

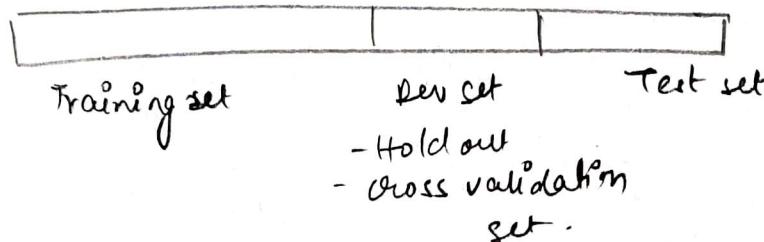
COURSERA - Improving Deep Neural Networks

⇒ Train | Dev | Test sets

Applied ML is highly iterative process

layers # hidden units learning-rates, activation functions

→ Data.



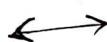
Mismatched train | test distribution

Training set:

cat pictures from
webpages }

Dev/test/sets:

cat pictures from your
app



→ make sure dev and test set come from same distribution.

→ Not having a test set might be okay (only dev set).

⇒ Bias / Variance

let's take cat classification

overfit

underfit

just
right.

Train set error: 1%

15%

15%

0.5%

Dev set error: 22%

16%

30%

1%

high variance

low bias

high bias

low variance

high bias

high variance

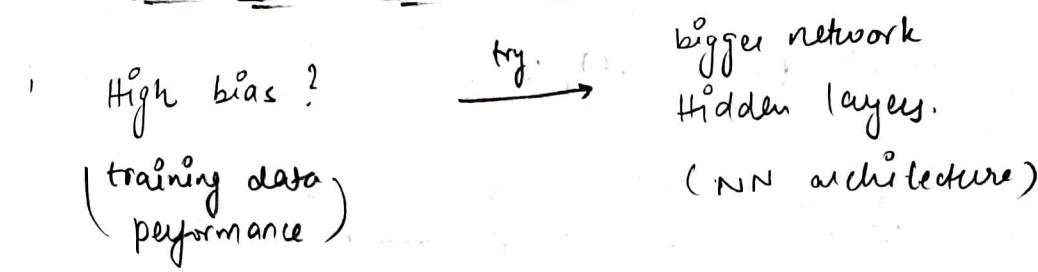
low bias

low variance

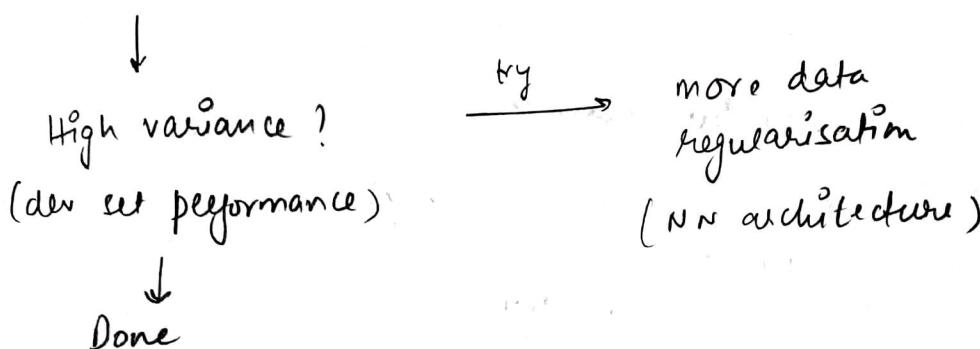
bias → due to inaccuracy of training set (high training error)

variance → due to difference between dev set error and
training set error

\Rightarrow Basic Recipe for Machine Learning



So, we firstly try to get rid of bias problem.



\Rightarrow Regularisation

Try when high variance is there \rightarrow overfitting.

Logistic regression

$$\min_{w, b} J(w, b) \quad w \in \mathbb{R}^n, b \in \mathbb{R}$$

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

$\lambda \rightarrow$ regularisation parameter.

$$L_2 \text{ regularisation} \quad \|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w$$

$$L_1 \text{ regularisation} \quad \frac{1}{2m} \sum_{j=1}^n |w_j| = \frac{\lambda}{2m} \|w\|_1$$

w will be sparse
 (lot of zeros)

Neural Network

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

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

"Frobenius norm"

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

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

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

called as weight decay.

⇒ Why regularisation reduces overfitting?

⇒ Dropout regularisation

↳ Dropout some nodes randomly during iterations

Implementing dropout ("Inverted dropout")

Illustrate with layer $l=3$, keep-prob = 0.8

$d_3 = np.random.rand(a_3.shape[0], a_3.shape[1]) < \text{keep_prob}$

$a_3 = np.multiply(a_3, d_3)$

$a_3 = a_3 / \text{keep_prob}$ (scaling factor)

For example, 50 units \rightarrow 10 units shut off.

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

↳ reduced by 20%.

thus, scale up a little by $a^{[3]} = \frac{a^{[3]}}{0.8}$.

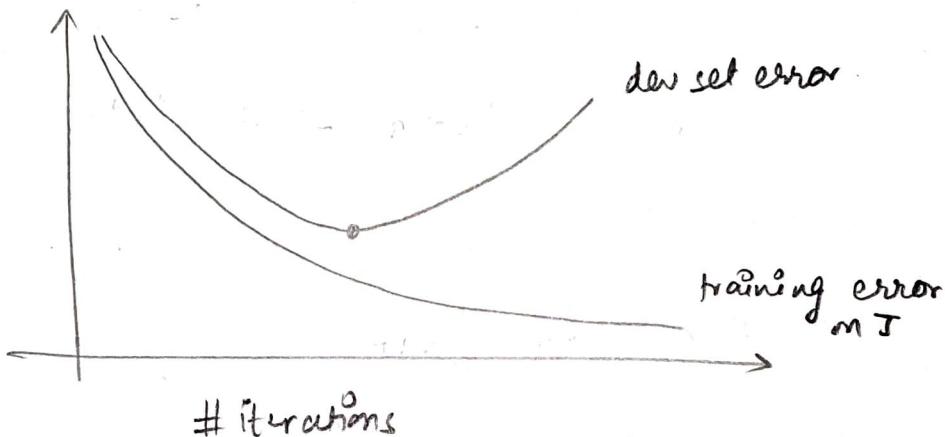
\rightarrow No dropout is done during test time.

\Rightarrow Other regularisation & techniques

(i) Data augmentation

- ↳ use flipped, tilted images
- ↳ distorted digits

(ii) Early stopping



Stop where dev-set error is minimum.

⇒ Normalising inputs

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

subtract mean

$$\mu = \frac{1}{m} \sum_{i=1}^m X^{(i)}$$

$$X = X - \mu$$

Normalise variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m X^{(i)} \star \star 2$$

element wise.

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

standard deviation

$$\text{So, in total. } X = \frac{X - \mu}{\sigma}$$

⇒ Vanishing | Exploding gradients

↳ This can be solved by careful initialisation of weights.

⇒ Weight initialisation for Deep NN

single neuron example

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

large $n \rightarrow$ small w_i

$$\text{Var}(w_i) = \frac{1}{n}, \frac{2}{n}$$

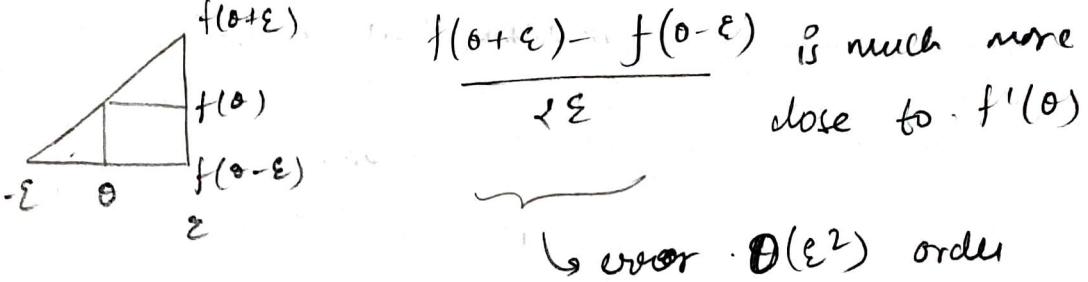
$$w^{(L)} = np \cdot \text{random_randn}(\text{shape}) * n \cdot \text{sqrt}\left(\frac{2}{n^{(L-1)}}\right)$$

This helps reduce the vanishing gradient problem
exploding

$$\text{For ReLU} \rightarrow \sqrt{\frac{2}{n}}, \text{ tanh} \rightarrow \sqrt{\frac{1}{n}}, \text{ other. } \sqrt{\frac{2}{n^{(L-1)} + n^{(L)}}}$$

⇒ Numerical approximation of gradients

$$\text{let } f(\theta) = \theta^3 \quad g(\theta) = \frac{f(\theta + \varepsilon) - f(\theta - \varepsilon)}{2\varepsilon}$$



$\frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon}$ is much more close to $f'(\theta)$
 \hookrightarrow error $\Theta(\epsilon^2)$ order

whereas $\frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} \rightarrow \Theta(\epsilon)$ order error.

gradient checking

Take $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, W^{[L]}, b^{[L]}$ and reshape into a vector θ

$$J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, b^{[L]}) = J(\theta)$$

Take $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$
 $\underbrace{\quad \quad \quad}_{\text{concatenate}}$

gradient checking

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

for each i :

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

$$\approx \frac{d\theta_i}{d\theta_i} = \frac{dJ}{d\theta_i}$$

check $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \approx 10^{-7}$ great
 10^{-5} aware
 10^{-3} concerned

$$\epsilon = 10^{-7}$$

\Rightarrow Gradient checking implementation, notes

- \rightarrow Don't use in training - only ^{to} debug.
- \rightarrow If algorithm fails grad check, look at components to try to identify bug.
- \rightarrow Remember regularization.
- \rightarrow Doesn't work with dropout
- \rightarrow Run at random initialisation again after training.

\Rightarrow Week-2 Optimization Algorithms

\Rightarrow Mini Batch Gradient Descent

If m is very large, like $m = 5,000,000$

minibatches of $m = 1000$

$$X = \left[\underbrace{x^{(1)}, x^{(2)}, \dots, x^{(1000)}}_{X^{(1)}}, \underbrace{u^{(1000)}, v^{(1000)}, \dots, x^{(2000)}}_{X^{(2)}}, \dots, \underbrace{x^{(m)}}_{X^{(5000)}} \right]$$

$$y = \left[\underbrace{y^{(1)}, y^{(2)}, \dots, y^{(1000)}}_{y^{(1)}}, \underbrace{\dots}_{y^{(2)}}, \dots, \underbrace{\dots}_{y^{(5000)}} \right]$$

$$X^{(1)} \rightarrow (n_x, 1000) \quad X^{(2)} \rightarrow (n_x, 1000)$$

$$y^{(1)} \rightarrow (1, 1000), \quad y^{(2)} \rightarrow (1, 1000)$$

process mini batches.

Mini batch gradient descent

1 step of GD using

for $t = 1 \dots 5000 \{$

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

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

$$z^{[1]} = w^{[1]} X^{\{t\}} + b^{[1]}$$

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

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

vectorize
1000 examples.

compute cost

$$J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^l \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_k \|w^{(k)}\|_F^2$$

Back prop to compute gradients - wrt $J^{\{t\}}$

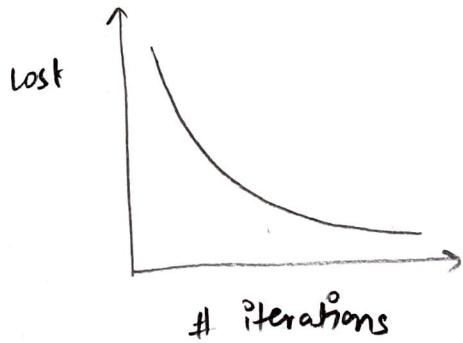
$$w^{[1]} = w^{[1]} - \alpha \nabla J^{\{t\}}$$

$$b^{[1]} = b^{[1]} - \alpha \nabla J^{\{t\}}$$

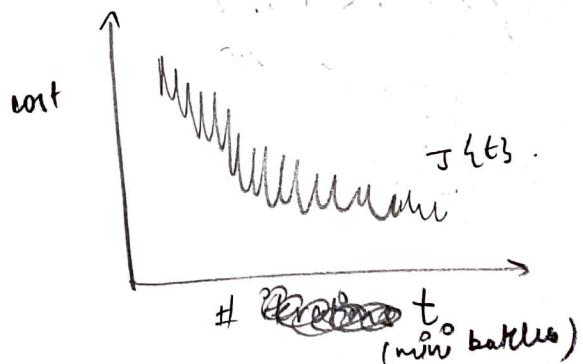
} one "epoch"

batch grad. descent

→ single pass through
training set.



Mini batch GD



Fluctuations due to different training sets, however,
overall trend is downward

Choosing your mini batch

mini-batch size = m : Batch grad descent $(x^{(1)}, y^{(1)})$
 $= (X, Y)$

mini-batch size = 1,

stochastic gradient

Descent

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

every example is
its own mini
batch.

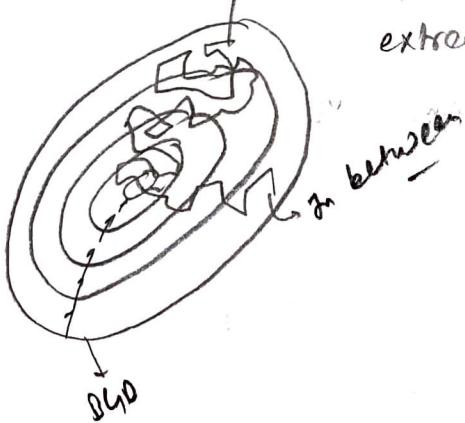
Practical mini-batch size is somewhere between those.

BGD \rightarrow Too long time per iteration.

SGD \rightarrow lose speedup from vectorisation.

In between \rightarrow Fast learning \hookrightarrow vectorisation

SGD makes progress without processing
extremely tiny set.



small training set : Use batch GD.
($m \leq 2000$)

Typical mini-batch sizes:

64, 128, 256, 512

Make sure x^{t+3}, y^{t+3} fit in CPU/GPU memory.

⇒ Exponentially weighted averages

Temperature in London

$$\theta_1 = 40^{\circ}\text{F}$$

$$\theta_2 = 49^{\circ}\text{F}$$

$$\theta_3 = 45^{\circ}\text{F}$$

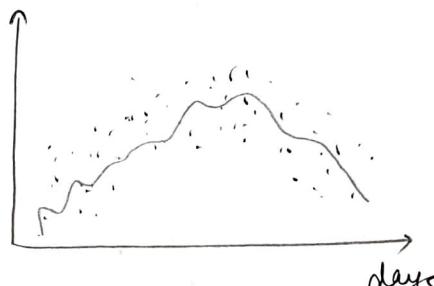
⋮

$$\theta^{180} = 60^{\circ}\text{F}$$

$$\theta^{181} = 58^{\circ}\text{F}$$

⋮

temperature



$$V_0 = 0$$

$$V_1 = 0.9V_0 + \theta_1$$

$$V_2 = 0.9V_1 + \theta_2$$

$$V_3 = 0.9V_2 + \theta_3$$

$$V_t = BV_{t-1} + (1-B)\theta_t$$

$$B = 0.9 \quad : \approx 10 \text{ days}$$

$$B = 0.98 \quad : \approx 50 \text{ days}$$

$$B = 0.5 \quad : \approx 2 \text{ days}$$

Too much noise.

V_t is interpreted as
approximate average
over \approx

$\frac{1}{1-B}$ days of
temperature.

Too smooth, but
latency.

\Rightarrow Exponentially weighted averages

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

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

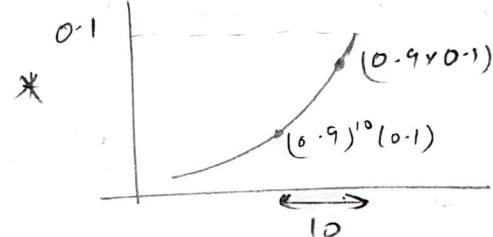
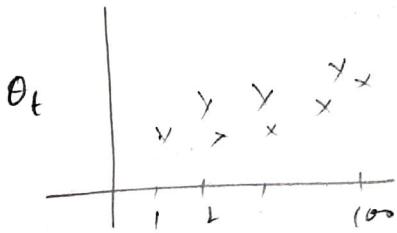
$$V_{99} = 0.9V_{98} + 0.1\theta_{99}$$

:

$$V_{100} = 0.1\theta_{100} + 0.9 \times (0.1\theta_{99} + V_{98} \times 0.9)$$

$$= 0.1\theta_{100} + 0.9 \times 0.1\theta_{99} + 0.9 \times 0.9 \times 0.1\theta_{98}$$

$$+ (0.9)^3 (0.1) \theta_{97} \dots$$



$$\epsilon = 1 - \beta.$$

$$(1-\epsilon)^{1/\epsilon} = \frac{1}{e}$$

$$(0.9)^{10} = \frac{1}{e} \approx 0.35 \approx \frac{1}{3} \text{rd}$$

\Rightarrow Understanding exponentially weighted averages

implementation

$$V_0 = 0$$

$$V_0 = \beta V_0 + (1-\beta)\theta_1$$

$$V_0 = \beta V_0 + (1-\beta)\theta_2$$

$$V_0 = 0$$

repeat {

get next θ_t

$$V_0 = \beta V_0 + (1-\beta)\theta_t$$

}

Just one
line of
code.

$\Rightarrow \beta^0$ as correction

During initial phase $v \approx 0$, we need to correct this

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

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

$$v_0 = 0$$

$$v_1 = 0.98 v_0 + 0.02 \theta_1$$

$$v_2 = 0.98 v_1 + 0.02 \theta_2$$

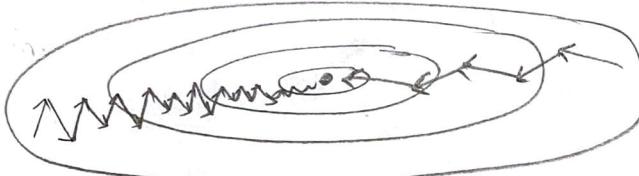
$$= 0.98 \times 0.02 \theta_1 + 0.02 \theta_2$$

$$= 0.0196 \theta_1 + 0.02 \theta_2$$

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

$$\frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

\Rightarrow Gradient Descent with momentum



Momentum:

on iteration t :

compute d_w, d_b in current minibatch.

$$Vd_w = \beta Vd_w + (1-\beta)d_w$$

$$Vd_b = \beta Vd_b + (1-\beta)d_b.$$

$$w = w - \alpha Vd_w.$$

Implementation details

Hyperparameters: $\alpha, \beta \sim 0.9$ generally

On iteration t :

compute dW, db on current
mini-batch.

$$\nabla_{dW} = 0 \quad \nabla_{db} = 0$$

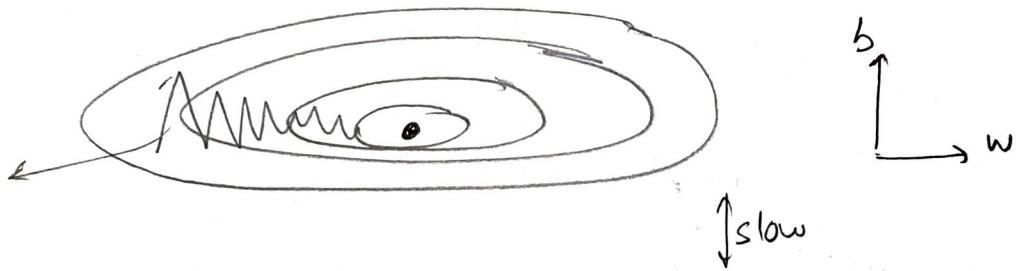
↓
some dim
as dW
as db

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

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

$$W = W - \alpha V_{dW}, \quad b = b - \alpha V_{db}$$

$\Rightarrow \underline{\text{RMS Prop}} \rightarrow \text{Root mean square prop}$



On iteration t

compute dW, db on current
mini-batch.

$$sdw = \frac{\beta}{2} sdw + (1-\frac{\beta}{2}) dW^2 \quad (\text{elementwise square})$$

$$sdb = \frac{\beta}{2} sdb + (1-\frac{\beta}{2}) db^2$$

$$W = W - \frac{\alpha dW}{\sqrt{sdw + \epsilon}} \quad b = b - \frac{\alpha db}{\sqrt{sdb + \epsilon}} \quad \begin{array}{l} \text{just to} \\ \text{ensure,} \\ \text{nothing} \\ \text{explodes.} \end{array}$$

Now, we want learning to go fast in w direction.

so, dW should be small and db should be greater and that is indeed the case with our contours.

Adam optimisation Algorithm

$$VdW = 0, \underbrace{sdw = 0}_{\text{initialisation}}, Vdb = 0, \underbrace{sdw^2 = 0}_{\text{initialisation}}$$

On iteration t :

compute dW, db using current mini batch.

$$VdW = \beta_1 VdW + (1 - \beta_1) dW, \quad Vdb = \beta_1 Vdb + (1 - \beta_1) db$$

$$sdw = \beta_2 sdw + (1 - \beta_2) dw^2 \quad \underbrace{\text{"Momentum" } \beta_1}_{\text{initialisation}}$$

$$sdw^2 = \beta_2 sdw^2 + (1 - \beta_2) dw^2 \quad \underbrace{\text{"RMS prop" } \beta_2}_{\text{initialisation}}$$

$$V_{dw}^{\text{corrected}} = \frac{VdW}{1 - \beta_1^t} \quad V_{db}^{\text{corrected}} = \frac{Vdb}{1 - \beta_1^t}$$

$$sdw^{\text{corrected}} = \frac{sdw}{1 - \beta_2^t} \quad sdw^2_{\text{corrected}} = \frac{sdw^2}{1 - \beta_2^t}$$

$$W = W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{sdw^{\text{corrected}}} + \epsilon} \quad b = b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{sdw^2_{\text{corrected}}} + \epsilon}$$

Hyperparameters choice

α : needs to be tuned

$$\left. \begin{array}{l} \beta_1 : 0.9 \quad (dw) \\ \beta_2 : 0.999 \quad (dw^2) \\ \epsilon : 10^{-8} \end{array} \right\} \text{generally accepted values.}$$

Adam : Adaptive moment estimation

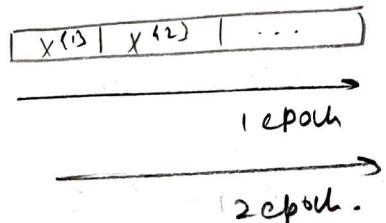
\Rightarrow Learning rate Decay

↳ decrease learning rate over iterations.

slowly reduce α .

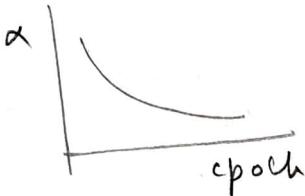
1 epoch = 1 pass through the data.

$$\alpha = \frac{1}{1 + \text{decay rate} \times \text{epoch num}} \alpha_0$$



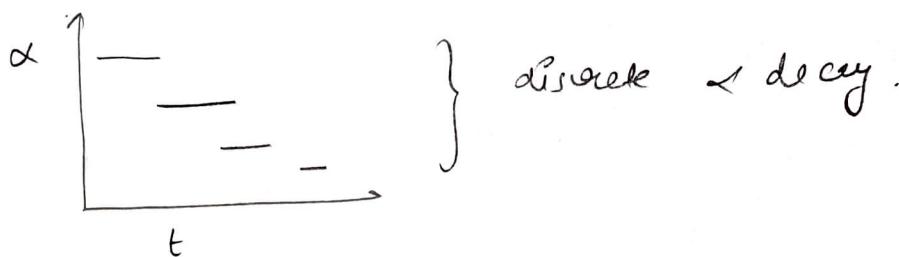
let $\alpha_0 = 0.2$ decay rate = 1

cepoch	α
1	0.1
2	0.067
3	0.05
4	0.04



Also, $\alpha = 0.95^{\text{epoch num}} \times \alpha_0$ (exponential decay)

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



⇒ The problem of local optima

→ In higher dimensions, most of the points are saddle points and there is a very less chance of an absolute minimum.

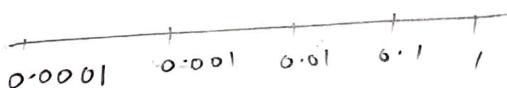
⇒ Week-3 Hyperparameter tuning

⇒ Hyperparameters

most important $\leftarrow \alpha$ # layers
 $B \sim 0.9$ # hidden units
 B_1, B_2, \dots learning rate decay.
 $(0.9) (0.999) (10^{-8})$
mini-batch size.

Try random values → Don't use a grid, try random values.

$\alpha \rightarrow \log$ scale



$$\alpha = -4 * \text{np.random.rand}() \quad \hookrightarrow \alpha \in [-4, 0]$$

$$\alpha = 10^{\alpha} \quad \rightarrow 10^{-4} \text{ to } 1$$

\Rightarrow Batch Normalisation

Normalising inputs works and speeds up GD.
(X)

So, can we normalise every $a^{[l-1]}$ to train every $w^{[l]}$,
 $b^{[l]}$ faster. We would actually normalise $z^{[l-1]}$.

\rightarrow Implementing batch norm.

Given some intermediate values in NN. $z^{[l-1]}(i)$

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

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

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

learnable parameters.

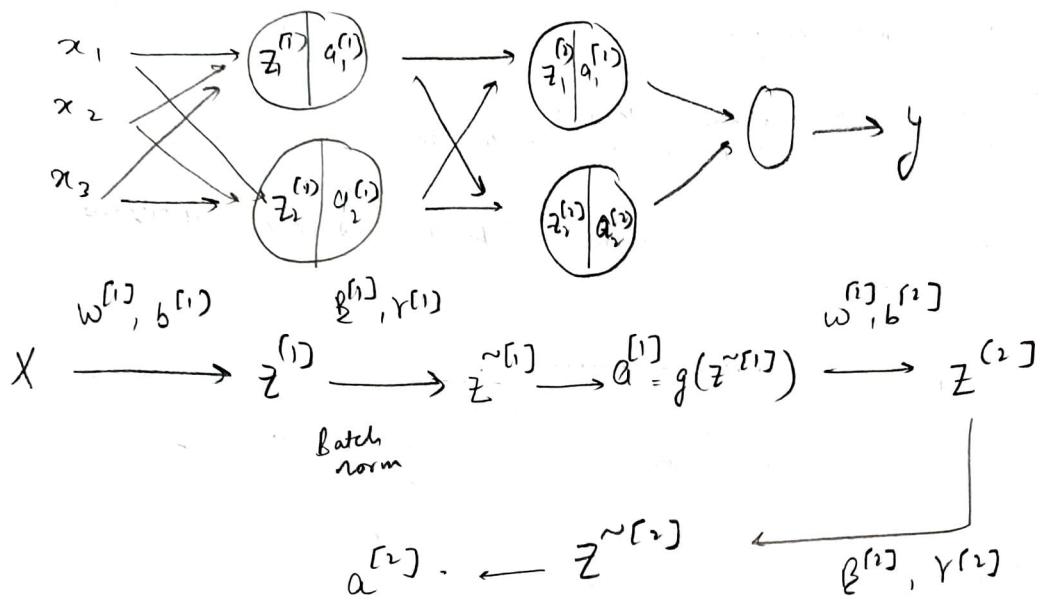
$$z^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

$$\left. \begin{array}{l} \text{if } \gamma = \sqrt{\sigma^2 + \epsilon} \quad \beta = \mu \\ z^{(i)} = z^{(i)} \end{array} \right\} \text{normalisation step is proceeded.}$$

So, using γ, β to make different mean & variance other than 0 and 1.

We do not always want mean and variance for $z^{(i)}$'s to be 0 and 1. With γ and β we can control them.

→ Fitting Batch Norm into a Neural Network

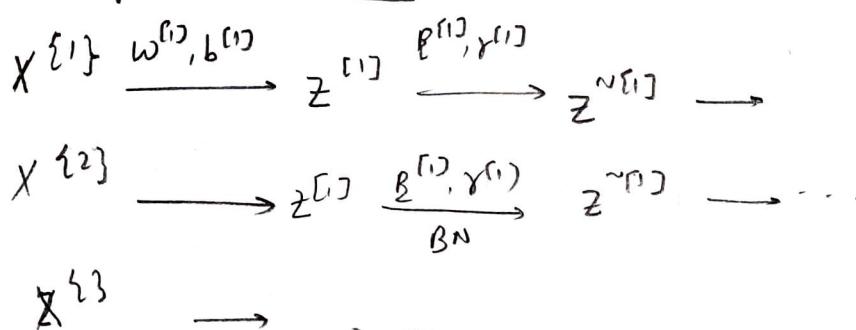


Parameters → $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$
 $\underbrace{\beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]}}$
 ↴ different β used in momentum & Adam.

We update β and γ similar to w and b

$$\text{GD.} \rightarrow \beta^{[t]} = \beta^{[t-1]} - \alpha d\beta^{[t]}$$

Working with mini batches



Parameters: $w^{[1]}, b^{[1]}, \beta^{[1]}, \gamma^{[1]}$

↳ no longer required as it gets cancelled in \tilde{z} norm.

$$\text{parameters } w^{(e)}, b^{(e)}, \gamma^{(e)}, \beta^{(e)}$$

$$(n^{(L)}, n^{(L-1)}) \quad (n^{(L)}, 1) \quad (n^{(e)}, 1)$$

$$z^{(e)} \rightarrow (n^{(e)}, 1)$$

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

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

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

$$z^{\sim(L)} = \underbrace{\gamma^{(L)} z_{\text{norm}}^{(L)}}_{\text{elementwise}} + \beta^{(L)}$$

multiplication

\Rightarrow Implementing GD

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

compute forward prop on $X^{(t)}$

in each hidden layer, BN to get $\tilde{z}^{(L)}$ from $z^{(L)}$

use backprop to compute $dW^{(L)}, d\gamma^{(L)}, dB^{(L)}$

$$\begin{aligned} \text{update parameters} \quad w^{(L)} &= w^{(L)} - \alpha dw^{(L)} \\ \beta^{(L)} &= \beta^{(L)} - \alpha dB^{(L)} \\ \gamma^{(L)} &= \gamma^{(L)} - \alpha d\gamma^{(L)} \end{aligned}$$

works same manner with Adam, RMSprop, momentum.

\Rightarrow Batch norm at test-time

$$\begin{aligned} \mu &= \frac{1}{m} \sum_i z^{(i)} \\ \sigma^2 &= \frac{1}{m} \sum_i (z^{(i)} - \mu)^2 \end{aligned} \quad \left[\begin{array}{l} m = \text{mini} \\ \text{batch} \\ \text{size} \\ \text{here.} \end{array} \right]$$

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

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

μ and σ used here cannot be applied at test time because they are different for different batches.

μ, σ^2 : estimate using exponentially weighted average
(across minibatches)

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

exponentially weighted.

$$\mu^{(1)}[t], \mu^{(2)}[t], \mu^{(3)}[t], \dots$$

$\mu^{[t]}$

$$\theta_1, \theta_2, \theta_3, \dots$$

exponentially weighted

$$\sigma^2 \{1\}[t], \sigma^2 \{2\}[t], \sigma^2 \{3\}[t] \rightarrow \sigma^2[t]$$

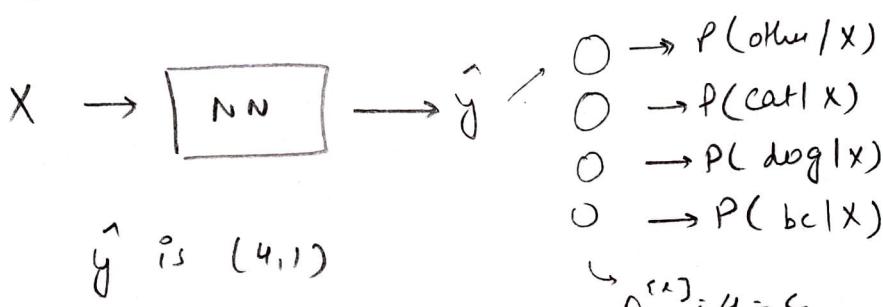
$$z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \left. \right\} \text{test time} \quad z^* = \gamma z_{\text{norm}} + \beta$$

\Rightarrow Softmax regression

↳ multiclass logistic regression

For example, we need to recognise, cats, dogs, baby chicks

$$c = \# \text{ classes} = (0, 1, 2, 3)$$



softmax layer

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

Activation function:

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

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum t_i}$$

$$\frac{t_i}{\sum t_i} = a_i^{[L]}$$

Loss function in softmax

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{cat} \quad \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad \text{say.}$$

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

$$\text{For this example } y_1 = y_3 = y_4 = 0$$

$\Rightarrow \mathcal{L}(\hat{y}, y) = -y_2 \log \hat{y}_2 = -\log \hat{y}_2$, make \hat{y}_2 big to minimize loss function.

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

$$Y = [y^{(1)}, y^{(2)} \dots y^{(m)}] = \begin{bmatrix} 0 & 0 & 1 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & \dots \end{bmatrix} (4 \times m)$$

\hat{y} = $(4 \times m)$ matrix.

\rightarrow go with softmax

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

$$\text{Backprop: } \frac{d\mathcal{L}^{[1]}}{dz^{[1]}} = \hat{y} - y \rightarrow \text{derive}$$

$$\frac{d\mathcal{L}}{dz^{[L]}}$$