

Improving Deep Neural Networks: Hyperparameters tuning, Regularisation, Optimisation:

* Train / Dev / Test sets:



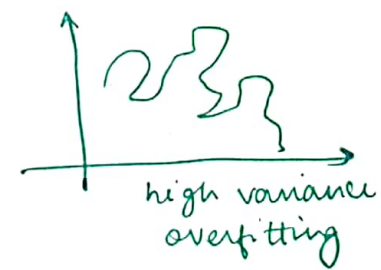
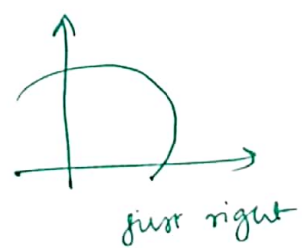
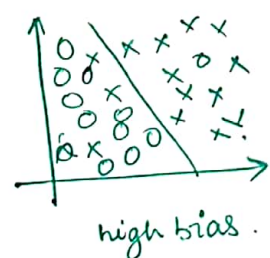
70/30 %
Train test

60/20/20
Train dev Test

Big data: 98/1/1
Train dev Test

- make sure dev and test set come from same ~~resolution~~ distribution.

* Bias / Variance:



Train set error 1% → very good
Dev set error 11% → relatively poor } High variance

Train set error 15%
Dev set error 16% } underfitting / High bias

15% } high bias & high variance:
30%

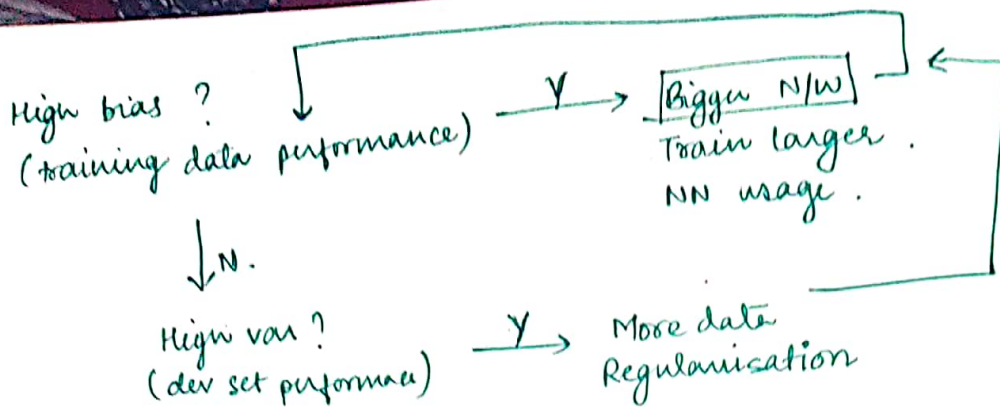
0.5% } low bias
1% } low variance.

Human ≈ 0%

Bayes / optimal error ≈ 0%
if 15%
high bias

HB HV





"Bias variance tradeoff"
 \rightarrow in new deep learning era, tradeoff has been ~~not~~ nullified

Regularisation of Neural NW

\rightarrow * logistic reg:

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

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

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

L_1 regularisation: $\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$

$\lambda \rightarrow$ regularisation parameter.

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

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

$$w: (n^{[L]}, n^{[L-1]})$$

$\uparrow \quad \uparrow$

F: "Frobenius norm"

$dw =$ (from backprop)

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

$$\frac{dJ}{dw^{[L]}} = dw^{[L]}$$

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

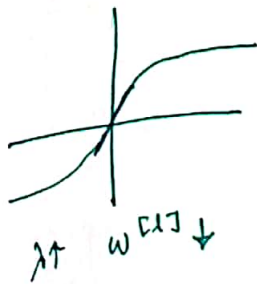
$$w^{[L]} = \frac{\alpha}{m} dw^{[L]} - \alpha (\text{backprop})$$

$$(1 - \frac{\alpha \lambda}{m})$$

* Regularisation vs overfitting

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

$w^{[L]} \approx 0$



$$\tanh(z) = g(z)$$

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

Every layer \approx linear
can lead to high bias.

* dropout regularisation

keep 0.5 chance of keeping 1 node, 0.5 chance of removing
then train a smaller network.

implementation:

• inverted dropout:

$$l=3:$$

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

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

$$\# a3 *= d3$$

$$a3 /= \text{keep_prob}$$

Ex: 50 units \rightarrow 10 unit shut off

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

\hookrightarrow reduce by 20%.

$$/= 0.8$$

• Making predictions at x time

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

no drop out

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

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

⋮

$$/= \text{keep_prob}$$

$\#$
dropout.

2 3 4 5 6

Q2

Smaller keep-probs makes more powerful drop ~~off~~ auto.

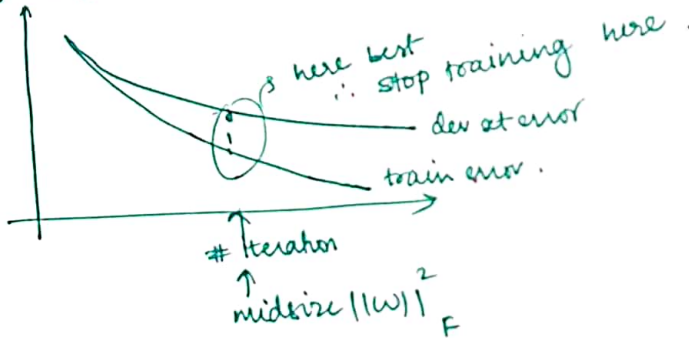
→ dropout helps prevent overfitting

→ downside: 'J' not well defined
∴ less calculatable.



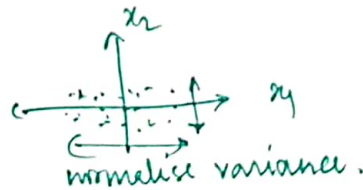
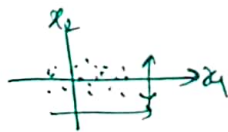
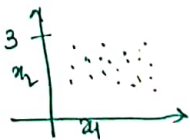
* other regularisation techniques:

- Data augmentation
 - make TS more redundant
 - random distortions
 - ...
- Early stopping



⊕ Setting up optimisation problem:

* Normalising Train sets.



Subtract mean:

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

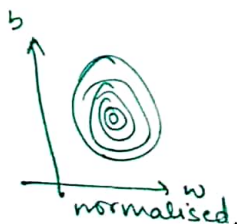
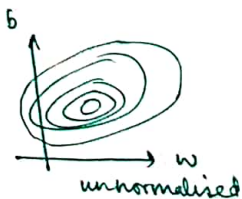
$$x' = x - \mu$$

Normalize var.

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

6 classes use

$$x / \sigma$$



* vanishing / Exploding gradients :

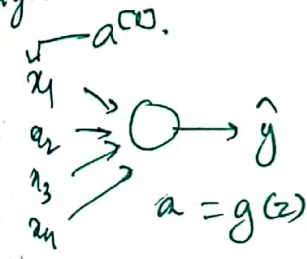
→ when training deep NN, grads. can become very large or very small.

if $g(z) = 2$ $b^{[L]} = 0$.

$$\hat{y} = w^{[L]} w^{[L-1]} \dots \underbrace{w^{[2]} w^{[1]} x_1}_{a_2} z_1^{[L]}$$



* weight initialisation :



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

larger $n \rightarrow$ smaller w_i

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

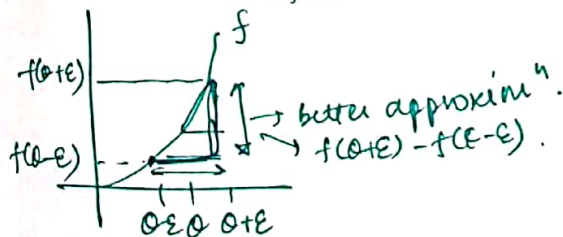
$$w^{[L]} = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}\left(\frac{2}{n^{[L-1]}}\right) \rightarrow \text{Relu}$$

tanh : $\sqrt{\frac{1}{n^{[L-1]}}}$

→ Xavier initialisation.

* Numerical approx of gradients :-

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



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

$$\frac{(1 - 0.1)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

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

$$\begin{matrix} 0(\epsilon) \\ 0.01 \\ 0.0001 \end{matrix}$$

$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \rightarrow \text{less accurate}$$

* Gradient checking:

$w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ reshape into big vector θ .

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

Take $dw^{[1]}, db^{[1]}, \dots, dw^{[L]}, db^{[L]}$ reshape into big vector $d\theta$.

is $d\theta$ is grad of 'J'?

grad check:

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

for each i :

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

$$\approx d\theta^{[i]} = \frac{\partial J}{\partial \theta} \quad \left| \quad d\theta_{\text{approx}} \approx d\theta \right.$$

$$\frac{\|d\theta_{\text{oppose}} - d\theta\|_2}{\|d\theta_{\text{oppose}}\|_2 + \|d\theta\|_2}$$

$\epsilon = 10^{-7}$ ✓. correct initialisations.

10^{-5} : check

10^{-3} : worried.

implementing:

- only to debug:

$$d\theta_{\text{opp}}[i] \quad d\theta$$

- if algo fails grad-check, look at comp. to find bugs.

$$\begin{matrix} db^{[L]} & dw^{[L]} \\ \nwarrow & \swarrow \end{matrix}$$

- Remember regularisation

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

$d\theta = \text{grad of } J \text{ wrt } \theta$.

- Doesn't work with dropouts. keep prob = 1.0

- Run at random initialisation; perhaps after some training.

1% NN seems to have high bias, then

- Inc no of units in each hidden layer.
- Make NN deeper

0.5% TS error
0.5% dev set error

→ Inc λ

→ get more training data.

• λ inc \Rightarrow weights $\rightarrow 0$

• using keep-prob param does regularisation effect
+ dec train set error.

• Dropout, Data augmentation & λ regularisⁿ + dec variance.

• normalise 'x' so that cost funcⁿ faster to initialise.

⊕ Optimisation Algos:

* Mini-batch gradient descent:-

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots \end{bmatrix}$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots \end{bmatrix}$$

if $n = 50,000,000$?

div to mini-batches.

1000 - each.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(1000)} & \left| \begin{array}{c} x^{(1001)} \dots x^{(2000)} \\ \dots \end{array} \right| \dots \end{bmatrix}$$

$\underbrace{\hspace{10em}}_{X^{\{1\}} \quad (n_x, 1000)}$
 $\underbrace{\hspace{10em}}_{X^{\{2\}} \quad (n_x, 1000)}$
 \dots
 $X^{\{5000\}}$

$$Y = \begin{bmatrix} y^{\{1\}} & y^{\{2\}} & y^{\{3\}} & \dots \end{bmatrix}$$

$(1, 1000) \quad (1, 1000)$

mini batch t : $x^{\{t\}}, y^{\{t\}}$.

mini batch q :

for $t = 1, 5000$

FWD propagⁿ on $x^{\{t\}}$

$$Z^{\{t\}} = w^{\{t\}} x^{\{t\}} + b^{\{t\}}$$

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

1 step of GD using $x^{\{t\}}, y^{\{t\}}$

$$\text{cost } J = \frac{1}{1000} \sum_{i=1}^p L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|w^{[l]}\|_F^2$$

Backprop to compute gradients wrt $J^{(t)}$

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

$$b^{[l]} = b^{[l]} - \alpha db^{[l]}$$

}

"1 epoch"

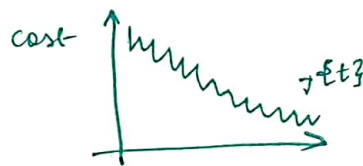
↳ pass through training set.

* MBGD understanding :

Batch GD.



MBGD



Plot $J^{(t)}$ using $x^{(t)}, y^{(t)}$.

Choosing size :

If minibatch size = m : Batch GD $(x^{(1)}, y^{(1)}) = (x, y)$

If minibatch size = 1 : Stochastic GD

every example is mini batch.

$(x^{(t)}, y^{(t)}) = (x^{(1)}, y^{(1)}) = \text{mini batch.}$



stochastic GD.

practice : 1 - m batch size

Batch GD

minibatch : m

↓
Too long.

Stochastic GD.

↓
lose sparsity from vectorization.

In between

↓
fastest learning
- Too much vectorization
- no waiting.

Small TS: Use batch GD.
($m \leq 2000$).

Typical mb size is 64, 128, ...
 $2^6 \quad 2^7 \quad 2^8$

29

make sure mini batch fit in CPU/GPU memory.
 $x^{(t)}, y^{(t)}$

power of 2 is faster.

* Exponentially weighted averages:

$$V_0 = 0$$

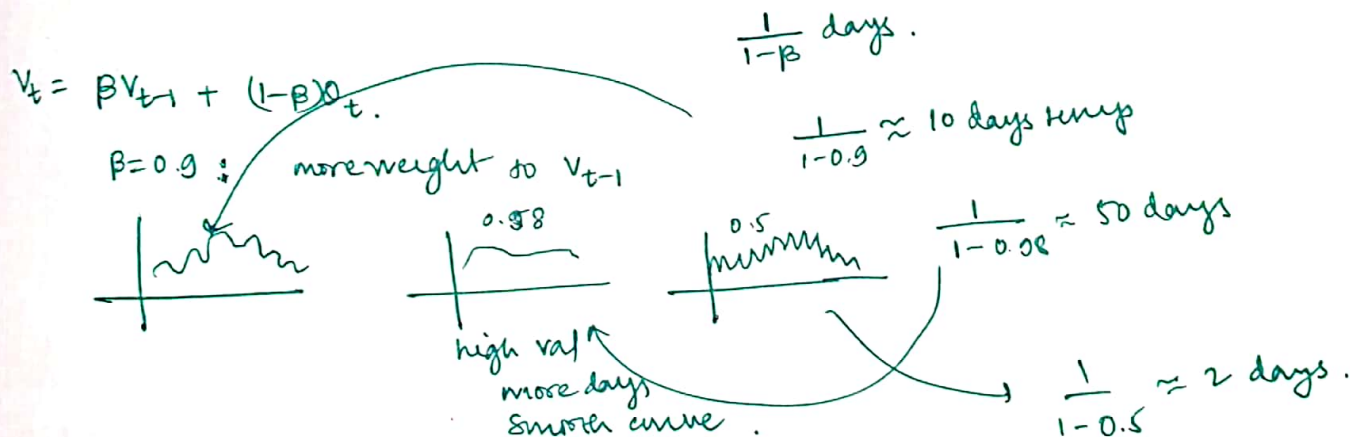
$$V_1 = 0.9V_0 + 0.1\theta_1$$

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

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

⋮

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



* Understanding exponentially weighted average.

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

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

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

⋮

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

$$= 0.1\theta_{100} + 0.1 \times 0.9\theta_{99} + (0.1 \times (0.9)^2)\theta_{98} + (0.1)(0.9)^3\theta_{97} + \dots$$



$$\approx \frac{1}{1-\beta} \text{ days}$$

$$\epsilon = \frac{1}{\beta} (1-\beta)$$

$$0.9^{10} \approx 0.35 \approx \frac{1}{e}$$

$$\frac{(1-\epsilon)^{1/\epsilon}}{0.9} = \frac{1}{e}$$

$$0.9^{50} \approx \frac{1}{e}$$

(for exponential decay)

implementation

$$V_0 = 0$$

$$V_1 = \beta V_0 + (1-\beta) \theta_1$$

$$V_2 = \beta V_1 + (1-\beta) \theta_2$$

$$V_3 = \beta V_2 + (1-\beta) \theta_3$$

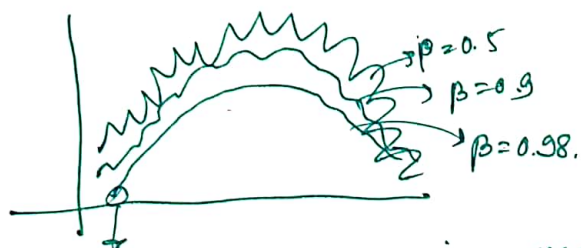
...

$$V_0 = 0$$

Repeat {
get next θ_t .

$$V_t := \beta V_t + (1-\beta) \theta_t \leftarrow$$

}



starts very low \therefore bias error.

$$V_1 = 0.9 V_0 + 0.02 \theta_1$$

$$V_0 = 0$$

$$V_1 = 0.02 \theta_1$$

$$V_2 = 0.98 V_1 + 0.02 \theta_2$$

$$V_2 = 0.0196 \theta_1 + 0.02 \theta_2$$

$$\frac{V_t}{1-\beta^t} \rightarrow \text{Bias correction.}$$

$t=2$:

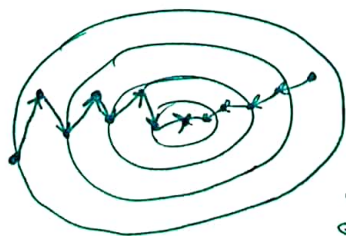
$$1-\beta^t = 1-(0.98)^2 = 0.0396$$

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

t large \rightarrow estimate corrected

* Gradient descent with momentum:

- works faster.



\downarrow slow learning
 \leftrightarrow faster learning.

momentum:
on iteration t :

compute dw & db on current mini batch.

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

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

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

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

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

vert dirⁿ: average out to zero.
 horizontal: quick osc. in horizontal dirⁿ.



A ball ~~refer~~ rolls down the bow

Hyperparam α, β . $\beta = 0.9$.

$$\text{avg} = \frac{1}{1-\beta}$$

Init:

$$V_{dw} = 0$$

$$V_{db} = 0$$

$$V_{dw} = \beta V_{dw} + dw \leftarrow$$

$$V_{db} = \beta V_{db} + db \leftarrow \text{leave out } \underline{(1-\beta)}$$

RMS Prop:

- not mean square prop.



on iteration t :

compute dw, db on current mini batch.

$$S_{dw} = \beta S_{dw} + (1-\beta)dw^2 \leftarrow \text{element wise. relatively small}$$

$$S_{db} = \beta S_{db} + (1-\beta)db^2 \leftarrow \text{relatively large.}$$

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

\rightarrow bigger $\rightarrow 10^{-8}$ mostly neglected.

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

\rightarrow (+E) $\rightarrow 10^{-8}$ mostly neglected.
 smaller.

after rms prop:



Adam optimization algorithm:

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

on iteration t :

compute dw, db using current MB.

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

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2)dw^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2)db^2 \leftarrow \text{RMS prop}$$

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

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

$$s_{dw}^{correct} = s_{dw} / (1 - \beta_2^t)$$

$$w := w - \alpha \frac{V_{dw}^{correct}}{\sqrt{s_{dw}^{correct}} + \epsilon}$$

$$s_{db}^{correct} = s_{db} / (1 - \beta_2^t)$$

$$b := b - \alpha \frac{V_{db}^{correct}}{\sqrt{s_{db}^{correct}} + \epsilon}$$

Hyperparameters choice:

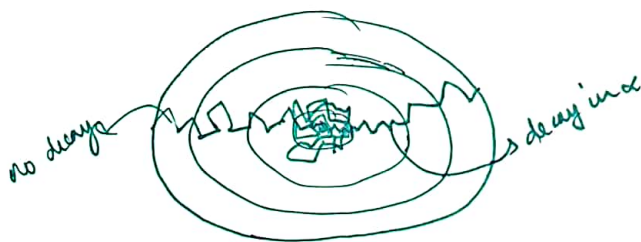
- α : needs to be tuned
- β_1 : 0.9 (dw)
- β_2 : 0.999 (dw²)
- ϵ : 10^{-8} .

Adam: Adaptive moment estimate

$\beta_1 \rightarrow 1^{st}$ moment

$\beta_2 \rightarrow 2^{nd}$ moment

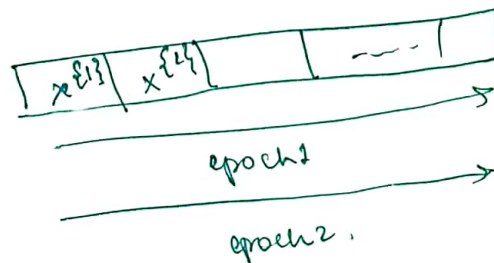
* learning rate decay:



1 epoch = 1 pass through data

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

Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4

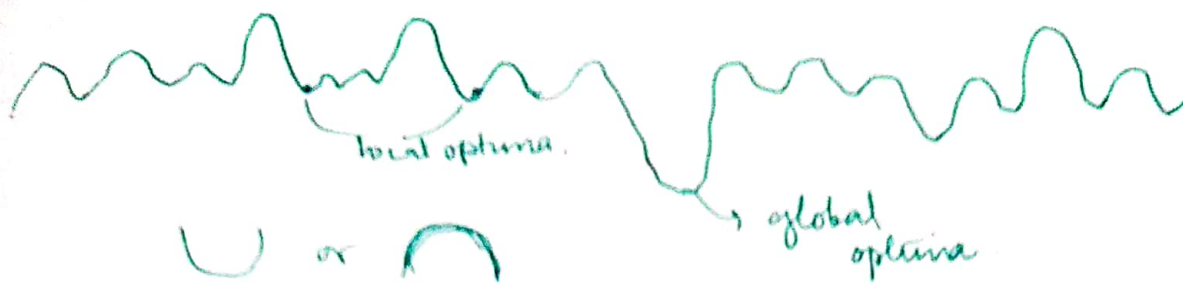


$$\alpha = 0.95^{\text{epoch number}} - \alpha_0 \Rightarrow \text{exp. decay}$$

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



the problem of local optima:



Saddle pt: pt where derivative = 0

problem of plateaus:

→ grad ≈ 0 for a long time

• Plateaus can make learning slow

Hyperparameter tuning :-

* Process :

$\alpha \rightarrow$ most imp H.P.

$\beta \sim 0.9$

$\beta_1, \beta_2, \epsilon \} \rightarrow$ Significant

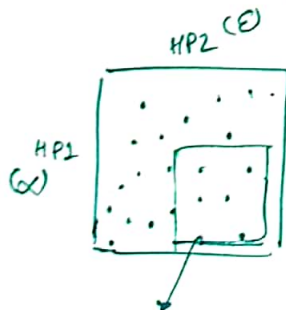
layers

hidden units $\} \rightarrow$ significant

learning rate decay

mini batch size $\} \rightarrow$ significant.

Try random val: no grid.

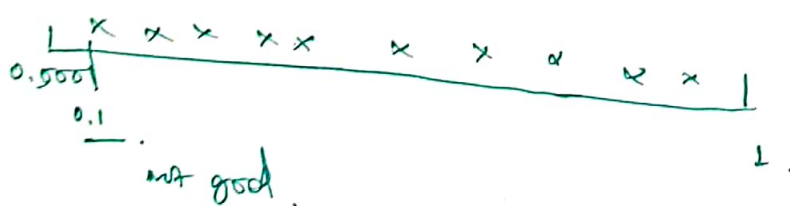


Coarse to fine searches

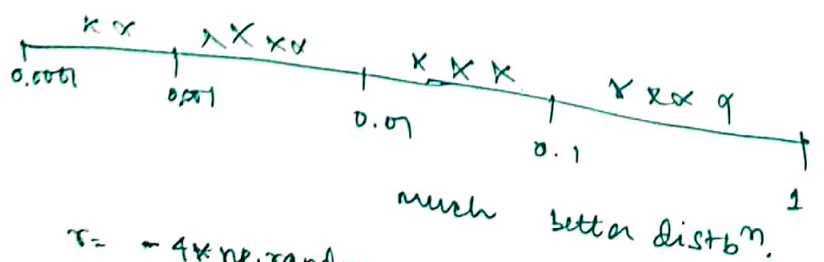
* Using appt scale to pick hyperparameters :-

- Picking HP at random

$\alpha = 0.0001, \dots, 1.$



(1)



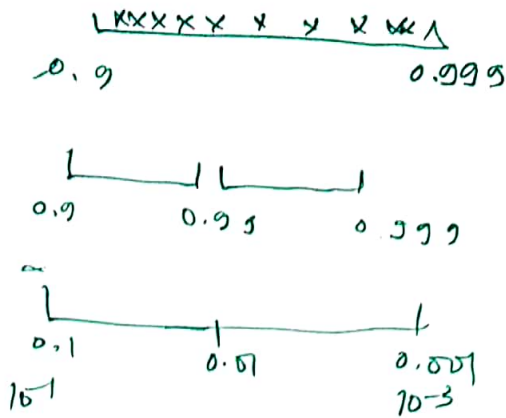
$r = -4 * np.random.rand()$
 $\alpha = 10^r$

$\leftarrow \sigma \in [-4, 0]$
 $\leftarrow 10^{-4} \sim 10^0$

$$\beta = 0.9 \quad \dots \quad 0.999$$

$$\frac{1}{10} \quad \downarrow \quad 1000$$

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

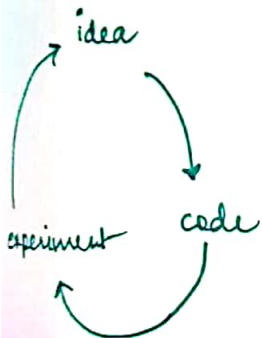


$$r \in [-3, +]$$

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

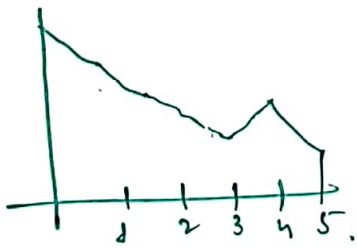
$$\beta = 1-10^r$$

Hyperparameters tuning in practice (Pandas vs Caviar)



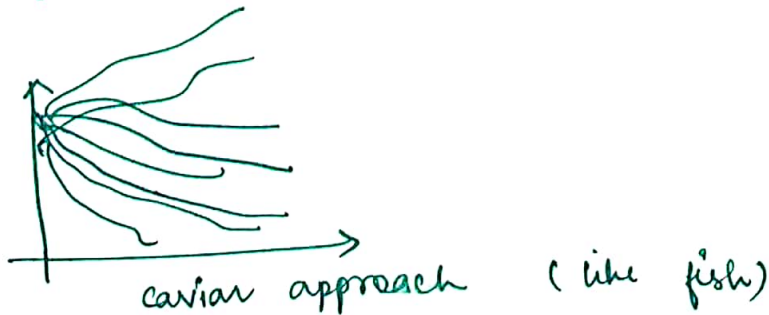
- NLP, Vision, Speech, ..
- Intuitions do get stale.

babysitting one model: monitoring learning rates one at a time.



pandas approach. (like panda)

Training many model in parallel.



⊕ Batch normalisation :

* Normalization of activations in a N/w.

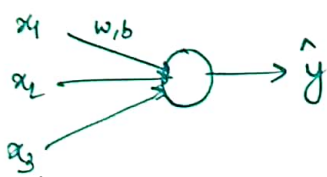
LR:

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

$$x = x - \mu$$

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

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



normalise $a^{[2]}$ so as to train $w^{[3]}, b^{[3]}$ faster.

implementing :

$$z^{[1]}, \dots, z^{[m]}$$

└──────────┘
 $z^{[1]}(i)$

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

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

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

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

γ, β : learnable params.

if

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

then

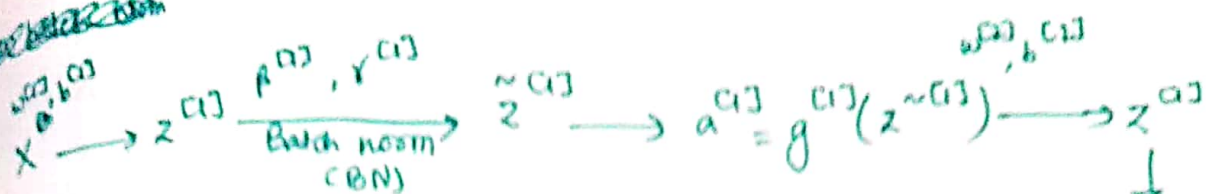
$$\beta = \mu$$

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

we $\tilde{z}^{(i)} = z^{(i)}$

Shifting Batch norm into a NN.

~~Batch norm~~



mean: $w^{(l)}, b^{(l)}, w^{(l)}, b^{(l)}, \dots, w^{(l)}, b^{(l)}$
 $\mu^{(l)}, \gamma^{(l)}, \mu^{(l)}, \gamma^{(l)}, \dots, \mu^{(l)}, \gamma^{(l)}$ } $d\mu^{(l)}$
 $\mu^{(l)} = \mu^{(l)} - \alpha d\mu^{(l)}$

$z^{(l)}$
 $(n^{(l)}, 1)$

$b^{(l)}$
 $(n^{(l)}, 1)$ $\mu^{(l)}$
 $(n^{(l)}, 1)$

implementing qd:

for $t = 1 \dots \text{no of MB}$:

Compute forwardprop $m \times \{t\}$

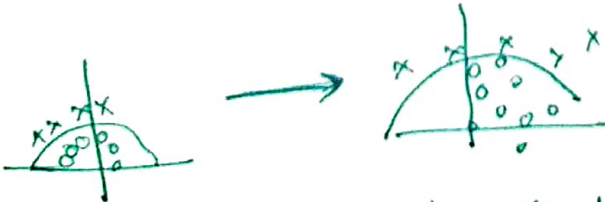
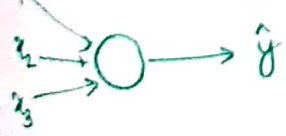
In each hidden layer use BN to replace $z^{(l)}$ with $\tilde{z}^{(l)}$
 use backprop to compute $d\mu^{(l)}, d\gamma^{(l)}, d\mu^{(l)}, d\gamma^{(l)}$

update params
 $w^{(l)} = w^{(l)} - \alpha dw^{(l)}$
 $\mu^{(l)} = \mu^{(l)} - \alpha d\mu^{(l)}$
 $\gamma^{(l)} = \gamma^{(l)} - \alpha d\gamma^{(l)}$

works in momentum, RMS prop, Adam.

Why Batch norm works?

Learning on shifting i/p distribution



data distribution changes through covariate shift
 if distbⁿ changes from $x \rightarrow y$, retrain algorithm

- batch norm \rightarrow the amt by which hidden values shift. Thus, $z_2^{[L]}$ & $z_2^{[L]}$ remain same.
- batch norm makes sure variance of $z_2^{[L]}$ & $z_2^{[L]}$ remain same.
- makes learning of later layers easier
- Also has a slight regularisation effect. (adds noise to hidden layers) \rightarrow unintended side effect.
- After training NN with batch norm, perform needed norm. using μ & σ^2 , estimate using exponentially weighted avg across minibatches at test time.

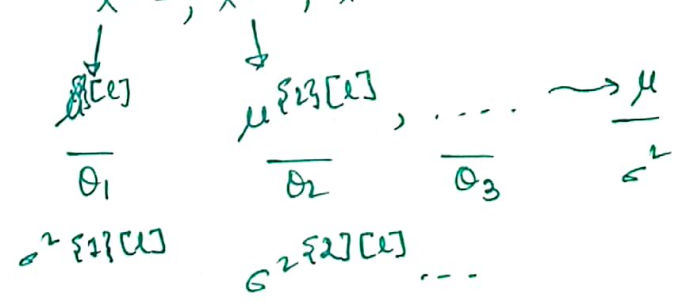
* Batch norm at test time:

\rightarrow minibatch at a time

but at ~~test~~ test time, we need to process examples one at a time.

\rightarrow estimate μ & σ^2 from a T.S.

μ, σ^2 : estimate using exponentially weighted average (across MB)
 $x^{[1]}, x^{[2]}, x^{[3]}$



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

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

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

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

* ϵ used to avoid division by zero.

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

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

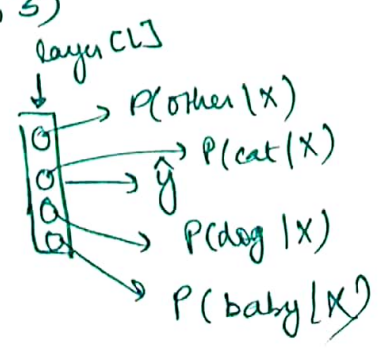
⊕ Multiclass classificⁿ:

* Softmax Regression

Ex:

$C = \# \text{ no of classes} = 4 \quad (0, \dots, 3)$

$X \rightarrow \dots$



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

Activⁿ funcⁿ

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

(4,1)

$$y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

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

$$z^{[L]}_{(4,1)} \rightarrow a^{[L]} = \hat{y} \rightarrow L(\hat{y}, y)$$

Backprop:

$$dz^{[L]}_{(4,1)} = \hat{y} - y$$

\downarrow
 $\frac{dJ}{dz^{[L]}}$

⊕ DL prog. frameworks :

* Tensorflow :

ex: $J = w^2 - 10w + 25 = (w-5)^2$
min at $w = 5$.

import tensorflow as tf

$w = tf.Variable(0, dtype=tf.float32)$

$cost = tf.add(tf.add(w**2, tf.multiply(-10, w)), 25)$

$train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)$

$init = tf.global_variables_initializer()$

$session = tf.Session()$

$session.run(init)$

$sess.run(w)$

$session.run(train)$ # runs 1 step of GD

$print(session.run(w))$ # 0.1

for i in range(1000):

$session.run(train)$

$print(session.run(w))$ # 4.999 → working.

coefficients = ...
 $x = \text{tf.placeholder}(\text{tf.float32}, [3, 1])$

$$\text{cost} = x[0][0] * w^{**2} + x[1][0] * w + x[2][0]$$

session.run(train, feed_dict={x: coefficients})

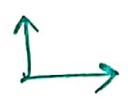
Structuring ML Projects:

Introduction to ML strategy.

* orthogonalisation

- what to tune in order to achieve one effect, people are clear about this.

tune ~~measure~~ different parameters differently

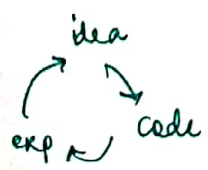


chain of assumptions in ML:

different
tuners
for
different
assumptions
independent
of each other
orthogonalisation.

- Fit training set well on cost funcⁿ.
 - bigger N/w
 - Better opt. algo
- ↓
- Fit dev set well on cost funcⁿ.
 - Regularization
 - bigger training set
- ↓
- Fit test set well on cost funcⁿ.
 - bigger dev set
- ↓
- Performs well in real world.
 - change dev set or cost funcⁿ.

* single number evaluation metric.



	precision	recall
A	95%	90%
B	98%	85%

→ 3 examples as cats, what % are cats
 what % of actual cats recognised.

F score = "Avg" of precision P & recall R

$$\left(\frac{2}{\frac{1}{P} + \frac{1}{R}} \right)$$
 "harmonic mean"