

Course II

* Train/dev/test sets

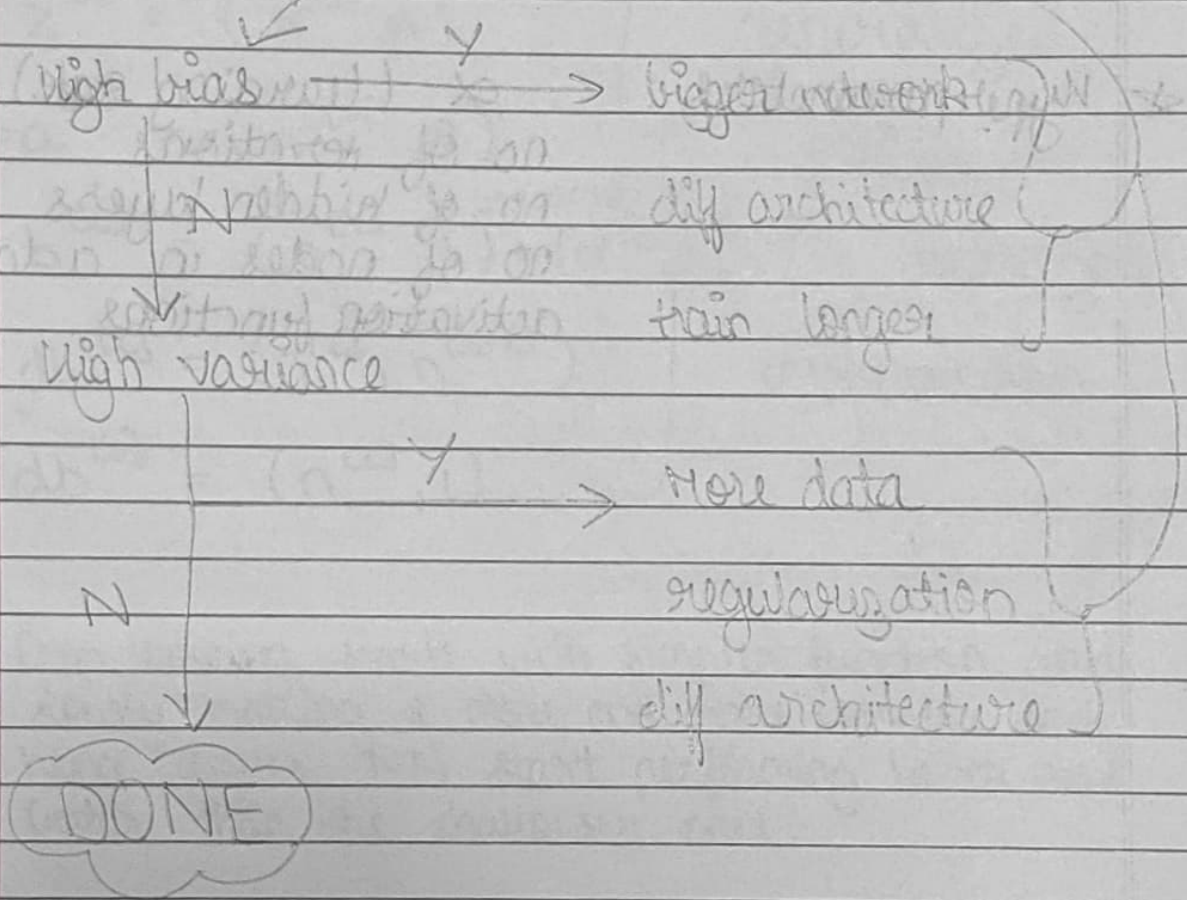
smaller dataset : 60/20/20

larger dataset : 98/1/1

Test and dev datasets should be from same distribution.

→ Train set error → bias

Dev set error → variance



★ Regularisation. (Logistic Regression)

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

$$\|w\|^2 = w^T w$$

L2 regularisation

★ Regularisation of neural network.

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

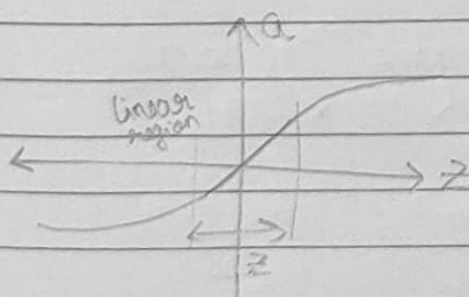
fröbénius norm

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

$$dw^{[l]} = \left(\text{from back} \right)_{\text{prop}} + \frac{\lambda}{m} w^{[l]}$$

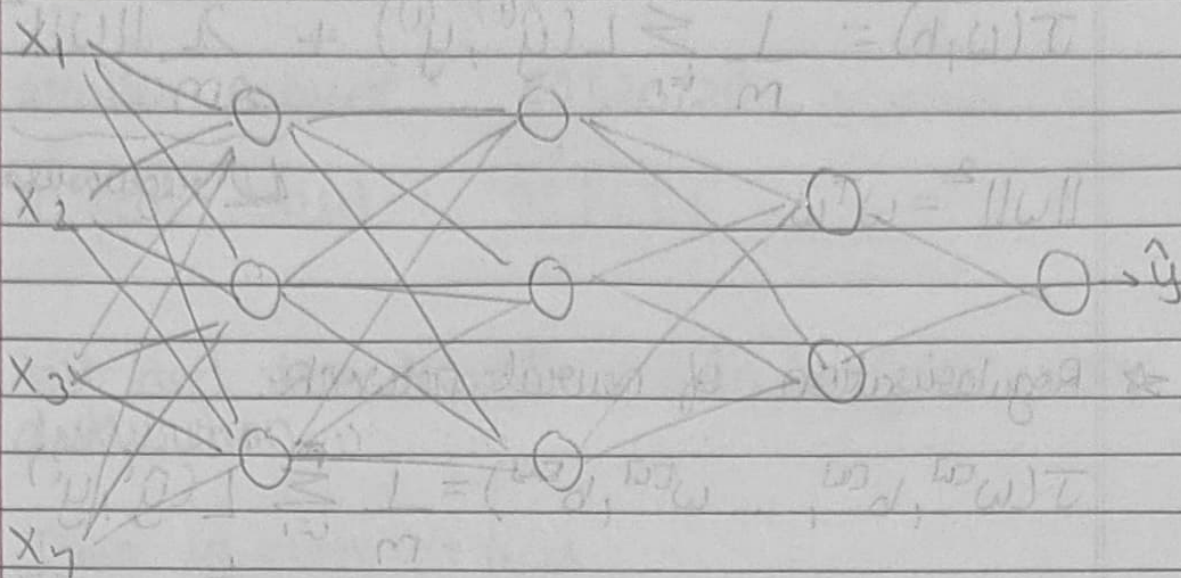
$$dw := dw - \alpha dw^{[l]}$$

Note: $\lambda \uparrow \quad w^{[l]} \downarrow \quad z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \downarrow$



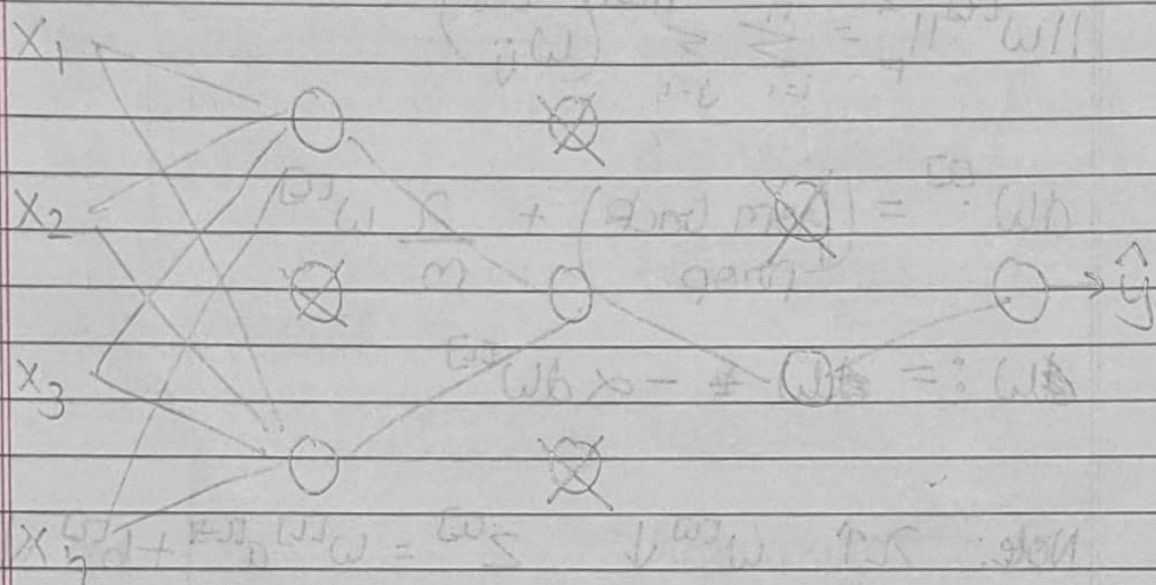
it become linear
and hence reduces
overfitting of
the function.

★ Dropout Regularisation



delete nodes at random

↓ dropout



★ Inverted dropout

neuron dropout

★

generates a matrix of given shape with values

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

< keep_prob

btw 0 and 1

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

$$a3 /= \text{keep_prob}$$
→ to not reduce the expected value of z

keep_prob → is the prob with which decide to keep or eliminate the neuron.

higher
→ ~~lower~~
issue

for layers where overfitting isn't a

lower ~~higher~~ for layers where overfitting is a
issue

downside is that J is not well defined

$$E[y] = \mu$$

$$y - \mu = x$$

$$E[y] = \mu$$

$$y - \mu = x$$

$$E[y] = \mu$$

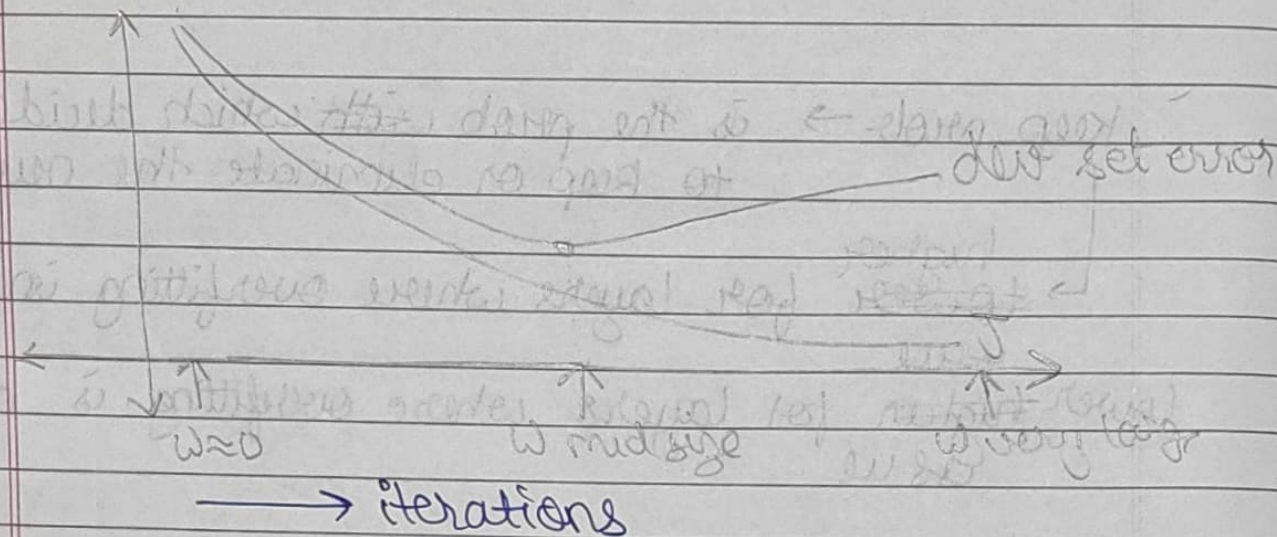
$$y - \mu = x$$

★ Data augmentation

- flip the image horizontally } don't add too much info and is
 → random crops of the image } cheap to do.

This way overfitting can be reduced.

★ Early stopping



orthogonalization → focus on one task at a time.

★ Normalize the data → σ_3 to optimize

Make mean 0

$$\mu = \frac{1}{m} \sum x_i$$

$$x := x - \mu$$

Normalise variance

$$\sigma^2 = \frac{1}{m} \sum (x_i)^2$$

$$x := x / \sigma$$

* Weight initialization:

ReLU: $\text{Var} := \sqrt{2/n}$

tanh: $\text{Var} := \sqrt{1/n}$

* grad check

θ = vector of $w_1^{[1]}, b_1^{[1]}, w_2^{[1]}, b_2^{[1]}, \dots$

$d\theta$ = vector of $dw_1^{[1]}, db_1^{[1]}, dw_2^{[1]}, db_2^{[1]}, \dots$

for each i :

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

$$\approx d\theta_i = \frac{\partial J}{\partial \theta_i}$$

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

$\epsilon = 10^{-7}$

use in debug only

NOT WILL TRAINING

★ Minibatch

→ split the huge training dataset into smaller parts

$$X = [X^{(1)} \quad X^{(2)} \quad X^{(3)} \quad \dots \quad X^{(1000)} \quad | \quad X^{(1001)} \quad \dots \quad X^{(2000)} \quad | \quad \dots \quad | \quad \dots \quad X^{(n)}]$$

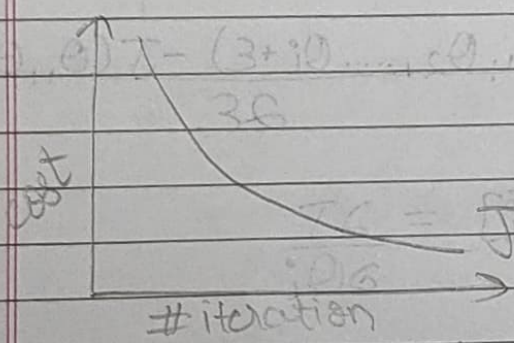
$(n, m) \leftarrow \xrightarrow{X^{(i)}}$

$$Y = [Y^{(1)} \quad Y^{(2)} \quad Y^{(3)} \quad \dots \quad Y^{(1000)} \quad | \quad Y^{(1001)} \quad \dots \quad Y^{(2000)} \quad | \quad \dots \quad | \quad \dots \quad Y^{(n)}]$$

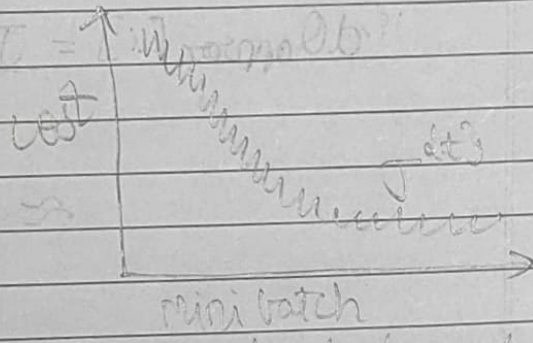
$(1, m) \leftarrow \xrightarrow{Y^{(i)}}$

$$Y^{(i)} : [Y^{(1)} \quad Y^{(2)} \quad \dots \quad Y^{(1000)}]$$

$$X^{(i)} : [X^{(1)} \quad X^{(2)} \quad \dots \quad X^{(1000)}]$$

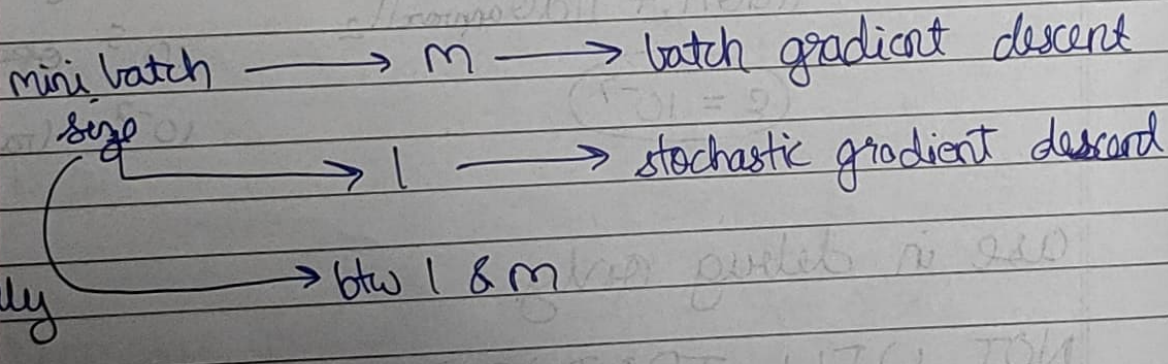


batch gradient



mini batch

gradient descent



* Exponentially weighted average

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

$$\frac{1}{1-\beta}$$

average

if $\beta = 0.9 \approx 10$ day average

if $\beta = 0.5 \approx 2$ day average

* Bias Correction

$\frac{V_t}{1-\beta^t}$ → initialize phase me help karta → as t gets bigger $1-\beta^t \approx 1$ and hence no effect.

* gradient descent with momentum

on iteration t :

$$dw, db$$

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

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

$$w = w - \alpha V_{dw}$$

$$b = b - \alpha V_{db}$$

$\alpha, \beta \rightarrow$ hyperparameter

★ RMS prop

On iteration t :Example hai
generally higher
dimensionalcompute dw, db

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2$$

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

small

large

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

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

to ensure it does
not blow up as S_{db}, S_{dw} can be 0

★ Adam optimization algorithm.

$$v_{dw} = 0, S_{dw} = 0 \quad v_{db} = 0, S_{db} = 0$$

On iteration t ,compute dw, db

$$v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dw$$

$$v_{db} = \beta_1 v_{db} + (1 - \beta_1) db$$

 $\beta_1 \rightarrow$ momentum $\beta_2 \rightarrow$ RMS prop.

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2$$

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

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{1 - \beta_1^t}$$

$$V_{db}^{\text{corrected}} = \frac{V_{db}}{1 - \beta_1^t}$$

$$S_{dw}^{\text{corrected}} = \frac{S_{dw}}{1 - \beta_2^t}$$

$$S_{db}^{\text{corrected}} = \frac{S_{db}}{1 - \beta_2^t}$$

$$w := w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}} + \epsilon}}$$

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

generally, $\beta_1 = 0.9$ $\beta_2 = 0.999$ $\epsilon = 10^{-8}$

* Learning rate decay

↳ initial make bigger step but as the network matures smaller steps.

∴ α is decreased with time.

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

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

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

use random values for hyperparameter tuning.

★ Batch Norm

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

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

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

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

$$\beta = \mu$$

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

learnable
parameter

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

$$\tilde{z}^{[2]} \xrightarrow[\beta, \gamma]{\text{BN}} \tilde{z}^{[2]} \rightarrow a^{[2]} = g(\tilde{z}^{[2]}) \rightarrow y$$

★ Batch norm has slight regularisation effect.

★ Batch norm at test time

↳ at test time we may not have enough examples to calculate the μ or σ^2

So, we estimate it using exponentially weighted moving average.

★ Softmax Layers.

let $z^{[L]}$ be $(n, 1)$

↳ activation func

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

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^n t_i} \quad (n, 1)$$

$$a_i^{[L]} = \frac{e^{z_i^{[L]}}}{\sum_{i=1}^n t_i} = \frac{t_i}{\sum t_i}$$

Eg: $y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ $\hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.4 \\ 0.1 \end{bmatrix}$

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

★ Plateaus can slow learning

local optima is very rare in higher dimensional space.