

## COURSE 2: Improving Deep Neural Networks.

Start date

29.12.18

### week 1: Practical Aspects of Deep Learning.

10.30 pm

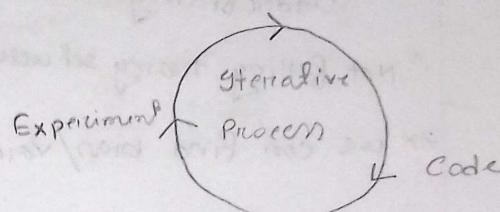
#### P1: Setting up your machine learning Application.

#### L1: Train / Dev / Test sets. [ 11:00 AM 30-12-2018 ]

For a machine learning algorithm, one has to define,

- i) Number of layers.
- ii) Number of hidden units.
- iii) Learning rates.
- iv) Activation functions.

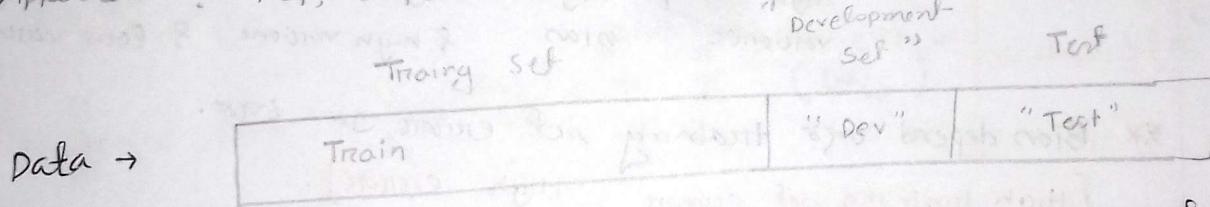
Idea



e.t.c and other hyperparameters.

[ Stochastic gradient descent code with constant learning rate value giving output better result ]

Application :- NLP, Computer vision, speech, structured data, [ Ads, search, securities ]



→ Training ↗ Parameter combination check ( result ) Development set  
→ Development set ↗ Model Test set Evaluate ( result )

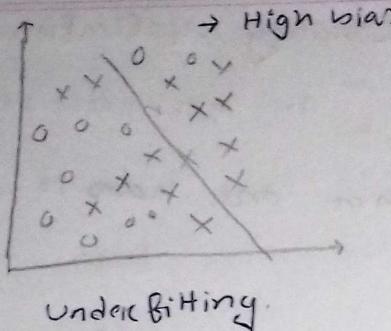
→ [ make sure that development and test set come from same distribution. ( otherwise model may fail in test set ) ]

\*\* [ Test set gives us the unbiased estimate of a model ]

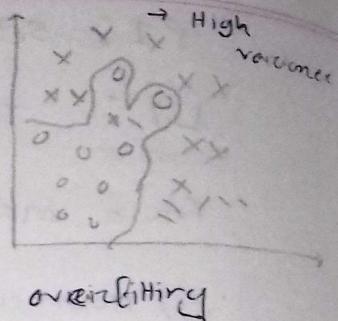
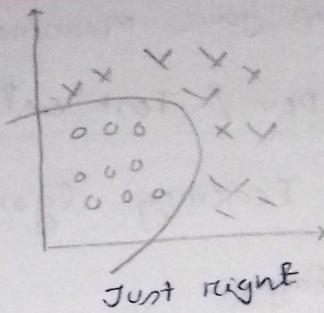
\*\* Not having a test set might be okay.

[ Train / dev / test sets help us to iterate more quickly and gives us good parameter estimation ].

## L2: Bias / Variance. [ 11:20 AM 30-12-18 ].



"Not fitting training set well"



"Not fitting test set well"

\*\* we can find bias/variance by taking train/ dev net error.

Train net error :	1%	15%	15%	0.5%
Dev net error :	11%	16%	30%	1%
	High variance	High bias	High bias & high variance	Low variance
				High variance & low variance.

\*\* Bias depend on training net error  $\Rightarrow$  Dev net error.

[ High training net errors, High error ].

\*\* Variance depend on training net error  $\Rightarrow$  Dev net error.

$\Rightarrow$  Low var error  $\Rightarrow$  Dev net error.

[ High difference in train and dev net errors, High variance ].

## L3 : Basic Recipe For Machine Learning

[ 11:40 AM 30-12-18 ]

High bias  
(training data performance)

→ Bigger network.  
Train longer.  
( NN architecture, search )

High variance  
( dev net performance )  
 $\downarrow N$   
Done

→ more data.  
regularization.  
( NN architecture, search )

P2: Regularizing your neural network. [Theoretical part but important: minimize L2]

L1: Regularization. [1:01 PM, 30.12.18]

Logistic regression :-  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$

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

Regularization parameter

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

$$\rightarrow \text{L}_2 \text{ regularization} \quad \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w. \quad \text{commonly used}$$

$$\text{L}_1 \text{ regularization} \quad \frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1.$$

Neural network :-

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

$$\text{"Frobenius norm"} \quad \|w^{[l]}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2 \quad \text{w.shape} = (n^{[l]}, n^{[l-1]})$$

Backpropagation.

$$\left\{ \begin{array}{l} dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \\ w^{[l]} = w^{[l]} - \alpha dw^{[l]} \end{array} \right. \quad \begin{array}{l} \text{Code J} \\ \text{implement} \\ \text{multi-} \end{array}$$

After expanding

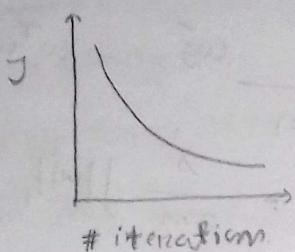
$$\left\{ \begin{array}{l} dw^{[l]} = w^{[l]} - \alpha \left[ (\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \right] \\ = w^{[l]} - \frac{\alpha \lambda}{m} w^{[l]} - \alpha (\text{from backprop}) \end{array} \right. \quad \begin{array}{l} \text{Code J} \\ \text{implement} \\ \text{multi-} \end{array}$$

$\rightarrow$  Hence weight is decreasing in each step. That is why L2 regularization is called "weight decay"

L2: why regularization reduces overfitting? [1:25 pm 30.12.18].

\*\* Regularization parameters penalizes weight matrix.

→ If regularization parameter ( $\lambda$ ) is very large then weights will be nearly zero(0) and our network become roughly linear.



[ If we use regularization const ( $\lambda$ ) will decrease monotonically]

L3: Dropout Regularization [1:35 pm 30.12.18]

[ A copy of our network is made. We will set a probability threshold. Neurons under this threshold will be dropped off from our network.]

Implementing dropout :- ("generated dropout")

Illustrate with layer  $l=3$ , keepprob = 0.8

$d_3 = \text{np.random.rand}(a_3.shape[0], a_3.shape[1]) < \text{keepprob}$ .  
[ Create random matrix of 0's and 1's with 80% 1's (20% 0's)]

$a_3 = \text{np.multiply}(a_3, d_3)$ . [ Activation Function (or 20% 0's)]

$a_3 /= \text{keepprob}$ . [ Generated dropout]

→ making predictions in test time:

$$a^{[0]} = x$$

No drop out,

$$z^{[1]} = w^{[1]} \cdot a^{[0]} + b^{[1]}$$

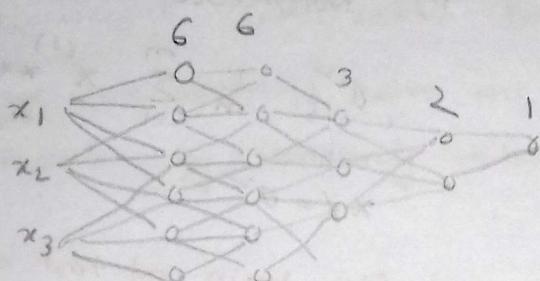
$$a^{[1]} = g^{[1]}(z^{[1]})$$

⋮  
g

training  
data

## L4: understanding Dropout [1.55 pm 30.12.18]

- \* Dropout randomly knocks out neurons of a network.
- \*\* Can not rely on any one feature, so have to spicce out weight. [shrink weight] similar as L2 regularization.



\*\* we can vary keep prob. in every layer  
\*\* input layer do feature Normalise 41.71, 22.5

- \*\* High dropout of neuron with lower keep prob.

- \*\* Low dropout of neuron with higher keep prob.

[Mainly used in computer vision] [overfit or not use drop out]

Problem [Cost function gradient sum 225 m.]

L5:- other regularization methods.

[2.05 pm 30.12.18]

[picture flip over training size about 225]

i) Data Augmentation

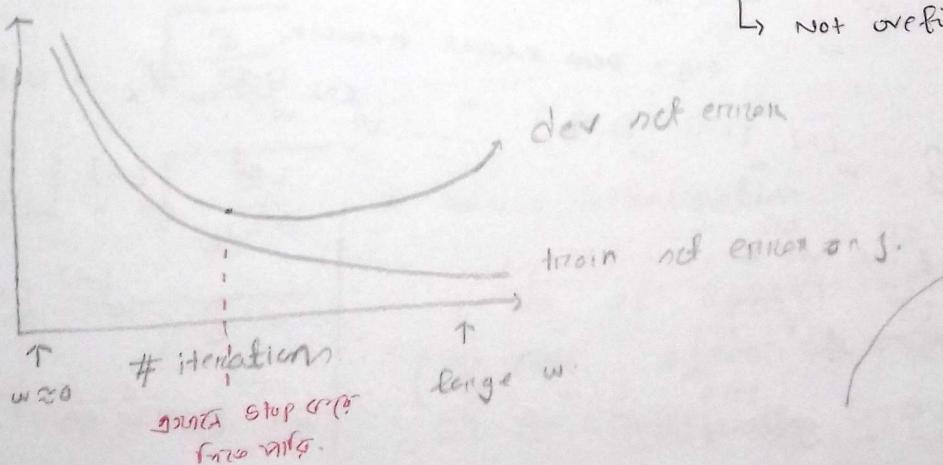
[Augmenting data change cost func.]

ii) Early stopping. (low prob 2018)

orthogonalization.

↳ optimization cost function.

↳ not overfit.



related material: Log + cosine  
and 2018

### P3: Setting up your optimization problem.

[ - solving optimization  
technique for eg. optimization  
problems ]

#### L1: Normalizing inputs

[ 2.15 PM 30.12.18 ]

\*\* Normalizing inputs speeds up our training process. [ put fractions in similar scales. ]

Normalization can be done in two ways.

i) subtract mean:

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

$$x' = x - \bar{x}$$

ii) Normalize variance.

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

+ elementwise

$$x' = \sigma^{-1}$$

\*\* Use same  $\bar{x}$  and  $\sigma^2$  to normalize train and test set.

[ bigger feature can give range of values different than others ]

#### L2: Vanishing or Exploding Gradients

[ 7.00 PM 05-01-19 ]

when training a very deep network your derivatives or slopes can sometimes get either very big or very small this is called (vanishing or exploding gradients.)

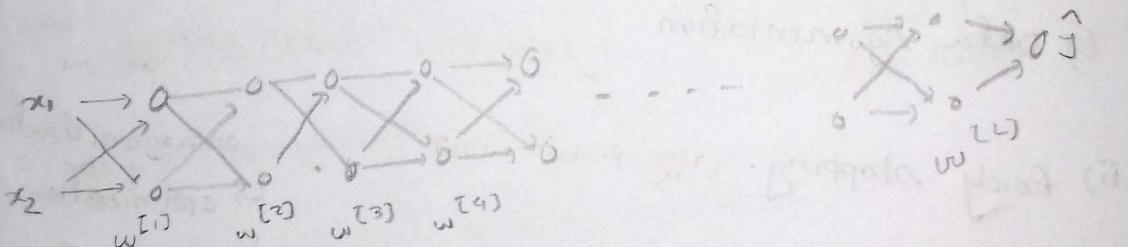


Fig: Deep neural network

$$y = w^{[L]} \cdot w^{[L-1]} \cdot w^{[L-2]} \cdots w^{[2]} \cdot w^{[1]} \cdot x$$

$$\text{Let } w^{[L]} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$z^{[1]} = w^{[1]} x$$

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

$$a^{[2]} = g(z^{[2]}) = g(w^{[2]}, a^{[1]})$$

$$\therefore \hat{y} = w^{[L]} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} x$$

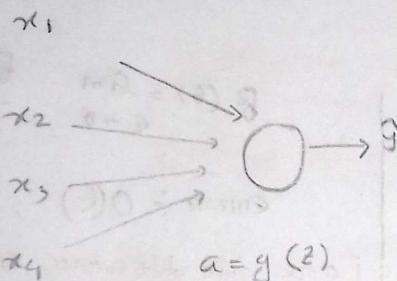
\*\* In this case value of  $\hat{y}$  will be very large and gradients will explode.

\*\* If  $w^{[l]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$  then value of  $g$  will be very small or vanish.

\*\* It makes training very difficult.

L3: weight initialization for Deep networks. [7.10 pm 05.01.19]

→ we can partially solve the vanishing/exploding gradients problem by randomly initialize the weights.



Let,

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

#  $n$  = number of features.

Large  $w_i \rightarrow$  Smaller  $w_i$

$$\text{var}(w) = \frac{2}{n} \quad \begin{array}{l} \text{For ReLU it} \\ \text{gives better} \\ \text{result.} \end{array}$$

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

→ If net the weight matrix not to much greater than 1 or not to much smaller than 1. So weight matrix will not explode or vanish too quickly.

\*\* Different initialization depending on variance

Variance is another hyperparameter that can be used to tune our algorithm.

$$1. \text{ReLU} = \sqrt{\frac{2}{n^{[l-1]}}}$$

$$2. \tanh = \sqrt{\frac{1}{n^{[l-1]}}} \quad * \text{ Xavier initialization.}$$

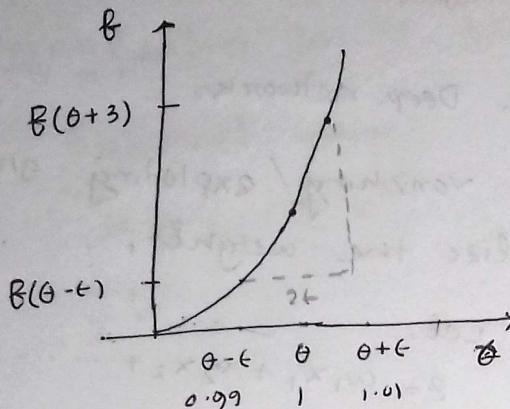
$$3. \sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

~~x.v.y.~~  
\*\* It's from a 2017 paper  
One, 2017 uses this 1/2 paper  
for references in 2017.

\* "He" initialization.

## 4: Numerical approximation of gradients [7:35 PM 05-01-19]

Checking your derivative computation.



$$f(\theta) = \theta^3$$

$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \xrightarrow{\text{width}} g(\theta)$$

$$\frac{(1.01)^3 - (0.99)^3}{2 \cdot (0.01)} = 3.0001 \approx 3.$$

Theory:-

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$$

$$\text{error} = O(\epsilon), 0.0001$$

[Two side differences give more accuracy].

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{\epsilon}$$

$$\text{error} : O(\epsilon), 0.01$$

[One side difference, gives less accuracy].

### \* LS: gradient checking [7:45 PM 05.01.19]

~~vvv~~ [It is useful to check whether our backpropagation implementation is correct or not].

Gradient check for a neural network,

- Take  $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$  and reshape into a big vector  $\theta$ .  $J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$ .

- Take  $d w^{[1]}, d b^{[1]}, \dots, d w^{[L]}, d b^{[L]}$  and reshape into a big vector  $d\theta$ .

\*\*  $d\theta$  is gradient of cost function ( $J$ )

## Gradient checking (Grad check)

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

For each  $i$ :

$$d\theta_{approx}[i] = \frac{f(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - f(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon} \quad \begin{array}{l} \text{Taking two} \\ \text{sided difference.} \end{array}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta_i} \quad \begin{array}{l} \rightarrow \text{for approximately some or 225} \\ \text{error error} \end{array}$$

$$\rightarrow \text{Now check } d\theta_{approx} \approx d\theta. \quad C = 10^{-7}$$

$$\text{Check: } \frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2} \quad \begin{array}{l} \text{If value of check is less than } 10^{-7} \\ \text{then it is good otherwise check.} \end{array}$$

L6: Gradient checking implementation notes. [8:00 PM 05.01.19].

- Don't use grad check in training - only to debug. [takes more time].
- If algorithm fails grad check, look at components to try to identify bug.
- Remember regularization in grad check
- grad check doesn't work with dropout.
- Grad check doesn't work after some training.
- Run at random initialization; perhaps again after some training.

## Programming Assignment

- 1. Setting ~~parameters~~ initialization for constant and ~~zeros~~.
- Neural network zero initialization after each ~~iteration~~.
- [zero initialization, Random initialization, "He" initialization].
- \*\* different initialization leads to different result.
- He initialization works well for networks with ReLU activation
- ~~✓ ✓ ✓~~ [Precision Boundary plot and track 2mn 92% error  
then for 27s]
- 2. Regularization helps to improve algorithm accuracy.
- 3. Gradient checking ensures that backpropagation is actually working.
- [I no longer use Backpropagation to monitor gradients - it's  
one 2mn]
- Gradient checking is only we will not use it in  
every iteration during training.

## Week-2 Optimization algorithms.

### L-1 Mini-batch gradient descent [4.15 pm 08.01.19]

Batch vs. mini-batch gradient descent.

→ faster computation as first we split whole training set into definite number of partitions. These partitions is called mini-batching.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(1000)} & | & x^{(1001)} & \dots & x^{(2000)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix}_{(n_x, m)}$$

$\underbrace{x^{(1)}}_{\text{min batch}} \quad \underbrace{x^{(2)}}_{\text{number}}$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(1000)} & | & y^{(1001)} & \dots & y^{(2000)} & | & \dots & | & \dots & y^{(m)} \end{bmatrix}_{(n_y, m)}$$

$\underbrace{y^{(1)}}_{\text{Y}} \quad \underbrace{y^{(2)}}_{\text{Y}}$

# number of training example

\* number of mini-batches  $\rightarrow$

# number of example in a mini-batch

mini-batch gradient descent. (implementation)

For  $t = 1, \dots, 5000$  <sup>number of mini-batches</sup>  $\underbrace{\text{for } t}_{t \in \{1, \dots, 5000\}}$ .

{ forward of on  $X$

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

$$A^{(1)} = g^{(1)}(z^{(1)})$$

$$A^{(2)} = g^{(2)}(z^{(2)})$$

we can vectorize it

(1000 example).

1 step of gradient descent using  $x^{(t)}, Y^{(t)}$ .

$$\text{Compute cost } J^{(t)} = \frac{1}{1000} \sum_{i=1}^l \ell(g^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_e \|w^{(e)}\|_F^2$$

Backprop to compute gradients of  $J^{(t)}$  (using  $x^{(t)}, Y^{(t)}$ )

$$w^{(e)} = w^{(e)} - \alpha \nabla w^{(e)}, b^{(e)} = b^{(e)} - \alpha \nabla b^{(e)}$$

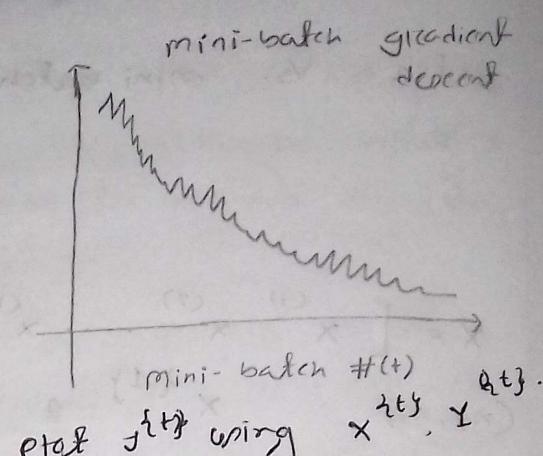
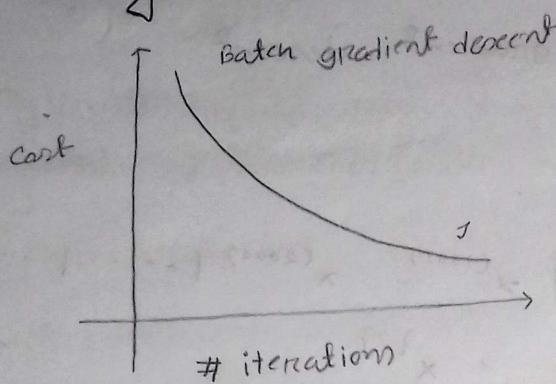
1-epoch -

mini-batch

\* for large training set runs, much faster than batch gradient descent.

## L-2 understanding mini-batch gradient descent [4.35 PM 08.01.19]

Training with mini-batch gradient descent.



choosing your mini-batch size. (one of hyperparameters)

if mini-batch size =  $m$  : Batch gradient descent  $(x^{(1)}, y^{(1)}) = (x, y)$

if mini-batch size = 1 : Stochastic gradient descent | Every example is its own ( $x^{(1)}, y^{(1)}$ ) mini-batch.

$$(x^{(1)}, y^{(1)}) = (x^{(1)}, y^{(1)})$$

→ In practice we choose mini-batch size between (1 to  $m$ )

Stochastic gradient descent  
(mini-batch size = 1)

↓  
large speed up  
from vectorization

(good)  
mini-batch size  
not too big not too small

- fastest learning
- vectorization
- make program
- without probm entire net.

Batch gradient descent  
(mini-batch size =  $m$ )

↓  
Too long for  
each iteration.  
(take long time)

1. If small training net : use batch gradient descent ( $m \leq 2000$ )

2. Typical minitbatch size :  $\underbrace{64, 128, 256, 512}_{\text{or larger}}$

3. Make sure that minibatch  $(x^{t+1}, y^{t+1})$  fits in your CPU/GPU memory.

### L3: Exponentially weighted averages. [4:57 pm 08.01.19].

Temperatur in London.

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

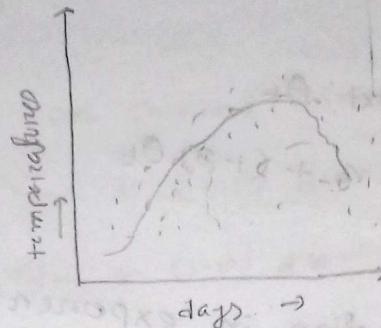
$$\theta_2 = 40^\circ F$$

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

$$\vdots$$

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

$$\vdots$$



$$v_0 = 0$$

$$v_1 = 0.9 v_0 + 0.1 \theta_1$$

$$v_2 = 0.9 v_1 + 0.1 \theta_2$$

$$\vdots$$

$$v_t = 0.9 v_{t-1} + 0.1 \theta_t$$

Exponentially weighted averages. formula for computing exponentially weighted moving averages)

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

↳ hyperparameter.

[ $v_t$  as approximately average over  $\approx \frac{1}{1-\beta}$  deg temperature].

$$\beta = 0.9 \quad : 10 \text{ deg temp}$$

$$\beta = 0.98 \quad : 50 \text{ days temp}$$

$$\beta = 0.5 \quad : 2 \text{ day temp}$$

Higher values of  $\beta$  in averages over larger values of days.]

↳ exponentially weighted averages. [5:10 pm 08.01.19]

L4: understanding exponentially

exponentially

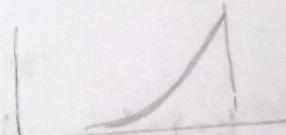
(from lecture 2022).

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

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

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

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



→ Hence we have a weight decay function.

→ Before weight decays to zero it has effect on our weight average equation.

→ we can find how many days we are averaging by using  $(\frac{1}{1-\beta})$  equation

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

↓  
current  
value.
↓  
previous  
value

Implementation.

$$v_0 = 0$$

repeat {

get next  $\theta_t$

$$v_t := \beta v_0 + (1-\beta) \theta_t$$

}

[gt takes really small amount of memory]

L-5 Bias Correction in exponential wei. Av [5.25 pm 08.01.2019]

(graph from 412T on 27/12)  
slide 2.

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

$$v_0 = 0$$

$$v_1 = 0.08 v_0 + 0.02 \theta_1$$

- if we take  $v_0=0$  then  $v_1 = 0.02 \theta_1$ . So we don't correctly estimate  $v_1$ .

→ we can correct this bias by using a simple equation,

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

$$t=2 : \quad 1-\beta^t = 1 - (0.08)^2 = 0.0306$$

$$\frac{v_2}{0.0306} = \frac{0.0106 \theta_1 + 0.02 \theta_2}{0.0306}$$

→ when value of  $t$  is high we can remove the  $\beta^t$  Hence we correctly remove bias.

[gt is used to build some other better optimization algorithm].

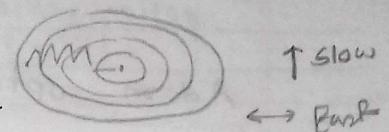
## L-6 Gradient descent with momentum [1:12 AM 09.01.19]

[Faster than traditional gradient descent] (slide ↴ graph incl. fig. adapt. note)

Momentum :-

on iteration t :

Compute  $d\omega, db$  on current mini-batch.



$$v_{dw} = \beta v_{dw} + (1-\beta) dw.$$

$$v_{db} = \beta v_{db} + (1-\beta) db.$$

$$\omega = \omega - \alpha v_{dw}, \quad b = b - \alpha v_{db}.$$

\* It smooths the steps of gradient descent. (Faster).

Here: Hyperparameter :  $\alpha, \beta$  - 2nd year parameter change  
2nd year student fastest output wrt.

## L-7 RMS prop. [1:25 AM 09.01.19]

[Root mean square prop speeds up gradient descent].

on iteration t :

Compute  $d\omega, db$  on current mini-batch.

$$s_{dw} = \beta_2 s_{dw} + (1-\beta_2) dw^2 \rightarrow \text{small}$$

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

$$\omega := \omega - \alpha \frac{dw}{\sqrt{s_{dw} + \epsilon}}$$

updates in horizontal direction divided by a small number  $\epsilon$   
it helps to converge faster.

$$b := b - \alpha \cdot \frac{db}{\sqrt{s_{db} + \epsilon}}$$

→ updates in vertical direction divided by a large number. So it slows down oscillations.

Ensures that it never becomes divided by zero.

L-8 Adam optimization algorithm [10:40 AM 09-01-19]

[Combination of momentum and RMS prop and a generalized optimization algorithm.]

Adam optimization algorithm.

Compute iteration t:

$$v_{dw} = 0, s_{dw} = 0, v_{db} = 0, s_{db} = 0$$

on iteration t:

compute  $dw, db$ , using current mini-batch

$$v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dw \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{momentum part}$$

$$v_{db} = \beta_1 v_{db} + (1 - \beta_1) db \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{"}\beta_1\text{"}$$

$$s_{dw} = \beta_2 s_{dw} + (1 - \beta_2) dw^2 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{RMS prop part}$$

$$s_{db} = \beta_2 s_{db} + (1 - \beta_2) db^2 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{"}\beta_2\text{"}$$

$$v_{dw}^{\text{corrected}} = v_{dw} / (1 - \beta_1^t), \quad v_{db}^{\text{corrected}} = v_{db} / (1 - \beta_1^t) \quad \left. \begin{array}{l} \text{corrected} \\ \text{bias} \end{array} \right\} \text{correction}$$

$$s_{dw}^{\text{corrected}} = s_{dw} / (1 - \beta_2^t), \quad s_{db}^{\text{corrected}} = s_{db} / (1 - \beta_2^t) \quad \left. \begin{array}{l} \text{corrected} \\ \text{part} \end{array} \right\}$$

$$w = w - \alpha \cdot \frac{v_{dw}^{\text{corrected}}}{\sqrt{s_{dw}^{\text{corr.}} + \epsilon}}, \quad b = b - \alpha \cdot \frac{v_{db}^{\text{corrected}}}{\sqrt{s_{db}^{\text{corr.}} + \epsilon}} \quad \left. \begin{array}{l} \text{corrected} \\ \text{update} \\ \text{of} \\ \text{gradient} \\ \text{descent} \end{array} \right\}$$

Hyperparameter choice :-

$\alpha$  : needs to be tuned.

$\beta_1$  : 0.9 (common choice) ( $dw$ ) [momentum]

$\beta_2$  : 0.999 (Adam's paper) ( $dw^2$ ) [RMSprop].

$\epsilon$  :  $10^{-8}$  (Adam's paper).

Adam : Adaptive moment estimation

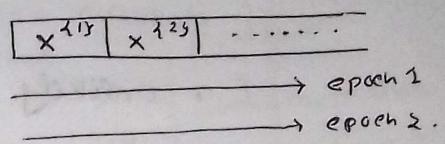
## L-9 Learning rate decay [1:00 PM 09.01.19]

[slowly reduce learning rate  $\alpha$ ]

Learning rate decay.

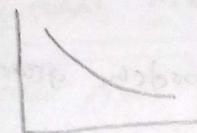
1 epoch = 1 pass through data

$$\alpha = \frac{1}{1 + \text{decay-rate} * \text{epoch-num}}$$



learning rate  
Hyperparameter we want to tune.

Epoch	$\alpha$
1	0.1
2	0.62
3	0.5
:	

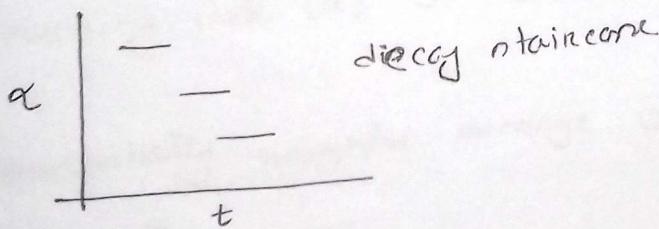


other learning rate decay methods.

$$\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \rightarrow \text{exponentially decay}$$

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

minibatch number



manual decay.

[Learning rate decay] helps to train neural network [Ranbir].

## L-10 The problem of local optima [ 1.10 PM 09-01-19 ].

[Saddle point]

[ plateaus are the points where derivatives are zero for a long time]

[ → unlikely to get stuck in bad local optima.

- plateaus can make learning slow

Programming Assignment :-

1. with a well-tuned mini-batch size, usually outperforms either gradient descent or stochastic gradient descent. (when training set is high).

→ 2. learn to implement mini-batch gradient descent.

3. Adam is the most effective optimization algorithm for training a neural network. It ~~descent~~ combines idea of RMSprop and momentum.

[ Adam outperforms mini-batch gradient descent and momentum].

[ Adam algorithm is best]

## Week 3: Hyperparameter tuning, Batch normalization and Programming frameworks.

### P1: Hyperparameter tuning.

#### L1: Tuning Process [2.20 pm 10.01.19]

V.V.G.

Hyperparameters for a neural network.

- 1 Learning-rate ( $\alpha$ ) → most important
- 2 Momentum ( $\beta$ )
- Adam optimization parameters  $(\beta_1, \beta_2, \epsilon)$
- 3 number of layers.
- 2 Hidden units.
- 3 Learning rate decay.
- 2 mini-batch size.

Theory type  
definition  
implementation  
↳ TensorFlow

→ In hyperparameter tuning of deep learning, we use random values.

→ Using an appropriate scale to pick hyperparameters [2.30 pm 10.01.19]

L2: Using an appropriate scale to pick hyperparameters

1. number of hidden unit & number of layers ↗ for linear search after 2nd.
2. learning rate ( $\alpha$ ) ↗ for random search
3. Exponentially weighted average ↗ for random search after 2nd.

mn 215

### L3: Hyperparameters tuning in practice: [2.42 pm 10.01.19]

[Computational models don't depend on my theoretical work]

→ Bandwidth one model.

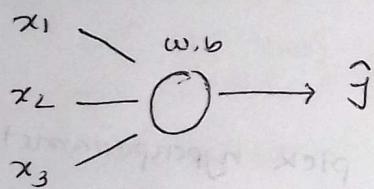
→ Training many models in parallel.

### P2: Batch Normalization.

#### L1: Normalizing activations in a network [2.50 pm 10.01.19]

**Note:** [Batch normalization makes hyperparameter search problem much easier, and makes neural network much robust].

Normalizing inputs to speed up learning.

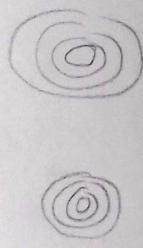


$$\bar{x} = \frac{1}{m} \sum_i x^{(i)}$$

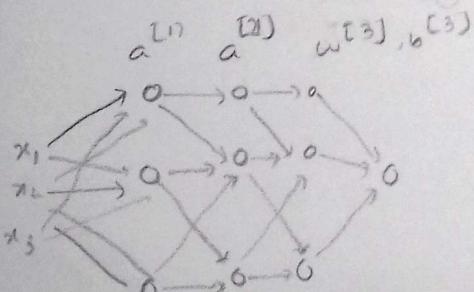
$$x = x - \bar{x}$$

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

$$x = x / \sigma^2$$



Here the question is for any hidden layer can we normalize the activations/inputs to speed up computation of next layer.



→ Here is a debate whether we normalize activation  $a^{[l]}$  or  $z^{[l]}$

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

→ normally we normalize  $z^{[l]}$  as a default.

## Implementing Batch Norm.

Given some intermediate values of a neural network  $[z^{(1)} \dots z^{(m)}]$ .

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

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

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

used for numerical stability  
if  $\sigma^2$  becomes zero.

$z \tilde{}$

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

gamma

Learnable parameters of a model. It allows us to set the mean of  $\tilde{z}$  ( $z \tilde{}$ ) to be whatever you want.

→ use  $\tilde{z}^{[l](i)}$  instead of  $z^{[l](i)}$  for the computation of later layers in the neural network.

~~very~~  
~~\*\*~~ Batch norm not only normalize the input layer but also the values of hidden layer in an NN.

→ By varying  $\gamma$  (gamma) and  $\beta$  (beta) we can control mean and variance. For example.

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

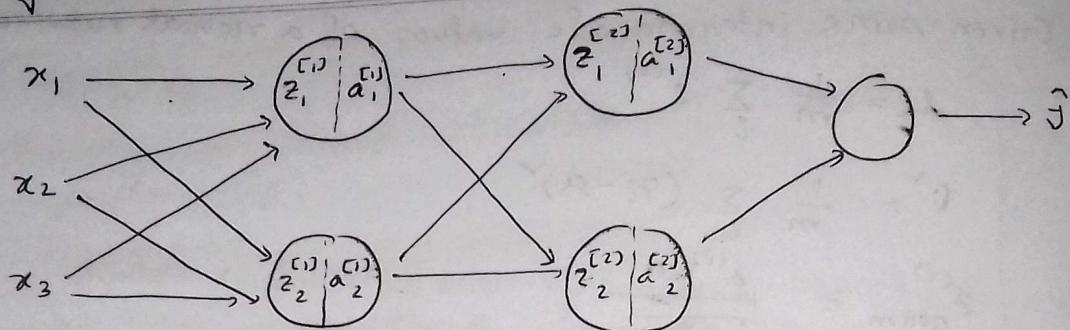
$$\text{and } \beta = \bar{m}$$

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

# Now our hidden units have standardize mean and variance controlled by two explicit parameters (gamma and beta).

## L2: Fitting Batch Norm into a NN [3:55 pm 10.01.19]

Adding Batch Norm to a network.



Steps of NN.

$$x \xrightarrow[w^{[1]}, b^{[1]}] z \xrightarrow[\text{BN}]{\hat{z}, \beta^{[1]}, \gamma^{[1]}} \tilde{z} \rightarrow a = g(\tilde{z}) \xrightarrow[w^{[2]}, b^{[2]}] z \xrightarrow[\text{BN}]{\hat{z}, \beta^{[2]}, \gamma^{[2]}} \tilde{z} = a.$$

Parameters:

$$\left. \begin{matrix} w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]} \\ \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]} \end{matrix} \right] \quad \left. \begin{matrix} d\beta^{[L]} \\ \beta^{[L]} = \beta^{[L]} - \alpha d\beta^{[L]} \end{matrix} \right)$$

→ we can implement batch Norm by only one line of code

working with mini-batches.

$$x^{(1)} \xrightarrow[w^{[1]}, b^{[1]}] z^{(1)} \xrightarrow{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{(1)} \rightarrow g(\tilde{z}^{(1)}) = a^{(1)} \quad \dots$$

$$x^{(2)} \xrightarrow[w^{[1]}, b^{[1]}] z^{(2)} \xrightarrow{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{(2)} \xrightarrow[w^{[2]}, b^{[2]}] z^{(3)} \xrightarrow{\beta^{[2]}, \gamma^{[2]}} \tilde{z}^{(3)} \xrightarrow[w^{[3]}, b^{[3]}] z^{(4)} \xrightarrow{\beta^{[3]}, \gamma^{[3]}} \tilde{z}^{(4)} \xrightarrow[w^{[4]}, b^{[4]}] z^{(5)} \xrightarrow{\beta^{[4]}, \gamma^{[4]}} \tilde{z}^{(5)} \xrightarrow[w^{[5]}, b^{[5]}] z^{(6)} \xrightarrow{\beta^{[5]}, \gamma^{[5]}} \tilde{z}^{(6)} \xrightarrow[w^{[6]}, b^{[6]}] z^{(7)} \xrightarrow{\beta^{[6]}, \gamma^{[6]}} \tilde{z}^{(7)} \xrightarrow[w^{[7]}, b^{[7]}] z^{(8)} \xrightarrow{\beta^{[7]}, \gamma^{[7]}} \tilde{z}^{(8)} \xrightarrow[w^{[8]}, b^{[8]}] z^{(9)} \xrightarrow{\beta^{[8]}, \gamma^{[8]}} \tilde{z}^{(9)} \xrightarrow[w^{[9]}, b^{[9]}] z^{(10)} \xrightarrow{\beta^{[9]}, \gamma^{[9]}} \tilde{z}^{(10)}$$

On batch norm parameters

subtraction step. So,

Parameters:  $w^{[l]}, \cancel{\beta^{[l]}}, \beta^{[l]}, \gamma^{[l]}$   
 $(n^{[l]}, 1) \quad (n^{[l]}, 1)$

$\beta^{[l]}$  is cancelled out in mean

$$z^{[l]} = w^{[l]} a^{[l-1]} + \cancel{\beta^{[l]}}$$

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

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

$$\tilde{z}^{[l]} = \gamma^{[l]} z_{\text{norm}}^{[l]} + \mu^{[l]}$$

Implementing gradient descent using Batch-Norm.

for  $t = 1 \dots$  number of mini-batches.

→ Compute forward prop on  $X^{(t)}$ .

→ In each hidden layer, use BN to replace  $Z^{(l)}$  with  $\hat{Z}^{(l)}$ .

→ Use backprop to compute  $d\omega^{(l)}, d\beta^{(l)}, d\gamma^{(l)}, d\delta^{(l)}$ .

→ Update parameters,  $w^{(l)} = w^{(l)} - \alpha d\omega^{(l)}$ ,  $\beta^{(l)} = \beta^{(l)} - \alpha d\beta^{(l)}$ ,  $\gamma^{(l)} = \gamma^{(l)} - \alpha d\gamma^{(l)}$

[we can also use momentum, RMSprop, Adam optimization]

→ Onward project to understand implementation

L-3 why does Batch Norm work? [4:20 PM, 10.01.19]

[Fast forward & far away]

1. normalizing features to mean zero and variance one  
speed up learning. Because features are ranging from zero to one rather than ranging from different range.

2. Batch-norm reduces the problem of input value-changing  
3. It allows each layer learn independently a little bit

then other layers.

4. It makes the job of learning in the later layers much easier.

Batch Norm as regularization.

1. Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
2. 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.
3. This has a slight regularization effect.

L-4 Batch norm at test time [ 8.20 pm 10.01.19 ]

Equation of Batch Norm at training time.

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

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

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

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

During testing we may not have enough example to form a mini-batch or may have only one example.

In that case:

$\bar{\mu}, \sigma^2$ : estimate using exponentially weighted average

(This keep track of latest weighted average of  $\bar{\mu}, \sigma^2$ )

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

$$\begin{matrix} \bar{\mu}^{[l]} \\ \sigma^2_{\text{norm}}^{[l]} \end{matrix}, \begin{matrix} \bar{\mu}^{[l]} \\ \sigma^2_{\text{norm}}^{[l]} \end{matrix}, \begin{matrix} \bar{\mu}^{[l]} \\ \sigma^2_{\text{norm}}^{[l]} \end{matrix}, \dots \rightarrow \left. \begin{matrix} \bar{\mu} \\ \sigma^2 \end{matrix} \right\}$$

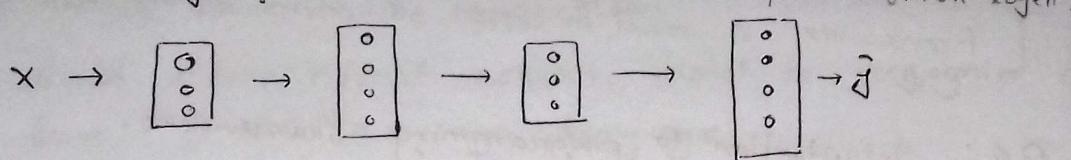
$$\left[ z_{\text{norm}} = \frac{z - \bar{\mu}}{\sqrt{\sigma^2 + \epsilon}}, \tilde{z} = \gamma z_{\text{norm}} + \beta \right].$$

Now apply exponentially weighted average and use it in test time.

### P3 : Multiclass classification

L1 :- Softmax Regression. [10.45 am 10.01.19]

Softmax Layer:-



→ Final output layer of a multiclass classifier is called softmax layer.

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

Activation function for softmax layer:

$$t_i = e^{(z_i^{[L]})} \rightarrow \text{ensures that probability lies between 0 to 1.}$$

$$a^{[L]} = \frac{e^{(z_i^{[L]})}}{\sum_{j=1}^{n^{[L]}} t_j} \text{ or } \frac{t_i}{\sum_{i=1}^{n^{[L]}} t_i}$$

[softmax for example has output [0.2 0.2 0.2]]

L2: Training a Softmax classifier. [11.00 am 10.01.09]

["soft-max" comes from "hard-max"]

Softmax regression generalizes logistic regression to C classes.

→ If C=2 softmax reduces to logistic regression.

Loss Function:  $y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$   $\hat{a}^{[L]} = \hat{j} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$

$$\ell(\hat{j}, j) = - \sum_{j=1}^C j_j \log \hat{j}_j \rightarrow \text{loss on one training example}$$

$$J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{j}^{(i)}, j^{(i)}) \rightarrow \text{loss in entire training set}$$

Input:- (C class, m examples)

$$Y = \begin{bmatrix} 0 & 1 & 0 & \dots & m \\ 0 & 0 & 1 & \dots & m \\ 0 & 0 & 0 & \dots & m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 0 & \dots & m \end{bmatrix}$$

$$\hat{Y} = \begin{bmatrix} 0.33 & \dots & m \\ 0.22 & \dots & m \\ 0.11 & \dots & m \\ 0.29 & \dots & m \end{bmatrix}$$

Gradient descent with softmax. [ slide 9 min 27 sec ]

[ framework: it built in code so often easier to run ]

P4: Introduction to programming frameworks.

L2: Deep learning frameworks [ 11:15 AM 10.01.19 ]

Choosing deep learning frameworks:

- Ease of programming (development and deployment)
- Running speed.
- Truly open (open source with good governance).

[ some frameworks ]

L2: Tensorflow [ 11:30 am 10.01.19 ]

[ Tensorflow knows automatically how to implement backprop ]

\*\* we have to define a cost function and tensorflow automatically finds the derivatives and a way to minimize cost function.

Running and writing program in tensorflow has following steps.

- i) Create tensors (variables) that are not yet executed.
- ii) write operations between those tensors.
- iii) Initialize your tensors.
- iv) Create a session.
- v) Run the session.

## Programming Assignment.

- Learn basic about tensorflow.
- Build a deep neural network model to recognize numbers from 0 to 5 in sign language.

will learn

- i) initialize variables.
- ii) start your own session.
- iii) Train algorithm.
- iv) Implement a neural network.

\* Remember to initialize your variables, create a session and run the ~~session~~ operations inside the session.

\* a placeholder is an object where value can be specified later

- one hot vectors. (One hot encoding)

- backpropagation automatically *homework due on 27-10-17*