

Deep
improving
Neural Networks
Neural Network Tuning
Hyper parameter tuning
Hyper optimization
help in training
Optimizing

Train set

"Dev set" "Validation Set" "Test set".

5 Train set: Keep on training on train set use "Dev set" to tune parameters to find which combination work "best".

10 Test set: Evaluate your best model on test set & report that accuracy.

Train / Test

70 / 30

(Train) > (Dev) > (Test)

Train do do

15 These ratios are pretty hold valid if your data set holds $\approx 10K$ samples.

In modern big data if you have $\approx 10,000,000$ then

so in case $\approx 10K$ dev fest is valid enough

$\approx 10K$ test data set is valid enough

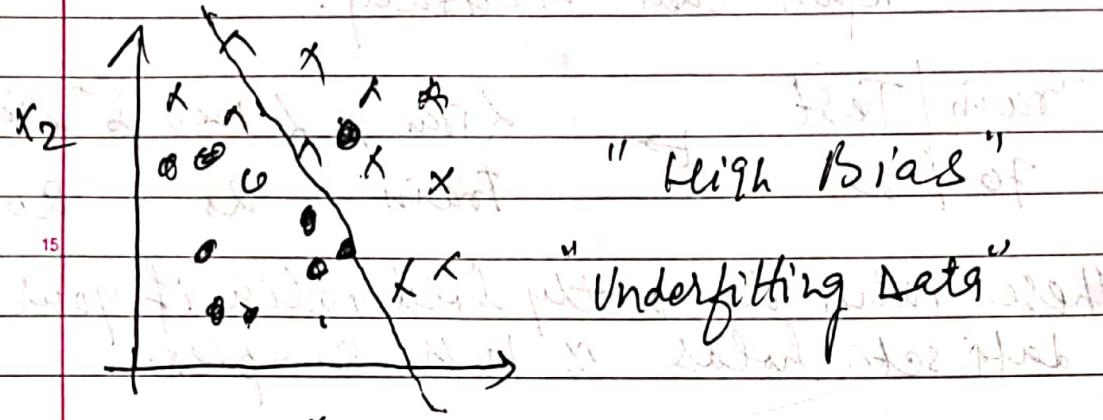
25 So in φ case Ratio =

| | | |
|-------|-----|------|
| 98% | 1% | 1% |
| Train | dev | test |

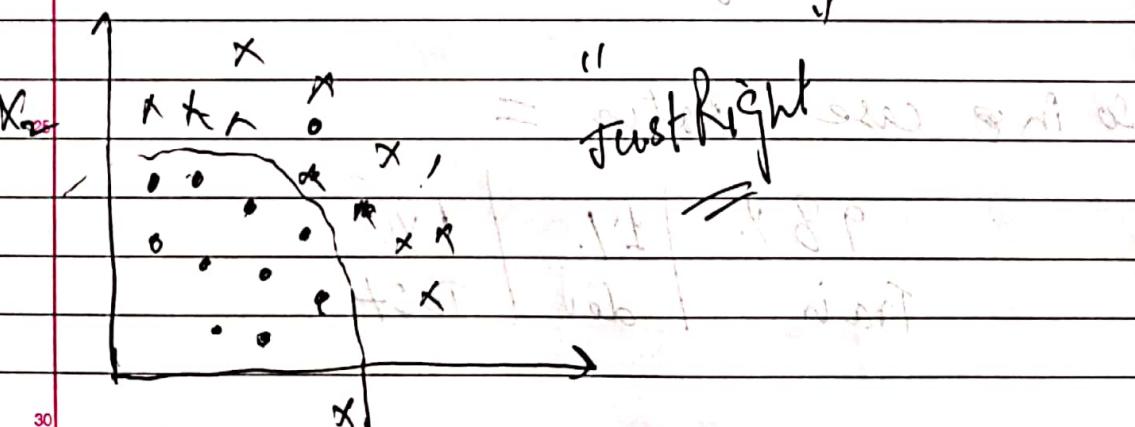
Mismatches Train/Test distribution

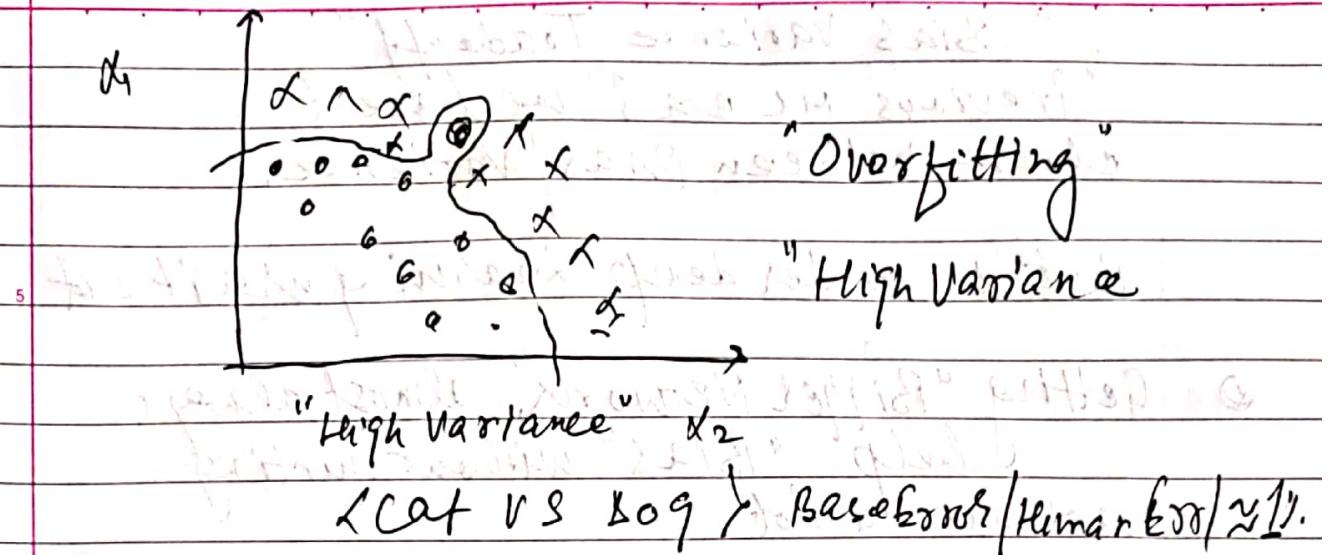
- ① Make sure "dev/test" & "test" come from same distribution.
- ② Not having test set is ok, We can report "dev" set accuracy.

of Bias & Variance



among all the fitted way fit stick pick is shown at





| | | | | | |
|----|-----------------|-------------|-------------|---------------------------|--------------------------|
| 10 | Train Set Error | 1% | 15% | 15% | 0.5% |
| 15 | Dev Set Error | 1% | 16% | 80% | 1% |
| | "High Variance" | "High Bias" | "High Bias" | High Bias Low Variance | Low Bias Low Variance |

The above analysis is based on Best Optimal error.

20 High bias? \rightarrow Bigger Network, Train longer
 ↓ " (NN Architecture search)"

25 "Training Bigger Network almost always help" to get rid of "bias" problem.

↓ High Variance? \rightarrow More data, regularization

30 "DONG"

" Bias Variance Trade-off "

"Previous MC era P we have to choose between bias/variance.

But in modern deep learning era it's not

- ④ Getting "Bigger Network" almost always help "Bias" without hurting variance a lot. (at 23 min)

2013-08-21 | 08:15 | 0.1 km south of the coast

100% of the time, the system will be able to correctly identify the target word.

2 hours before you travel

ANSWER 15

Consequently, no valid responses were available.

20

Constitutive Regulation (RSP16) \rightarrow 2nd phase

21 [View Details](#)

~~Aspetto~~ Salvo tutti dalla Asmashinjasi scivolti

25

• numbers "2nd" is big

~~Ames Beach, Berkeley~~, X 25 mm. *Mytilus californianus*

30

卷之三

Regularization

logistic regression: - $w \in \mathbb{R}^n$ $b \in \mathbb{R}$

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

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w = L_2 \text{ regularization}$$

$$L_1 \text{ regularization} \quad \frac{1}{2m} \sum_{i=1}^{n_x} |w_i| = \frac{\lambda}{2m} \|w\|_1$$

$$\lambda = \underline{\text{lambd}}$$

Neural Network

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

$$+ \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|^2$$

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

"Frobenius Norm" of a matrix, $\| \cdot \|$ L_2 regularization

$$dw^{[l]} = (\text{from backprop}) + \frac{1}{m} w^{[l]} \quad = \quad \text{Weight Decay}$$

$$w^{[l]} = w^{[l]} - dw^{[l]}$$

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

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

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

It's like taking the making the $w^{[l]}$ matrix a little smaller.

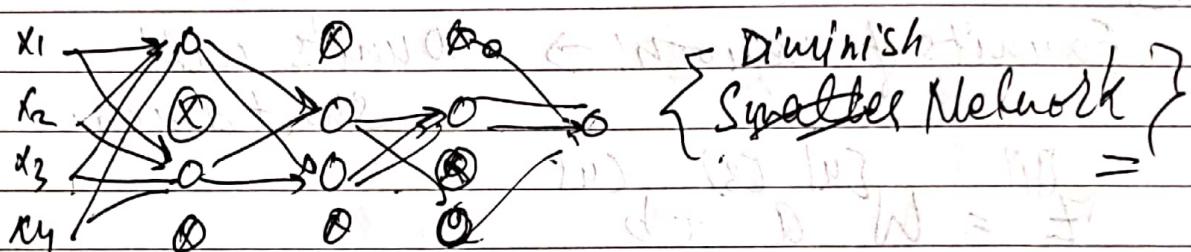
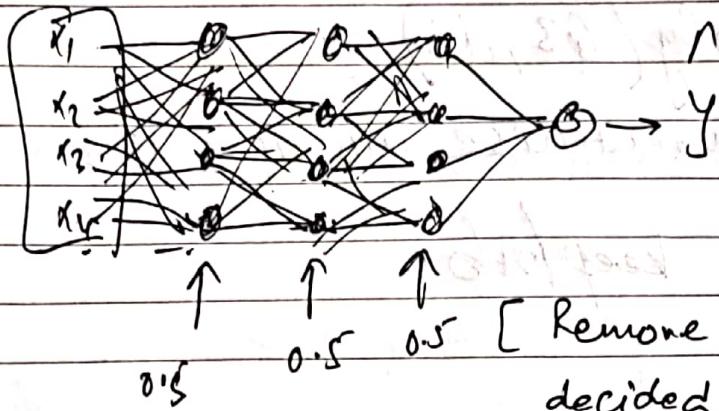
$$\begin{aligned} & \text{(Forward pass)} \\ & \text{(Backward pass)} \\ & \text{...} \end{aligned}$$

the bigger the α the more "aggressive" it is.

$w^{[l+1]} = w^{[l]} - \alpha (backward\ pass) \cdot \frac{1}{m} \sum_{i=1}^m w^{[l+1]}$

$$w^{[l+1]} = w^{[l]} - \alpha \frac{1}{m} \sum_{i=1}^m w^{[l+1]}$$

(Dropout Regularization)



We do this for each Training Example.

Implementing Dropout ("Inverted Dropout")

Illustrate with layer $l=3$;

$d3 = np.random.rand(a3.shape[0], a3.shape[1])$

$random_arr = np.random.rand(3, 1)$

$d3 = random_arr \leq keepProb | keepProb = 0.6$

$random_arr = [[0.476060], [0.617575], [0.449407]]$

$[False]$

$[True]$

$[True]$

$]$

~~$a_3 = a_3$~~

$$a_3 = \text{np.multiply}(q_3, d_3)$$

element-wise multiplication

$$a_3' = a_3 / 0.8 \Rightarrow \text{keep frb.}$$

for example

50 units of neurons \rightarrow 10 units will shut down.

$$z^{(l)} = w^{(l)} a + b$$

↑ reduced by 20%.

so we to reduce it $a_3' = a_3 / 0.8$

(PumpIt)

At test time: given X

$$a^{(l)} = X \quad \{\text{No Dropout Test Time}\}$$

$$z^{(l)} = w^{(l)} a^{(l)} + b$$

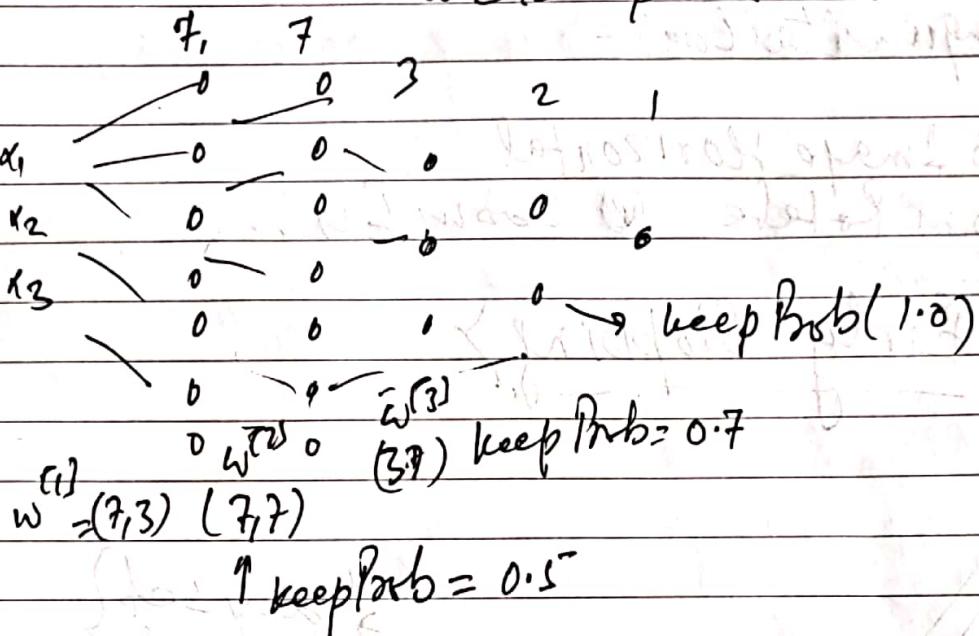
$$a^{(l)} = g^{(l)}(z^{(l)})$$

$$y^{(l)} = w^{(l)} a^{(l)} + b$$

$y =$ fill

Q Why Does Dropout Work?

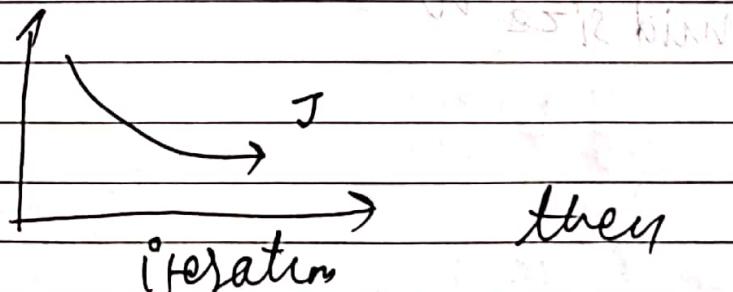
Intuition: Can't rely on any one feature, so have to spread out weights.



In practice we don't apply dropout at input layer hence for that layer keepProb = 1.0 or ~~0.5~~

Downside: Cost function is not well defined.

So in order to check the code validity we FIRST keep [keepProb=1] plot the



they enable dropout.

Other regularization

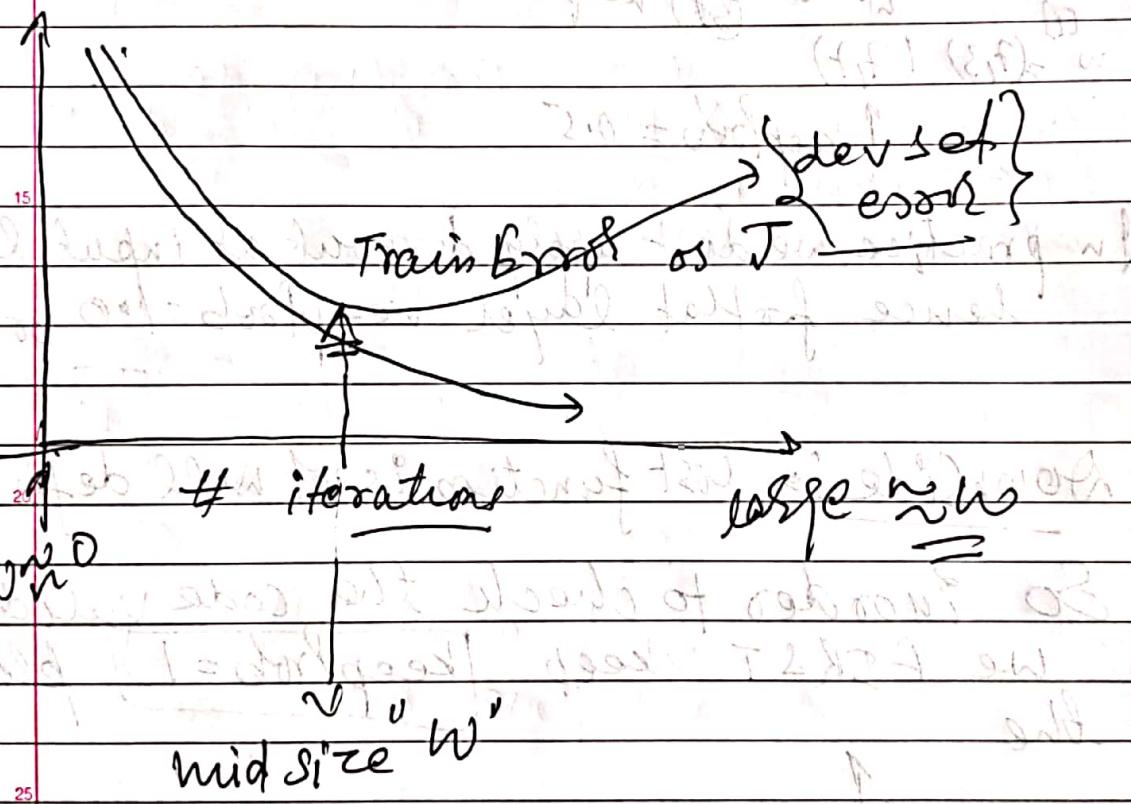
In order to get more training data:

data Augmentation:-

(1) flip Image Horizontal

(2) Random Rotate (3) zoom In :-

or Early Stopping



Q)

Ans width

Orthogonalization :- { Segregation Task }

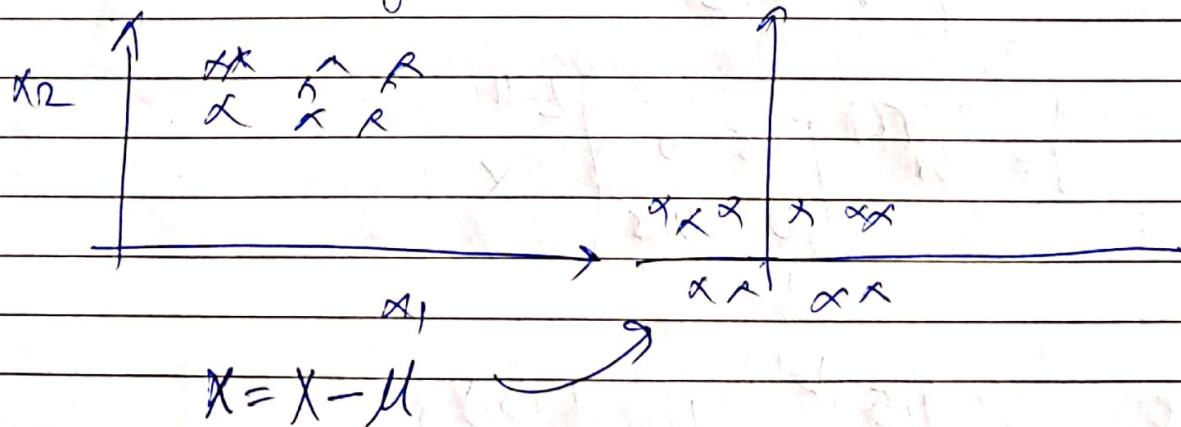
= Optimize Cost function J }
 - Gradient Descent, RM } $J(w, b)$

= Noise filtering }
 - Regularization }

Early Stopping couple these two hence no overfitting.

Normalizing Input }

1. Subtracting mean



Normalize variance

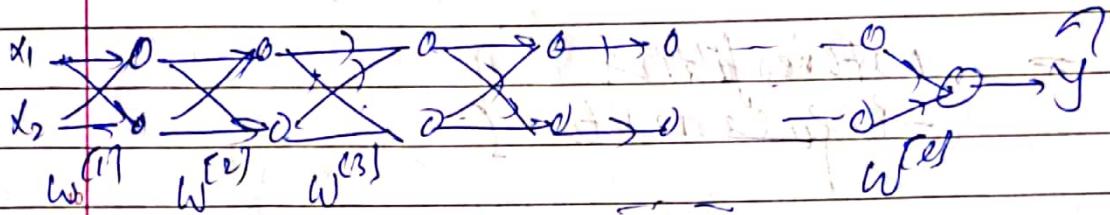
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})^2 \quad | \quad X' = \frac{X}{\sigma}$$

(mean) (variance)

| | |
|--------|------|
| Camlin | Page |
| Date | / / |

use same μ & σ^2 to normalize
the TRAIN & TEST Data

Vanishing & Exploding Gradients



$$g(z) = \tanh(z) \quad \text{activation function}$$

$$y = w^{(L)} w^{(L-1)} \cdots w^{(1)} x - b$$

$a^{(L)} = w^{(L)} x + b \Rightarrow 0$

$a^{(L-1)} = g(a^{(L)}) = z^{(L)}$

$$w^{(L)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

$$y = w^{(L)} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x$$

$$(1.5)^{L-1} \quad \text{if } L \gg \infty$$

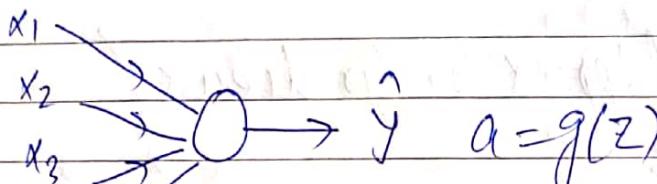
so if very deep network

y will explode conversely

$(1.5)^{L-1}$ is very less

Weight Initialization for Deep Networks

Partial solution for Vanishing & Exploding problem
in Deep Neural Network is careful init of
random weight:



$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n$$

We need z not to be very small nor to be very big. What you need is:-

large $n \rightarrow$ smaller w_i

Set Variance(w_i) = $\frac{1}{n}$

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

for ReLU activation

$$\text{Variance}(w_i) = \frac{2}{n} \quad] \text{ To do this}$$

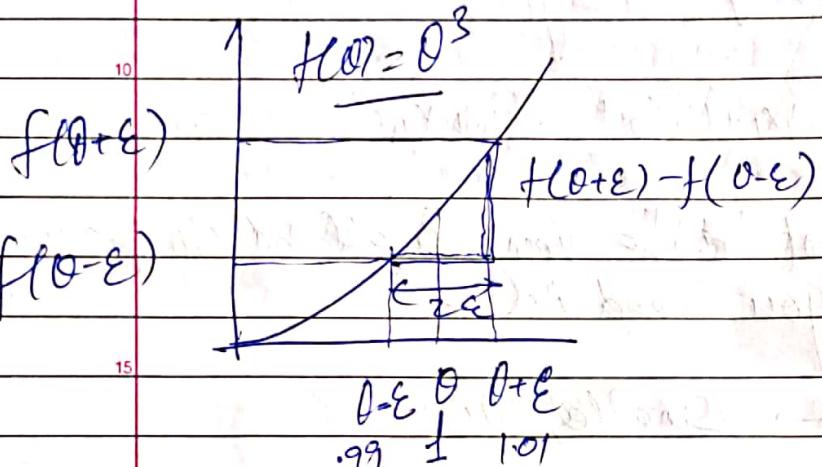
$$w^{[l]} = np.random.randn(shape) * np.sqrt\left(\frac{2}{n^{[l]}}\right)$$

Other variance

tanh activation $\sqrt{\frac{1}{n(\sigma-1)}}$

*These are not the first parameters to tune.

[Numerical Approx for Gradients]



$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta)$$

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

SS

$$g(\theta) = 3\theta^2 = 3 \text{ very close}$$

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

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

error Order
 $O(\epsilon)$

\downarrow
error dimension
 $O(\epsilon^2)$

L Gradient Check for Neural Network

Take $W^{(1)}, b^{(1)}, -W^{(1)}, b^{(1)}$ and reshape into big vector θ

Take $d\theta^{(1)}, db^{(1)} \rightarrow d\theta^{(1)}, db^{(1)}$ reshape big vector $d\theta$

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

for each i :

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

$$\approx \frac{\partial J}{\partial \theta_i} \quad \downarrow \quad d\theta_{approx} \approx d\theta$$

check $\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta\|_2} \approx 10^{-7}$ or smaller
 $\|d\theta_{approx}\|_2 + \|d\theta\|_2$ if $\epsilon = 10^{-7}$ [great]

10^{-5} It's ok you should check

10^{-3} Wrong

Grading Checking Implementation

④ Grad check doesn't work with dropout

⑤ If algorithm fails grad check, look at the components to try to identify bug.

Ex. for some layer values of $db^{(l)}$ are very far off 0. If this is the case then the bug is in how you find out " db "

15

20

25

30

Optimization Algorithm

Batch vs Mini-batch gradient descent

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

$n = 50000$; what if n is too large.

We create mini-batch let say of 1000 each

$$X = \{x^{(1)}, x^{(2)}, \dots, x^{(1000)}, x^{(1001)}, \dots, x^{(2000)}, \dots, \underbrace{\dots}_{\{x^{(5000)}\}}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, \underbrace{y^{(1000)}, y^{(1001)}, \dots, y^{(2000)}}_{\{y^{(815)}, y^{(825)}, \dots, y^{(8005)}\}}]$$

Minibatch t : $X^{\{t\}}, Y^{\{t\}}$

$$X^{\{t\}} = (n_x, 1000) \quad Y^{\{t\}} = \{1, 1000\}$$

for $t=1 \rightarrow 5000$
forward prop $X^{\{t\}}$

$$Z^{(t)} = W^{(t)} X^{\{t\}} + b^{(t)}$$

$$A^{(t)} = g^{(t)}(Z^{(t)})$$

Compute cost: - $\frac{1}{J} \sum Z^{\{t\}}$

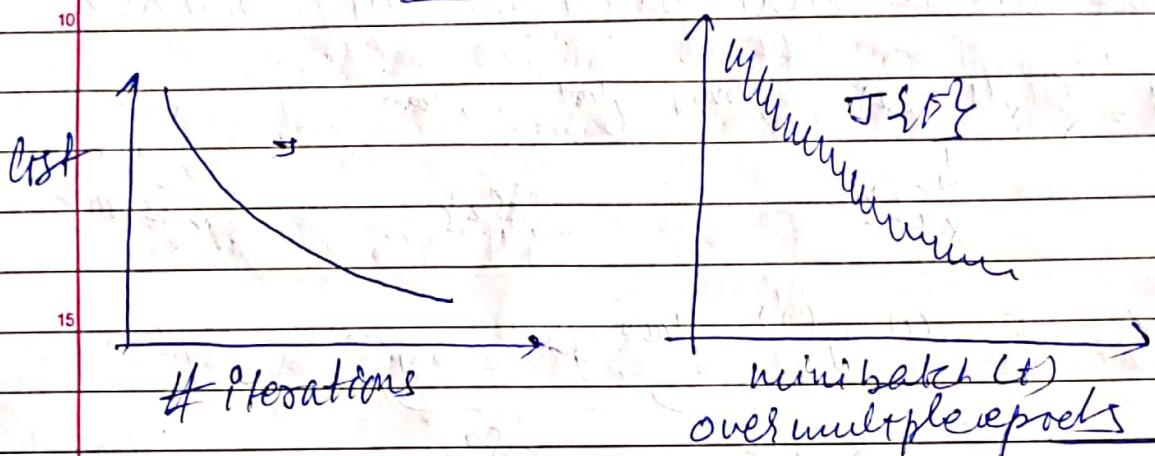
$$J = \frac{1}{1000} \sum_{i=1}^l \hat{L}(Y^{(i)}, \hat{Y}) + \frac{\lambda}{2 \times 1000} \sum_{l} \|W^{(l)}\|_F^2$$

Backprop to compute gradients w.r.t
JST using $x^{[t]}, y^{[t]}$

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

$$b^{[l]} = b^{[l]} - \alpha \frac{\partial J}{\partial b^{[l]}}$$

This one pass to training set this is
called "epoch"



↑
JST

mini-batch (t)
over multiple epochs

- if mini-batch size = m = Gradient descent
if mini-batch size = 1 = Stochastic gradient
descent

if small train set $m \approx n \leq 2000$

else mini batch size is anything
of 64, 128, 256, 512 ...

[Exponentially Weighted Averages]

$$\theta_1 = 40^{\circ}\text{F} \quad 4^{\circ}\text{C}$$

$$\theta_2 = 60^{\circ}\text{F} \quad ;$$

$$\theta_{MA} = 56^{\circ}\text{F} \quad 15^{\circ}\text{C}$$



$$V_0 = 0$$

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

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

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

If plot V_t we will get the RFD curve

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t \quad \beta = 0.9$$

V_t can be considered as average over

$\approx \frac{1}{1-\beta}$ days temperature if $\beta = 0.9$ days

so it is average over $\approx \frac{1}{1-0.9} = 10$ days temp.

if $\beta = 0.5$ ≈ 2 days avg. we get GREEX curve.

Bias Correction in Weighted Avg

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

5 we put $V_0 = 0$

$v_1 = .98 V_0 + .02 v_1$ so even if
first day temp $o_1 = 40^{\circ}\text{F}$ we get

10 $V_1 = .98 \times 0 + .02 \times 40$

= 8°F only hence curve
starts much lower initially.

so bias correction can be -

15 $V_t = \beta V_{t-1} + (1-\beta) v_t$

$$V_t = \frac{v_t}{(1-\beta)}$$

20 $t=2 \rightarrow 1-\beta^2 \Rightarrow 1-(.98)^2 = 0.0396$

so estimate

25 $\underline{\underline{v_2}} = \frac{.0196 o_1 + .02 v_2}{0.0396}$

$(-\beta)$ → in a long term get vanish

30 $(\text{long}) 1 - (.98) \rightarrow \approx 1$

[Gradient Descent With Momentum]

Gradient Descent Example



What we want | vertical axis slower &
horizontal axis faster -

Momentum :-

→ Compute dW, db on current mini-batch

→ compute $V_{dw} = \beta V_{dw} + (1-\beta) dw$ (similar to
analogy to $V_0 = \beta V_{0t} + (1-\beta) \delta t$)

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

$$w = w - \alpha V_{dw} \quad | \quad b = b - \alpha V_{db}$$

by doing this we will find "GREEN" above.

most of the time $\beta = 0.9$ works
very well.

In Adam optimization
this is $\beta_1 = 0.9$

RMS Prop

On iteration t:-

① compute dW, dB

$$sdw = \beta_2 Sdw + (1-\beta_2) dw^2 \quad (\text{elementwise square})$$

$$sdb = \beta_2 Sdb + (1-\beta_2) db^2$$

$$w = w - \alpha \frac{dw}{\sqrt{sdw}}$$

$$b = b - \alpha \frac{db}{\sqrt{sdb}}$$

initially sdw will be very small, sdb very large. impulsive to remove divide by zero we do

$$w = w - \alpha \frac{dw}{\sqrt{sdw} + \epsilon}$$

$$b = b - \alpha \frac{db}{\sqrt{sdb} + \epsilon}$$

every very
small number

$$\epsilon \approx 10^{-8}$$

RMS prop given by Zeffry Hilton
on course bog

ADAM Optimization [Run/ST Momentum]

$$V_{dw} = 0, S_{dw} = 0 \quad | \quad V_{db} = 0 \quad | \quad S_{db} = 0$$

On iteration t :

compute \hat{d}_{dw} , \hat{d}_{db} with mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) d_{dw} \quad \} \text{momentum}$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) d_{db} \quad \} \leftarrow \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) d_{dw}^2 \quad \} \text{RMS}$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) d_{db}^2 \quad \} \varepsilon \beta_2$$

$$\begin{array}{c} \text{bias} \\ \text{correct} \\ \text{Corr} \\ \text{on} \end{array} \left[\begin{array}{c} V_{dw} = \frac{V_{dw}}{1 - \beta_1^t} \\ V_{db} = \frac{V_{db}}{1 - \beta_1^t} \end{array} \right]$$

$$\left[\begin{array}{c} S_{dw}^{\text{correct}} = \frac{S_{dw}}{1 - \beta_2^t} \\ S_{db}^{\text{correct}} = \frac{S_{db}}{1 - \beta_2^t} \end{array} \right]$$

$$w = w - d \left(\frac{V_{dw}^{\text{correct}}}{\sqrt{S_{dw}^{\text{correct}}} + \epsilon} \right)$$

$$b = b - d \left(\frac{V_{db}^{\text{correct}}}{\sqrt{S_{db}^{\text{correct}}} + \epsilon} \right)$$

$$\left. \begin{array}{c} \alpha \text{ needs to tune} \\ \beta_1 = 0.9 \end{array} \right| \left. \begin{array}{c} \beta_2: 0.999 \\ \epsilon \approx 10^{-8} \end{array} \right| \left. \begin{array}{c} \text{common} \\ \text{values} \end{array} \right| =$$

Learning Rate Decay.

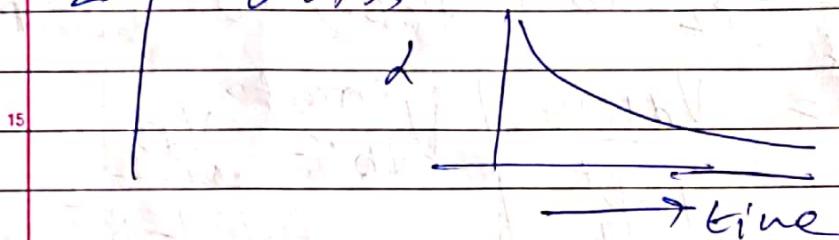
$$\frac{\alpha_{t+1}}{\alpha_t} = \frac{\alpha_0}{\alpha_0 \cdot \text{decay-rate}}$$

→ epoch

epoch = 1 pass to data

$$\alpha = \frac{\alpha_0}{1 + (\text{decay-rate}) * \text{epoch_num}}$$

| | | |
|-------|------------------|--|
| 10 | $\alpha_0 = 0.2$ | $\frac{0.2}{1+1} = 0.1$ |
| epoch | 1 | |
| 1 | 0.1 | $\frac{0.2}{1+2} = \frac{0.2}{3} = 0.0667$ |
| 2 | 0.0667 | |



$$\alpha = (0.95)^{\text{epoch_num}} \times \alpha_0 \quad [\text{exp-decay}]$$

20

25

30

"Hyperparameter Tuning Process"

Hyperparameters are :- α (Learning Rate)

β_1 (momentum) (β_1, β_2) Adam params ϵ

layers # hidden units
learning rate decay | mini batch size.

α [Alpha Learning Rate]

~~$\beta_1 = 0.9$~~

hidden units

mini batch size

Decreasing Order
of Priority to
tune hyperparams
as per

"Andrew Ng"

layers

learning rate decay

$$\beta_1 = .9 \quad \beta_2 = .999 \quad \epsilon = 10^{-8}$$

Adam params

25

30

H2

Grid Search Method

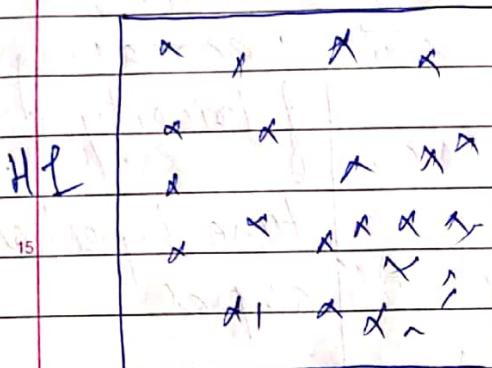
Try the combination

of H1 & H2. [25 Combos]



This old method we should not use
in "Deep Learning"

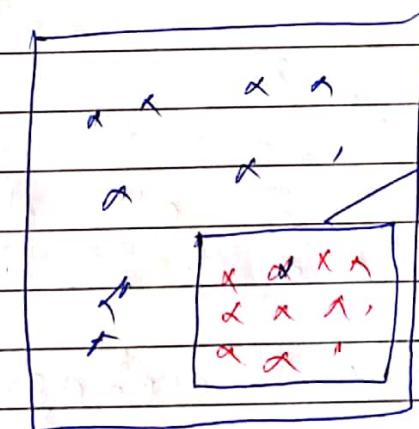
H2



choose the point at random

L Course to fine p

H1



This area works better so zoom in at take more dense samples —

L Appropriate Scale to pick hyperparameters

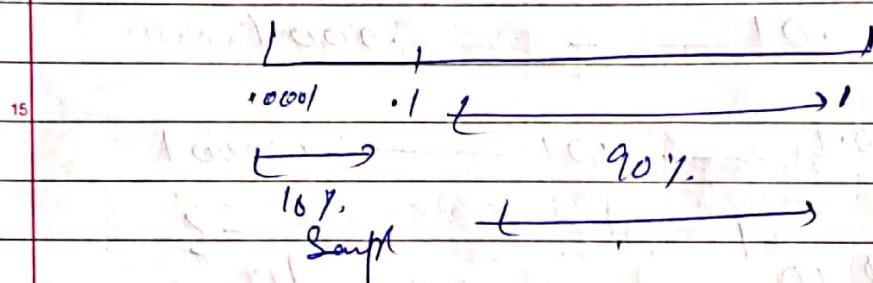
$$h^{(l)} = 50 \text{ ————— } 100$$

$$50 \text{ } \underline{x} \text{ } \underline{\times 0.1 \times 0.01 \times 1/100}$$

layer L: 2-4 (no do the coarse fine)

linear search is pretty ok here but not true
for other.

$$d = .0001 \text{ ————— } 1$$



So taking the samples on linear scale will be biased
hence it is good to take the steps on log "scale".

$$\begin{aligned} & \text{A logarithmic scale from } 10^{-4} \text{ to } 10^0 \text{ is shown.} \\ & \text{The values } 0.0001, 0.001, 0.01, 0.1, 1 \text{ are marked.} \\ & \text{A bracket labeled } 90\% \text{ spans from } 0.0001 \text{ to } 0.01. \\ & \text{A bracket labeled } 10\% \text{ spans from } 0.01 \text{ to } 0.1. \\ & \text{The word "Sample" is written below the } 0.01 \text{ mark.} \\ & \text{The formula } d = 10^r \text{ is given, where } r \in (-4, 0) \\ & \text{and } d \in (10^{-4}, 10^0). \end{aligned}$$

Now what we do sample between $[a, b]$
in our case $[-4, 0]$ if put

$$r \in [a, b]$$

$$d = 10^r$$

↳ Hyperparameter for exponentially weighted moving averages

$$\beta = .9 - .999$$

5

$$\begin{array}{c} 1 \\ \hline .9 & .999 \\ \hline \end{array}$$

→ Doesn't make sense to sample the value on linear scale

10

The way we can think about this is to we have to explore the

$$1-\beta = 1.01 - .0001$$

15

$$0.1 - .01 - .0001$$

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

20

Logscale logscale

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

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

25

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

β randomly chosen hyperparameter value:

30

Mathematical Justification why to sample

on "log" :-

$$\beta : -0.9 \rightarrow .9005 \text{ (This doesn't much impact)}$$

$$\approx \frac{1}{1-0.9} \times \frac{1}{0.1} = 10 \text{ day Avg.}$$

$$\frac{1}{1-0.9005} = \frac{1}{0.0995} \approx 10.05 \text{ days Avg.} \xrightarrow{\text{data}}$$

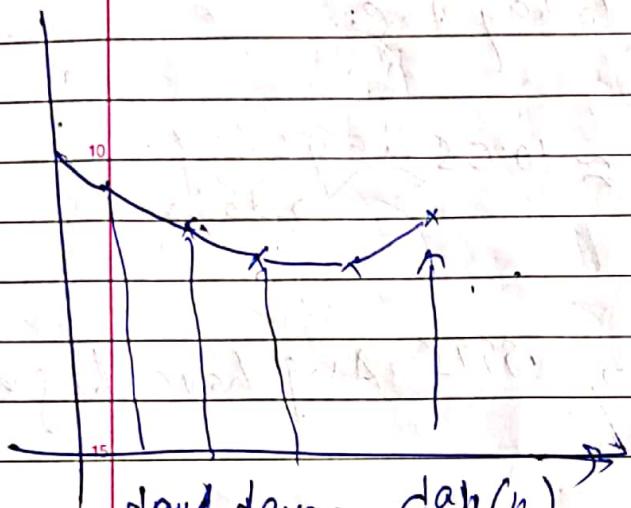
but $\beta = .889 \rightarrow .9995$ Avg have huge impact

$$\frac{1}{1-0.889} = \frac{1}{0.01} \approx 1000 \text{ days Avg} \xrightarrow{\text{past data}}$$

$$\frac{1}{1-0.9995} = \frac{1}{0.0005} \approx 2000 \text{ days Avg} \xrightarrow{\text{date}}$$

Hyperparameter Tuning in practice: Panda Vs Cavia!

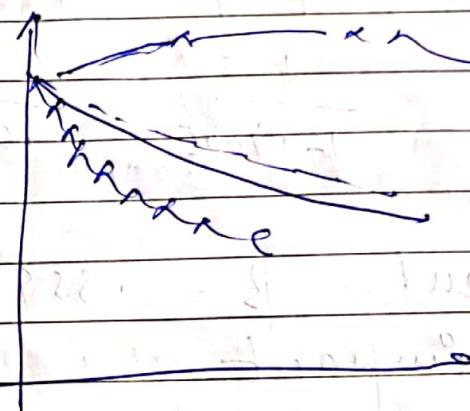
Panda! One model
at a time



day 0 → day(n)
day 0 → Reverse
batch

Patiently Training One
model at a time

Cavia! Multiple
model in parallel

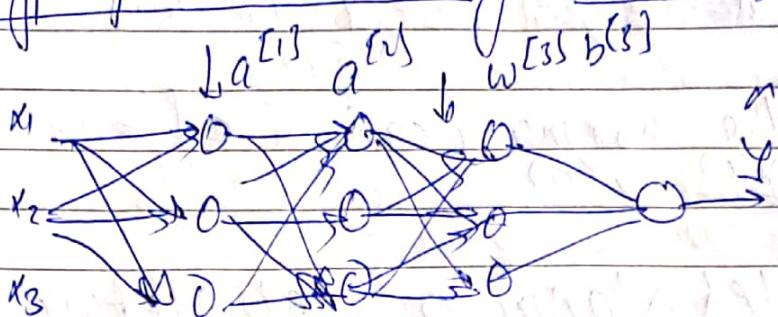


Try multiple model
& choose one.

25

30

Hyperparameter Tuning: Batch Normalization



Q: for any hidden layer can we normalize the values $a^{(l)}$ to next layer weight $w^{(l+1)}$ faster.

But Technically we normalize $z^{(l)}$ no $a^{(l)}$

Implementing Batch Norm

Given some intermediate values in NN $\underbrace{z^{(1)} \dots z^{(n)}}_{\text{hidden layer No.}}$

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

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

i^{th} unit $z^{Fe3(i)}$
hidden layer No.:-

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

Just in case $\sigma^2 = 0$
 $\epsilon = 10^{-6}$

but we don't want my hidden unit to be 0

$$\sigma^2 = 1 \text{ always.}$$

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

Gama Beta

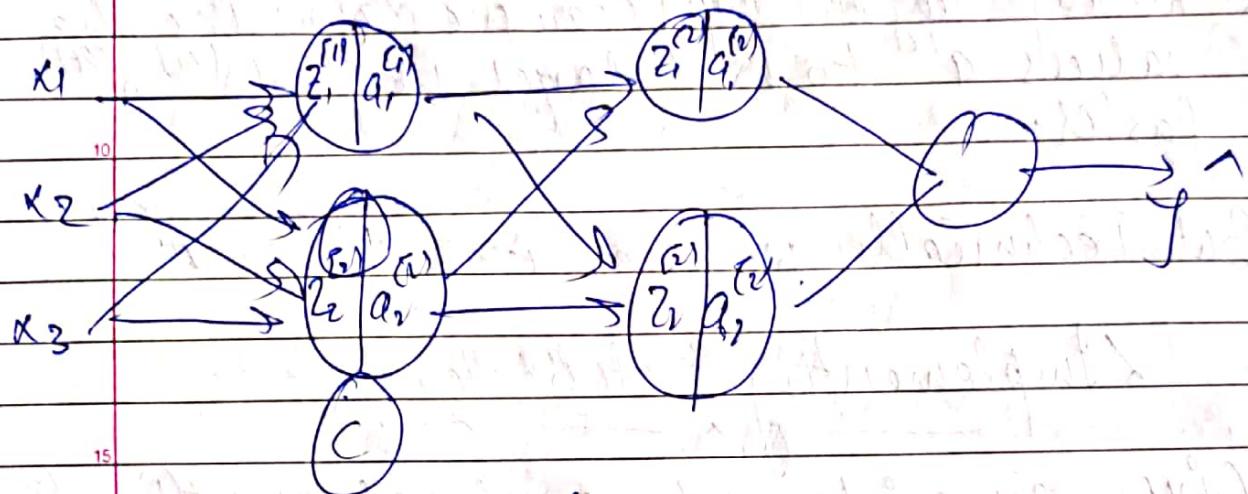
γ, β is learnable param using Gradient Descent / RMS / ADA / Momentum - etc }

if $\gamma = (\sigma^2 + \epsilon) / \beta = 1$ it will

nullify the normalization effect.

so $\hat{z}^{(l)} = z^{(l)}$

Applying Batch Norm to Deep Networks



$$x \xrightarrow{W, b} z^{(l)} \xrightarrow{\text{BatchNorm(BN)}} \hat{z}^{(l)} \rightarrow a^{(l)} = g(\hat{z}^{(l)})$$

$$a^{(l)} \xrightarrow{W, b} z^{(l)} \xrightarrow{\text{BN}} \hat{z}^{(l)} \rightarrow a^{(l+1)}$$

So params are now to be learnt:-

$$\begin{aligned} & w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)} \\ & \beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)}, \dots, \beta^{(L)}, \gamma^{(L)} \end{aligned}$$

$$\frac{d\beta}{d\delta^{(L)}} \quad \frac{d\beta}{d\delta^{(L)}} \quad \frac{d\beta}{d\delta^{(L)}}$$

{ BatchNorm + MultiBatches }

$$x^{(1)} \xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow{\text{BN}} \tilde{z}^{(1)} \xrightarrow{g^{(1)}(\tilde{z}^{(1)})} q^{(1)}$$

$$x^{(2)} \xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow{\text{BN}} \tilde{z}^{(1)} \xrightarrow{g^{(1)}} \hat{x} \rightarrow \dots$$

Params: $w^{(l)}, b^{(l)}, \beta^{(l)}, \gamma^{(l)}$

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

↓ while BN $b^{(l)}$ doesn't change anything while taking mean so $b^{(l)}$ can be scrapped so the params are now:-

Param :- $w^{(l)}, \beta^{(l)}, \gamma^{(l)}$

$$\tilde{z}^{(l)} = \frac{z^{(l)}}{\sqrt{\text{var}_{\text{norm}}}} = \gamma^{(l)} z_{\text{norm}} + \beta^{(l)}$$

$$q^{(l)} = (n^{(l)}, 1) \quad | \quad \beta = (h^{(l)}, 1) \quad \left. \begin{array}{l} \text{Dimension} \\ \gamma \end{array} \right\} =$$

25

30

for $t=1$ — num mini Batch
compute forward pass on $x^{[t]}$

In each hidden layer, use BN to
replace $\tilde{x}^{[t]}$ $\rightarrow \hat{x}^{[t]}$

Use backprob to compute

$$dw^{[t]}, d\beta^{[t]}, dr^{[t]}$$

Update Params -

$$w^{[t+1]} = w^{[t]} - \alpha dw^{[t]}$$

$$d\beta^{[t+1]} = \beta^{[t]} - \alpha dp^{[t]}$$

$$dr^{[t+1]} = r^{[t]} - \alpha dy^{[t]}$$

This works with ADAM/momentum/RMSprop

BATCH NORM TEST TIME?

At a test to predict μ, σ^2 is used
exponentially weighted avg of
mini batches:-

$$x^{[1]} \quad x^{[2]} \\ x \quad x$$

$$\mu^{[1]} \quad \mu^{[2]} \quad \dots \quad \mu^{[3]} \\ \text{expwt. avg.} \rightarrow \text{Meaning.}$$

Same way

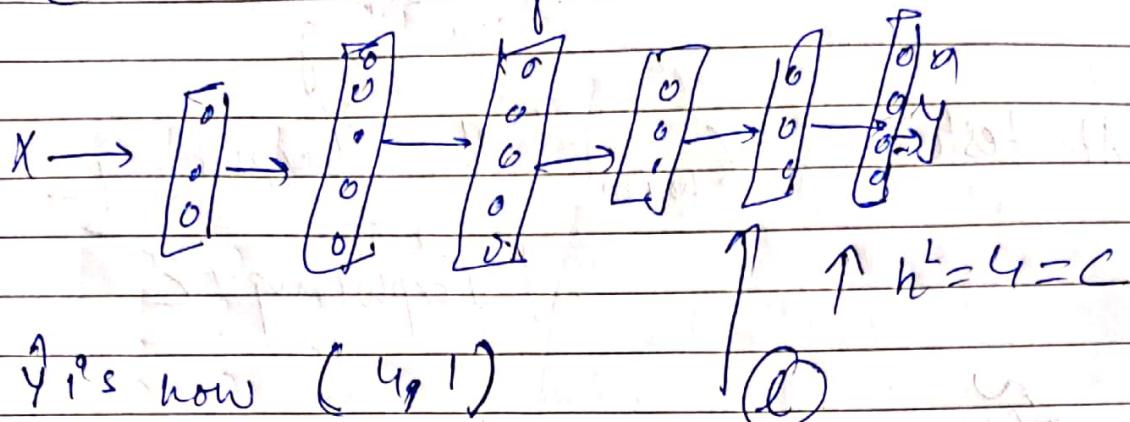
$$\sigma^2_{\text{fixes}}, \sigma^2_{\text{fit}(\alpha)} - \xrightarrow{\text{expwt ang}} \sigma^2_{\text{expwt ang}}$$

At testime $Z_{\text{norm}} = \frac{Z - \bar{M}_{\text{expwt ang}}}{\sqrt{\sigma^2_{\text{expwt ang}} + E}}$

$$Z = Z_{\text{norm}} + \beta$$

Softmax Regression

$C = 4 = \text{Number of Classes:-}$



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

Affiliation Softmax :-

$$t = e^{z^{[L]}} = (4, 1) \text{ vector}$$

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^{C=4} t_i} = (4, 1) \text{ vector.}$$

$$\text{Ex:- } z = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 0.4 \\ 0.1 \\ 20.1 \end{bmatrix}, \sum_{i=1}^4 t_i = 176.3$$

$$a^{[L]} = \frac{t}{176.3} = \begin{bmatrix} e^5/176.3 \\ e^2/176.3 \\ e^{-1}/176.3 \\ e^3/176.3 \end{bmatrix}$$

2 Training Softmax Classifier

$$z^{(L)} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$a^{(L)} = g^{(L)}(z^{(L)}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix}$$

Softmax

$$\text{Softmax} \rightarrow \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \xrightarrow{\text{hard max.}} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \text{"hard max"}$$

$$\text{Loss function} : - \sum_{j=1}^C y_j \log(\hat{y}_j)$$

$$y = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad \hat{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \text{GT}$$

$$\hat{y}_2 = 0.2 \quad \text{GT} = 1$$

make it small is to make \hat{y}_2 small and to make we have to make \hat{y}_2 as big as possible

This Loss on single train example

So on loss on entire train set:

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

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

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

$(4, m)$ dimension

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

$$\hat{\mathbf{y}} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$(4, m)$ dimension
sign matrix

{ Gradient descent with Softmax }

$$\text{Backprop} - d\hat{z}^{(t)} = \hat{y} - y$$

$$(4, 1) = (4, 1) - (4, 1)$$

$$\frac{\partial J}{\partial z^{(t)}} =$$