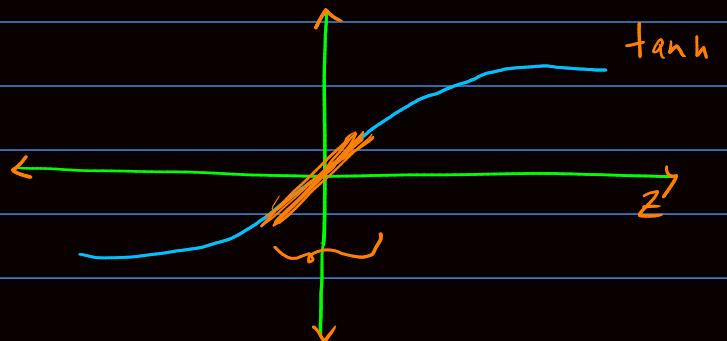
penalizes the wt. matrices from being too large

① If  $\lambda$  very large.,  $\mathbf{w}^{[l]} \approx 0$

reducing the impact of wt. in hidden units like a logistic Reg unit

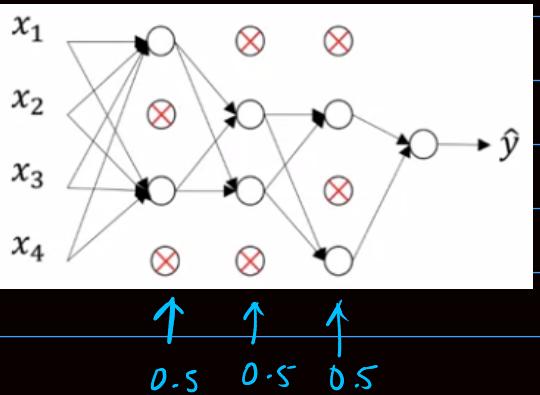
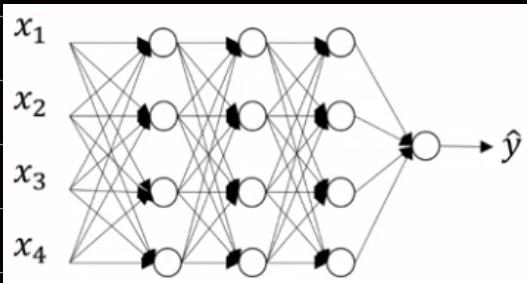


②



If  $\lambda \uparrow$ ,  $\mathbf{w}^{[l]} \downarrow$ ,  $z \downarrow$   $\therefore$   $\lim g(x)$  will end up in the linear regime which is bad.

## DROPOUT REGULARIZATION {\* Mostly used in C.V Problems}



### Implementation

At layer 3

$$\text{keep\_prob} = 0.8$$

Generates Boolean Value  $\rightarrow d3 = \text{np.random.rand}(a3.\text{shape}[0], a3.\text{shape}[1]) < \text{keep\_prob}$

$$a3 = \text{np.multiply}(a3, d3) \quad \# a3 \neq d3$$

$$a3 /= \text{keep\_prob}. \quad \left\{ \begin{array}{l} \text{50 units} \leftrightarrow 10 \text{ units shut off} \end{array} \right\}$$

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

Example:

```
In [46]: keep_prob = 0.8
a3 = np.random.rand(50, 1)
print(a3[:5])
d3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep_prob
print(d3[:5])
a3 = np.multiply(a3, d3)
a3 /= 0.8
print(a3[:5])
```

```
[[0.31604719]
 [0.94878238]
 [0.32860215]
 [0.18180209]
 [0.3065882 ]]

[[False]
 [ True]
 [ True]
 [ True]
 [ True]]
```

```
[[0.
 [1.18597798]
 [0.41075269]
 [0.22725262]
 [0.38323525]]]
```

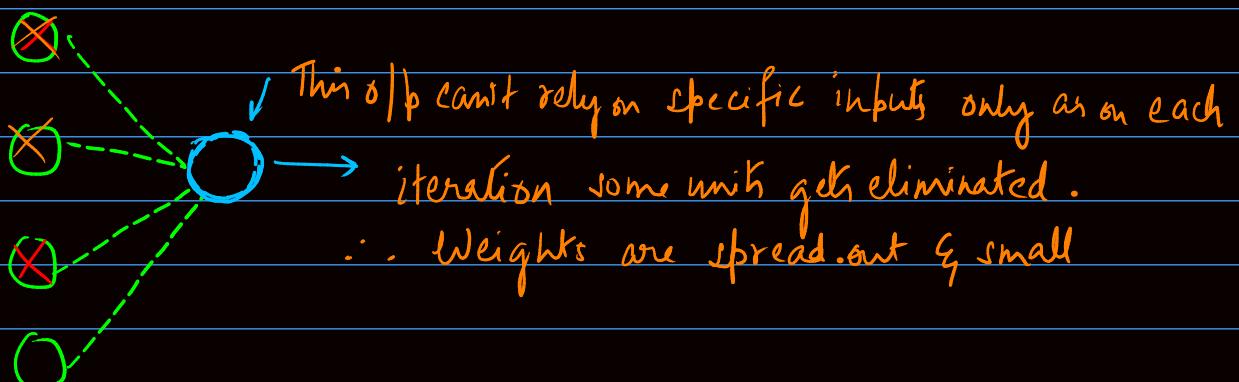
↑  
reduced by 20%.

## Making Predictions at Test Time

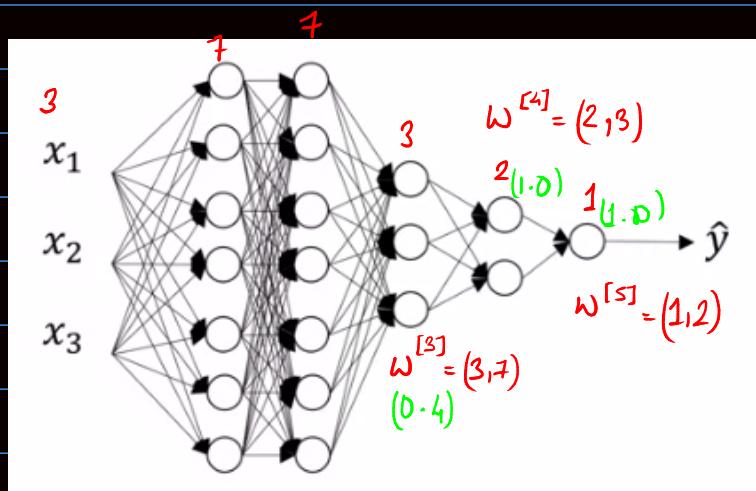
- \* Don't use .drop out { Add noise to our prediction & also we don't want anything random at test time }

## Why does Dropout Work?

Intuition: Can't rely on any one feature, so have to spread out weights  $\rightarrow$  shrink cells. (L2)



- \* Vary keep\_prob in each layer



\* keep\_prob values

$$\begin{matrix} W^{[1]} & W^{[2]} \\ (1,3) & (1,1) \\ (0,1) & (0,5) \end{matrix}$$

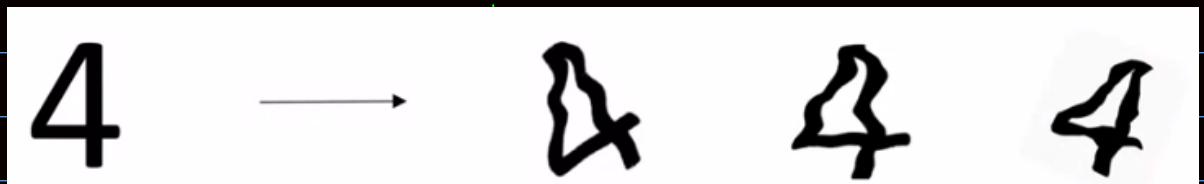
\* Don't generally add dropout to the first layer { keep it close to 1 }

\* Drawback  $\rightarrow$  Not well defined. I think.

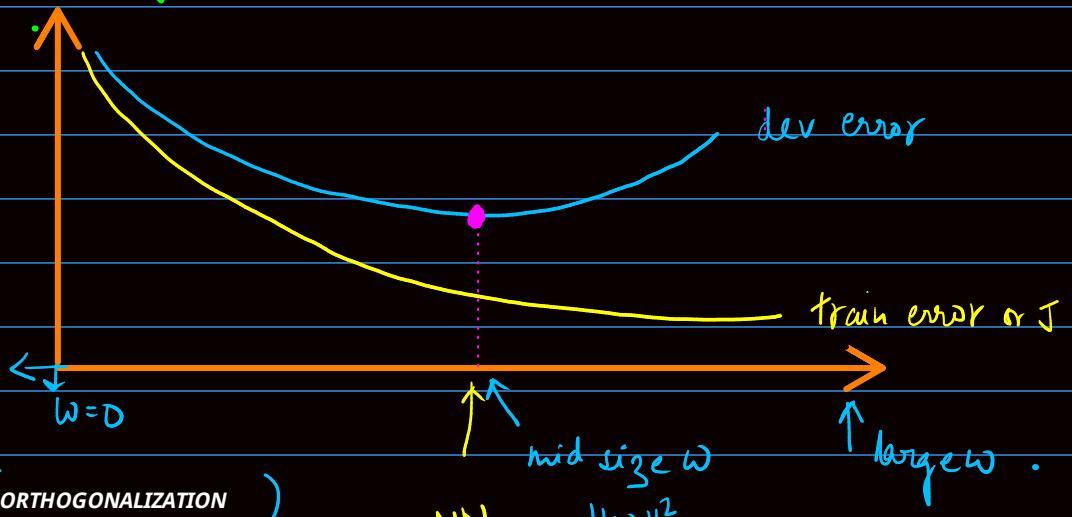
## Other Regularization Techniques

### ① Data Augmentation

Example:



### ② Early Stopping



Any ML Problem (ORTHOGONALIZATION)

One task at a time

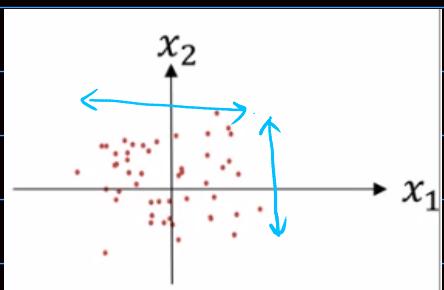
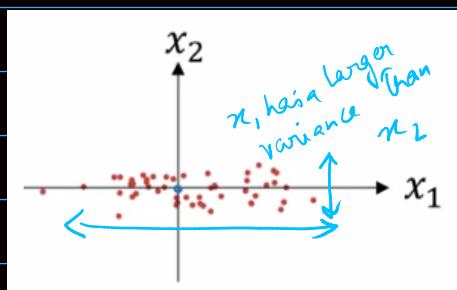
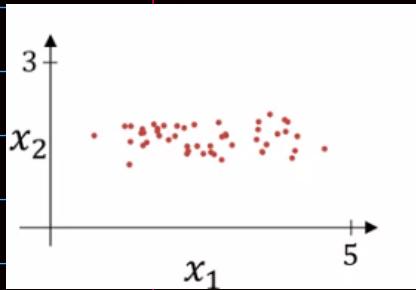
$$\|w\|_F^2 = \text{small norm}$$

1. Optimize Cost fun  $J$ .  
Gradient Descent...  $\Rightarrow J(w, b)$  till here all to dev error
2. Not Overfit  
Regularization  $\Rightarrow$  Separate task

\* But with Early Stopping I have to take care of both task simultaneously

Instead L2 Regularization is better though it requires much more expensive computation

## NORMALIZING TRAINING SETS



$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Subtract  $\mu$

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

$$X := X - \bar{x}$$

Normalize Variance

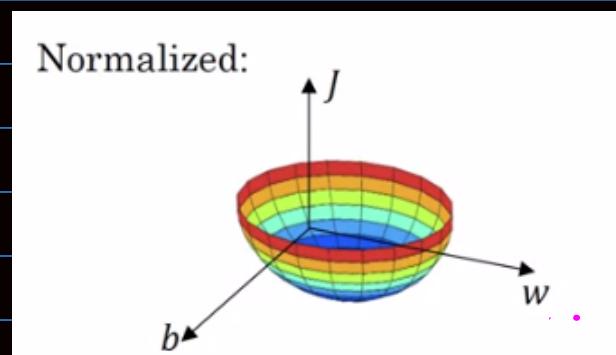
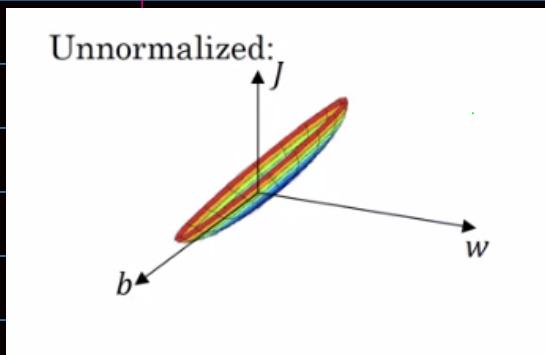
$$\sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \bar{x})^2$$

elementwise

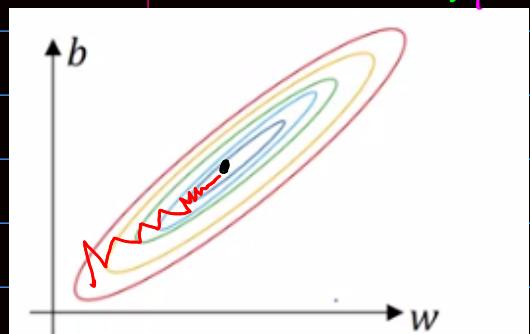
$$X := \frac{X}{\sigma}$$

\* Use the same  $\bar{x}$  &  $\sigma$  for normalizing test set.

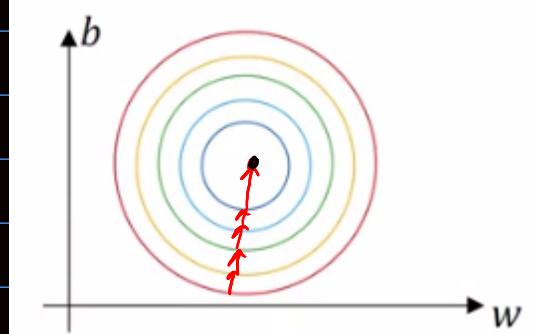
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$



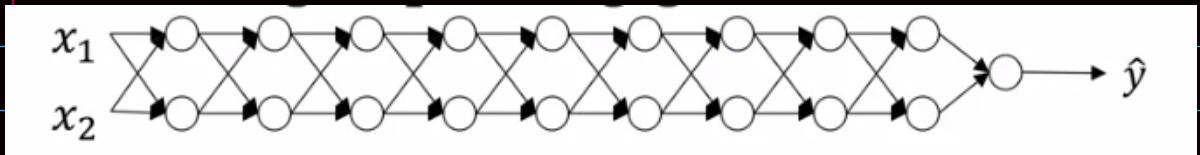
If input's are  $x_1 = (1, \dots, 1000)$   $\uparrow$   
 $x_2 = (0, \dots, 1)$  for similar ranges



not much of an issue!

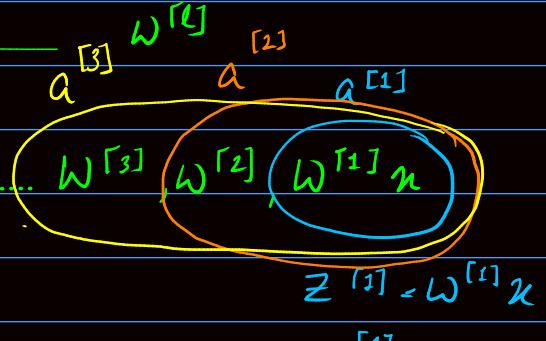


# VANISHING & EXPLODING GRADIENTS



lets tried  $g(x) = \text{identity} \therefore b^{[L]} = 0$

$$\hat{y} = \omega^{[L]} \omega^{[L-1]}$$



$$\omega^{[L]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \quad \hat{y} = \omega^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x$$

$$= 1.5^{(L-1)} x \quad \text{very large.}$$

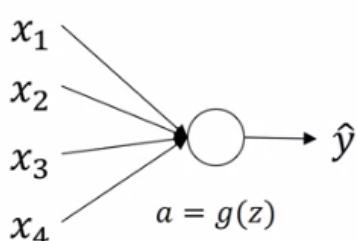
$$y. \quad \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \quad = (0.5)^{L-1} x \quad \text{very small}$$

$\omega^{[L]} > I \rightarrow$  exp. ↑ gradient descent & activation

$\omega^{[L]} < I \rightarrow$  " ↓ " " " "

\* Careful weight Initialization can solve this problem upto some extent

## Single neuron example



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

larger n → smaller  $w_i$

$$\text{Var}(w_i) = \frac{1}{n}$$

∴ In code;

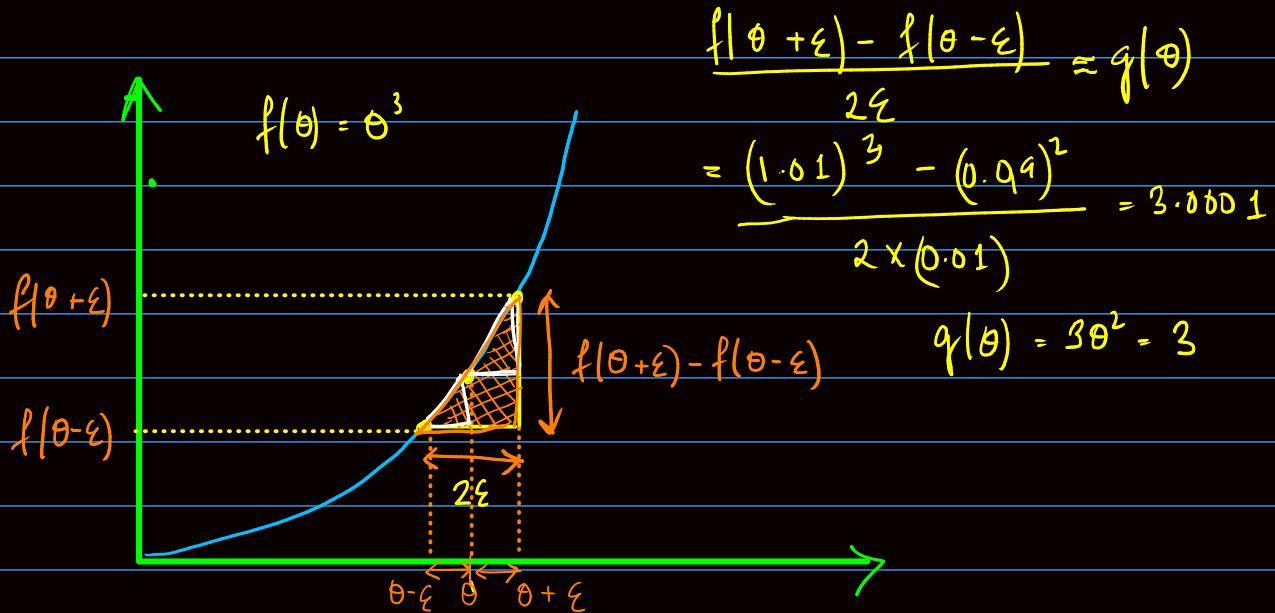
tanh. (Xavier's Initialization)

$$W^{[l]} = np.random.randn(shape) * np.sqrt\left(\frac{1}{n^{[l-1]}}\right)$$

If ReLU is used as activation,

$$W^{[l]} = np.random.randn(shape) * np.sqrt\left(\frac{2}{n^{[l-1]}}\right)$$

### Numerical Approximation of gradients



approx error = 0.0001 { two sided diff way i.e. considering the bigger Δ with both sides from θ

One-sided error = 0.03

$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$	$\frac{\mathcal{O}(\epsilon^2)}{0.01} = 0.0001$	$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}$ $\uparrow \quad \uparrow$ $\text{err: } \mathcal{O}(\frac{1}{\epsilon}) = 0.01$
--	---	---

more accurate

less accurate

## GRADIENT CHECKING { Debugging technique }

Step 1: Take  $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$  and reshape into a big vector  $\theta$ .

Concat

Step 2: Take  $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$  and reshape into a big vector  $d\theta$ .

Concat

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

for each  $i$ :

$$d\theta_{approx}^{[i]} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_i - \epsilon)}{2\epsilon}$$

$$\approx d\theta^{[i]} = \frac{\partial J}{\partial \theta_i}$$



Check  $d\theta_{approx} \stackrel{?}{\approx} d\theta$

$$\rightarrow \frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2}$$

$10^{-7} \rightarrow$  great  
 $10^{-5} \rightarrow$  Not bad but  
 have a look  
 again  
 $10^{-3} \rightarrow$  must be  
 some error  
 somewhere.

\* Don't use in training — only to debug { very slow operation }

\* If algorithm fails grad check, look at the components to identify bug

- \* Remember to include the regularization term if used.
- \* Doesn't work with dropout. { Turn off dropout when using grad check }
- \* Run grad check at random initialization perhaps again after some training of ai sometime backprop works better when  $w, b \approx 0$  }

## WEEK 2

### Batch vs mini-batch gradient descent

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & & x^{(10^6)} \end{bmatrix}_{(n_x \times m)}$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(10^6)} \end{bmatrix}$$

In such big data we split the data into mini batches of 1000 (say)

$$X = \left[ \underbrace{x^{(1)} x^{(2)} \dots x^{(1000)}}_0 \mid \underbrace{x^{(1001)} \dots x^{(2000)}}_{\{2\}} \mid \dots \mid \underbrace{\dots}_{\{5000\}} \right]_{(n_x \times 1000)}$$

$$Y = \left[ \underbrace{y^{(1)} y^{(2)} \dots y^{(1000)}}_0 \mid \dots \mid \dots \mid \underbrace{y^{(1m)}}_{\{5000\}} \right]_{(1 \times 1000)}$$

for  $t = 1 \dots$  range :

One step of gradient descent

Forward Prop on  $X^{[t]}$

$$Z^{[1]} = \omega^{[1]} X^{[t]} + b$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$A^{[L]} = g^{[L]}(Z^{[L]})$$

Over 1000 examples

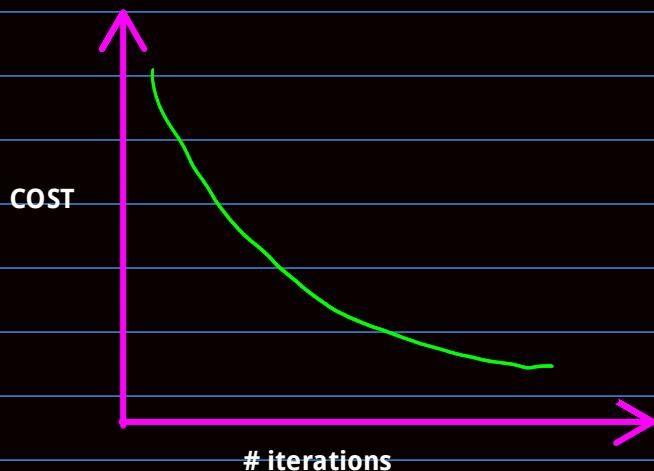
$$\text{Compute cost} = \frac{1}{1000} \sum_{i=1}^0 L(\hat{y}^{(i)}, y^{(i)}) := J^{[t]}$$

Back Prop w.r.t  $J^{[t]}$

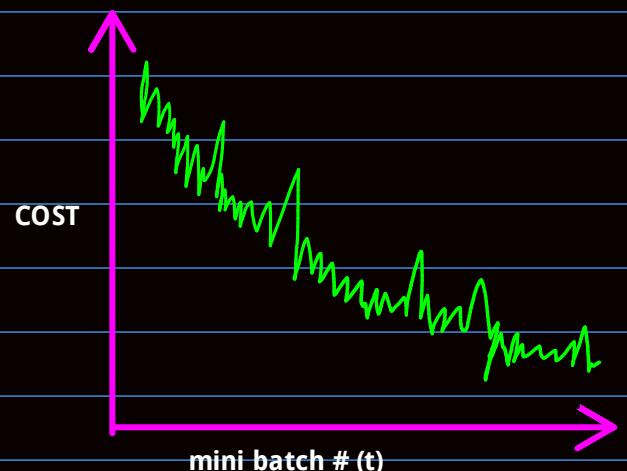
$$\omega^{[t]} = \omega^{[t]} - \alpha d\omega^{[t]}, b^{[t]} = b^{[t]} - \alpha db^{[t]}$$

]  
]

Batch Gradient Descent



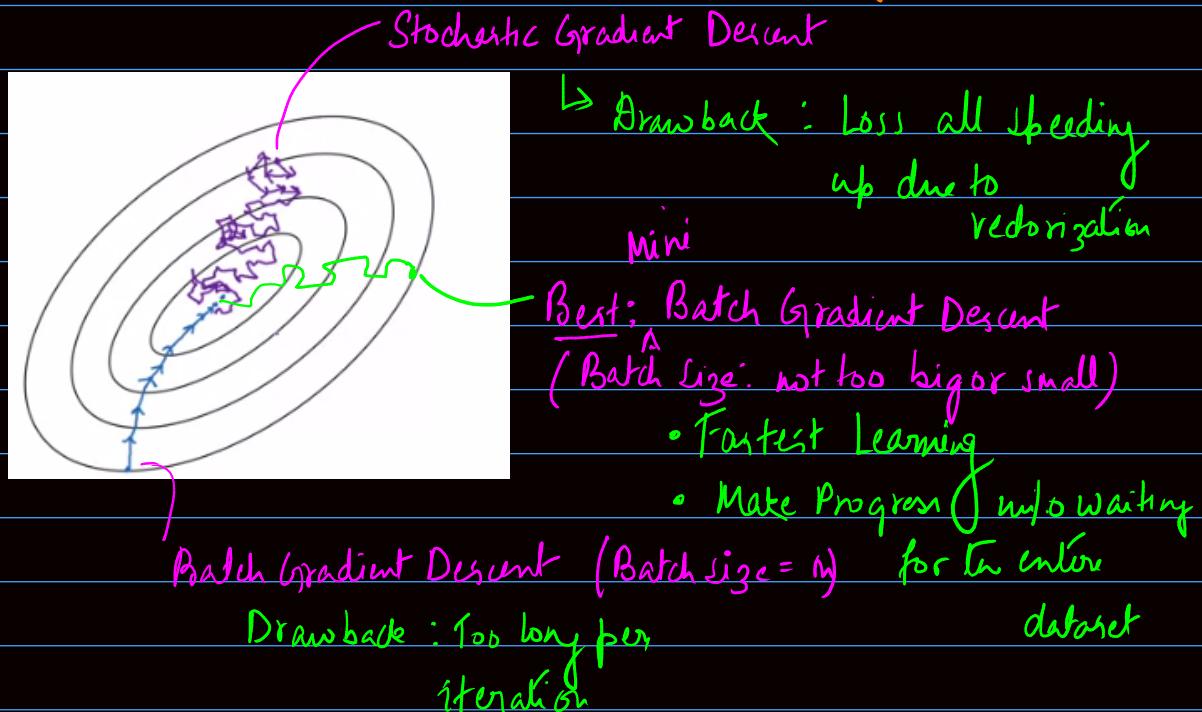
Mini Batch Gradient Descent



plot  $J^{[t]}$  captured using  $x^{[t]}, y^{[t]}$

If mini batch size = m : Batch Gradient Descent  $(X^{[1]}, y^{[1]}) = (x, y)$   
(M.B)

- If M-B size = 1 : Stochastic Gradient Descent  
Each batch is one training example.



\* If small dataset: Batch Gradient Descent ( $m \leq 2000$ )

\* Typical MiniBatch Size: 64, 128, 256, 512

\* Make sure minibatch size in CPU/GPU Memory

### EXponential WEIGHTING AVERAGES

#### Temperature Plot

$$\theta_1 = 40^\circ F$$

$$\theta_2 = 40^\circ F$$

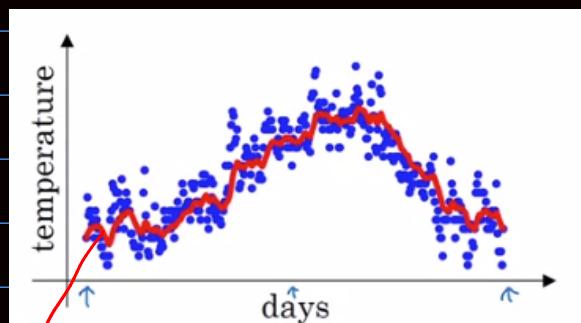
$$\theta_3 = 45^\circ F$$

:

:

$$\theta_{180} = 60^\circ F$$

$$\theta_{181} = 56^\circ F$$



Moving average Plot

Let  $V_0 = 0$

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

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

$$V_3 = 0.9V_2 + 0.1\theta_3$$

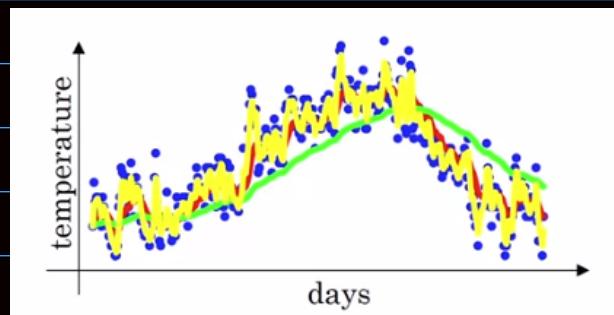
$$V_t = 0.9V_{t-1} + 0.1\theta_t$$

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

When  $\beta = 0.9$

$\beta = 0.98$

$\beta = 0.5$



$$v_t = \beta v_{t-1} + (1-\beta)\theta_t$$

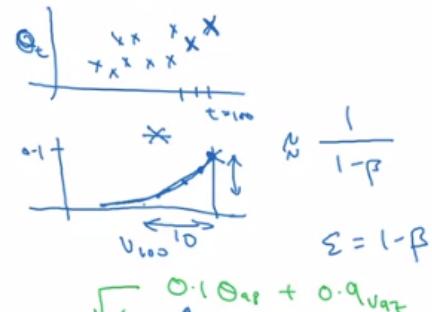
$$v_{100} = 0.9v_{99} + 0.1\theta_{100}$$

$$v_{99} = 0.9v_{98} + 0.1\theta_{99}$$

$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$

...

$$\begin{aligned} \underbrace{v_{100}}_{\text{Example}} &= 0.1\theta_{100} + 0.9 \underbrace{(0.1\theta_{99})}_{\text{green}} + 0.9 \underbrace{(0.1\theta_{98})}_{\text{green}} \\ &= 0.1\theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1(0.9)^2 \theta_{98} + 0.1(0.9)^3 \theta_{97} + 0.1(0.9)^4 \theta_{96} \\ &\approx 0.9^{10} \theta_{100} \approx 0.35 \approx \frac{1}{e} \end{aligned}$$



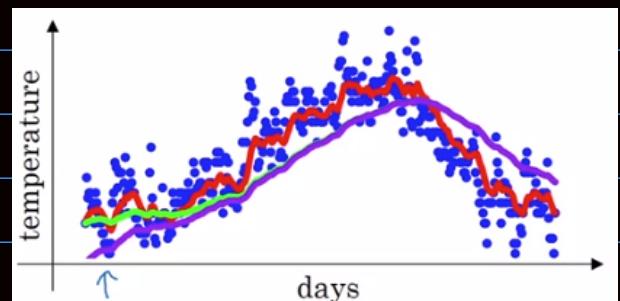
$$\frac{(1-\varepsilon)^{\frac{1}{\varepsilon}}}{\alpha} = \frac{1}{e}$$

$$\varepsilon = 0.02 \rightarrow 0.98^{\frac{1}{0.02}} \approx \frac{1}{e}$$

Andrew Ng

## Bias Correction

When  $\beta = 0.98$  we should get the green line but we get the purple line which starts off really low.



$$V_0 = 0$$

$$V_1 = 0.98 V_0 + 0.02 \theta_1$$

$$V_2 = 0.98 V_1 + 0.02 \theta_2 = 0.9 (0.98 V_0 + 0.02 \theta_1) + 0.02 \theta_2$$

$$= \underbrace{0.96 \theta_1}_{\text{much small}} + \underbrace{0.02 \theta_2}$$

Rather,

Correction

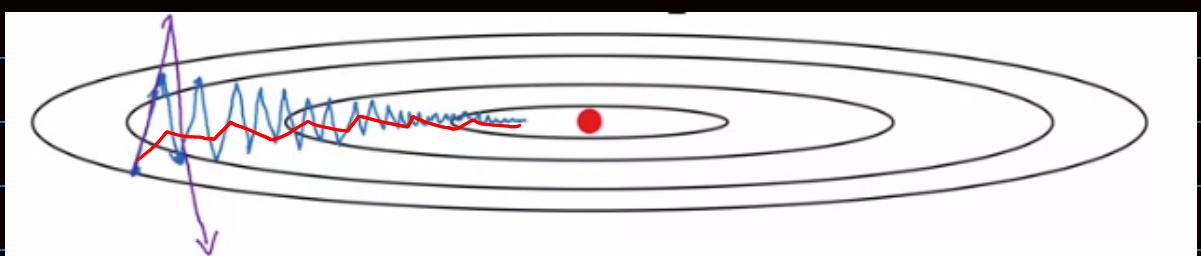
$$\frac{V_t}{1 - \beta^t}$$

$$t=2 : 1 - \beta^t = 1 - (0.98)^2 = 0.0396$$

$$\frac{V_2}{0.039} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

When  $t$  is large  $1 - \beta^t \approx 1$ , so the purple & green line becomes one

### GRADIENT DESCENT WITH MOMENTUM



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

On iteration  $t$ :

Compute  $dw, db$  on current mini batch

$$\begin{cases} V_{dw} = \beta V_{dw} + (1-\beta) dw \\ V_{db} = \beta V_{db} + (1-\beta) db \end{cases} \quad \left\{ \text{similar } V_\theta = \beta V_\theta + (1-\beta) \right.$$

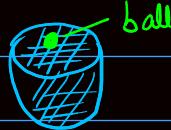
Slow learning  $\uparrow$

faster learning  $\leftrightarrow$

$$\omega = \omega - \alpha V_{dw}; b = b - \alpha V_{db}$$

\* In practice bias correction is not really necessary

Another analogy,



$$V_{dw} = \beta V_{dw} + (1-\beta) dw \quad \text{provides acceleration to the ball.}$$

$$V_{db} = \beta V_{db} + (1-\beta) db \quad \begin{matrix} \text{Friction} \\ \uparrow \end{matrix} \quad \begin{matrix} \text{provides velocity} \end{matrix}$$

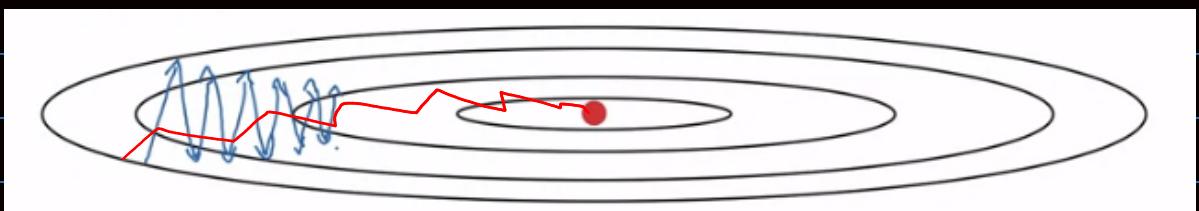
So here we have 2 hyperparameters  $\alpha, \beta$

common  $\beta = 0.9 \rightarrow$  Iterating over the last 10 gradients

\* Sometimes in practice  $(1-\beta)$  term is omitted from the equation.

$$V_{dw} = \beta V_{dw} + dw \rightarrow \text{But not much preferred}$$

### Root Mean Square Prop. (RMP Prop)



let,



On iteration  $t$ :

Compute  $dw, db$  on current minibatch

$$S_{dw} = \beta S_{dw} + (1-\beta) dw^2 \quad \begin{matrix} \leftarrow \text{element wise} \\ \downarrow \text{small} \end{matrix}$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2 \quad \begin{matrix} \leftarrow \text{large} \end{matrix}$$

(b)

$$\omega = \omega - \frac{\alpha \text{dw}}{\sqrt{S_{\text{dw}}}} \quad b = b - \frac{\alpha \text{db}}{\sqrt{S_{\text{db}}}}$$

*wrong.*

\*  $\sqrt{S_{\text{dw}}} + \epsilon$  or  $\sqrt{S_{\text{db}}} + \epsilon$  in which where  $\epsilon = 10^{-8}$

\* Adaptive Momentum

ADAM OPTIMIZATION (Momentum & RMS Prop Combination)

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

On iterate t:

compute  $\text{dw}, \text{db}$  using current minibatch

Momentum  $\rightarrow V_{\text{dw}} = \beta_1 V_{\text{dw}} + (1-\beta_1) \text{dw}^2, V_{\text{db}} = \beta_1 V_{\text{db}} + (1-\beta_1) \text{db}^2$

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

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

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

$$\omega = \omega - \frac{\alpha V_{\text{dw}}^{\text{corrected}}}{\sqrt{S_{\text{dw}}^{\text{corrected}}} + \epsilon}, b = b - \frac{\alpha V_{\text{db}}^{\text{corrected}}}{\sqrt{S_{\text{db}}^{\text{corrected}}} + \epsilon}$$

Hyperparameters:

$\alpha$  = needs to be tune

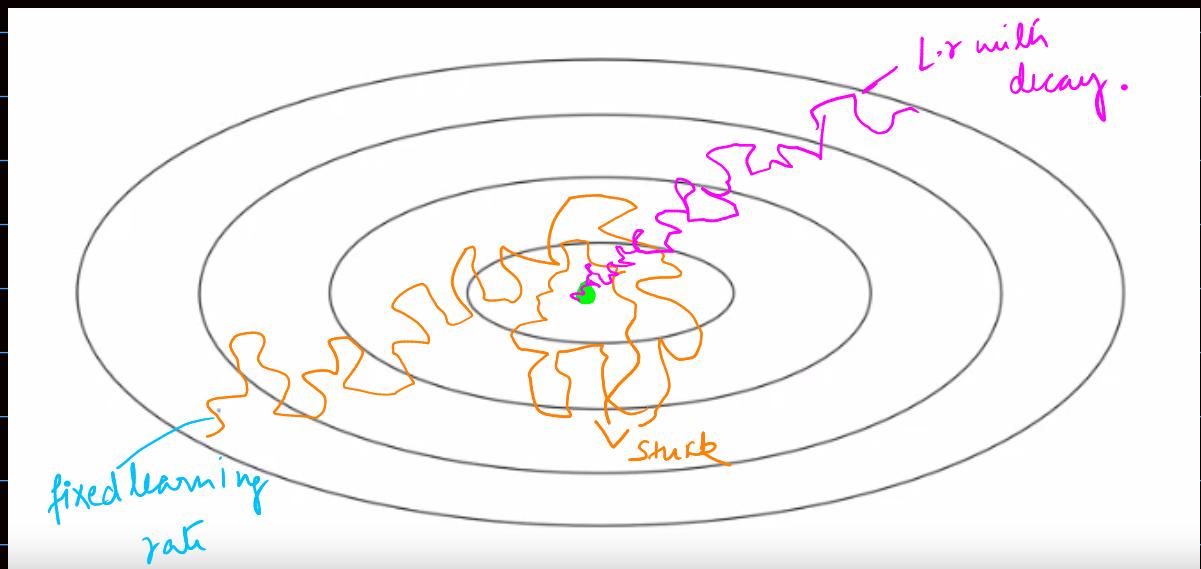
$\beta_1 = 0.9$  (dw)

$\beta_2 = 0.999$  (dw<sup>2</sup>)

$\epsilon = 10^{-8}$

} Usually not tweaked.

## Learning Rate Decay



1 epoch  $\rightarrow$  1 pass through the entire data

hyperparameter  $\alpha = \frac{1}{1 + \text{decay rate} * \text{epoch}} \alpha_0$

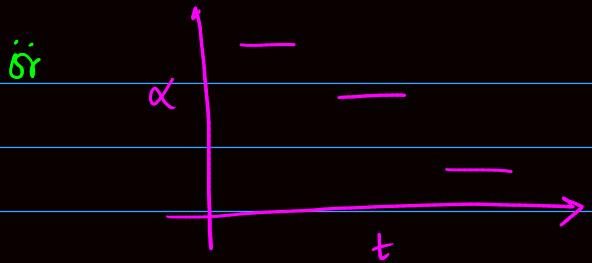
Let  $\alpha_0 = 0.2$ , decay-rate = 1.

Epoch	$\alpha$
1	0.1
2	0.067
3	0.05
4	0.04

Other L.r. decay

Exponential Decay,  $\alpha = 0.95^{\text{epoch num}} \alpha_0$

or  $\alpha = \frac{k}{\sqrt{\text{epoch num}}} \alpha_0$  or  $\frac{k}{\sqrt{t}} \alpha_0$

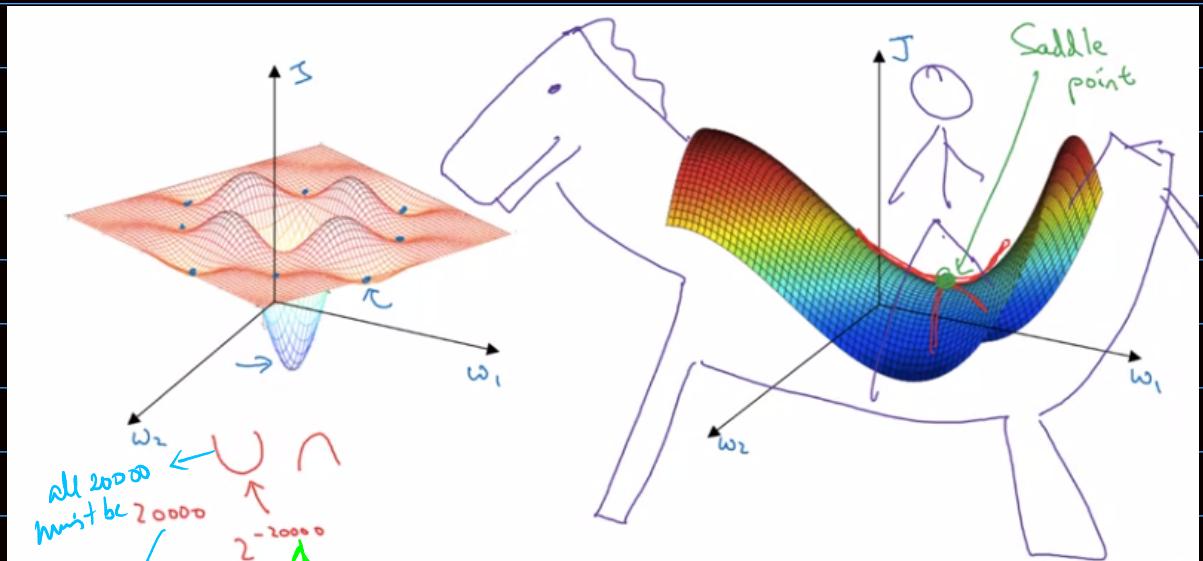


discrete stepwise decay

or

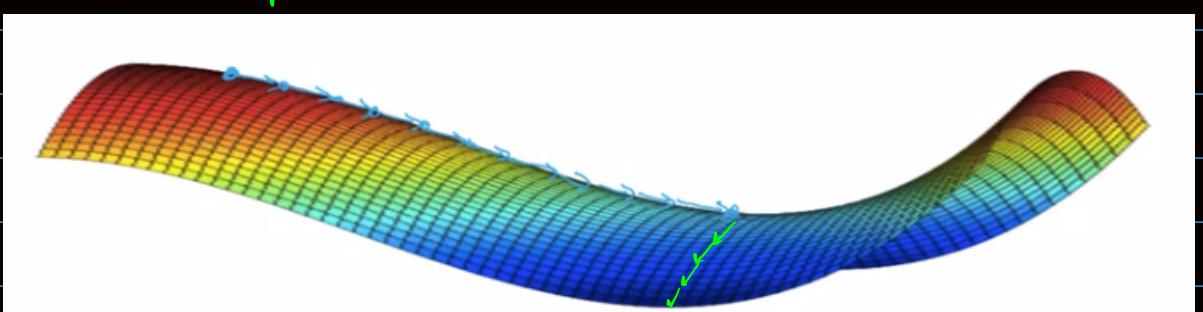
\* Manual Decay

### LOCAL OPTIMA IN NEURAL NETWORKS



Andrew Ng

### Problem of Plateau



- Unlikely to get stuck in a bad optima.
- Plateau can make learning very slow (Use Adam, RMS Prop)

## WEEK 3

### Tuning Process

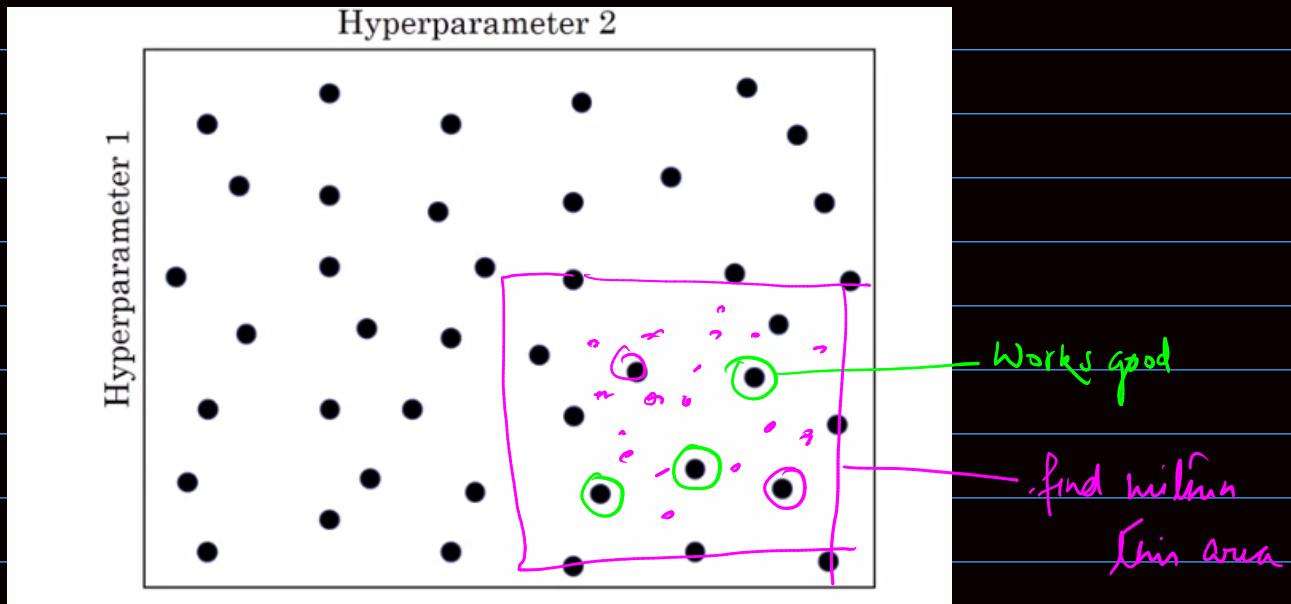
In order  
{Tip}

$$\alpha \xrightarrow{\text{decay}} \beta \quad | \quad \begin{matrix} \text{hidden units} \\ \text{mini batch size} \end{matrix} \quad | \quad \begin{matrix} \# \text{ layers} \\ \text{learning rate} \end{matrix}$$

$$\beta_1, \beta_2, \epsilon$$

\* Try Random Values: Don't follow a pattern strictly

\* Converge fine



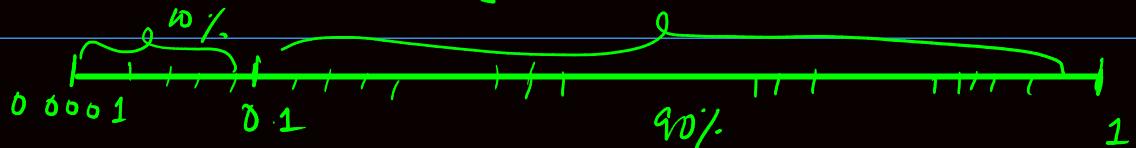
Using an appropriate scale to pick hyperparameters

For eg

$$① \quad n^{[l]} = [50, 100]$$

$$② \quad \# \text{ layers} = [2, 6]$$

$$③ \quad \alpha = [0.0001, 1]$$



- Use log scale here, which will be more reasonable



$$\textcircled{4} \quad \beta = 0.9 \dots 0.999$$

$\downarrow$

Avg over  $\sqrt{1000}$  values

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

Don't use linear scale



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

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

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

## \* To organize hyperparameters

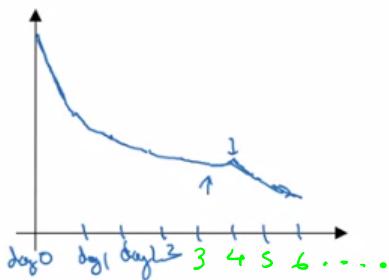
Intuition

\* Hyperparameters do get stale overtime so re-evaluate them

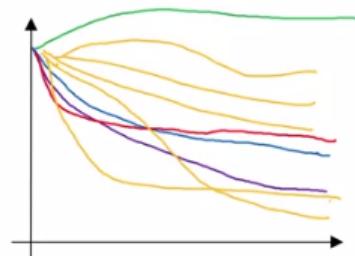
occasionally

Approach

Babysitting one model

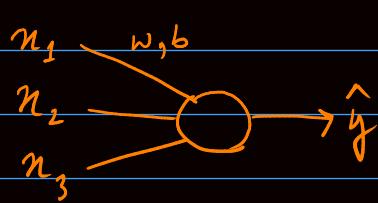


Training many models in parallel



Choose the approach according to the computational expenses

## BATCH NORMALIZATION

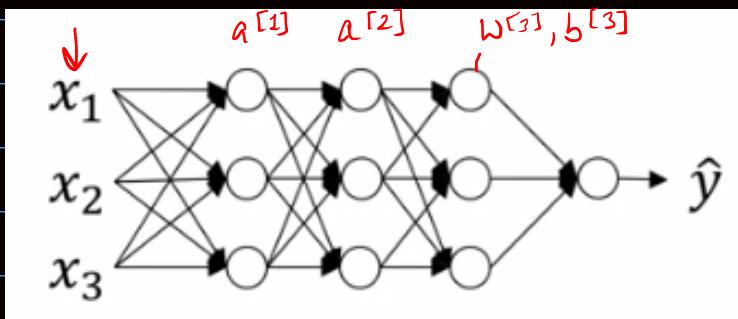


$$y = \frac{1}{m} \sum x^{(i)}$$

$$X = X - y$$

$$\sigma^2 = \frac{1}{m} \sum_i x^{(i)2}$$

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



In bigger model we normalize after every step

\* Normalizing  $z^{[2]}$  or  $a^{[2]}$  to train  $w^{[3]}$  &  $b^{[3]}$  faster is a debate but  $z^{[2]}$  is used more often

\* Given some intermediate values in NN  $z^{(1)}, z^{(2)} \dots \dots z^{(n)}$

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

$\downarrow$   
 $z^{(L)}(i)$

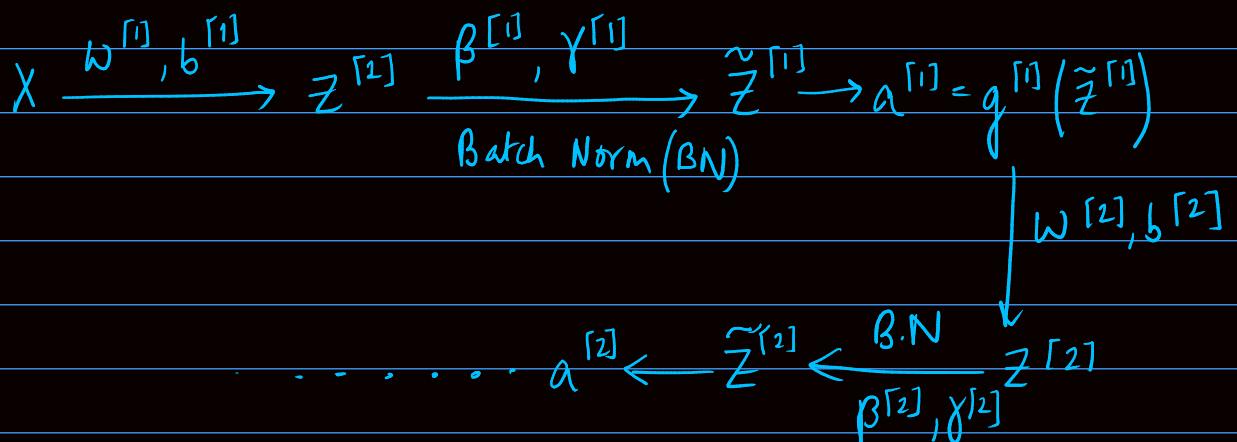
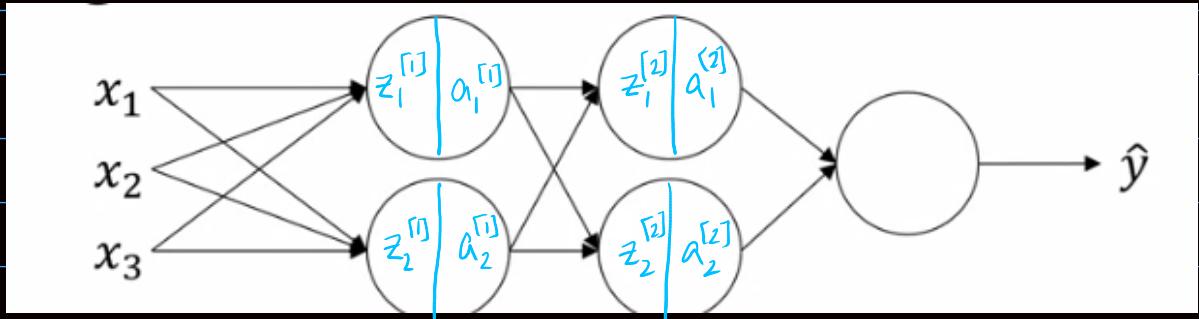
$$\sigma^2 = \frac{1}{m} \sum_i (z_i - y)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - y}{\sqrt{\sigma^2 + \epsilon}} \approx \text{Mean } \approx 0 \quad \text{Var } = 1$$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta \quad \begin{cases} \text{learning params of the model.} \\ \text{If we want diff mean \& var} \end{cases}$$

If  $\gamma = \sqrt{\sigma^2 + \epsilon}$  &  $\beta = y$  then  $\tilde{z}^{(i)} = z^{(i)}$

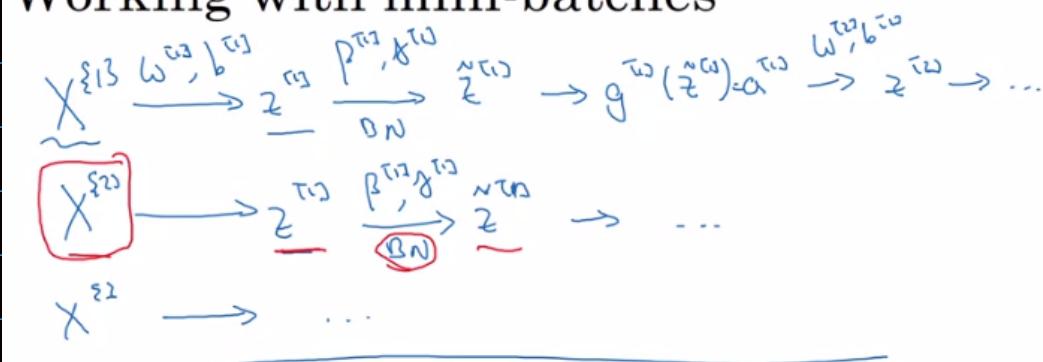
So, in batchnorm we use  $\tilde{z}^{[l](i)}$  instead of  $z^{[l](i)}$



Parametros :  $\omega^{[1]}, b^{[1]}, \omega^{[2]}, b^{[2]}, \dots \dots$   
 $\beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots \dots$

B-N Paper  
X Momentum  $\beta$

## Working with mini-batches



Parametros:  $\omega^{[1]}, \cancel{b^{[1]}}, \beta^{[1]}, \gamma^{[1]} \rightarrow \tilde{Z}^{[1]} = \omega^{[1]} a^{[1]} + \boxed{b^{[1]}}$

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

This gets cancelled in the mean subtraction step

$b^{[l]}$  is not required if BN is used.

Implementing gradient descent.

for  $t = 1 \dots \text{num of minibatches}$  :

{ Compute Forward Prop on  $X^{[t]}$

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

Use back Prop to compute.

$dW^{[l]}$ ,  $dB^{[l]}$ ,  $dY^{[l]}$

Update

$$W^{[l]} = W^{[l]} - \alpha dW^{[l]}$$

$$B^{[l]} = B^{[l]} - \alpha dB^{[l]}$$

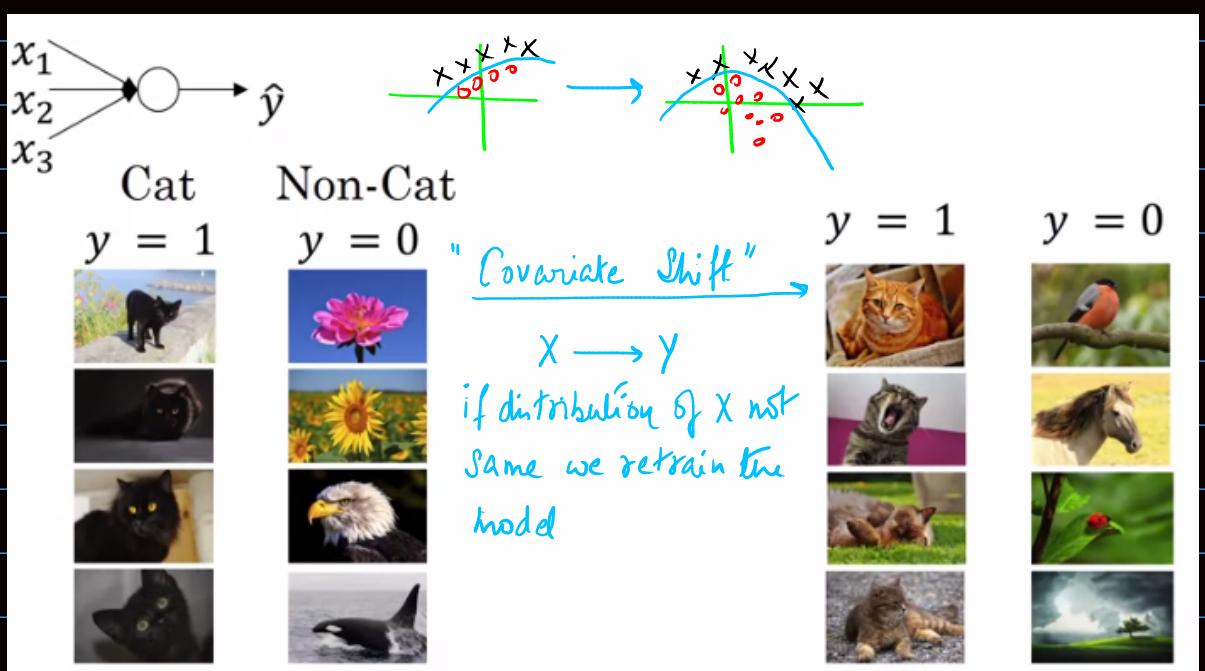
$$Y^{[l]} = Y^{[l]} - \alpha dY^{[l]}$$

Work with Momentum / Adam / RMS Prop

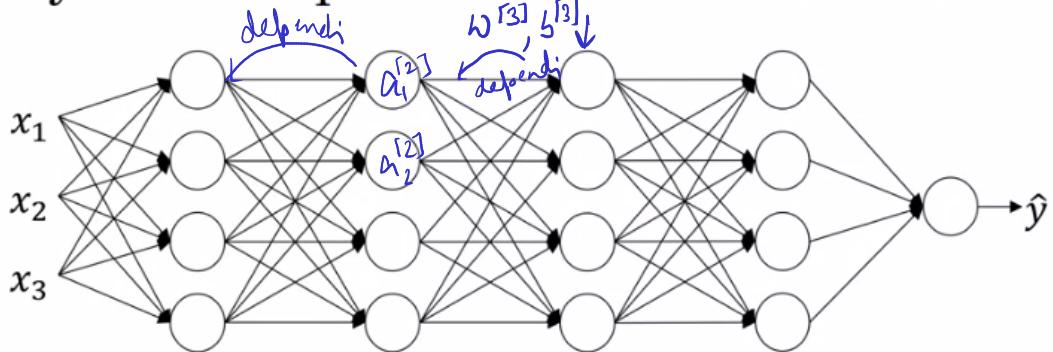
}

Why Batch Norm works?

Learning on shifting input distribution

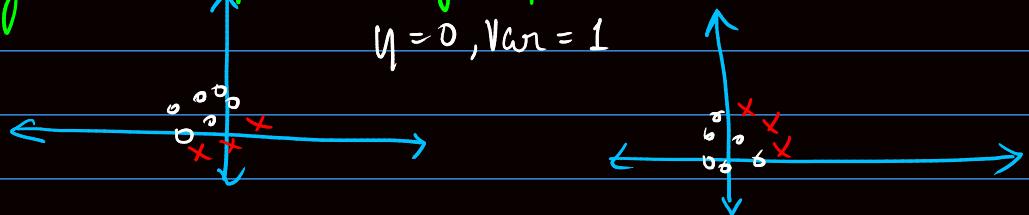


Why this is a problem with neural networks?



Each layer deep into the network is directly dependent on the parameter changes in the previous layer

- \* So, at least to make the distribution uniform with same mean & variance B.N is used. & also amount of adaptability of the later layers to the previous layers ↑



## 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. (Both multiplicative & additive noise)
- This has a slight regularization effect.

*The bigger minibatch size we reduce The regularization effect*

Batch Norm at test time:

Regular  
B.N. on  
Train set →

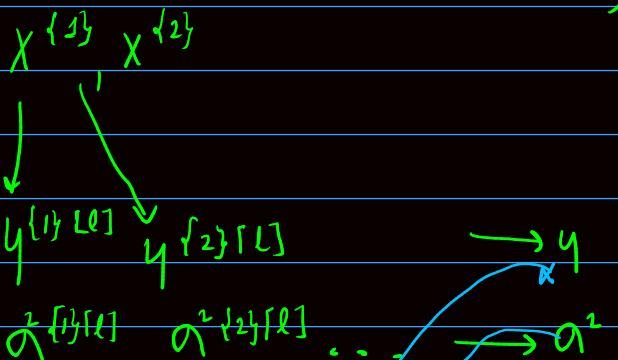
$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

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

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

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

$\hat{y}, \hat{\sigma}^2$ : estimate using exponentially weighted average (across mini batch)



Test time

$$z_{\text{norm}} = \frac{z - \hat{y}}{\sqrt{\hat{\sigma}^2 + \epsilon}}$$

$$\hat{z} = \gamma z_{\text{norm}} + \beta$$

\* ML Frameworks will do this automatically.

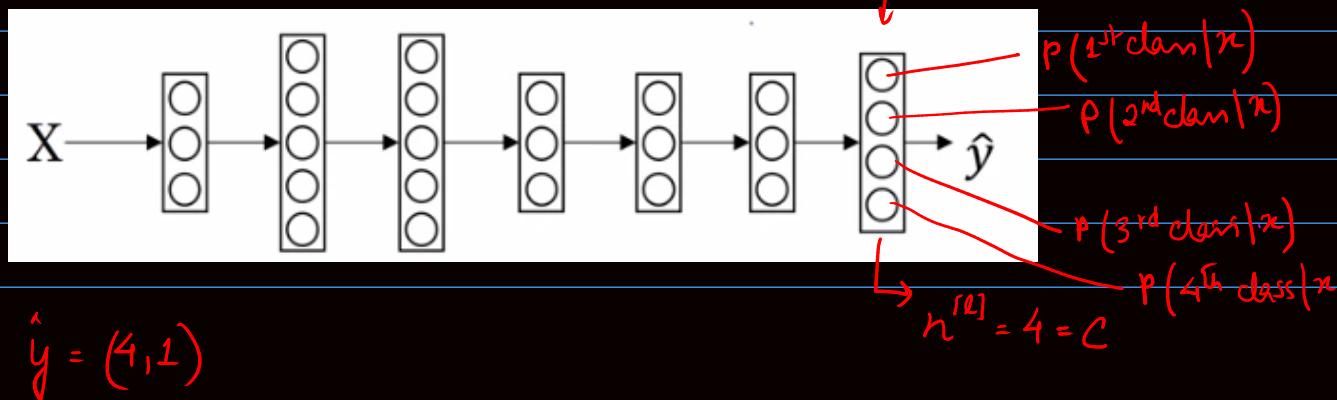
## SOFTMAX REGRESSION

Recognizing cats, dogs, and baby chicks



$$C = \# \text{classes} = 4 \quad (0, \dots, 3)$$

layer L



$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$$

Activation:

$$t = e^{z^{[L]}} \quad \{ \text{element wise exponentiation}$$

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^4 t_i} \quad \rightarrow a_i^{[L]} = \frac{t_i}{\sum_{i=1}^4 t_i}$$

Ex

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

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

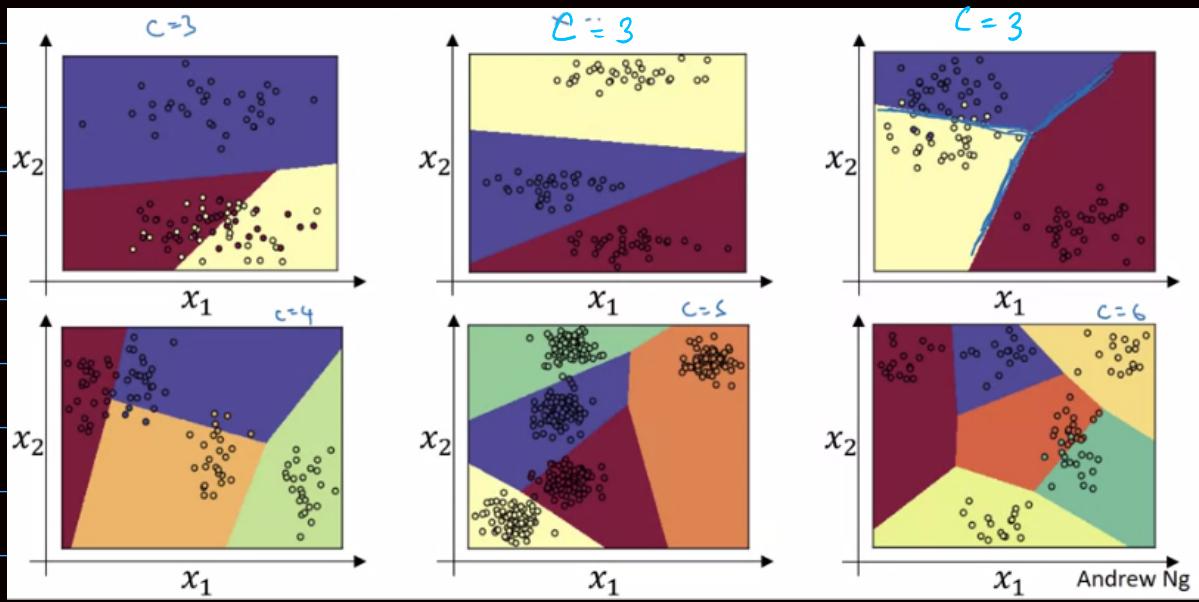
$$a^{[2]} = \frac{148.4}{176.3} = 0.842 \quad a^{[2]} = \frac{7.4}{176.3} = 0.042.$$

$$a^{[3]} = \frac{0.4}{176.3} = 0.002 \quad a^{[4]} = \frac{20.1}{176.3} = 0.114$$

Softmax Activation:

$$a^{[L]} = \frac{t}{\sum_{i=1}^4 t_i}$$

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



Training a softmax classifier.

(4,1)

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

softmax

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

"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  $a^{[L]} = \begin{bmatrix} 0.842 \\ 0.157 \end{bmatrix}$

Loss function

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \text{ cat}$$

but,

$$a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$y_1 = y_3 = y_4 = 0; y_2 = 1$$

$$L(\hat{y}, y) = - \sum_{i=1}^4 y_i \log \hat{y}_i$$

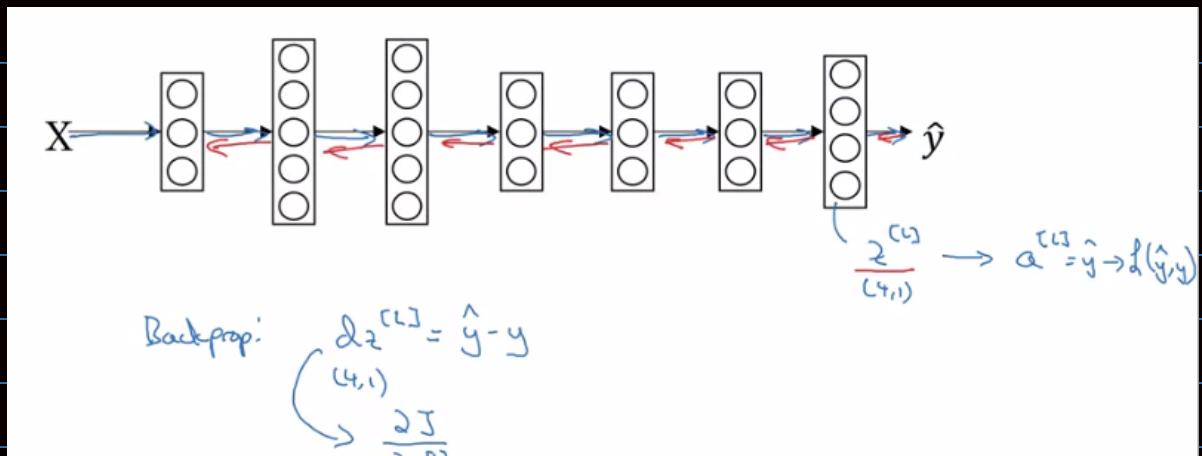
$$= -y_2 \log \hat{y}_2 = -1 \log \hat{y}_2 \quad \therefore \hat{y}_2 \uparrow$$

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

$$Y = [y_1^{(1)} \ y_2^{(2)} \ y_3^{(3)}, \dots]$$

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

## Gradient Descent



# 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)