

A quality product by
Santi Paper Products
Plt no 358/477, Tamando,
Bhubaneswar, 751030
Feedbacks and queries :
santipaperproducts16@gmail.com



Size : 297X210 mm
MRP : ₹ 45.00
(Incl. of all taxes)
Exercise Book
140 Pages
(Total Pages Include Index & Printed Information)



8906060291933

Type of Ruling:
Single Line



Notebook

INDEX

MC

NAME : Kush Pandya. SUB.: Science
STD.: 7 DIV. A ROLL NO. 10

S. NO.	DATE	TITLE	PAGE NO.	SIGN./REMARKS
--------	------	-------	----------	---------------

Page:
Date:

S. NO. DATE TITLE PAGE NO. SIGN./REMARKS

NAME: NAME: NAME:
SDR: SDR: SDR:
ROLL NO.: ROLL NO.: ROLL NO.

PAGE NO. PAGE NO. PAGE NO.

$(x, y) \in \mathbb{R}^n \rightarrow$ a single column vector of n features.
 $y \in \{0, 1\}$

m training example: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots \\ | & | & | & | \end{bmatrix}_{n \times m}$$

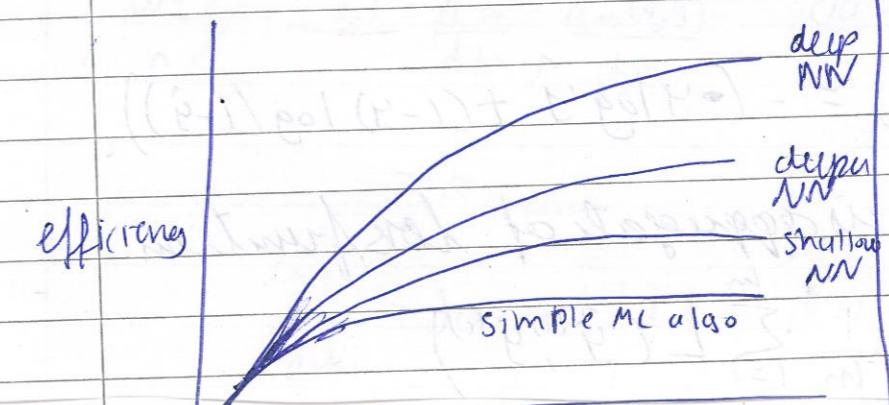
$$X \in \mathbb{R}^{n \times m}$$

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

$$y \in \mathbb{R}^{1 \times m}$$

$$y.\text{shape} = (1, m)$$

Why deep neural network work better



Logistic Regression

Given x want $y = P(y=1|x)$.

$x \in \mathbb{R}^{n_x}$ $b \in \mathbb{R}$.

Parameter: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$. Parameters $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

for linear regression $\rightarrow w^T x + b$.

for logistic regression.

Output $y = \sigma(w^T x + b) \rightarrow$

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$y' = \sigma(w^T x + b) \text{ where } \sigma(z) = \frac{1}{1+e^{-z}}$$

In logistic regression we do not use R^c method.

Since in gradient descent it gives multiple minima.
which becomes hard to find the least minimum.

The loss function used for logistic regression

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

Cost function is aggregate of loss function.

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

Gradient descent method.

Gradient descent method.

$$w = w - \alpha \frac{d J(w)}{d w} \Rightarrow dw$$

$$\Rightarrow w = w - \alpha dw.$$

and another one for b .

Naming convention.

d Final outvar \rightarrow dvar

d var

any var previous node variable.

Logistic regression derivatives.

$$\begin{aligned}
 & \begin{array}{c} x_1 \\ w_1 \\ x_2 \\ w_2 \\ b \end{array} \rightarrow z = w_1 x_1 + w_2 x_2 + b \rightarrow a = \sigma(z) \rightarrow L(a, y) \\
 & \frac{d z}{d t} = \frac{d L}{d z} = \frac{d L(a, y)}{d z} \quad \frac{d a}{d t} = \frac{d L(a, y)}{d a} \\
 & = -\frac{y}{a} + \frac{1-y}{1-a} \\
 & = a - y.
 \end{aligned}$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a} = x_i \frac{\partial a}{\partial z} \quad \frac{\partial a}{\partial z} = \sigma'(z) \quad \frac{\partial z}{\partial w_i} = x_i$$

logistic regression gradient descent.

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

$$a^{(i)} = g^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial J(w, b)}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial L(a^{(i)}, y^{(i)})}{\partial w_i}}_{dw_1^{(i)}}$$

$$dw_1^{(i)}$$

lets initialize

$$J=0, dw_1=0, dw_2=0, db=0$$

for $i=1$ to m . // It's a single iteration of a gradient descent.

~~$$z^{(i)} = w^T x^{(i)} + b$$~~

$$a^{(i)} = \sigma(z^{(i)})$$

$$J+ = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

}

$$J+=m$$

$$dw_1 / m$$

$$dw_2 / m$$

$$db / m$$

So here we have done this recursion.

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

Since to implement the above code we will need so many loop's that will make our program less efficient.

So we use something called Vectorization.

Vectorization

Just get rid of for loop.

$$z = \text{np.dot}(w, x) + b$$

$$w = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad x = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad b = []$$

avoid explicit for loop.

$$\begin{aligned} z^{(1)} &= w^T x^{(1)} + b \\ a^{(1)} &= \sigma(z^{(1)}) \end{aligned}$$

$$\begin{aligned} z^{(2)} &= w^T x^{(2)} + b \\ a^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

$$\begin{aligned} z^{(3)} &= w^T x^{(3)} + b \\ a^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

$$X = \begin{bmatrix} 1 & x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}^T \quad n_x$$

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

$$w^T = [w_1, w_2, \dots, w_m]$$

$$\begin{aligned} z &= [z^{(1)}, z^{(2)}, \dots, z^{(m)}] = w^T X + [b, b, \dots, b] \\ &= [w^T x^{(1)} + b, w^T x^{(2)} + b, \dots, w^T x^{(m)} + b] \end{aligned}$$

$$z = \text{np.dot}(w^T, x) + b \rightarrow \text{a simple no}$$

$$A = [a^{(1)}, a^{(2)}, \dots, a^{(m)}] = \sigma(z)$$

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

$$dz = A - Y = [a^{(1)} - y^{(1)}, a^{(2)} - y^{(2)}, \dots]$$

$$dt = [d z^{(1)}, d z^{(2)}, \dots, d z^{(m)}]$$

$$db = \frac{1}{m} \sum_{i=1}^m d z^{(i)} \equiv \frac{1}{m} \text{np.sum}(dz)$$

$$dw = \frac{1}{m} \times dz^T$$

The final implementation of logistic regression.

Put this in a for loop if you want to carry multiple iteration.

$$\begin{aligned} z &= w^T x + b \\ &= \text{np.dot}(w^T, x) + b. \end{aligned}$$

Note: this is for 1 iteration or gradient descent

$$A = \sigma(z)$$

$$dz = A - Y$$

$$dw = \frac{1}{m} \times dz^T$$

$$db = \frac{1}{m} \text{np.sum}(dz)$$

$$w = w - \alpha dw$$

$$b = b - \alpha db.$$

$A^{[0]}$ is also called ~~Input~~ Input.

Broadcasting in python.

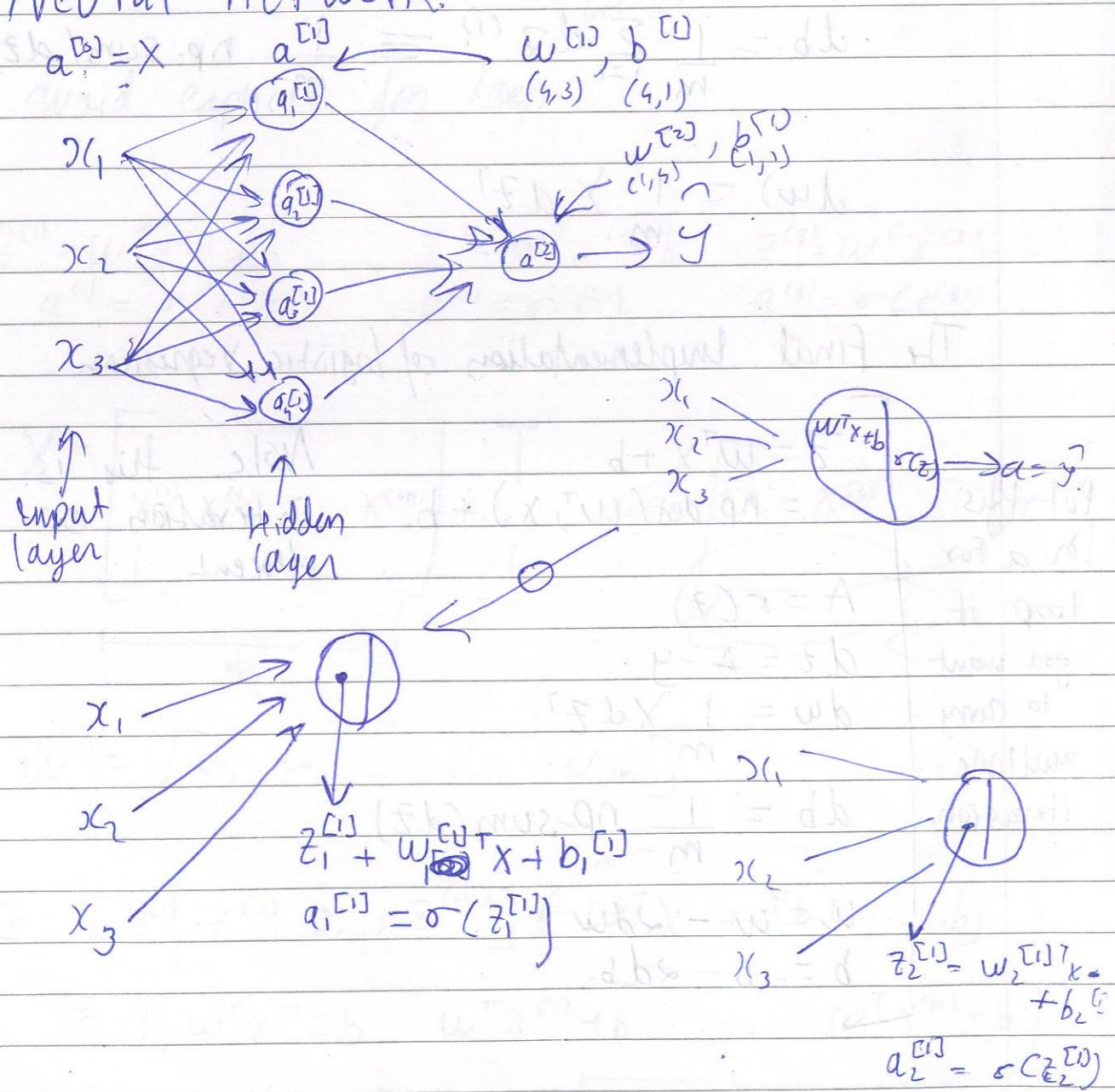
$$cal = A.sum(axis=0)$$

$$\text{Percentage} = 100 * A /$$

Vectorisation.

$$z^{[0]} = \begin{bmatrix} w_1^{[0]T} \\ w_2^{[0]T} \\ w_3^{[0]T} \\ w_4^{[0]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_1^{[0]} \\ b_2^{[0]} \\ b_3^{[0]} \\ b_4^{[0]} \end{bmatrix} = \begin{bmatrix} w_1^{[0]T} + b_1^{[0]} \\ w_2^{[0]T} + b_2^{[0]} \\ w_3^{[0]T} + b_3^{[0]} \\ w_4^{[0]T} + b_4^{[0]} \end{bmatrix}$$

Neural Network.



$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]})$$

(i.e) $z^{[1]} = W^{[1]}x + b^{[1]}$ (In short - For the σNN)
 $a^{[1]} = \sigma(z^{[1]})$
 $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$
 $a^{[2]} = \sigma(z^{[2]})$

$X \rightarrow a^{[2]} = y$
 $X^{(1)} \rightarrow a^{[2](1)} = y^{(1)}$
 $X^{(2)} \rightarrow a^{2} = y^{(2)}$

$X^{(m)} \rightarrow a^{[2](m)} = y^{(m)}$

example.
layer

On vectorized method

for $i=1$ to m :

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

$$a^{[1]i} = \sigma(z^{[1]i})$$

$$z^{[2]i} = W^{[2]} a^{[1]i} + b^{[2]}$$

$$a^{[2]i} = \sigma(z^{[2]i})$$

Vectorization

$$X = \begin{bmatrix} | & | & | & | \\ x^{[1]} & x^{[2]} & \dots & x^{[m]} \\ | & | & | & | \end{bmatrix}_{(n \times m)}$$

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

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

$$z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

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

You may also call it forward Propagation.

$$z^{[1]} = \begin{bmatrix} | & | & | & | \\ z^{[1]i_1} & z^{[1]i_2} & \dots & z^{[1]i_m} \\ | & | & | & | \end{bmatrix}$$

$$= \begin{bmatrix} | & | & | & | \\ a^{[1]i_1} & a^{[1]i_2} & \dots & a^{[1]i_m} \\ | & | & | & | \end{bmatrix}$$

If you need to understand more in details see video week 3 Justification for vectorized implementation.

Activation Func.

$$S = \text{Sigmoid} = \frac{1}{1+e^{-z}} \quad 0 \leq S \leq 1$$

$$t = \tanh = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad -1 \leq t \leq 1$$

\downarrow
It's better than sigmoid \rightarrow The means of the t^h comes more closer to 0.

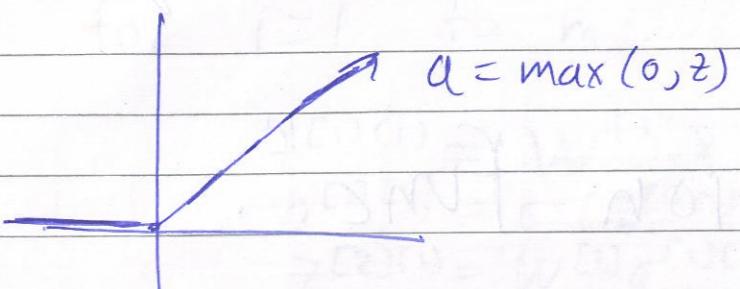
It centers your data.

It mostly strictly superior than sigmoid.

Use Sigmoid in output layer of binary classification.

The disadvantage of sigmoid and tanh is if large the input is large or small then the gradient derivative or slope of this function becomes very small. Which takes time getting the perfect descent.

$R \rightarrow ReLU \rightarrow$ Rectified linear unit.



It's derivative is zero so it takes less time.

Gradient descent for neural Network.

Parameters $w^{[1]} \in \mathbb{R}^{(n^{[1]}, n^{[0]})}$, $b^{[1]} \in \mathbb{R}^{(n^{[1]}, 1)}$, $w^{[2]} \in \mathbb{R}^{(n^{[2]}, n^{[1]})}$, $b^{[2]} \in \mathbb{R}^{(n^{[2]}, 1)}$.

$n_x = n_o$.

$$\text{cost function: } J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}, y)$$

Gradient descent.

Repeat {

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

$$dw^{[1]} = \frac{dJ}{dw^{[1]}}$$

$$db^{[1]} = \frac{dJ}{db^{[1]}}$$

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

$$b^{[1]} = b^{[1]} - \alpha db^{[1]}$$

} backpropagation.

Back propagation. (For binary classification)

$$dz^{[2]} = A^{[2]} - y$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]} +$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims=True})$$

$$dz^{[1]} = w^{[2]T} dz^{[2]} \times (g^{[1]'}(z^{[1]}))$$

any activation
fn derivative

↑ element wise product.

$$dw^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

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

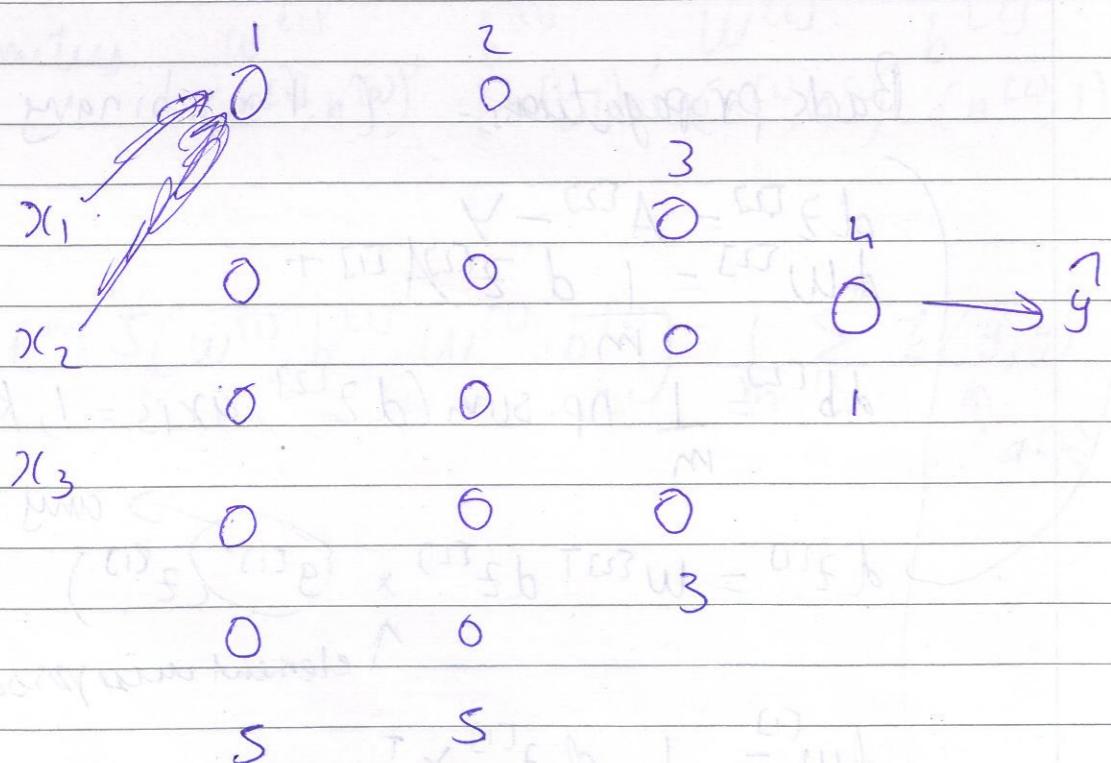
The Basic Deep Neural Network

$$\begin{cases} z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \\ a^{[L]} = g^{[L]}(z^{[L]}) \end{cases}$$

→ The basic of forward propagation

Deepness of NN.

More the layer more deep the NN.



$L =$ no. of layer = 4

$n_{[L]} =$ no. of units in layer $L =$

Page :
Date :

Page :
Date :

$a^{[L]}$ = activation in layer L.

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

$w^{[L]}$ = weight for $z^{[L]}$

$$x: z^{[1]} = w^{[1]} x + b^{[1]}$$

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

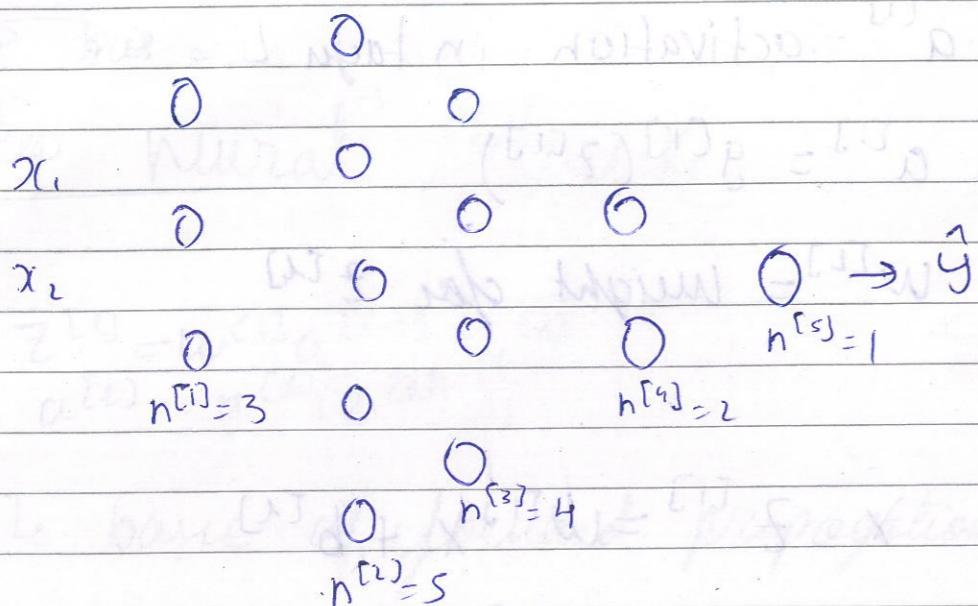
$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

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

$$z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

$$a^{[3]} = g^{[3]}(z^{[3]}) = y.$$

e.g.



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

$$\begin{matrix} (3, 1) \\ (n^{[1]}, 1) \end{matrix} \quad \begin{matrix} (2, 1) \\ (n^{[2]}, 1) \end{matrix}$$

$$\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$\text{dimension}(w^{[1]}) = (n^{[1]}, n^{[1-1]})$$

$$\text{dimension}(b^{[1]}) = (n^{[1]}, 1)$$

$$\dim(dw^{[1]}) = (n^{[1]}, n^{[1-1]})$$

$$\dim(db^{[1]}) = (n^{[1]}, 1)$$

One reason why we can't have a single layer of NN which has too many units instead of ~~large~~ deep NN, is because, if we want both of this model to be equally efficient the shallow NN model should have exponentially large units. Compared to input efficiency. And deep neural networks would have a logarithmic efficiency (not sure).

Forward function and backward

layer $\rightarrow l$.

Forward f^l will require $w^{[l]}, b^{[l]}$

input $\rightarrow a^{[l-1]}$

output $\rightarrow a^{[l]}$

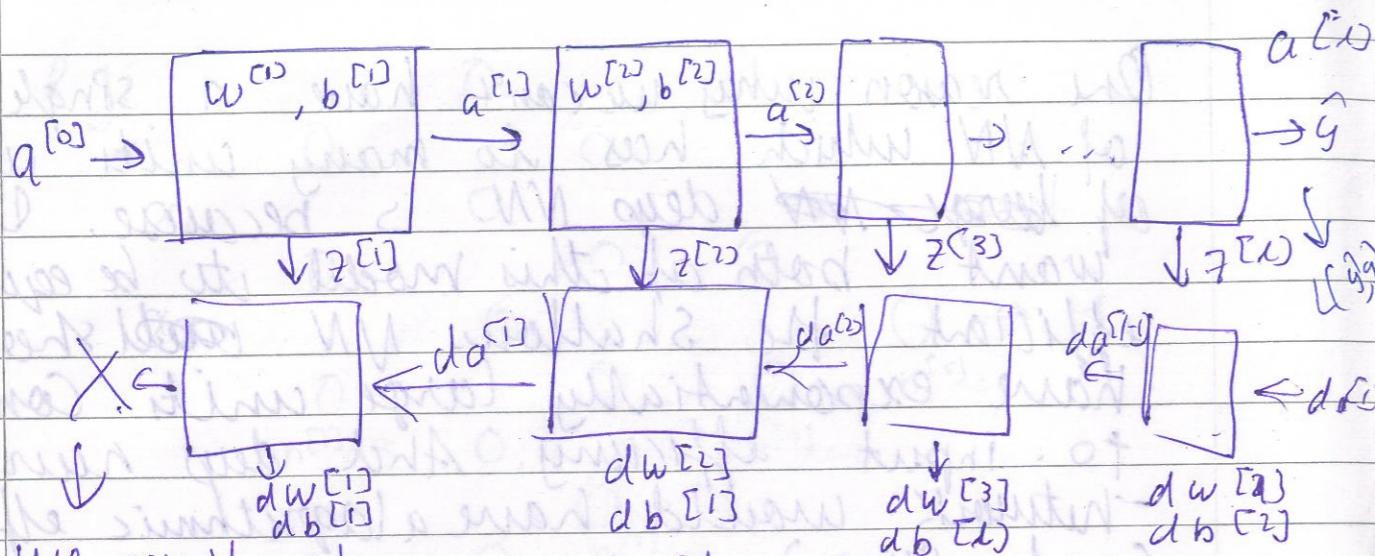
It will also cache(store) $z^{[l]}$

Backward f^l will require

input $\rightarrow d a^{[l]}$

output $\rightarrow d a^{[l-1]}$, ~~output~~ after this $d w^{[l]}, d b^{[l]}$

It will



We can't change input \mathbf{f}_h

Forward \mathbf{f}_h

$$\begin{aligned} z^{[1]} &= w^{[1]} a^{[1-1]} + b^{[1]} \\ a^{[1]} &= g^{[1]}(z^{[1]}) \end{aligned}$$

backward \mathbf{f}_h

$$\begin{aligned} w^{[1]} &= w^{[1]} - \alpha dw^{[1]} \\ b^{[1]} &= b^{[1]} - \alpha db^{[1]} \end{aligned}$$

Backward propagation for layer 1

steps.

$$dz^{[1]} = da^{[1]} * g'^{[1]}(z^{[1]})$$

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

$$db^{[1]} = dz^{[1]}$$

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

$$\begin{aligned} dz^{[L]} &= d\hat{y} * g'^{[L]}(z^{[L]}) \\ dw^{[L]} &= \frac{1}{m} dz^{[L]} \cdot a^{[L-1]T} \end{aligned}$$

$$\begin{aligned} db^{[L]} &= \frac{1}{m} np.sum(dz^{[L]}, axis=1, keepdims=True) \\ dA^{[L-1]} &= w^{[L]T} \cdot dz^{[L]} \end{aligned}$$

Parameters and hyperparameters

Parameter $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, w^{[3]}, b^{[3]}$

Hyper parameterisation \rightarrow learning rate
 \rightarrow iteration

\rightarrow hidden layer

\rightarrow hidden units $n^{[1]}, n^{[2]}$

\rightarrow choice of activation f_h .

i.e. \rightarrow momentum, minibatch, regularization

COURSE - 2

We have to identify the this thing before making neural networks.

- layers
- hidden units
- learning rates
- activation function

Idea
Experiment
Code

If you have small dataset try 70/30 train test split.

But if you have larger data set \rightarrow The size of test can be reduced drastically.

Basic Recipe.

High bias $\xrightarrow{\text{vs}}$ Biger Network
(training data)

↓ N
High variance $\xrightarrow{\text{vs}}$ More data regularization
(dev set perform)
↓ N
Done

Train set error		Dev set error
high Var	low	high
high bias	high	high (comparable to train set)
both bias		Super

Regularization

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

omit.

$$\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w.$$

L₂ regularization (Lasso Regres.)

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

L₂ is used most often.

λ = regularization parameter.

In Neural Network.

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

$$\|w^{[l]}\|_2^2 = \sum_{i=1}^{n^{[l+1]}} \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2$$

$w: (n^{[l]}, n^{[l+1]})$

"Frobenius norm"

Changes in backprop due to $dW^{[L]}$

$$dW^{[L]} = (\text{from backprop}) + \frac{\lambda}{m} W^{[L]}$$

$$\tilde{W}^{[L]} = W^{[L]}$$

$$W^{[L]} = W^{[L]} - \alpha dW^{[L]}$$

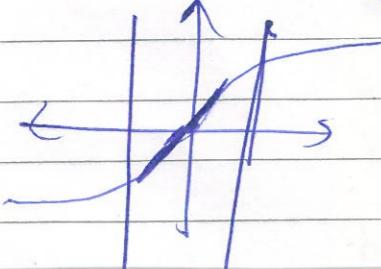
$$\Rightarrow W^{[L]} = W^{[L]} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{m} W^{[L]} \right]$$

$$= W^{[L]} - \frac{\alpha \lambda}{m} W^{[L]} - \alpha (\text{From backprop})$$

How does regularization help in reducing overfitting in NN.

$$J(\dots) = \frac{1}{m} \sum I(Y^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum \|W^{[l]}\|^2$$

So if $\lambda \uparrow \therefore W^{[L]}$ which is close to zero. In starting does not move much as λ amplifies it effect. We can see in graph below that λ would remain in the linear region and \therefore The NN will not be able make complicated decision boundary.



Dropout regularization

introduction drop out

Here illustrate with layer $L=3$ keep prob = 0.8

$$d_3 = \text{np.random.rand}(a_3.shape[0], a_3, \text{shape}[l]) < \text{keep-prob}$$

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

$$a_3 /= \text{keep-prob.}$$

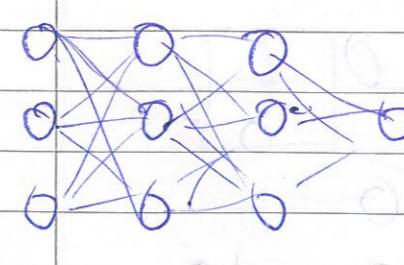
V.

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

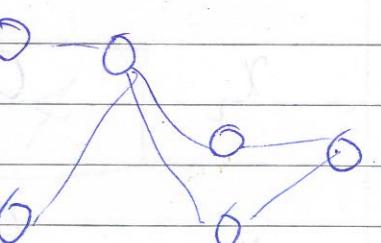
→ reduced by 20% (assumed by above operation so to make it expected to be by what we bump it)

At keep a_3 as expected value same

What is dropout.



remove certain node of a layer with some probability



iteration during

* For different training examples you zero out different hidden unit.

For making prediction

We don't use drop out during that time.
Since it will add noise to your predictions.

Keep prob ~~so~~ can be different for different layers just remember that with larger layers (more nodes) we keep the "Keep prob" low ~~(0.3)~~, while with small layers we keep them high.

We can put keep prob in input layer but its value should be very high.

Disadvantage.

~~Intuition~~

If forward prop a particular weight is large compared to other weights ^{in that layer} The dropout will have a effect on it and reduce its influence. increasing others. But it will have a overall effect of shrinking it.

$$0 \quad 0$$

$$x_1 \quad 0 \quad 0 \quad 0$$

$$x_2 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0$$

$$x_3 \quad 0 \quad 0 \quad 0 \quad 0.10 \quad 1.0$$

$$0 \quad 0 \quad 0.9$$

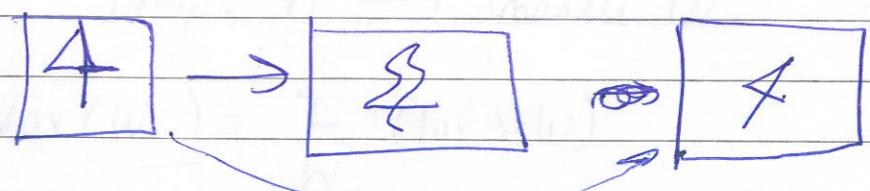
$$0 \quad 0$$

S-loss function. (loosely defined)

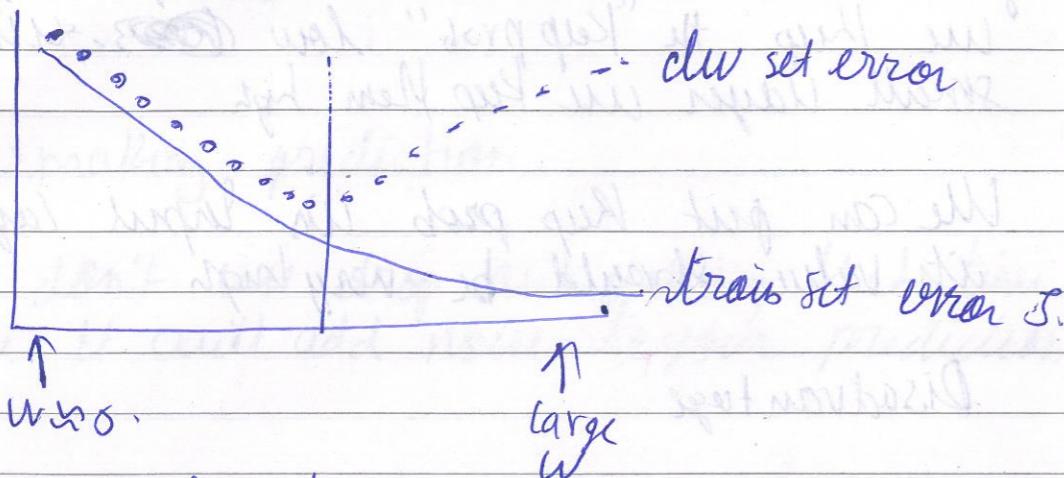
Data augmentation

Suppose we want to make cat classifier but we don't have enough data. One thing we can do is we can flip the image. Take some random zoom ins. etc.

For digits we can put random rotations and distortions to it.



Early stopping

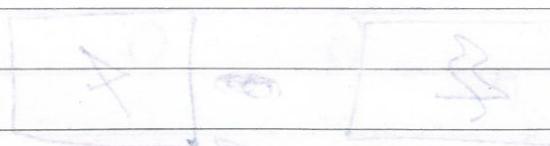


We just stop training the NN after some point. We can check whether the dev set error is decreasing or increasing.

Normalising.

Importance \rightarrow help your learning algorithm learn faster.

Meaning converting mean and variance of each of input attribute to same order of magnitude

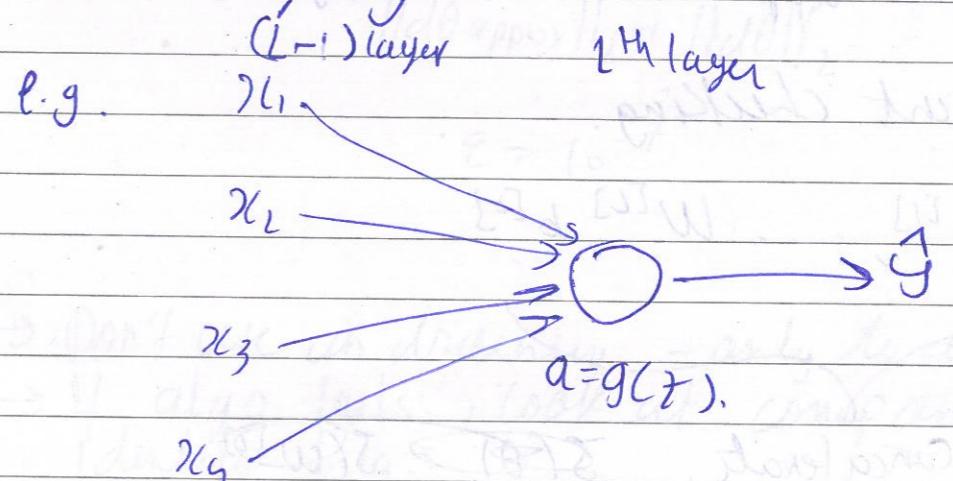


Vanishing/exploding gradient

See video [\[link\]](#) for further explanation.

We see that if a deep neural network starts working. Then we ~~have~~^{very} see that sometimes there is a vanishing and exploding gradients that really slows down the learning rate. ~~for the~~ complementation

To avoid this problem we try to initialize weight more carefully.



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

larger $n \rightarrow$ smaller w_i (\therefore to reduce the z based on previous layer size)

$$\text{Var}(w_i) = \frac{1}{n} \text{ (For relu)}$$

$$W^{(l)} = \text{np.random.randn(slope)} * \text{np.sqrt}\left(\frac{2}{n^{1-l}}\right)$$

Other Activation

tanh use $\sqrt{\frac{1}{n^{(l-1)}}}$

Xavier initialization

also $\sqrt{\frac{2}{n^{(l-1)} + n^{(l)}}}$

Numeric approximation of gradient.

Gradient checking.

Take $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$

concatenate. $J(\theta) \rightarrow J(W)$

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

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

to do the grade

$$J(dW^{[1]}, db^{[1]}, \dots) = \partial J(d\theta)$$

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

$$d\theta[i] = \frac{\partial J}{\partial \theta_i} \quad | \quad d\theta_{approx}$$

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

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

10^{-7} - great
 10^{-5} - OK (converge)
 10^{-3} - wrong.

- Don't use in training - only to debug.
- If algo fails., look at component by component to identify bug.
- Remember regularization.
- Doesn't work with dropout
- Run at random initialization: perhaps again after some training

Mini batch gradient descent.

$$X = \underbrace{\begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \end{bmatrix}}_{(n, m)} \quad y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

Size of minibatch $\rightarrow m \rightarrow$ batch learning (as previous)
 $\rightarrow 1 \rightarrow$ stochastic gradient descent.

Stochastic
 ↓
 less speedup
 lose vectorization
 advantage.
 As 1 is being
 processed at
 a time.

Minibatch
 (mini batch size
 not too small / large)

fault learning
 vectorization
 Make progress without
 processing the whole
 training data.

Batch
 Gradient
 Descent
 (mini batch grad.)

↓
 high time.

choosing minibatch size.

If small trainset: use batch gradient descent.
 $(m \leq 2000)$

Typical minibatch size

$$64, 128, 256, 512$$

$$2^6 \quad 2^7 \quad 2^8 \quad 2^9$$

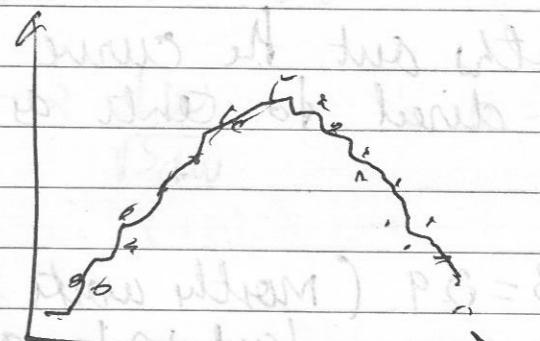
Make sure your minibatch fits in (CPU/GPU) memory

Exponentially weighted averages.

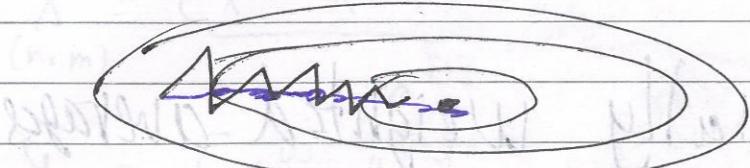
$$\theta_1 = 40^\circ F$$

$$\theta_2 = 60^\circ F$$

$$V_t = \beta V_{t-1} + (1 - \beta) \theta_t \quad \text{with } V_0 = 0 \quad \text{const}$$



Gradient descent with momentum



↑ Slower learning.
← Faster learning

Momentum
on iteration t

compute dw, db on wrt mini-batch

initialise \rightarrow

$$V_{db}, V_{dw} = 0 \quad V_{dw} = \beta V_{dw} + (1-\beta) dw \quad \left. \begin{array}{l} \text{changing} \\ \text{of } V_{dw}, V_{db} \end{array} \right\}$$

$$V_{db} = \beta V_{db} + (1-\beta) db \quad W = W - \alpha V_{dw}, b = b - \alpha V_{db} \quad \left. \begin{array}{l} \text{remaining} \\ \text{same.} \end{array} \right\}$$

If smooths out the curve and takes a smooth direct to center as shown in blue pen.

$\beta = 0.9$. (mostly works well)
average on last 10 gradients

Sometime ~~the~~ the above 2 equations
is sometimes written

$$\begin{aligned} V_{dw} &= \beta V_{dw} + dw \\ V_{db} &= \beta V_{db} + db \end{aligned}$$

but don't use this one.

gradient descent with momentum always works better than regular gradient descent

RMSprop

on iteration t

compute dw, db on the current minibatch.

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

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

$$W = W - \frac{\alpha dw}{\sqrt{S_{dw}}} \quad b = b - \frac{\alpha db}{\sqrt{S_{db}}}$$

Adam optimization Algo.

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

On iteration t

Compute dw, db using ^{current} minibatch

$$\left. \begin{array}{l} V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw \\ V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \end{array} \right\} \text{momentum}$$

$$\left. \begin{array}{l} S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2 \\ S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \end{array} \right\} \text{RMS prop.}$$

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

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

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^t)$$

$$S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^t)$$

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

α : needs to be fine

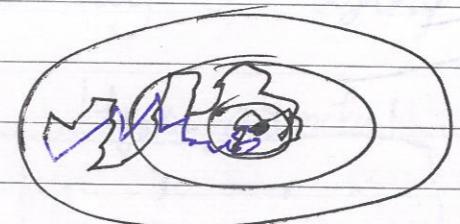
$\beta_1: 0.9$ (dw) recommended

$\beta_2: 0.999$ (dw^2) recommended

$\epsilon: 10^{-8}$ recommended

ADAM \rightarrow Adaptive moment estimation.

Learning rate decay.



mini batch

Slowly reduce α so as you get near to
through you take small steps.

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

Epoch	α	$\alpha_0 = 0.2$
1	0.1	decay rate = 1.
2	0.67	
3	0.05	
.	.	
.	.	

Other decay method

$\alpha = 0.95^{\text{epoch_num}} \alpha_0$ - exponentially decay

$$\alpha = \frac{R}{\sqrt{\text{epoch_num}}} \alpha_0 \quad \text{or} \quad \frac{R}{\sqrt{t}} \alpha_0$$

Formula

$$\begin{array}{c} - \\ - \\ - \\ - \end{array}$$

Manual decay.

Hyperparameters tuning.

most important priority

Priority given



α

B , #hidden units, mini batch size
layers, learning rate decay

least priority.

One way is using random value.

Hyperparameters

HP 1	.	,	.	,	.
.
.
.

Try ~~each~~ each of the value indicated by dots, equispace between them. And choose the one which gives best result.

But this works when no of HP are less.
When more HP are there to tune, then use random values given in next page

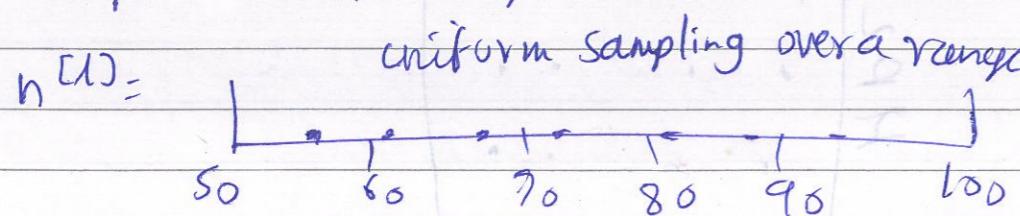
10 sets $\Rightarrow 11$ hyper parameters + 1 = 12
Y01 = b

HP1

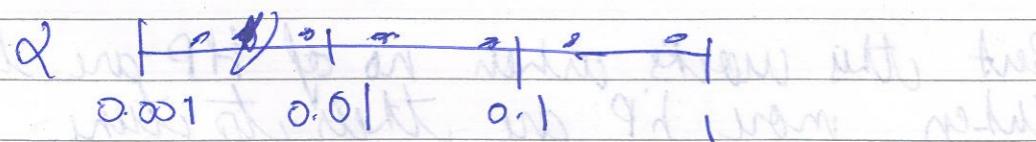
Random

Using appropriate scale for choosing random hyperparameters.

for hyperparameters like no of units, no of layers. choosing in an uniform sampling in a particular range would be suitable.



But for alpha try for exponential log scale.



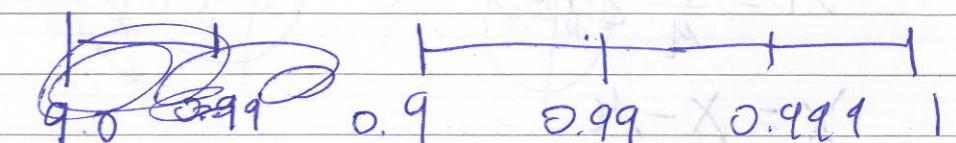
To implement in python.

$r = -5 * \text{np.random.rand}() \leftarrow \sigma \in [-5, 0]$

$$\alpha = 10^r$$

For B. (exponentially weighted average)

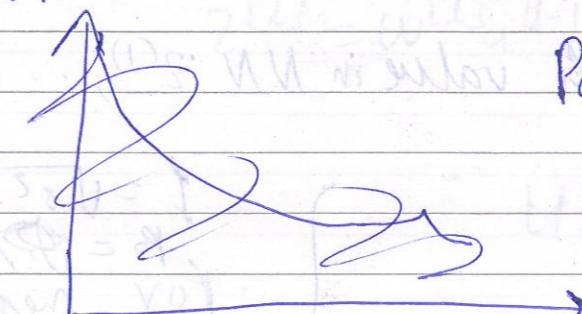
It's the same as before. But it starts with some bias.



$\beta: 0.900 \rightarrow 0.9005$: averaging to last 10 value

$\beta: 0.999 \rightarrow 0.9995$ averaging to around 1000 value

Hyperparameter tuning practice



Pandas \rightarrow manually done to tune it

Batch normalization.

In logistic regression.

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

$$X = X - \mu$$

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

$$X = X / \sigma^2$$

Now we can do this for every layer after applying activation fn.

Give some intermediate value in NN $z^{(1)}, \dots, z^{(n)}$

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

$$\left. \begin{array}{l} f = \sqrt{\sigma^2 + \epsilon} \\ \beta = \frac{\mu}{f} \\ \text{FOR mean=0} \\ \text{Variance=1} \end{array} \right\}$$

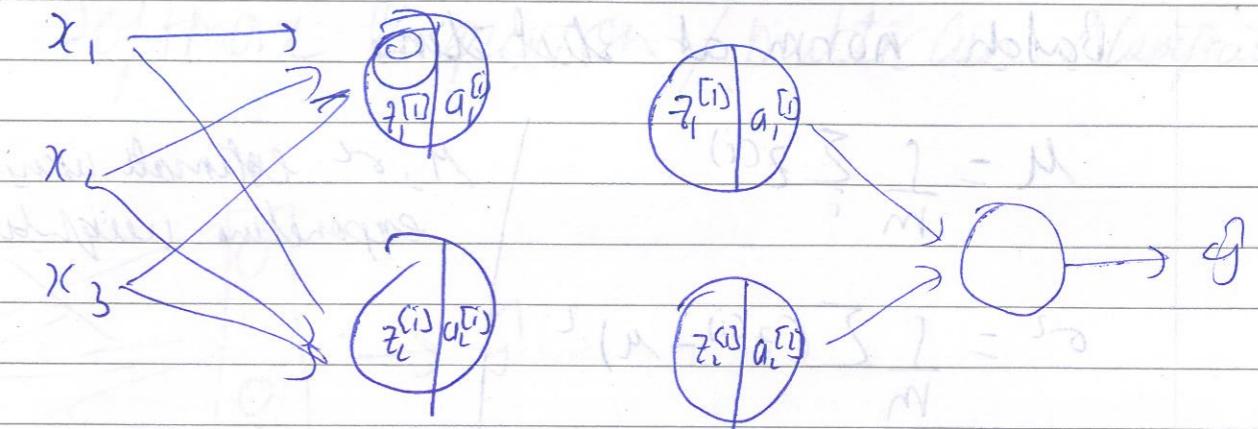
$$\sigma^2 = \frac{1}{m} \sum (z_i - \mu)^2$$

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

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

learnable parameters
or v-model.

Use $z^{(\text{norm})}$ itself if $z^{(\text{norm})}$



$$\begin{aligned} X &\xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow{\beta^{[1]}, \gamma^{[1]}} z^{[1]} \xrightarrow{a^{[1]} = g^{[1]}(z^{[1]})} \\ &\xrightarrow{w^{[2]}, b^{[2]}} z^{[2]} \xrightarrow{\beta^{[2]}, \gamma^{[2]}} z^{[2]} \xrightarrow{a^{[2]}} \end{aligned}$$

In batch normalization make the mean zero.

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

$\therefore b^{[t]}$ is not needed.

For $t=1 \dots \text{num Mini Batches}$

compute forward pass on $X^{[t]}$ (minibatch)
in each hidden layer use BN to repair
 $z^{[t]}$ with $\bar{z}^{[t]}$

Use backprop & update $dW^{[t]}, dB^{[t]}, d\beta^{[t]}$

(you can also momentum, Adam, RMSprop.)

Batch norm at test time

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

μ , or estimate using exponential weighted

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

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

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

(watch the last video of Batch normalization)

Softmax Regression / multiclass Classification.

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 6 \end{bmatrix}$$

$$\rightarrow y$$

y

last layer.

Now suppose we want to classify 4 things
4 nodes in last layer.

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

Activation fn.

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

$$a^{(L)} = e^{z^{(L)}}$$

$$(y, i) = \frac{e^{z^{(L)}}}{\sum_j e^{z^{(L)}}}, a_i^{(L)} = \frac{f_i}{\sum_j f_i}$$

This gives probabilities out a vector of shape.

The softmax gives out a vector (not a value like ReLU, sigmoid, tanh) of probabilities for each layer L .

Orthogonalization.

chain of assumption in ML

→ fit training set well. In human level
y train bigger algo, like better training system.

Performance

→ fit dev set well

g. regularization. ↓
Bigger training set

→ fit test set well.

g get bigger dev set. ↓

→ Perform well in real world.

g change the
dev set & test set.

As per author early stopping is not used by him
because it affects both the 1 & 2 chain
in above figure.

Does not mean you can't use it. But it takes
away orthogonalization.

Setting up your goal.

→ single number evaluation metric →

classifier	Precision	Recall
A	95%	90%
B	98%	85%

We know that there is a tradeoff between precision and recall. But what if the circled one above (i.e. Precision of B & Recall of A) are good. Then how to know which is better

i. Use F₁ score. (useful iterating over other ML algo).

Now suppose your algo

Algo	wrong → US	Canada	India	China	Avg
A	3%	7%	5%	9%	6%
B	-	-	-	-	-
D	-	-	-	-	-

In above just take Avg. of all. Assuming Avg is a good way.

Now you can also combine time as attribute too.

→ Structuring train / dev / test set.

Choose dev & test set such that it should have same distribution.

→ Going away from certain attributes

$$\text{Cost}_{\text{Fn}} = \frac{1}{m_{\text{dev}}} \sum_{i=0}^{m_{\text{dev}}} w^{(i)} \text{loss}(y_{\text{pred}}^{(i)}, y^{(i)})$$

$$\text{where } w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is not } g \\ 100 & \text{if } x^{(i)} \text{ is } g \text{ attribute} \end{cases}$$

where g attribute is something we want to ~~not~~ avoid.

Footnote → If your model is not doing well on test set real world application then change the test set.

Comparing to human level performance

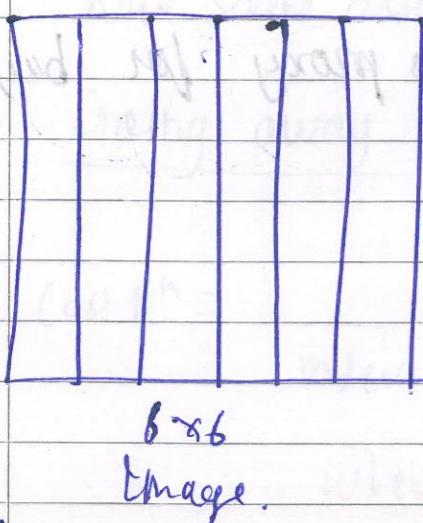
→ ~~Model~~

→ The best possible error (e.g. least possible error) on a set is called bayes optimal error.

Human level error is used as proxy for bayes error.

Convolutional Neural Network

Edge detection.



Vertical edge detection

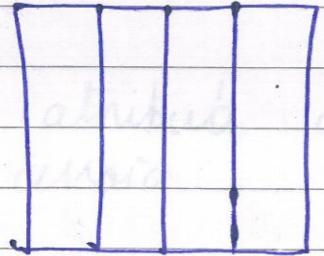
$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

filter

6×6

image.

This is
called
convolutional.



=

~~tf. nn. conv2d~~.

Page :
Date :

Page :
Date :

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

→ Vertical

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

→ Horizontal

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

→ Sobel Filter (Vertical)

$$\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$$

→ Sobel Filter

Sometimes for a complicated image
we can use all the elements of a filter as
a parameter

~~train other filters~~

$$\begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{bmatrix}$$

Padding

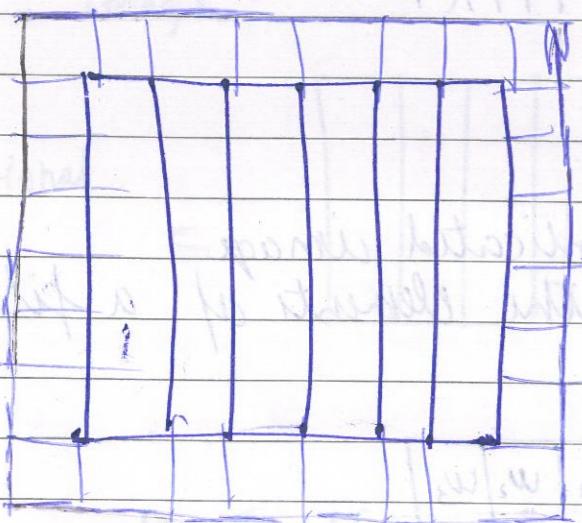
Suppose you have image of $n \times n$ & filter of $f \times f$

Then you will have a output of

$$(n-f+1) \times (n-f+1)$$

∴ Every time you apply convolution ~~normal~~ ~~it~~ to your image start shrink.

Also many edge number in image gets multiplied by very filter.



$P=1$ (Padding layer)

$$P = \frac{f-1}{2}$$

use mostly initialize padding to zero

6x6 image padded 1 block from edge making multiplication output with filter as 6x6 matrix

Page :
Date :

Page :
Date :

Valid and same convolution.

"Valid" → no padding.

"same" → output matrix same as input

f is mostly usually odd

→ Since it means it has a central filter which can be useful.

→ And also due to the formula $P = \frac{f-1}{2}$ which give whole number when $f = \text{odd}$.

Strided Convolution

Just see the video. Hard to explain in words

$n \times n \rightarrow$ image

$f \times f \rightarrow$ filter

padding → p

stride → s

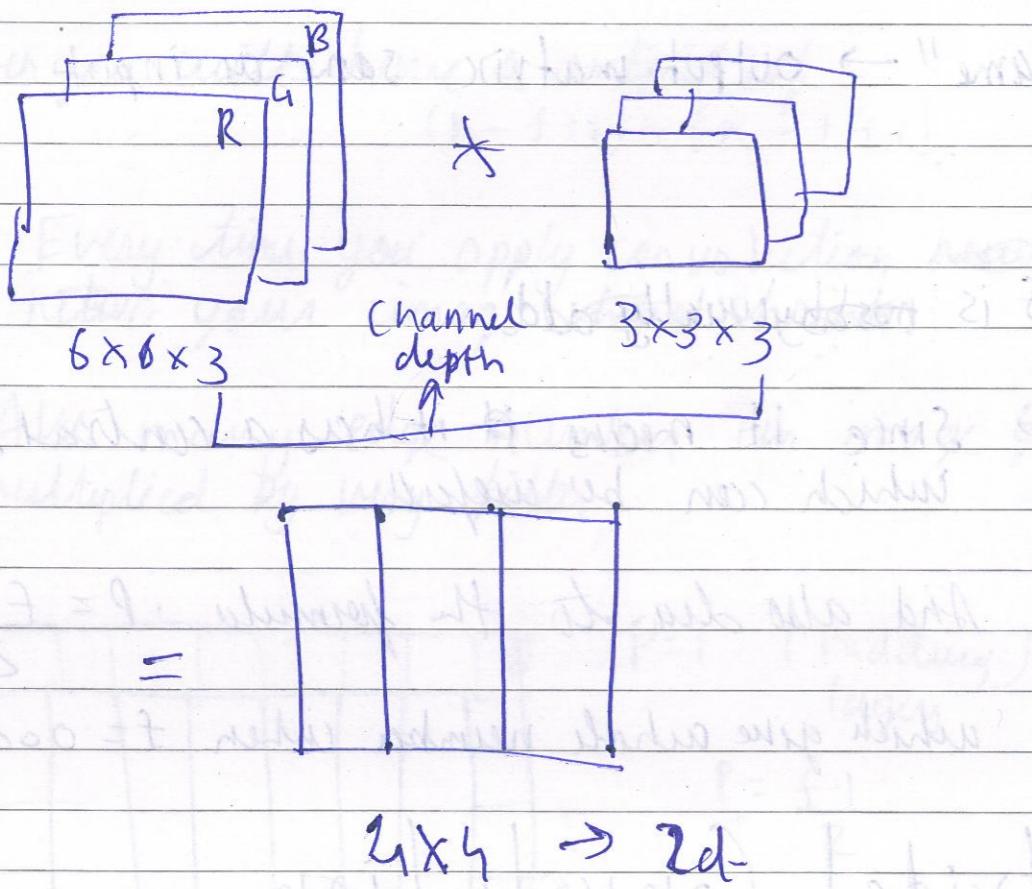
~~REPS~~

$$\text{output matrix} = \left[\frac{n+2p-f+1}{s}, \frac{n+2p-f+1}{s} \right]$$

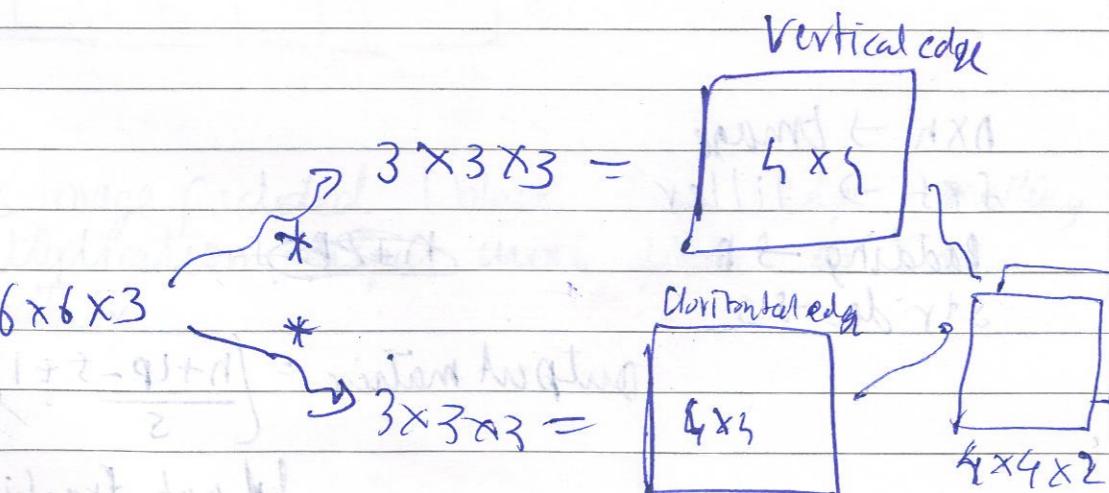
by not fraction than use

For RGB Image we have three channels R, G, B.

We use $3 \times 3 \times 3$ filter.



Multiple vector filter

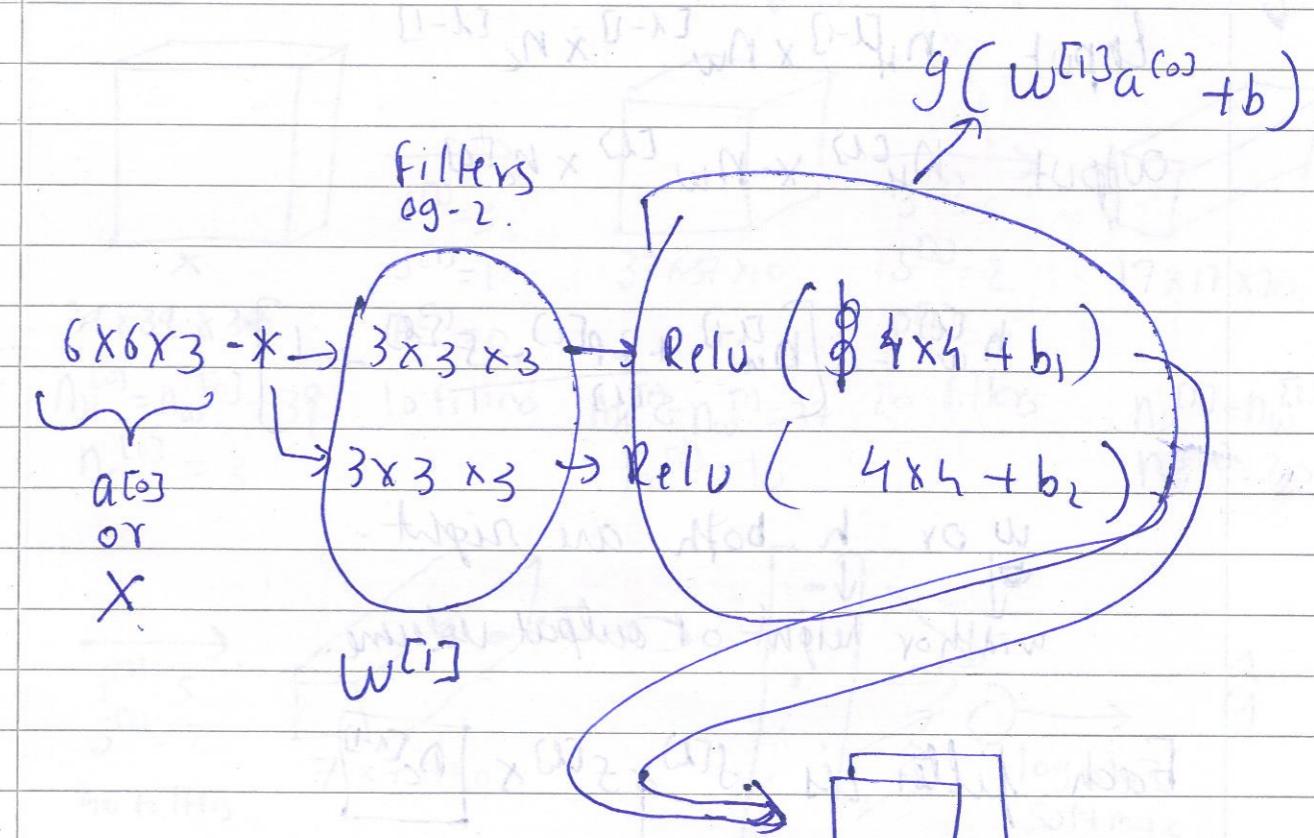


Summary.

$$n \times n \times h_c * F \times r \times h_c \rightarrow n - k + 1 \times n - t + 1 \times h_c$$

$$6 \times 6 \times 3 \quad 3 \times 3 \times 3 \quad 4 \times 4 \times 3$$

One layer of CNN



Notation

If layer l is a convolution layer $\times \times \times$

$f^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

$n_c^{[l]}$ = num of filters

Input $n_h^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]}$

Output $n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

$$h_w^{[l]} = \frac{[h_w^{[l-1]} + 2p^{[l]} - f^{[l]}] + 1}{s^{[l]}}$$

w or h both are right.

width or height of output volume.

Each filter is $f^{[l]} \times s^{[l]} \times n_c^{[l-1]}$

activation : $a^{[l]} \rightarrow n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

(same as output dimensions above)

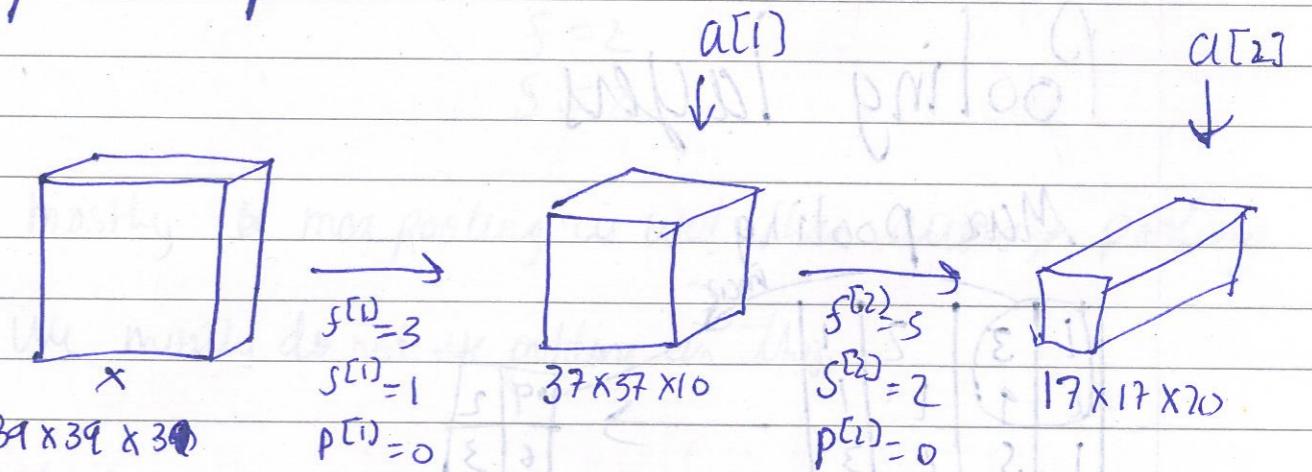
for ~~vector~~ batch/minibatch gradient descent
your output

$A^{[l]} \rightarrow m \times n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]} \rightarrow$ (not
univarsal)
no of training example

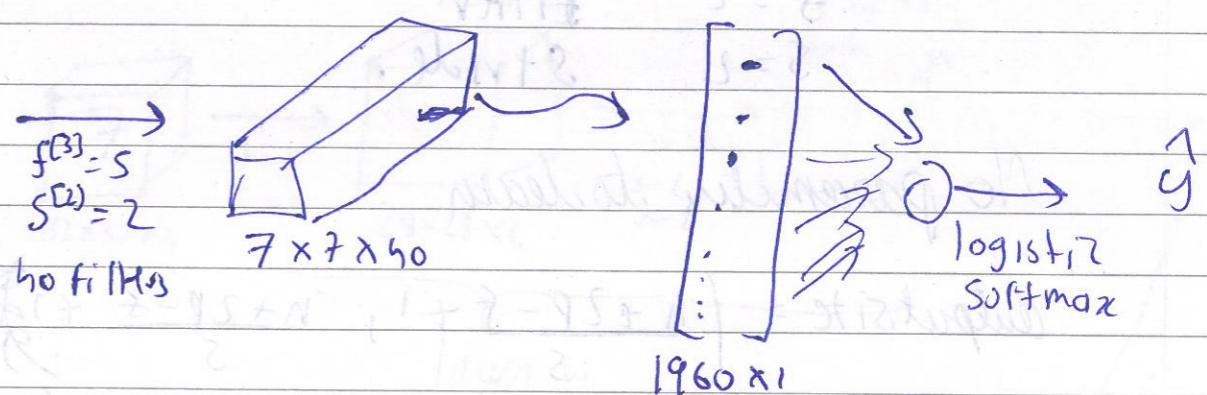
Weights : $f^{[l]} \times n_c^{[l-1]} \times n_h^{[l-1]} \times n_w^{[l-1]}$

bias : $n_c^{[l]}$ (i.e. $1, 1, 1, n_c^{[l]}$)

Example



$$\begin{aligned} n_h^{[0]} &= n_w^{[0]} = 39 & 10 \text{ filters} & n_h^{[0]} = n_w^{[1]} = 37 & 20 \text{ filters} & n_h^{[1]} = n_w^{[2]} = 17 \\ n_c^{[0]} &= 3 & & n_c^{[1]} &= 10 & n_c^{[2]} = 20 \end{aligned}$$



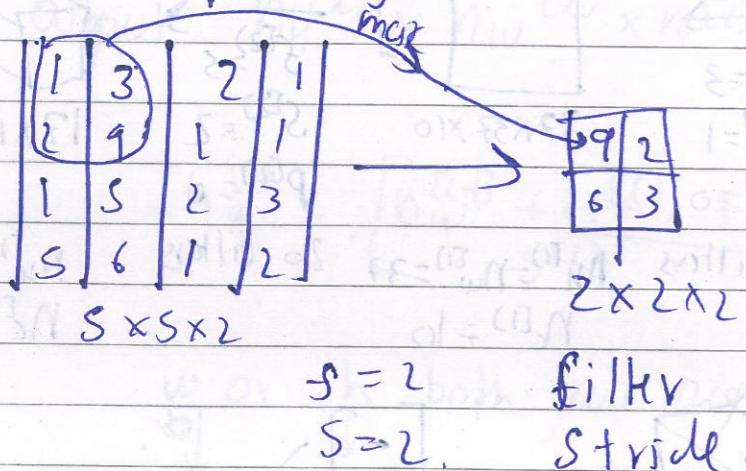
As you see height & width of decreases layer by layer and depth increases.

Types of layers in Convolution network.

- convolution (Conv) ←
- Pooling (Pool) ←
- Fully connected (FC) ←

Pooling layers

Max pooling



No parameters to learn

$$\text{output size} = \left(\frac{h + 2p - f + 1}{s}, \frac{w + 2p - f + 1}{s} \right)$$

Same as previous

Average pooling.

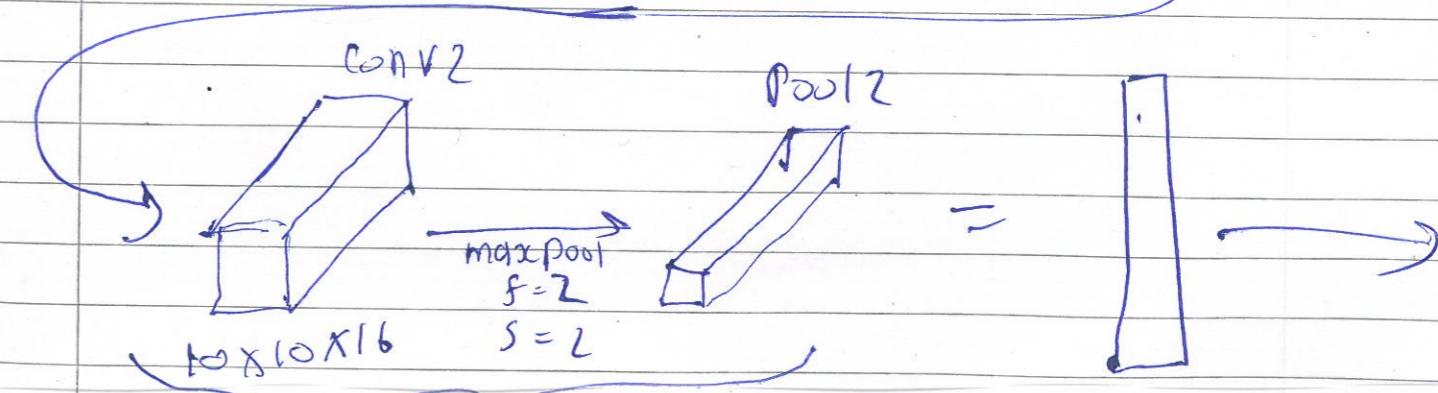
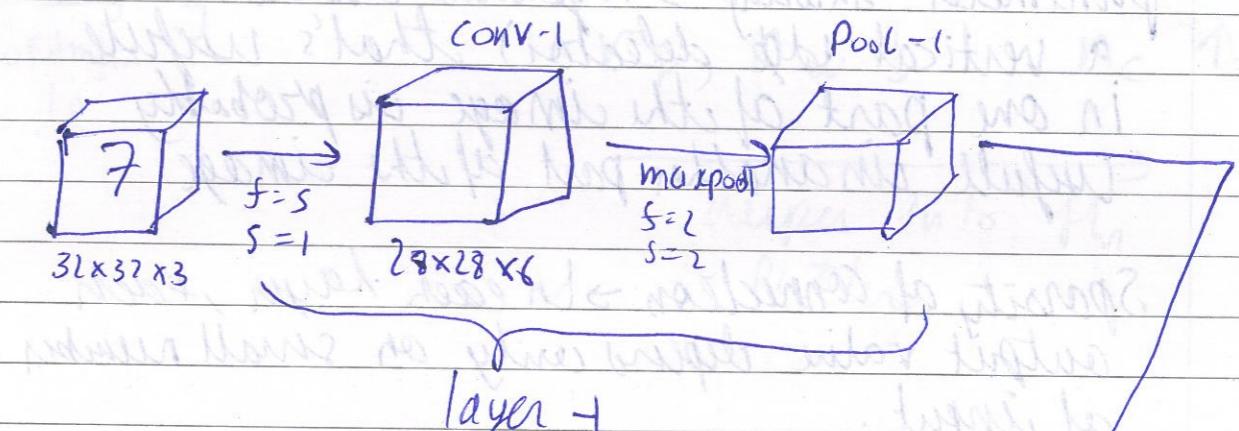
$$\begin{array}{|c|c|c|c|c|} \hline i & 3 & 2 & 1 & \\ \hline 2 & 9 & 1 & 1 & \\ \hline 1 & 4 & 2 & 3 & \\ \hline 5 & 6 & 1 & 2 & \\ \hline \end{array} \rightarrow \begin{array}{|c|c|} \hline 3.75 & 1.25 \\ \hline 4 & 2 \\ \hline \end{array}$$

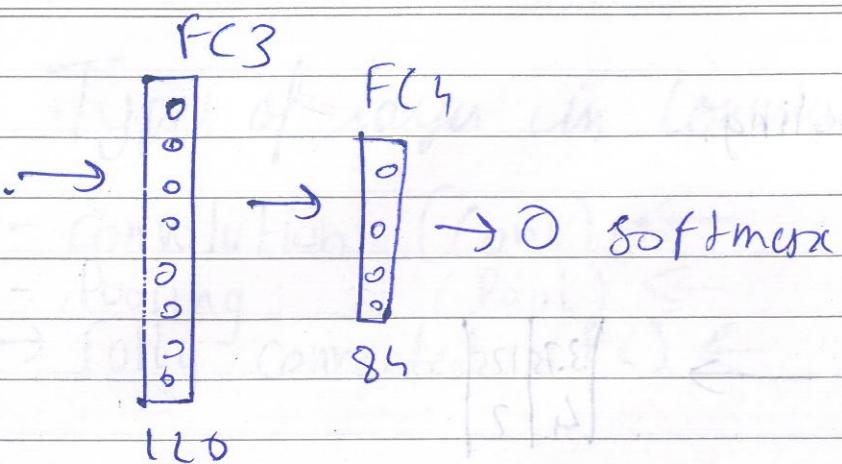
$f = 2$
 $s = 2$

mostly max pooling is used than average pooling

We mostly do not use padding in this.

LeNet-5 (example).





$n_h, n_w \downarrow$ $n_c \uparrow$ trend

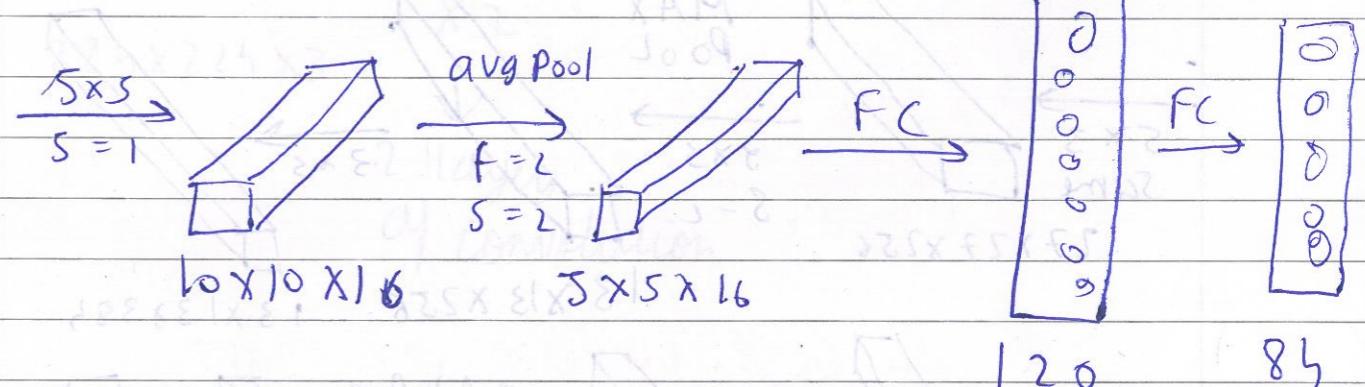
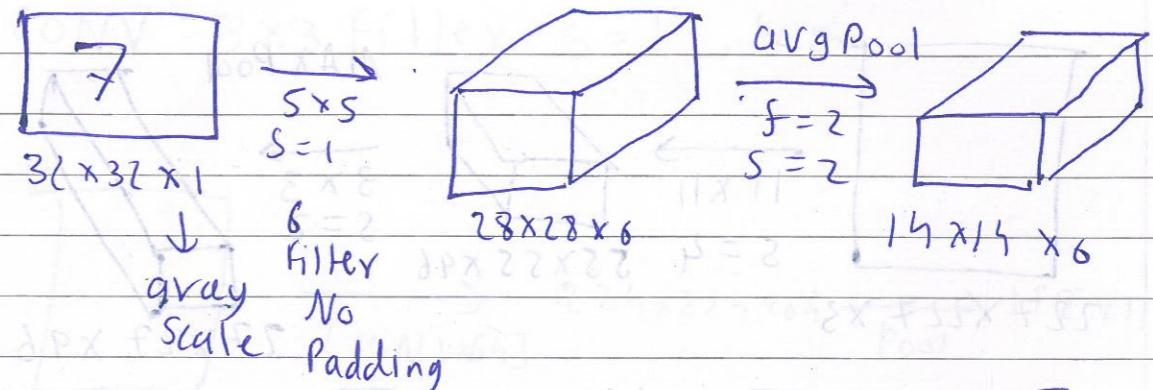
conv pool - conv - pool - FC - FC - softmax

Why Convolution.

Parameter sharing \rightarrow A feature detected such as a vertical edge detector that's useful in one part of the image is probably useful in another part of the image.

Sparcity of connection \rightarrow In each layer, each output value depends only on small number of inputs.

Lenet-5



\rightarrow softmax
10

$n_h, n_w \downarrow$ $n_c \uparrow$

deeper into ph
Pattern