

$$w^{[l]} = [n^{[l]}, n^{[l+1]}]$$

$$b^{[l]} = [n^{[l]}, 1]$$

Deep Neural Networks Representations

Forward propagation:-

Input $a^{[L-1]}$

output $a^{[l]}$, cache ($z^{[l]}$)

compute intermediate
use to store ~~result~~ of
one layer in back
propagation

Back propagation

Input $da^{[l]}$

output $da^{[l-1]}, dw^{[l]}, db^{[l]}$

$$\rightarrow dz^{[l]} = da^{[l]} * g^{[l]}(z)$$

$$\rightarrow dw^{[l]} = dz^{[l]} \cdot a^{[l-1]T}$$

$$\rightarrow db^{[l]} = dz^{[l]}$$

$$\rightarrow da^{[l-1]} = w^{[l]T} \cdot dz^{[l]}$$

$$\rightarrow dz^{[l]} = w^{[l+1]T} \cdot dz^{[l+1]} * g^{[l]}(z^{[l]})$$

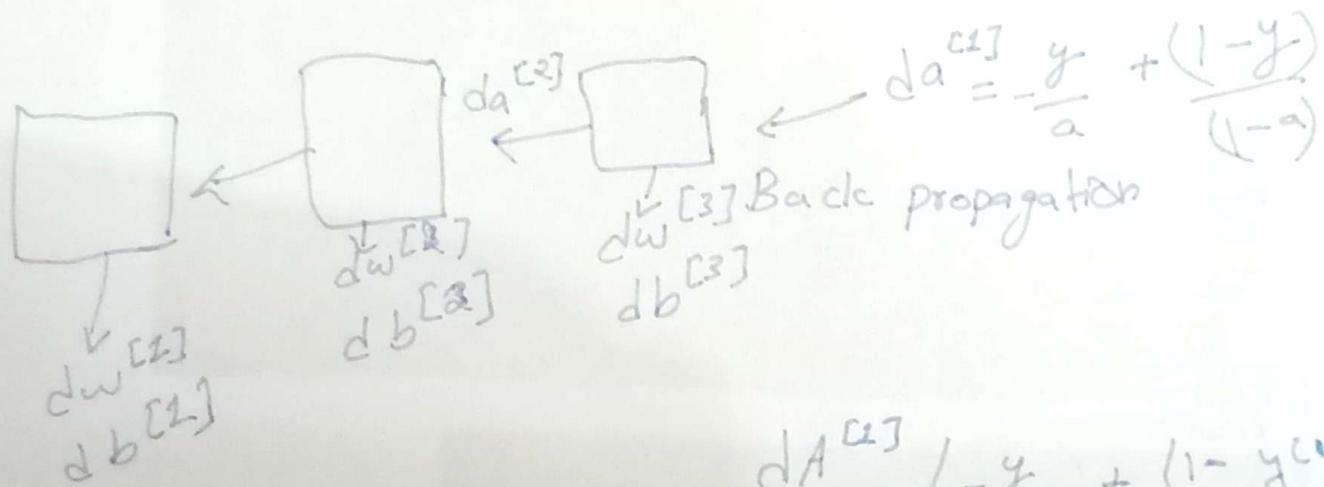
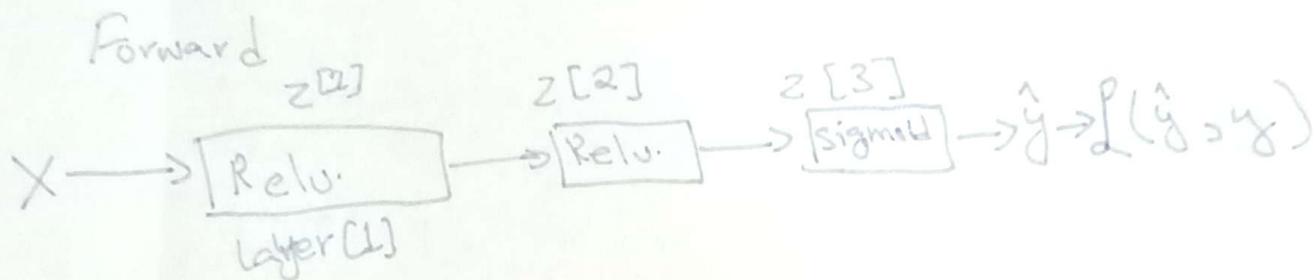
$$dz^{[l]} = dA^{[l-1]} \star g^{[l]}(z^{[l]})$$

$$dw^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{keepdims=True})$$

$$dA^{[l+1]} = w^{[l]T} \cdot dz^{[l]}$$

Summary:



$$dA^{[2]} = \left(\frac{-y}{a^{[2]}} + \frac{(1-y^{[1]})}{1-a^{[2]}} \right)$$

~~Hyperparameters~~

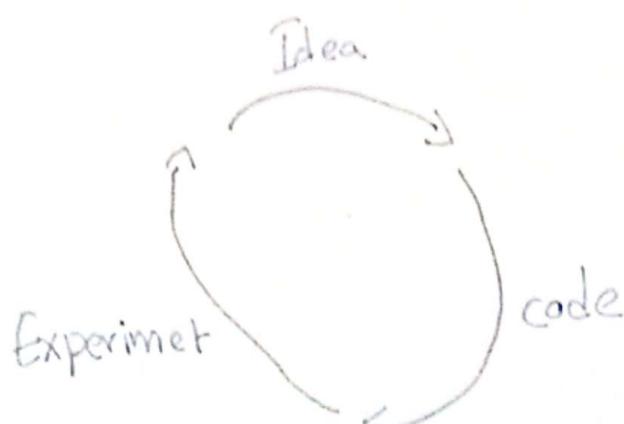
$w^{[2]}, b^{[2]}, w^{[3]}, b^{[3]}, w^{[4]}, b^{[4]}, \dots$

Hyperparameters: Defines what u want to train net. This includes

- Learning rate & iterations.
- Iterations
- Hidden layer l
- Hidden units (n_1, n_2, n_3)
- Choice of activation functions (Relu, Tan, Sigmoid)

Many other

: Momentum Term, mini batch size, and other are parameters in



to get best model we have to change hyperparameters or try different hyperparameters like

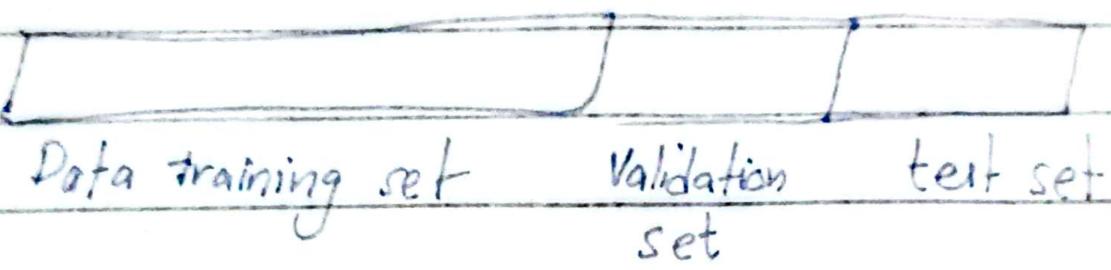
LR, α , Relu, tan, trying activation to get the best model

Deep learning by Andrew Ng CA

Train / Dev / Test sets Idea

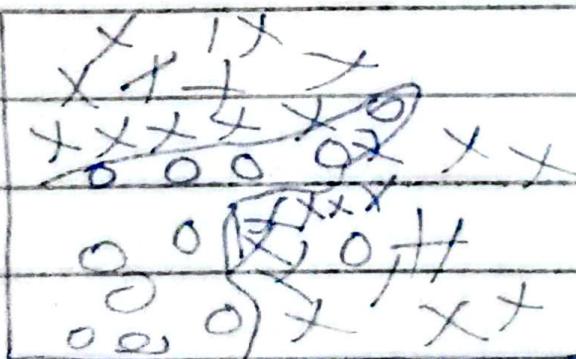
Experiment) code

ML is iterative model.



Bias and Variance trade off

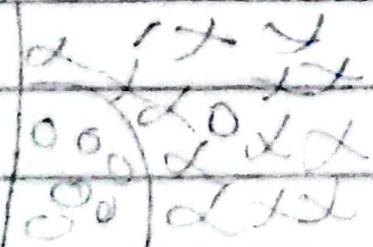
Sq bias = $\frac{1}{m} \sum \hat{y} - y$, Sq error absolute mean err



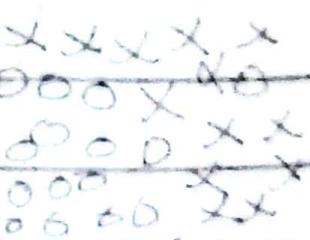
high bias

overfitting

overfit



overfit



Just right

high variance

Bias and Variance Cat classification

$y=1$

$y=0$



Train set error: 1% of error 15% error

Dev set error: 11% of error 16% error

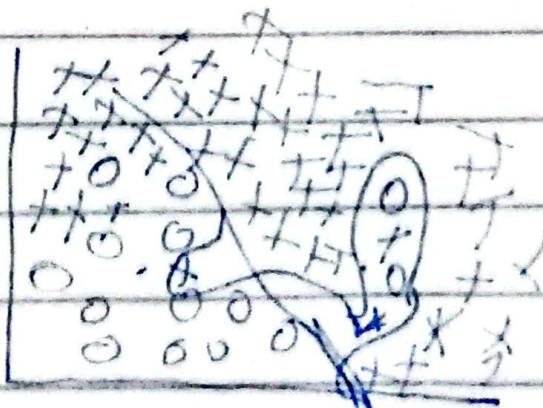
high variance error

underfitting error

high bias problem

Train set error	15 %	0.5 %
Dev set error	30 %	1 %
high bias	low bias	problem
high variance	low variance	problem

0% optimal error:-



high bias, high variance

topic) Basic recipe for ml

High bias? \rightarrow Bigger network
(training data performance) try train longer.

If we have to check high variance or bias problem first we have to check

High bias to reduce the error. to get rid of High bias problem.

Now for High variance

High Variance → More data
(data set performance) Regularization
Sometimes (NN architecture)

By C. GPT

Symptoms false

Bias : Error due to simple model
 that can't capture patterns well (training set) High training and dev error Underfitting

Variance:

- Errors due to overly complex low training but Derivative model that fits noise in training high dev/valid data error

Step 2:

Observation	Interpretation	Problem Type
High training error	can't fit training data well	High Bias (Underfitting)
Low training error → High test error	Model fits training but fails generalization	High Variance (Overfitting)
Low training and test error	Model is doing well	Good fit

If High bias (Underfitting)

Use complex model / add layers / neurons in NN
(deeper tree in decision tree)

+ Train longer / better optimization.

+ Add more relevant feature in training

- Decrease regularization (L1 and L2)

High Variance

- Use simple model (reduce layers, limit tree depth)
- Reduce number of x
- + Use more training data
- Use cross-validation to tune parameters

Data Cleaning.



Regularization.

When model learns very well from training data, including noise or irrelevant patterns it perform poorly on unseen data, this is called overfitting.

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

Frobenius norm $\|w\|_2^2 = \sum_{j=1}^{n_k} \sum_{i=1}^{n_{k-1}} w_{i,j}^2$

L_2 regularization is more often use than regularization

$$\frac{1}{2m} \sum (w_j) \sum_{j=1}^n |w_j| = \frac{1}{2m} \|w\|_1$$

where λ is a parameter.

Neural Network

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

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

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

L2 regularization

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

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

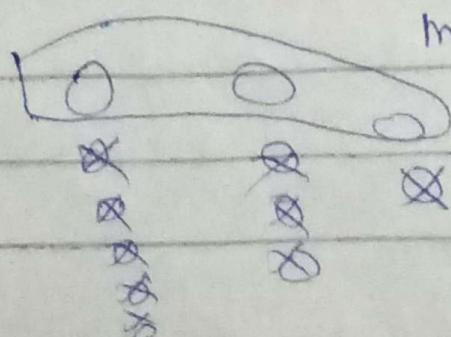
$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}$$

Why regularization reduce overfitting

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

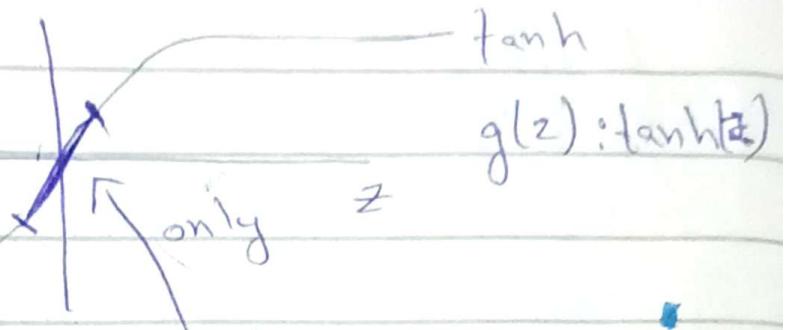
high variance or bias

In regularization $w^{[l]} = 0$ cause 0 in neural networks



make nn simple

cause high variance



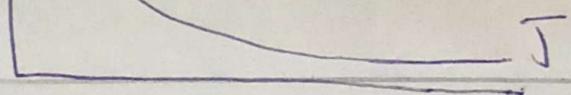
$$\lambda \hat{w}^l, \downarrow w^l, z^{(l)} = w^{(l)}_0 z^{(l-1)} + b^{(l)}$$

when $w^{(l)}$ decreases, $z^{(l)}$ also increases causing only linear performance which cause not a good performance

i.e Every layer \approx linear.

$$J(L) = \sum_i L(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2} \sum_l \|w^{(l)}\|_F^2$$

$J_{(l)}$ for debugging



Dropout Regularization Technique

It causes in the drop of neural network cause result in more simple and generalize NN causing in good performance as it learns patterns not memorize everything.

How regularization happens:-

In cause the complex neural network into simple or general which help network to learn general or simple pattern by

- dropping some nn connections

Implement dropout.

Layer 3

$d_3 = \text{np.random}(a_3, \text{shape}(0), a_3, \text{shape})$

→ keep_prob = 0.8

0.8 probability of
dropping a neuron

0.2 percent to eliminated

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

If anyone has value zero remove it from network.

probability $a_3 \geq \text{keep_prob}$ (if $NN=50$)

0.8 this cause 10 drop of NN from network

How regularization happens:-

In case the complex neural network into simple or general which help network to learn general or simple pattern by

- dropping some nn connections

Implement dropout -

Layer 3

$d_3 = \text{np.random}(a_3, \text{shape}(a_3), \text{shape}(d))$

→ keep_prob = 0.8
0.8 probability of dropping a neuron

0.2 percent to elim

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

? if ~~zero~~ anyone has value - zero remove it from network.

probability $a_3) = \text{keep_prob}$ (if $NN=50$
0.8 this cause 10 drop of NN from network

$$w^{[3]} = 3 \times 7$$

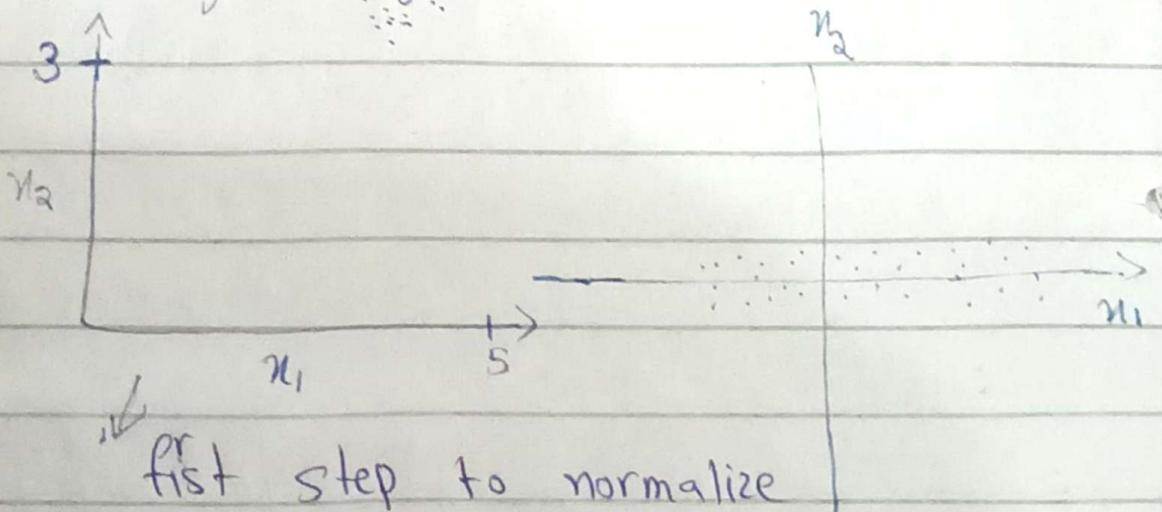
~~Redundant~~ Understanding Propert-

Data augmentation

Suppose we have one

Photo of cat we can use it by tilting, decrease
increase brightness, by zooming, and we
can train these photo to train our models.

Normalizing training sets



Subtract mean

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

$$X := X - \bar{X}$$

2nd step

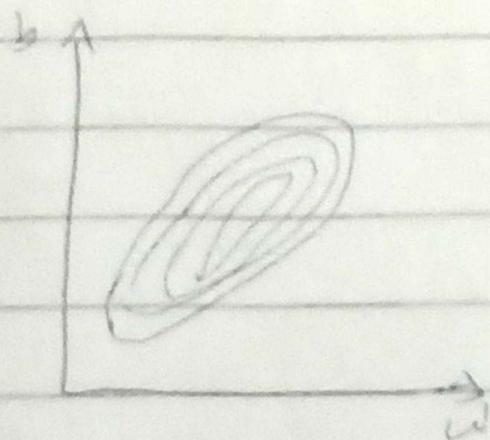
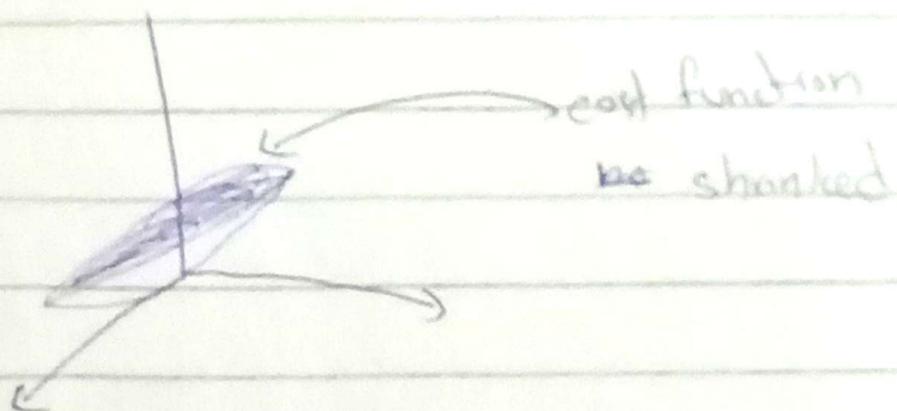
Normal variance

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

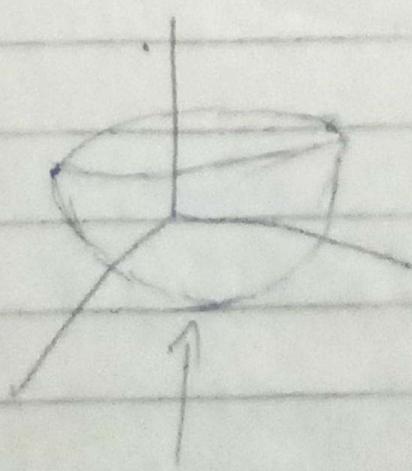
$$\bar{X} = 6$$

cost function (no normalization) - It tells him wrong system is

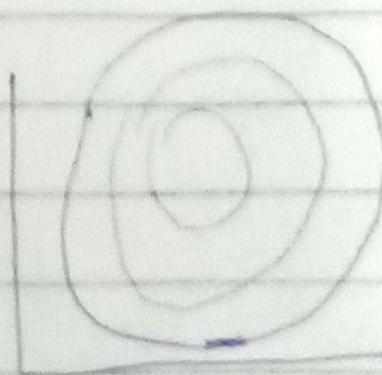
Why we use normalize inputs
unnormalized:



normalize data



cost function



cost function

Vanishing - Exploding Gradients -

When neural network increases the activation function decrease like

$$w = \begin{bmatrix} 0.5 & & \\ & 1.5 & 0 \\ 0 & & 0.5 \end{bmatrix}$$

$$y^{[l]} = w^{[l]} w^{[l-1]} w^{[l-2]} \dots$$

Weight initialize for Deep learning -

$$z = w_1 n_1 + w_2 n_2 + \dots + w_m n_m + b$$

Bigger n is we want smaller w , will be. One way is to

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

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

$$\text{Relu} \quad g^{[l]}(z) = \text{Relu}(z)$$

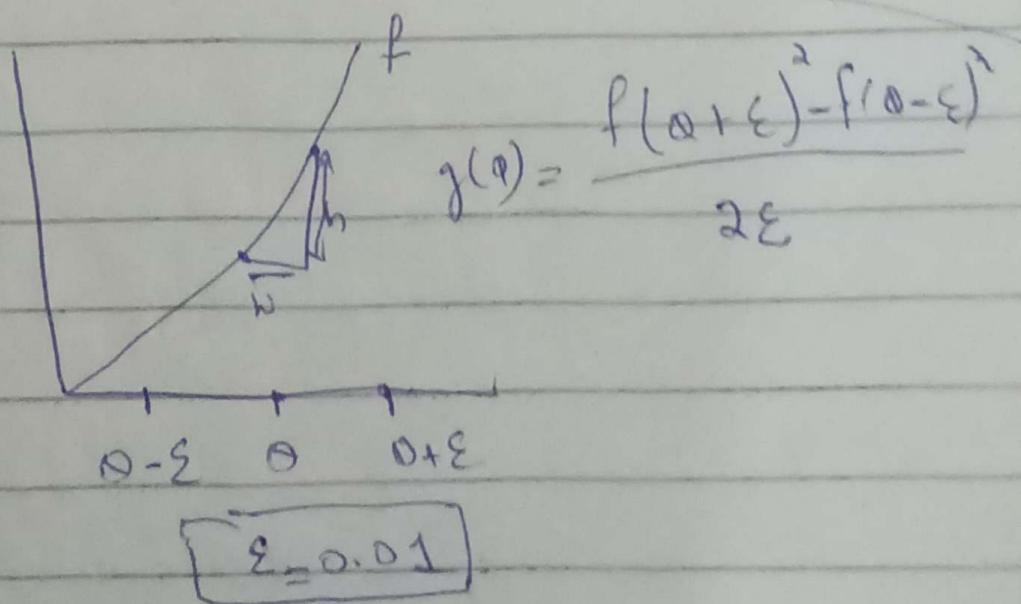
other variance

$$1) \text{ tan } h = \sqrt{\frac{1}{n(l-1)}}$$

a) Xavier initiation

$$\sqrt{\frac{2}{n(l-1)} + n[l]}$$

They are all used to decrease w_i cost function

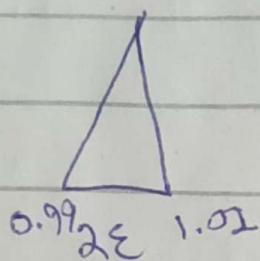


to

$$\underline{f(a+\epsilon) - f(a-\epsilon)}$$

2ϵ

it is more accurate formula for derivative.



We can use this value
 $\boxed{\epsilon=0.01}$ for own

but will be small

$\xrightarrow[3]{(1.01)} \text{always built in library/creator}$

$$\underline{(1.01) - f(0.99)}$$

$$\frac{2(0.01)}{= 3.001}$$

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

Gradient checking (apply)

for each i :

$$d\theta_i [i] = \frac{\overbrace{J(0, \theta_0, \dots, \theta_{i-1} + \epsilon, \dots)}^r - \overbrace{J(0, \theta_0, \dots, \theta_{i-1} - \epsilon)}^r}{2\epsilon}$$

this is used to check errors in back propagation

$$\text{do}([]) = \frac{25}{20}$$

Check the loop if

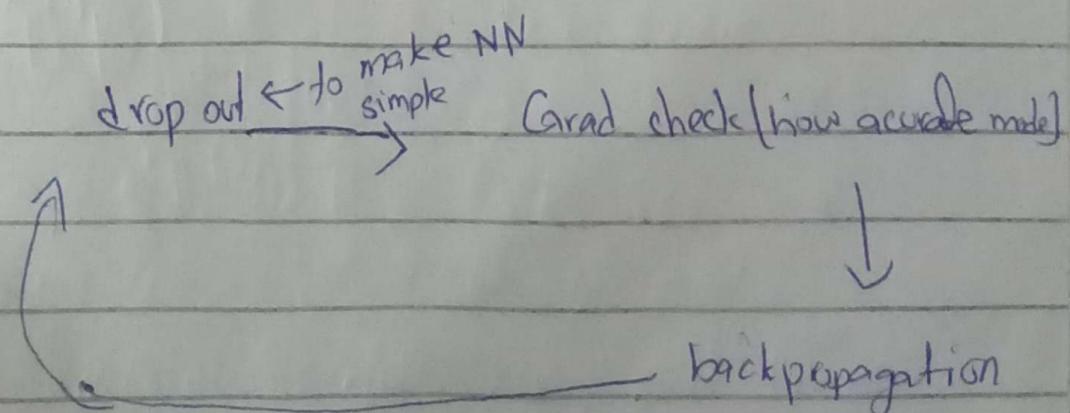
answer 10^{-7} u did a good job.

if 10^{-5} u have to look into model but ok

10^{-3} something is in error

- Don't use in training only to debug.
- If algorithm fails grad, check at component
- Remember regularization.

* We cannot use drop out with Grad check
we vs in this sequence.



Note: Grad check is only use we we have
to check back propagation accuracy at once.

Week 2 (Use case: backpropagation + optimization algorithms)

Optimization algorithms

if we have 500000 examples to train

Batch

- ⇒ Use all training gradient data
- ⇒ Slower, takes high memory.
- ⇒ More accurate gradients.

Mini Batch

for $t=1..5000$

forward prop on $X^{[t]}$

$$Z^{[1]} = W^{[1]} X^{[1]} + b^{[1]}$$

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

compute cost $J = \frac{1}{1000} \sum_{i=1}^n \ell(y^{(i)}, \hat{y}^{(i)})$

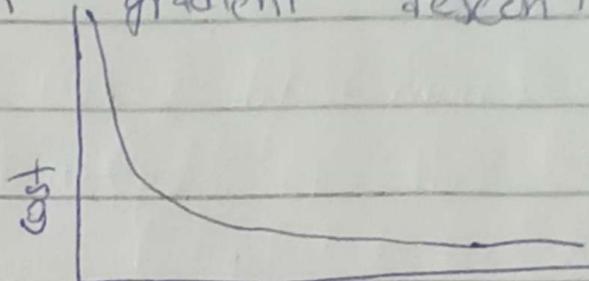
$$\frac{1}{2 \cdot 1000} \sum_k \|W^{[k]}\|_F^2$$

convert large data set into small batches result in efficiency, faster working or training

Mini-batch and batch model

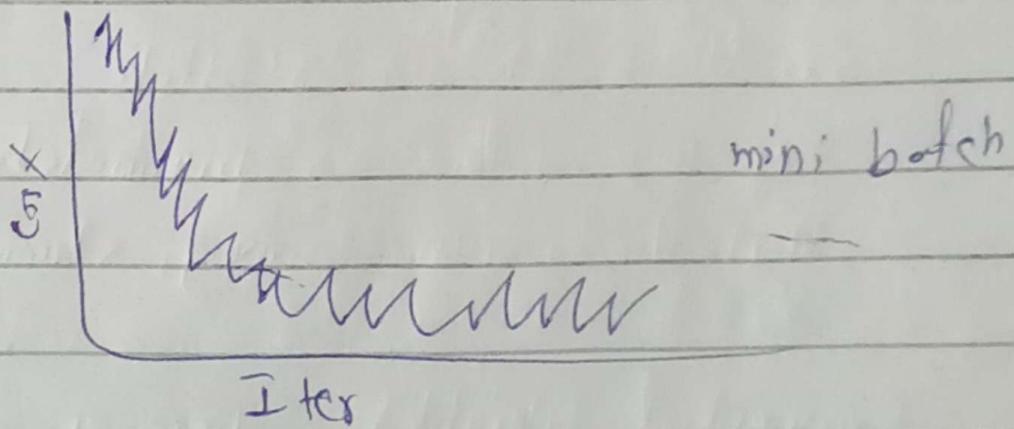
Create small batches of 32, 64, 128 samples

Batch gradient descent



iterations

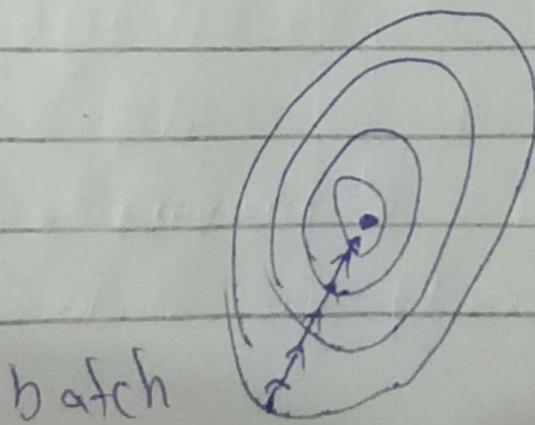
Iteration ↑ cost Increase



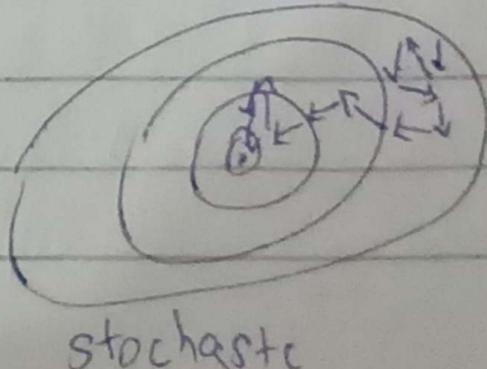
If minibatch size = 1 -

stochastic gradient descent

If minibatch size = m = Batch gradient descent.



batch



stochastic

Stochastic gradient descent (Use for small, faster
streaming, data)

lose speed up from vectorization

(mini-batch size
not to by (small))

↓
faster learning

• Vectoring (1000)

small ~~batch~~^{data} size \Rightarrow batch

very large data \Rightarrow minibatch

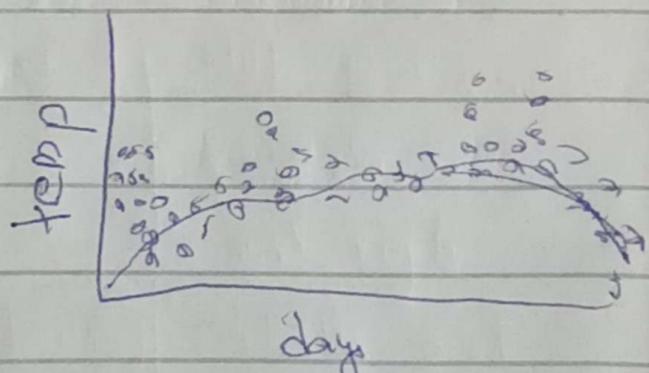
Stochastic \Rightarrow online learning, streaming, IoTs

Exponential Weighted averages

$$\bar{x} = 40^{\circ}\text{F } 4^{\circ}\text{C}$$

$$\bar{x}_2 = 49^{\circ}\text{F } 9^{\circ}\text{C}$$

$$\bar{x}_{180} = 60^{\circ}\text{F}$$



They create line & moving example
of daily data (any interval) like stockmarket

Key components EWA

$$v_t = B v_{t-1} + (1-B) \varphi_t \quad \because \varphi_t = \varphi - \alpha^* \sqrt{t}$$

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

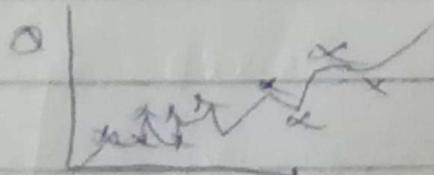
$$v_{100} = 0.1 \varphi_{100} + 0.1 \varphi_{99} + 0.9$$

$$B = 0.9, \varphi_0 = 0$$

$$\text{raw gradient } \nabla \varphi = [4, 3, 2, 1] \\ = 0.9 \times 0 + (1-0.9) 4$$

$$v_{t+1} = (0.1) 4 = 0.4$$

$$v_{t+2} = 0.9 \times 0.4 + 0.1 \times 3 \\ = 0.66$$



Bias correction in this case:-

$$t=2, 1-B^t = 1-(0.98)^2 = 0.0396$$

$$\underline{v_2} >$$

$$0.0396$$

Bias correction

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

$$\eta = [10, 10, 20, 10]$$

$$v_0 = 0$$

$$\beta = 0.9$$

with out bias correction we get

$$1) 0.9 \cdot 0 + 0.1 \cdot 10 = 1$$

$$2) 0.9 \cdot 1 + 1 \cdot 0 = 1.9$$

$$= 2.9 \Rightarrow 4.9$$

this cause increase only

Now with Bias correction

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

$$v_t = v_t (1 - \beta^t)$$

$$v_t = 1 \quad 1 - 0.9^4 = 0.1$$

$$= \frac{1}{0.1} = 10.0$$

Concept	Without Bias	With Bias correction
v_1	1.0	10.0
Variate scaling	Too large	Just right
Real effect	Unstable learning	Stable and accurate

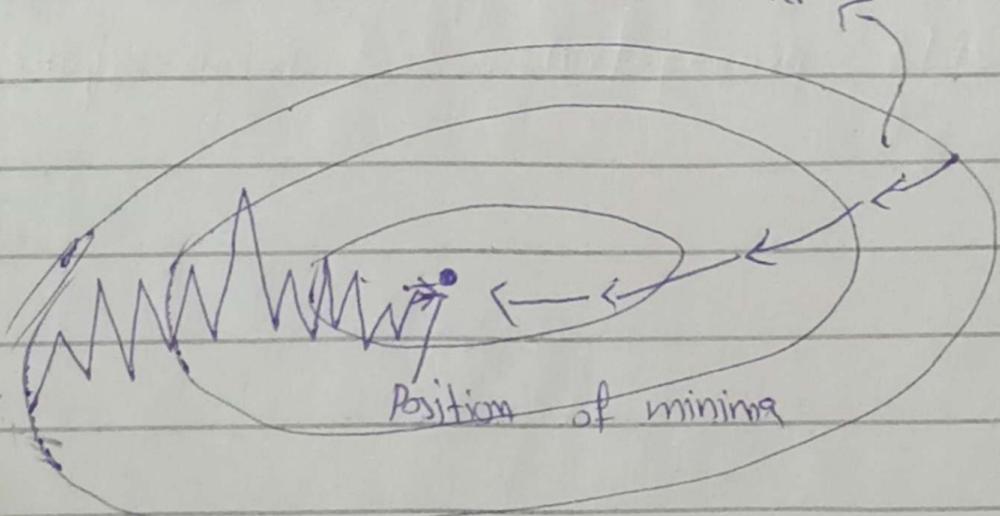
it corrects the initial value of v_f as
 v_f initial may not be zero.

Reasons-

We cannot believe every data start with zero-

- * Bias don't change initial value but it corrects or impacts the early steps.

Gradient descent with momentum



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

$$v_b = \beta v_b + (1-\beta)\partial_c$$

this cause in increasing speed

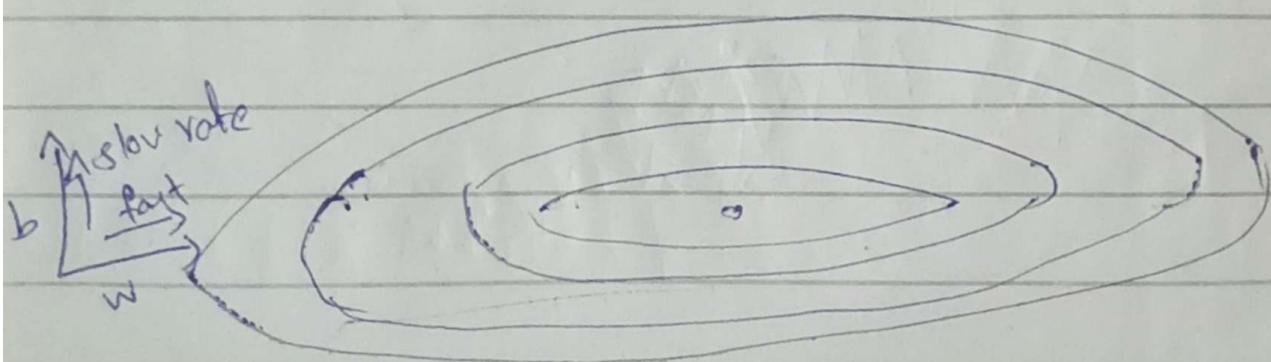
GD Momentum

- Used for speeding convergence
- Reduce oscillation

GD FWA

- Smoothing noisy data
- Detecting trends
- e.g.
SGD, Adams,

RMS propagation: Root mean square propagation



Compute dw, db on count mini-batch

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

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

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

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

$$\sqrt{S_{dw}}$$

$$\sqrt{S_{db}}$$

In next step, we will use RMS with Momentum

Adam Optimization

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

This algorithm takes GD Momentum and put it with RMS

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) d_w, V_{db} = \beta_1 V_{db} + (1 - \beta_1) d_b$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) d_w^2, S_{db} = \beta_2 S_{db} + (1 - \beta_2) d_b^2$$

Bias correcting-

$$V_{dw} = V_{dw} / (1 - \beta_1^t), V_{db} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw} = S_{dw} / (1 - \beta_2^t), S_{db} = S_{db} / (1 - \beta_2^t)$$

Hyperparameter choice

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999 \quad (\text{dw}) \cdot \text{default adam}$$

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

Adam: Adaptive moment estimation

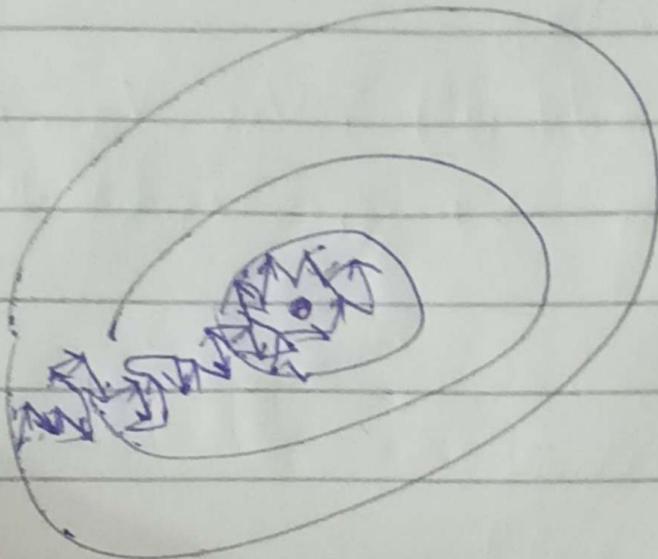
Clarification about learning rate decay.

$$\alpha = \frac{1}{\alpha_0}$$

It decayRate \times epoch num.

epoch = one cycle/or 1 pass training cycle

Learning rate decay:-



Suppose we are using mini-batch, due to small value (α) it will not reach minima, instead it will rotate around minima.

$$\alpha = \frac{1}{1 + \text{decayRate} \times \text{epoch}}$$

Epoch

1

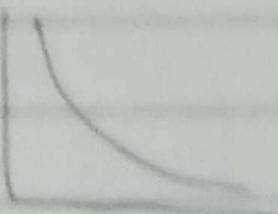
.

α

0.1

$\alpha = 0.01$

decay $\alpha \times 1$



Learning rate
decrease

other learning rate decay.

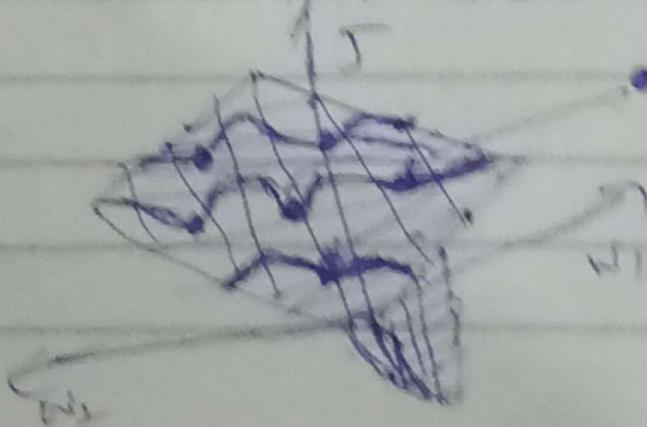
$$\alpha = 0.95^{\text{epoch_num}}$$

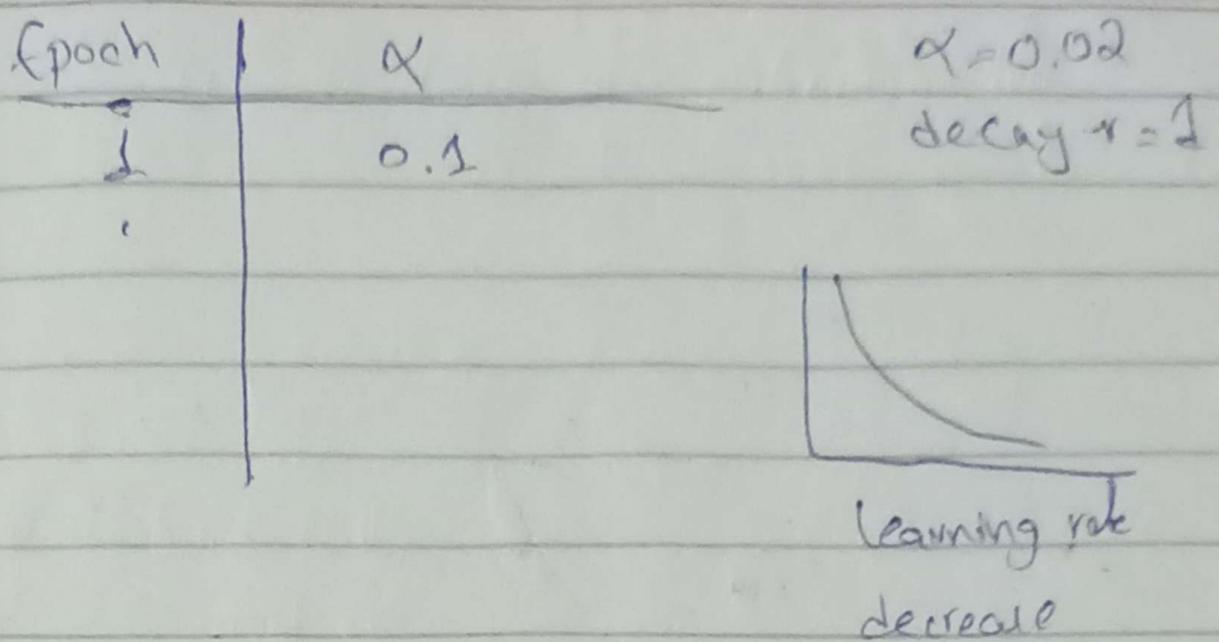
$$\alpha = \frac{1}{\sqrt{\text{epoch_num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \cdot \alpha_0$$

Manual decay

Problem with local optima

• local optima





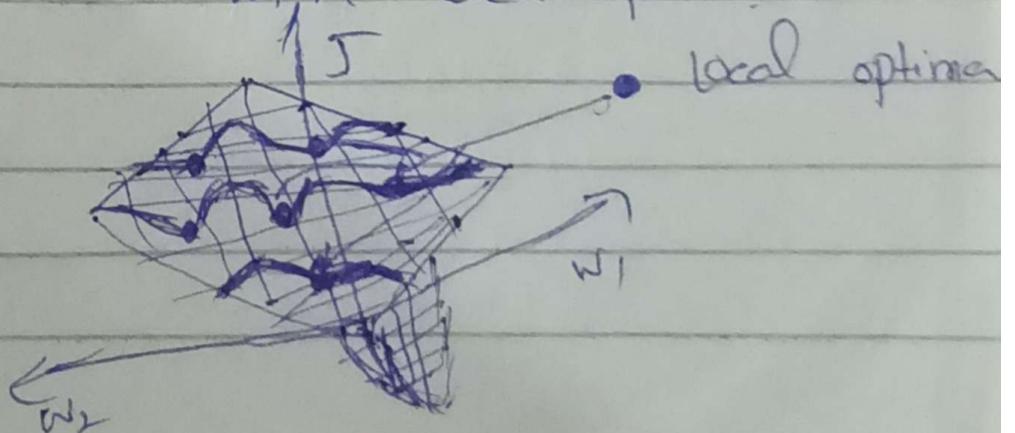
other learning rate decay.

$$\alpha = 0.95^{\text{epoch_num}}$$

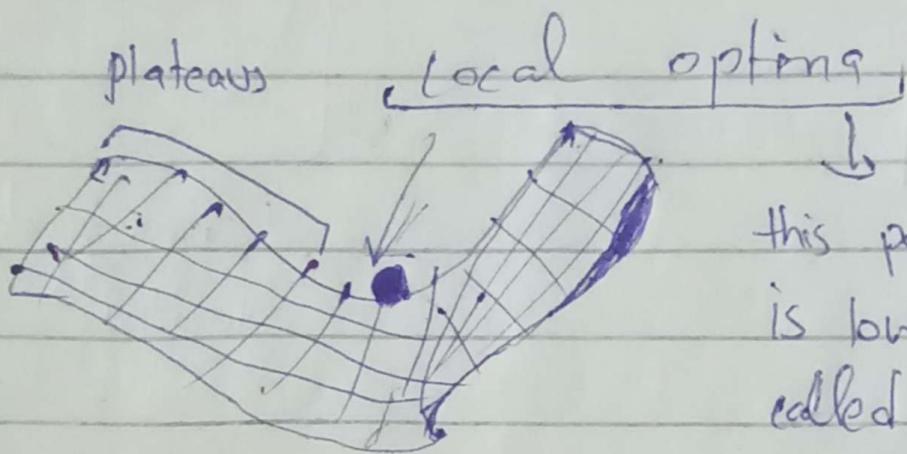
$$\alpha = \frac{1}{\sqrt{\text{epoch num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \cdot \alpha_0$$

Manual decay.

Problem with local optima.



but according to graph it is
not always true it is



local optima
this point loss
is low
called saddle
point

saddle ~~in 3d~~

⇒ plateau cause in slow learning

WEEK THREE : 3

Hyperparameter tuning,

$\beta = 0.9$ α , B_1, B_2, ϵ ,
most important hyperparameter

layers

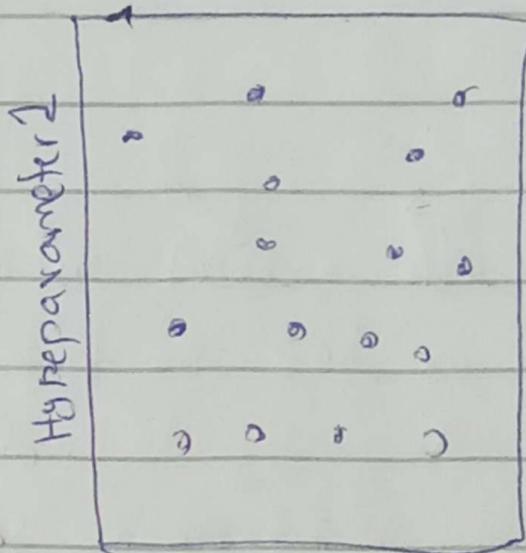
1 hidden unit and unit import

learning rate decay 3rd import
mini batch size

How can make it in work

Hyperparameter 2

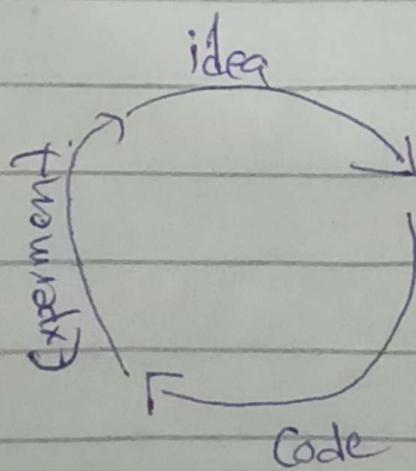
Search Best hyper
parameters



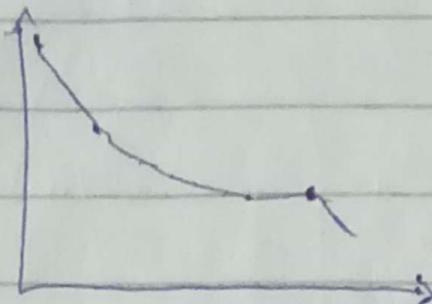
Grid Box

Using appropriate scale to pick hyperparameters
Layers (2-4)

Tips to choose hyperparameters.



Babysitting on model



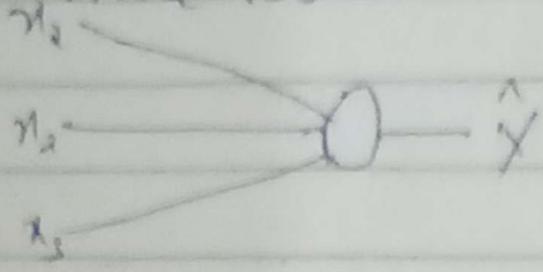
We can train ~~the~~ a model on different hyperparameters if we have no large capacity we can do it day by day, on each we test ~~one~~ hyperparameter specific hyperparameter known as panda approach one baby.

Caviar approach (many babies at a time)

We can try many hyperparameters if we have efficient resources -

Batch Normalization. ~~$\times = \gamma/\sigma$~~ $\times = \gamma/\sigma$

parameters work well and gradients make robust.



$$y = \frac{1}{m} \sum r^{(i)}$$

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

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

$$x = x / \sigma$$

Can we normalize $a^{[2]}$ so that we can train $w^{[3]}, b^{[2]} \dots$ faster. Yet this is the way

Standardization: Model perform better

Normalization: If you don't normalize algorithm KNN, SVM, NN will give more importance with features with large range

Implementing - Batch Norm

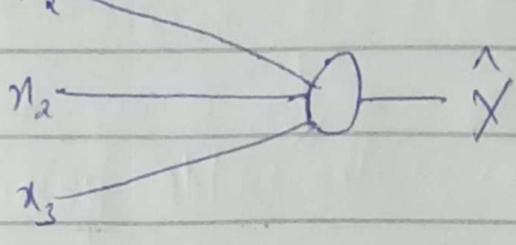
Given some intermediate values in NN

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

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

Batch Normalization ~~\rightarrow~~ $x = \bar{x}/\sigma$

\Rightarrow parameters work well and parameters make robust.



$$\bar{y} = \frac{1}{m} \sum n^{(i)}$$

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

$$\text{variance} = \sigma^2 = \frac{1}{m} \sum (x^{(i)} - \bar{y})^2$$

$$x = x'/\sigma$$

Can we normalize $a^{[2]}$ so that we can train $w^{[3]}, b^{[2]}$ faster. Yes this is the way

Standardization: Model perform better

Normalizations: If you don't normalize algorithm KNN, SVM, NN will give more importance with features with large range

Implementing - Batch Norm

Given some intermediate values in NN

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

$y = \text{mean}$

$$z^{(i)} = \frac{\bar{z} - y}{\sqrt{\sigma^2 + \epsilon}} \quad \text{for stability } \epsilon = \text{small value}$$

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

learnable parameters
of model.

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

$$\beta = y$$

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

Use $\tilde{z}^{(i)}$ instead of $z^{(i)}$

Fitting Batch to a network

$$x \xrightarrow{w^{(i)}, b^{(i)}} z^{(i)} \xrightarrow{\gamma^{(i)}, g^{(i)}, \text{Norm}} \tilde{z}^{(i)} \xrightarrow{\text{normalize}} q^{(i)} = g^{(i)} \tilde{z}^{(i)}$$

$$\xrightarrow{} z^{(1)} \xrightarrow[B_n]{B^{(1)}, f^{(2)}} \tilde{z}^{(2)} \xrightarrow{} a^{(2)} \xrightarrow{} \text{layer 2 normalized}$$

Parameters we get

$$\Rightarrow \beta^{(l)}, \gamma^{(l)}, \beta^{(a)}, \gamma^{(a)}, w^{(l)}, b^{(l)}, w^{(a)}, b^{(a)}$$

Working with minibatch

$$x^{(l)} \xrightarrow{w, b} z^{(l)} \xrightarrow{\beta, \gamma} \tilde{z}^{(l)} \xrightarrow{\text{normalize}} g^{(l)}(\tilde{z}^{(l)})$$

Parameter

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

Implementation gradient descent

\Rightarrow for $t < 1$ num Mini Batch -

Compute forward propagation, Use BN
to repa. $z^{(l)}$ with $\tilde{z}^{(l)}$.

Use back prop & compute

Update parameters

Work w/ GD momentum, RMS, Adam.

Batch Norm:

Summary: It is a process in which parameters w, b are normalized into $B, \gamma(\text{gamma})$, such that it increase the working of model by preventing the importance to features with large ranges.

Why does ~~no~~ Batch normalize work:-

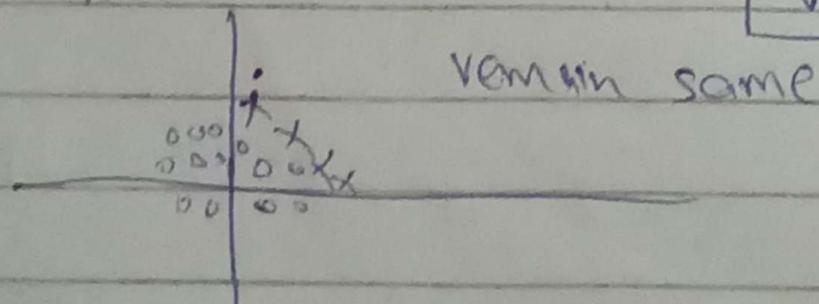
2nd reason:-

Let train cat detection model, if we compare colored or black cats, without Batch Norm.

our model or cost func don't work well.

What batch Norm do

mean = 0
variance = 1



it limits the amount to which updating the parameters in early layers which can effect distribution of layers in later layers, and batch layer effect reduce the changes less

* allow each layer layer independently
* It allows to add some noise to values $z^{[l]}$, within that minibatch.
~~for similar of hidden layer's activation~~

* This has a slight regularization effect.

* Adversing don't use Batch norm as regularizer but it has extra effect on learning algorithm.

* don't use it as a regularizer

* Use it only to regularize hidden layers parameters only -

Use Batch Norm at test time

Batch Norm process data 1 mini batch at a time but at test ^{time} we want process all at a time.

$$\begin{aligned} \bar{z} &= \frac{1}{m} \sum_i z^{(i)} \\ \sigma^2 &= \frac{1}{m} \sum_i (z^{(i)} - \bar{z})^2 \end{aligned}$$

~~to~~ \bar{z}, σ^2 = estimate this using Weighted average (across mini-batch)

$x^{(1)}, x^{(2)}, x^{(3)}$
minibatches

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

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

in this same way we get some

$$\gamma^{[1]}, \beta^{[1]}, \gamma^{[2](l)}, \beta^{[2](l)}, \gamma^{[3](l)}, \beta^{[3](l)}$$

after getting this we find this

Softmax Regression

$$a^{[l]} = \frac{e^{z^{[l]}}}{\sum_{i=1}^4 t_i}$$

Recognizing cats, dogs, and baby chicks

$$z^{[l]} = w^{[l]} q^{[l-1]} + b^{[l]}$$

Activation function:-

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

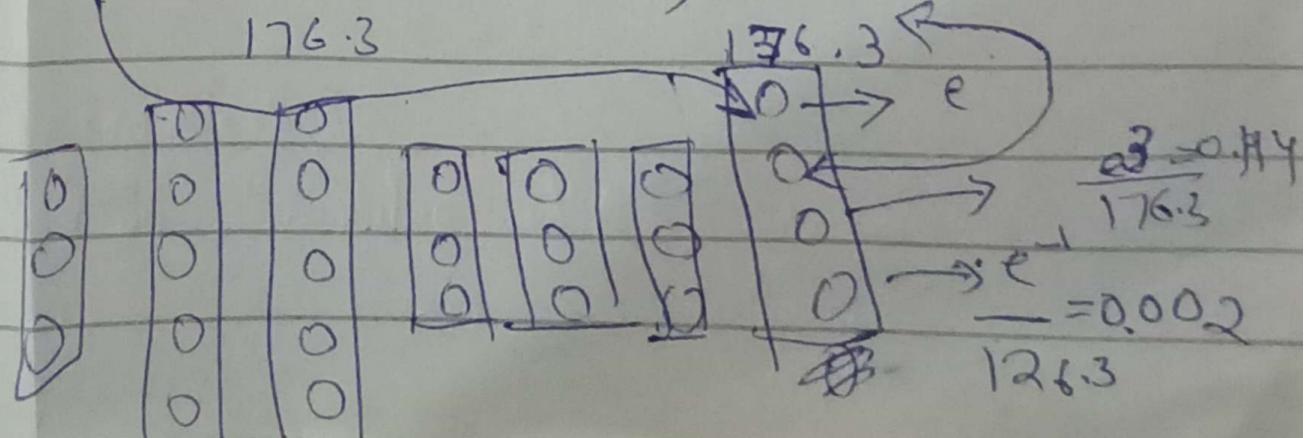
$$a^{[l]} = \frac{e^{z^{[l]}}}{\sum_{j=1}^n t_j}, a^{[l]} = \frac{t_j}{\sum t_i} [4, 1]$$

Example:

$$\begin{bmatrix} t_1 \\ t_2 \\ t_3 \\ t_4 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} = \sum_{j=1}^4 t_j = 176.3$$

$$a^{[l]} = \frac{1}{176.3} =$$

For $\frac{e^5}{176.3} = 0.842$, $\frac{e^2}{176.3} = 0.042$



foreach MN

loss function

Softmax Regression \Rightarrow extension of logistic Regression.

Example :-

Classifying images of handwritten digits (0-9)
predicting fruit {apple, banana, orange}

Training a Softmax classifier

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

$$g^{(i)}(z^{(i)}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

Introduction

Frameworks

Caffe/Caffe

CNTK

DL4J

Keras

TensorFlow

mxnet

PaddlePaddle

Torch

Probability

0.84 it means 84% confidence of being familiar or detecting apple.

foreach MN

Softmax Regression \Rightarrow extension of logistic Regression.

Example :-

classifying images of handwritten digits (0-9)
predicting fruit {apple, banana, orange}

Training a Softmax classifier

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

$$g^{(1)}(z^{(1)}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

7

Probability

0.84 it means 84% confidence of being familiar or detecting apple.

loss function

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \hat{y} = \begin{bmatrix} 0.5 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

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

$$= -y_1 \log \hat{y}_1 - \log \hat{y}_2$$

Cost

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

Introduction to Programming

Frameworks:

Caffe/Caffe2

CNTK

DL4J

Keras

Lasagne

mxnet

PaddlePaddle

TensorFlow

Theano

Torch

Tensor flow

If I want to learn tensor flow I get through the link of Andrew Ng course on Improving NN in coursera. At last there are link of further course on tensor flow.

Course 3:-

Structuring Machine learning strategy-

When training a model we have ideal, we will use it in a structure way, finding and creating strategies.