

~~week 4~~

Course 4 :- Convolutional Neural Networks. [23.01.19].

week 1 : Foundation of convolutional Neural Networks.

L-1 : Computer vision. [10.00 Am . 23.01.19]

Computer vision problems.

- image classification.
- object detection.
- neural style transfer. etc.

Main challenge of computer vision is higher input space.

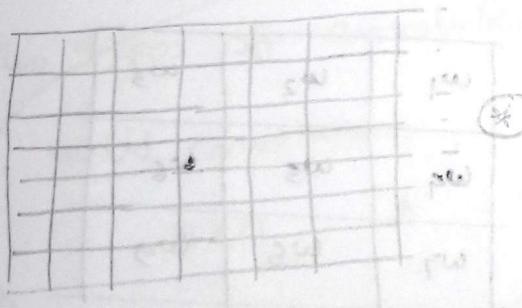
[convolution operation is the fundamental building block of convolutional neural network]

L-2 : Edge detection example. [10.10 Am 23.01.19]

vertical edge detection :-

6x6 image

python :	Conv. Backward
tensorflow :	tf.nn.conv2d
keras :	conv2D



$$\begin{matrix} 1 & 3 & -1 \\ 1 & 2 & -1 \\ 1 & 0 & -1 \end{matrix}$$

filter stride = 2
basic output

$$\begin{matrix} 8 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$$

☞ image processing \hookrightarrow Convolution \Rightarrow basic output

* A vertical edge is a region where bright region on the left and dark region on the right.

(convert 2x2 simple output for output)

L3: More Edge detection. [10.25 Am 23.01.19]

Matrix for
detecting
vertical
edge

1	6	-1
1	0	-1
1	0	-1

(vertical edge detector)

matrix for

detecting
horizontal
edges

1	1	1
0	0	0
-1	-1	-1

(Horizontal edge detector)

other filters to detect edges.

1. Sobel filter.

** Add extra weight to the center
of the image.

1	0	-1
2	0	-2
1	0	-1

2. Scharr Filter.

3	0	-3
10	0	-10
3	0	-3

w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9

** Our goal is to find best choices of this nine parameters.

→ use backpropagation to find out best parameters.

[g+ is one of the most powerful ideas of computer vision].

If we convolve a $(n \times n)$ matrix with a $(f \times f)$ matrix then size of resultant matrix will be $(n-f+1) \times (n-f+1)$.

$$\left\{ \begin{array}{l} n \times n \quad \otimes \quad f \times f \quad \rightarrow \quad (n-f+1) \times (n-f+1) \\ 6 \times 6 \quad * \quad 3 \times 3 \quad \rightarrow \quad 4 \times 4 \end{array} \right.$$

Disadvantages:-

- 1. It shrinks the image after each convolution.
 - 2. Corner pixels of an image touched only once. [A lot of information on the edges are thrown away].
- * To solve these problems we use padding.

** Adding extra pixels around the image is called padding.
Normally we add 0 as padding pixel.

Two types of convolutions,

1. valid convolution :- (No padding)
 $n \times n \quad * \quad f \times f \quad \rightarrow \quad n-f+1 \quad * \quad n-f+1$.

2. same convolution :- Pad so that output size is the same as the input size.

$$\text{input size} \rightarrow n+2P-f+1 \quad \times \quad n+2P-f+1. \quad \boxed{P = \text{Padding amount}}$$

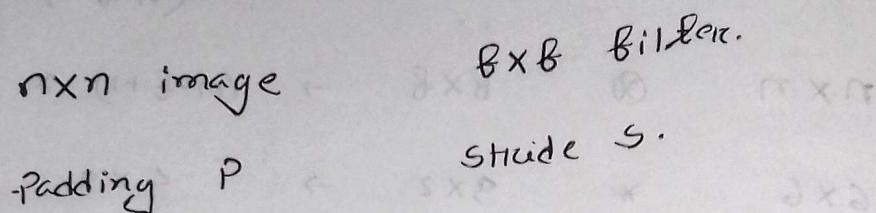
If we want input and output size same,

$$\left[\begin{array}{l} n+2P-f+1 = n \\ \therefore P = \frac{f-1}{2} \end{array} \right] \quad \begin{array}{l} \text{By using this equation we} \\ \text{can find suitable padding size.} \end{array}$$

** filter size is usually odd.

L-5 Strided Convolutions. [11:10 Am 23.01.19]

[we will skip/jump over $\frac{s}{s}$ number of rows and columns.]



Output size:

$$\left\lfloor \frac{n+2P-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2P-f}{s} + 1 \right\rfloor$$

Floor sign
if it is not an integer then we will take floor.

→ Cross-correlation vs. Convolution.

[Convolution is filter matrix for normally double flip over the input image and output is sum of product.

L-6 Convolutions Over Volume. [11:25 Am 23.01.19]

→ Convolution on RGB images:

→ channel of image and filter are equal.

[Computation part]

* we can use multiple filters on the image.

$$n \times n \times n_c * f \times f \times n_e \rightarrow n-f+1 \times n-f+1 \times n'_e$$

$n \times n \times n_c$ $f \times f \times n_e$ $n-f+1 \times n-f+1 \times n'_e$

number of channels. ↓ ↓

$6 \times 6 \times 3$ $3 \times 3 \times 3 \rightarrow 4 \times 4 \times 2$ number of filters.

* depending on the number of filters we can detect/multiple features. (number of features we are detecting)

L-7 one layer of convolutional Network

[11:45 AM 24.01.19]

Example of a layer :-

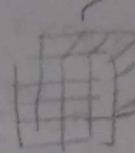


$6 \times 6 \times 3$

Billets
 $\star 3 \times 3 \times 3 \rightarrow \text{Relu}$

$(4 \times 4 + b) \rightarrow 4 \times 4$

$3 \times 3 \times 3 \rightarrow \text{Relu}(4 \times 1 + b) \rightarrow 4 \times 1$



$4 \times 1 \times 2$

because of 2 filters

** If you have 10 filters that are $3 \times 3 \times 3$ in one layer of neural network, how many parameters does that layer have?

$$3 \times 3 \times 3 = 27 \text{ parameters.}$$

$$+ \text{bias} = 28 \quad "$$

$$(28 \times 10) = 280 \text{ parameters.}$$

** One property of convolutional neural net is that number of parameters is independent of the (input size) of the image.

Summary of notation:-

If layer l is a convolutional layer:-

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

$P^{[l]}$ = padding.

$s^{[l]}$ = stride.

$n_c^{[l]}$ = number of filters.

number of filters in previous layer

Each filter size = $f^{[l]} \times f^{[l]} \times n_c^{[l]}$

weights : $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$.

bias : $1 \times 1 \times 1 \times n_c^{[l]}$

Activation : $a^{[l]} \rightarrow n_{ht}^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

Activations from Previous Layer

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

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

number of filters in layer l

size : $n_{ht}^{[l]} = \left\lfloor \frac{n_h^{[l-1]} + 2P^{[l]} - f^{[l]}}{s^{[l]}} \right\rfloor + 1$

$n_w^{[l]} = \left\lfloor \frac{n_w^{[l-1]} + 2P^{[l]} - f^{[l]}}{s^{[l]}} \right\rfloor + 1$

$A^{[l]} = m \times n_{ht}^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

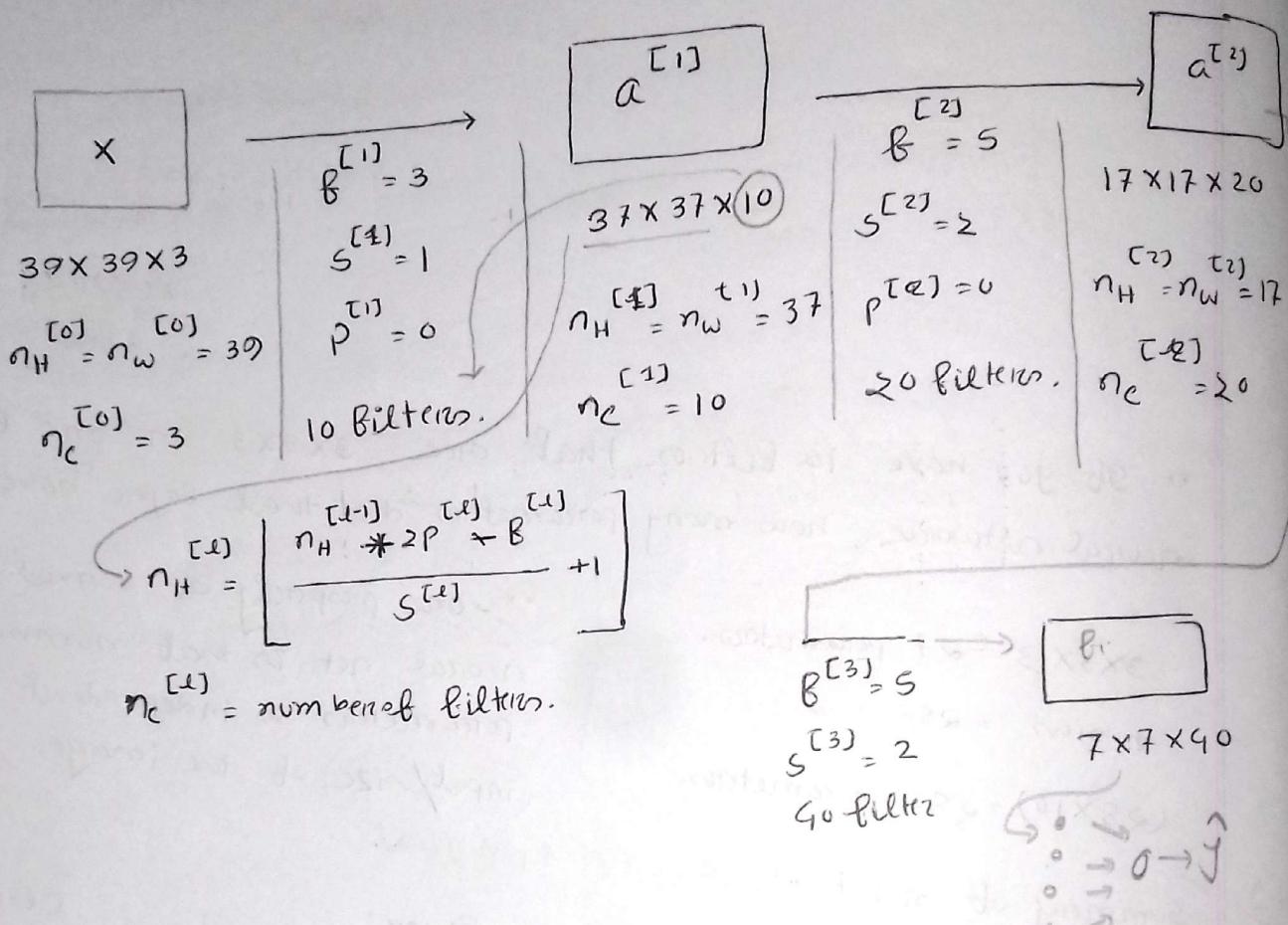
vectorized.

L-8 Simple Convolutional Network Example.

[12:20 PM 29.01.19]

Example ConvNet :- (Vector Notation Diet Example for 2nd L2)

Q



** As we go deeper of the network size decreases but the number of channel increases.

** Some of the hyperparameters of the convolutional neural network are.

1. Filter size.

2. Number of filters.

3. Padding size.

4. Stride size.

Types of layers in a convolutional network:-

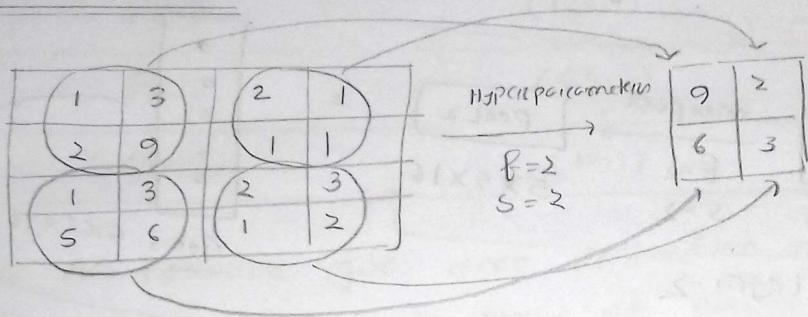
- Convolution. (conv)
- Pooling (pool)
- Fully connected. (fc)

L-9 pooling layers. [12.40 PM 24.01.19]

Max pooling :- (used frequently).

[max pooling works better in different applications
works on Theory once]

** No parameters to learn.



$$\left\lceil \frac{n+2P-p}{s} + 1 \right\rceil$$

** Size of the max pooled matrix

** Matrix size, stride & padding ensure one of region \rightarrow maximum

first 275

Average pooling :- In average pooling we do not take the maximum but take average of the numbers inside a window

Summary of pooling.

Hyperparameters :-

1. p : Filter size. } usually $p=2$

2. s : stride. }

$s=2$.

max or average pooling.

[max pooling usually do not use padding]

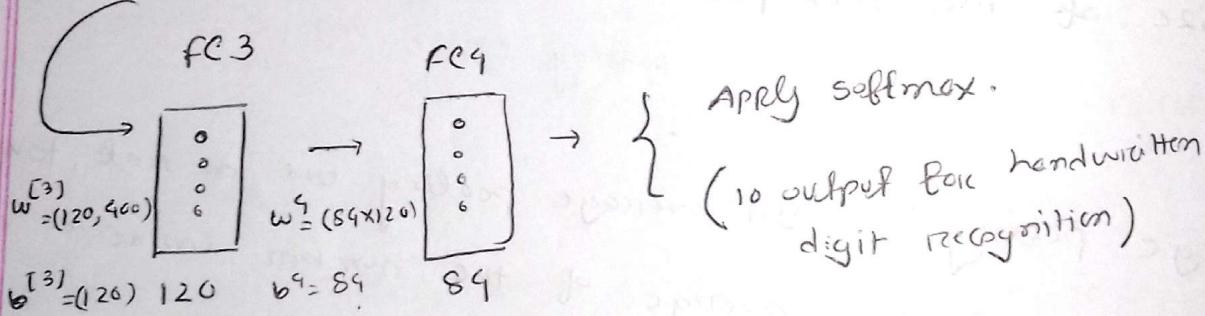
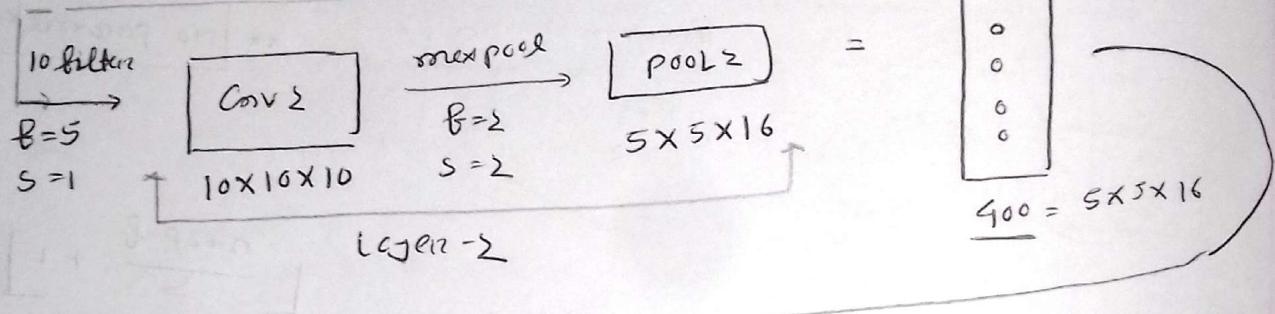
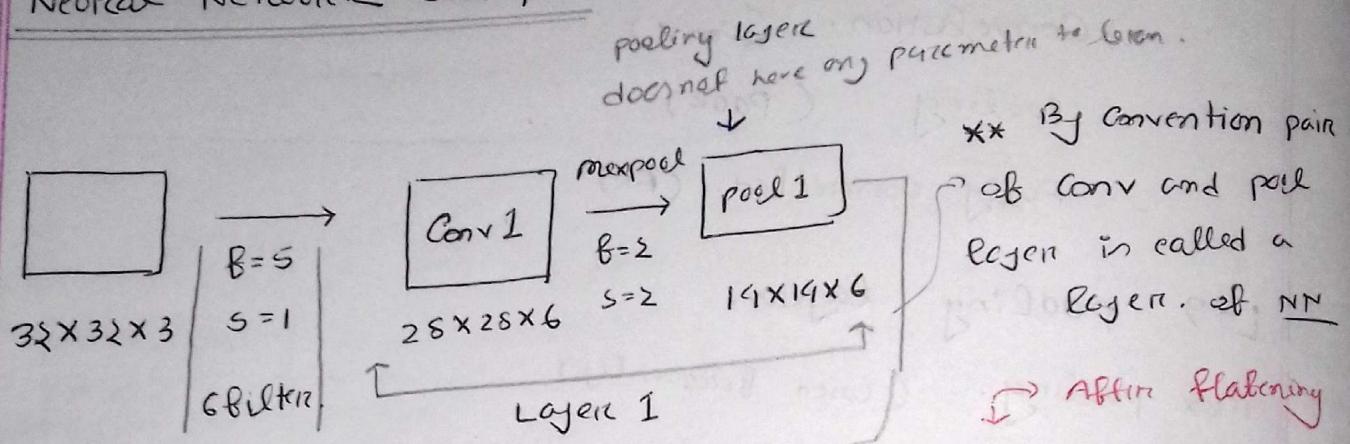
$$n_H \times n_W \times n_C$$

\downarrow (assume $P=0$)

$$\left\lceil \frac{n_H-p}{s} + 1 \right\rceil \times \left\lceil \frac{n_W-p}{s} + 1 \right\rceil \times n_C$$

** [pooling is applied in each channel independently]

Neural Network example. [LeNet-5].



Common pattern :- $\boxed{\text{Conv-pool, Conv-pool, fc-fc-fe - softmax}}$.

25x25
Paper (32x32)
2x4x4 (16x16)
2x2x8 (8x8)

** See the concrete example of how others do it and try to use it in your application.

	Activation Shape	Activation Size	# parameters
input:	(32, 32, 3)	3072	0
Conv1 ($B=5, S=1$)	(28, 28, 8) 8 filters	6272	208 $(25+1) \times 8$
POOL1 ($B=2, S=2$)	(14, 14, 8)	1568	0
Conv2 ($B=5, S=1$)	(10, 10, 16) 16 filters	1600	416 $(25+1) \times 16$
POOL2 $B=2, S=2$	(5, 5, 16)	400	0
FC3	(120, 1)	120	98001 - $(400 \times 120) + 1$
FC4	(84, 1)	84	10,081 - $(120 \times 84) + 1$
Softmax	(10, 1)	10	841 $(84 \times 10) + 1$

Observation

- ** As network goes deep activation size shrinks.
- ** No parameters to learn at pooling layers.
- ** Maximum parameters are in fully connected layers.

L-11 why convolutions? [6.45 pm 29.01.10].

Main advantages of convolution layers over fully connected layers acc.

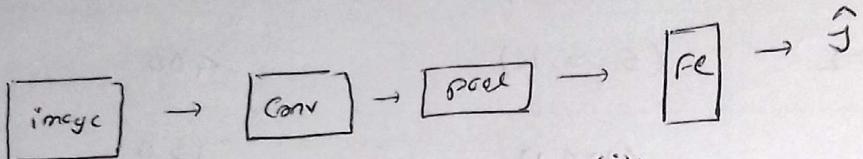
1. Parameters sharing :- A feature detector (such as vertical edge detector) that's useful in one part of image is probably useful in another part of the image.

2. Sparsity of connections :- In each layer, each output value depends only on a small number of inputs;

** Convolutional layer requires small number of parameters.
Compare to fully connected layer.

→ Translation invariance.

→ Training set $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$



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

use gradient descent to optimize parameters to reduce J .

→ Programming Assignment.

Conv-Forward() → \rightarrow \rightarrow calculate and implement convolution

In this assignment, we will implement

1. helper functions that will use when implementing tensorflow model.

2. Implement a fully connected convnet using tensorflow.

[Build and train a Convnet in Tensorflow for a classification problem]

** built-in functions for implementing backward propagation in tf

1. tf.nn.conv2d (x, w1, strides = [1, s,s,1], padding = 'SAME') :-

→ given an input x and group of filters $w1$, this function convolves $w1$'s filter on ' x '. Third input $[1, s,s,1]$ represents stride for each dimension of the input ($m, n_H_prev, n_W_prev, n_C_prev$).

2. tf.nn.max_pool (A, ksize = [1, k,k,1], strides = [1, s,s,1], padding = 'Same')
→ given an input A , this function uses a window of size (k, k) and strides of size (s, s) to carry out max pooling over each window.

3. tf.nn.relu(z1).

→ Computes elementwise relu of $z1$, which can be any shape.

4. tf.contrib.layers.flatten (P).

→ Given an input P , this function flattens each example into a 1D vector while maintaining batch size. It returns a flattened tensor with shape [batch-size, K] \rightarrow After flattening size of the vector.

5. tf.contrib.layers.baby_connected (F, num_outputs).

→ give the flattened input F , it returns the output computed by fully connected layer.

** (In fully connected layer we do not need to initialize weights, tensorflow model automatically initializes it).

6. `tf.nn.softmax_cross_entropy_with_logits` (`logits=23, labels=Y`) .

→ It calculates softmax entropy loss. This function both computes softmax activation function and resulting loss.

7. `tf.reduce_mean`

→ Computes mean of elements across dimension of a tensor.
use this to sum the losses over all the examples to get overall cost.

[open source Community to [long gone era] to get noticed]

P1: Case studies.

L-1: why look at case studies. [4.00 pm 26.01.19]

outline:

1. classic networks:

- LeNet-5.
 - Alex Net
 - VGG
- You can read research papers for the neural network architecture which may be useful for your task also.

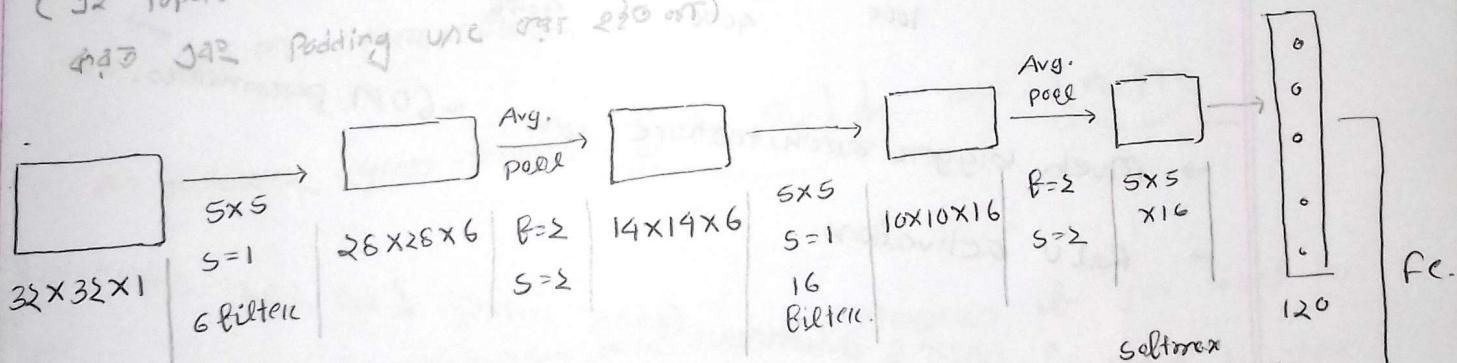
2. ResNet (Residual Network) [it has 152 layers]

3. Inception Neural Network.

L-2: classic networks.

→ paper year (1998)

1. LeNet-5. (useful for recognizing handwritten digit)

(paper year implemented our project with matrix Avg. pool use
and JAR Padding use our code)

** 60K parameters.

** As we go deeper ($(n_H, n_w) \rightarrow (n_C)$)

** Conv pool conv pool fc fc output.

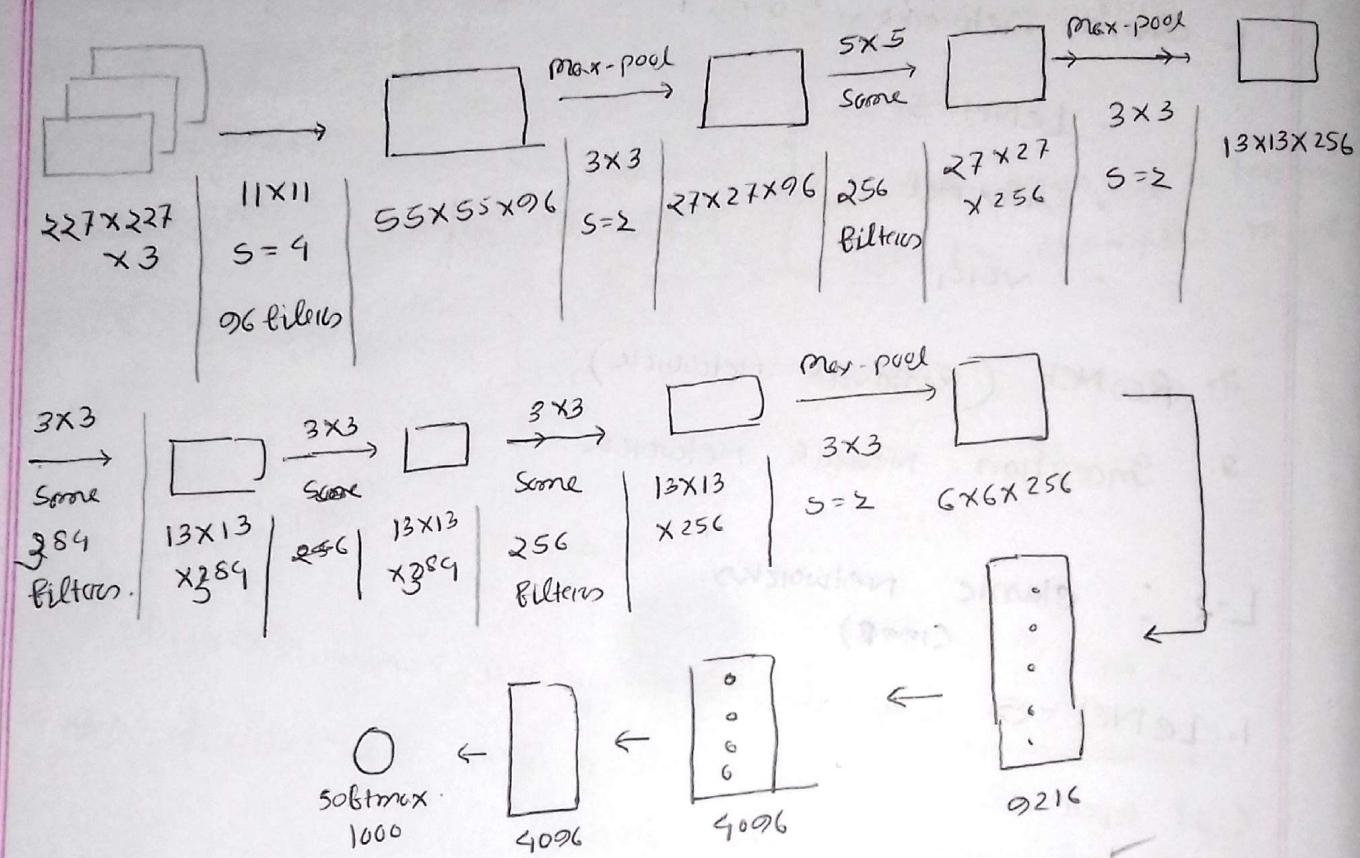
** Conv pool conv pool fc fc output.

[

→ ImageNet classification with deep Convolutional
neural networks.

2. AlexNet.

Architecture of AlexNet.



- Much bigger architecture
- ReLU activation.

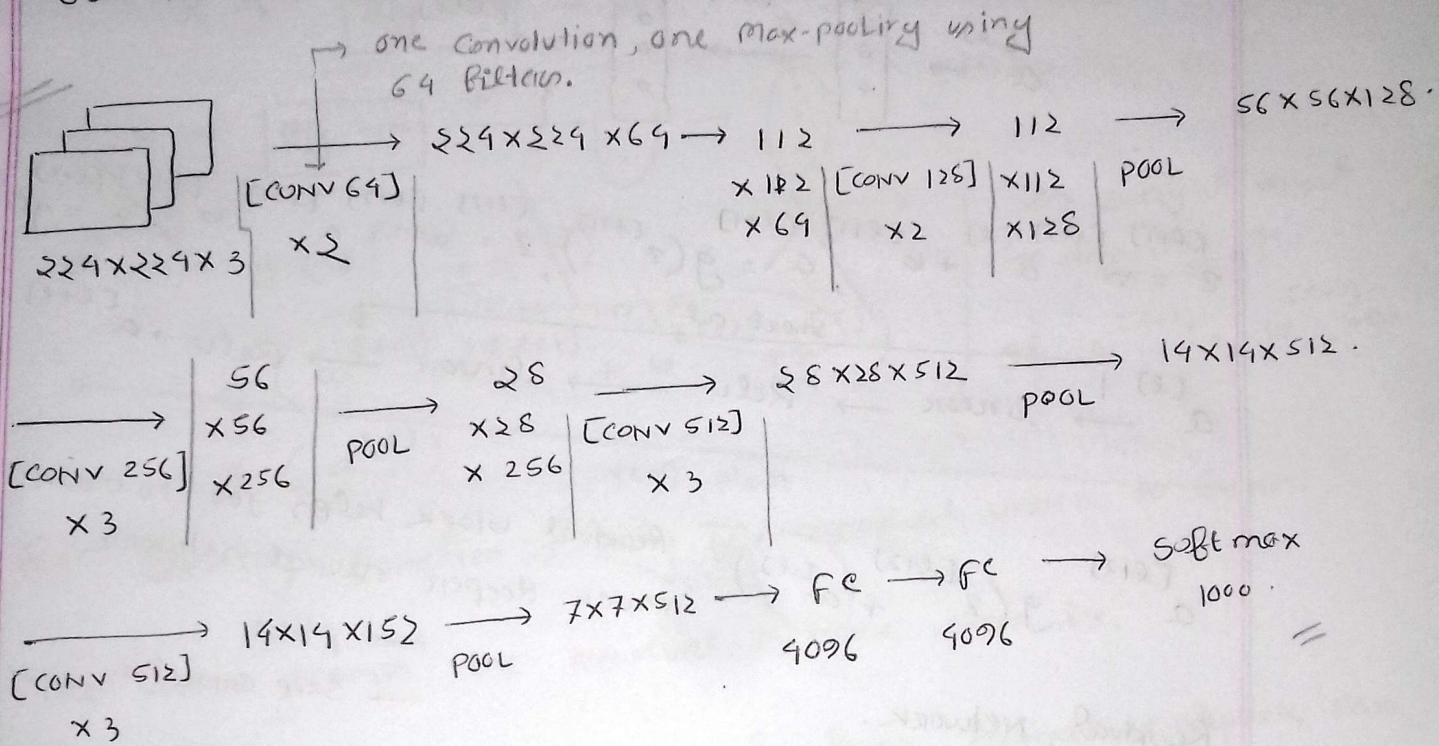
[Paper 2 uses 3 advanced features [ReLU, Dropout, L2 Norm]]

[This is a quite easy paper if you want to look at it
start with]

(Very Deep Convolutional Networks for Large-Scale Image Recognition).

3. VGG-16

CONV = 3×3 filters, $s=1$, same MAX-POOL = 2×2 , $s=2$.



* $\times 138M$ parameters.

* As network goes deeper

$[n_H, n_W] \downarrow$ and $[n_e] \uparrow$

[These three networks are classic networks if you want to get an idea about deep neural network architecture you can go through papers of these networks].

[useful to train a very deep neural networks]

$$a^{[l]} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \xrightarrow{a^{[l+1]}} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow a^{[l+2]}$$

$$z^{[l+1]} = w^{[l+1]} \cdot a^{[l]} + b^{[l+1]}, a^{[l+1]} = g(z^{[l+1]})$$

$$z^{[l+2]} = w^{[l+2]} \cdot a^{[l+1]} + b^{[l+2]}, a^{[l+2]} = g(z^{[l+2]})$$

~~short-cut / skip connection.~~

Main Path:

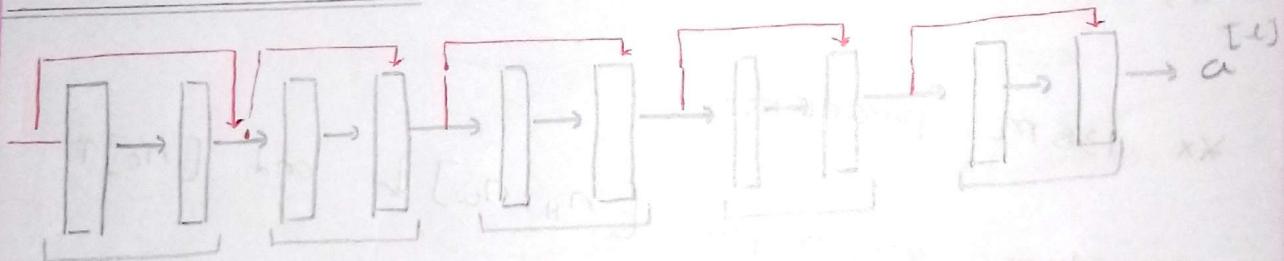
$$a^{[l]} \xrightarrow{\text{linear}} \text{ReLU} \xrightarrow{\text{linear}} \text{ReLU} \rightarrow a^{[l+2]}$$

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

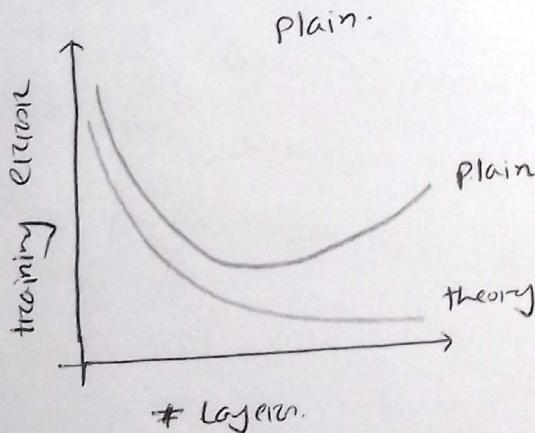
Residual block helps you to train much deeper neural network.

Residual Network.

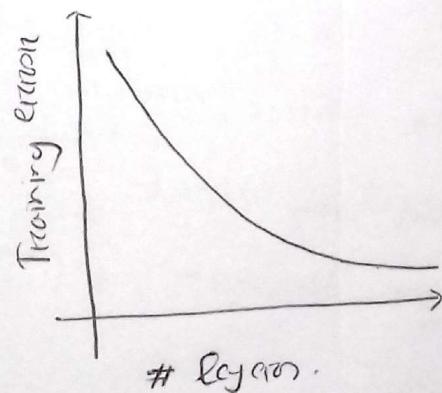
— plain
— skip connection



** Residual blocks.



ResNet



[This really helps with vanishing and exploding gradient problems and allows to train much deeper neural network]

why do feed-forward network work?

$$\begin{aligned}
 X &\rightarrow \boxed{\text{Big NN}} \rightarrow a^{[l]} \\
 X &\rightarrow \boxed{\text{Big NN}} \xrightarrow{a^{[l]}} \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \rightarrow \left(\begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right) \rightarrow a^{[l+2]}.
 \end{aligned}$$

$$\begin{aligned}
 a^{[l+2]} &= g(z^{[l+2]} + a^{[l]}) \\
 &= g(w^{[l+2]} \cdot a^{[l+1]} + b^{[l+2]} + a^{[l]}) \\
 &= g(a^{[l]})
 \end{aligned}$$

As weight must
be decreasing
so.
if $w^{[l+2]} = 0$, $b^{[l+2]} = 0$.

④ Identity function is easy for residual block to learn.

The main reason the residual network work is,

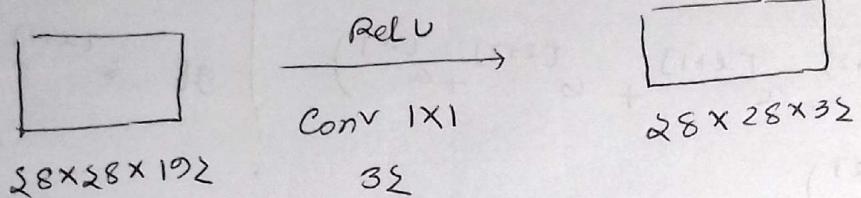
1. Extra layers of the identity function. residual network can learn identity function easily.
2. Do not hurt the performance of NN, also helps to improve.

[residual network has same input and output dimension 256
Some layer of size 224x224x3]

L-5 Networks in Networks and 1×1 Convolution. [3.50 pm 27.01.19]

** 1×1 convolutions are used to shrink number of channels. i.e.

** Pooling layers are used to shrink height (n_H) and width (n_W).

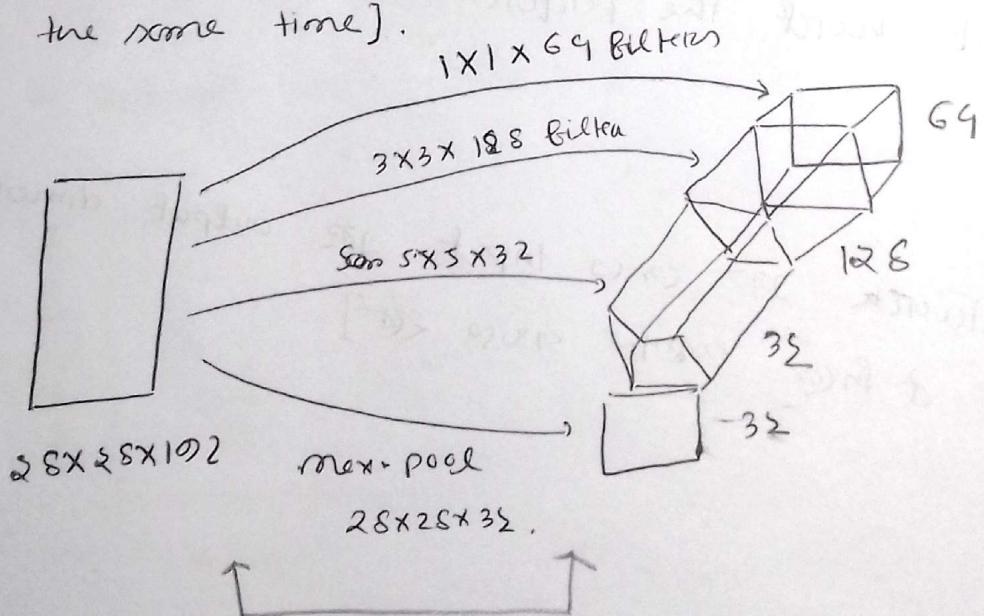


shrink the number of channel from 102 to 32.

→ Allows to increase, decrease or keep same number of channels of the volume

L-6 Inception Network motivation. [4.05 pm - 27.01.19]

The heart of inception network is that do convolutions with multiple size filters, max-pooling in one layer at the same time.

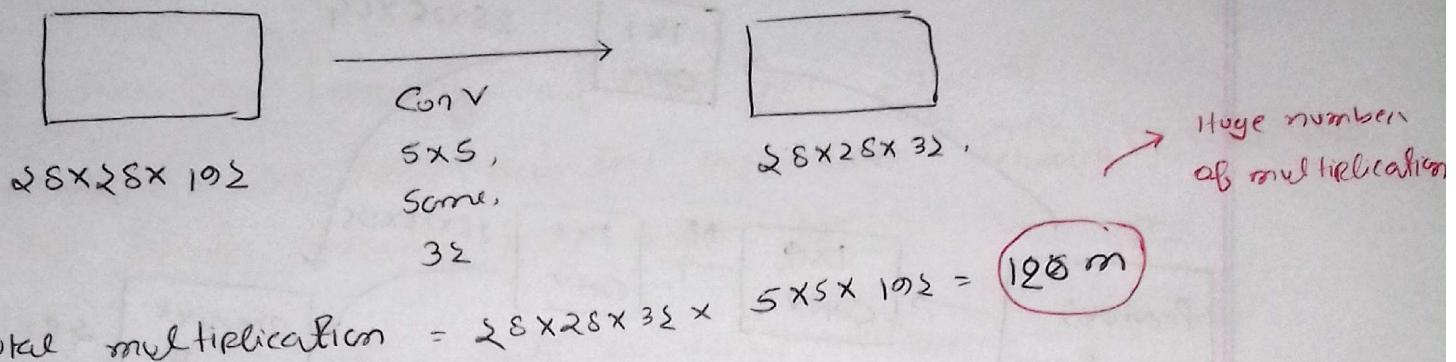


** To keep dimension same in input and output use same convolution.

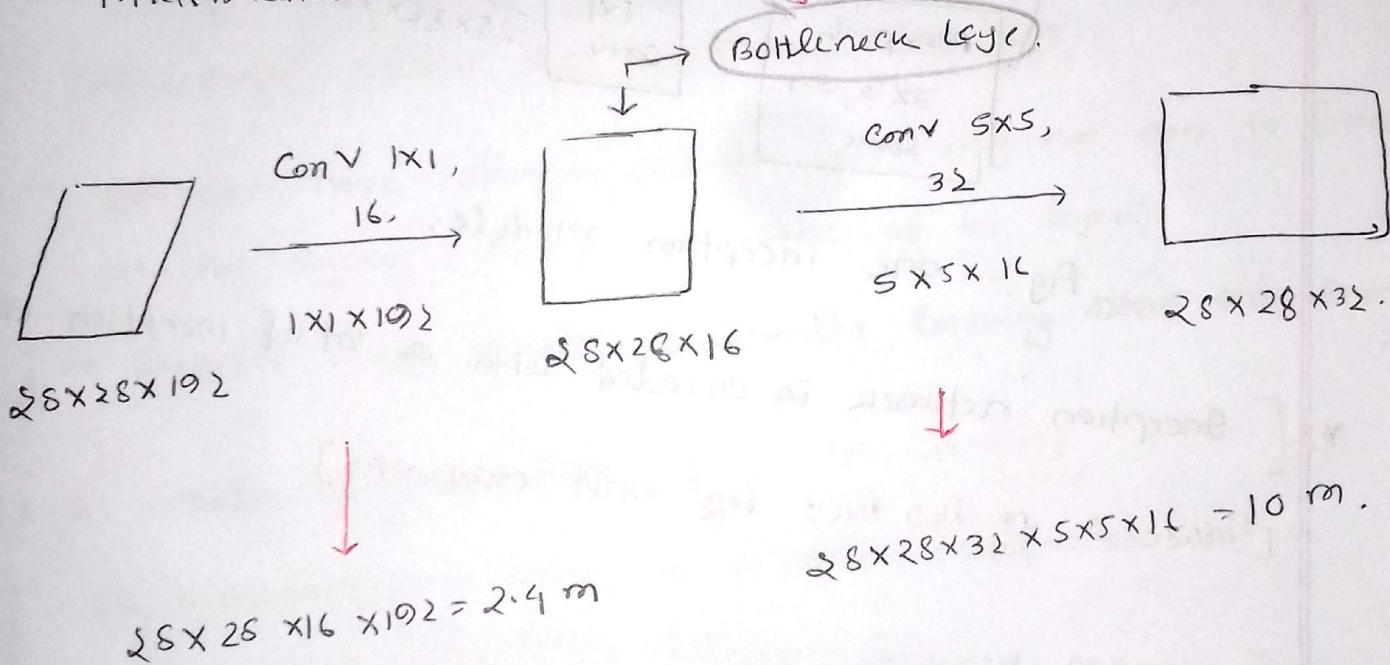
* Here we are combining and pooling

with multiple size of filters at the same time

[Computational cost is much higher].



** [This computation can be reduced by factor of 10 by using intermediate 1x1 convolutions.]



$$\text{Total multiplication} = 12.4 m.$$

[Shrinking the number of channel in bottleneck layer does not hurt the performance much].

Inception module.

(2017 01 26)

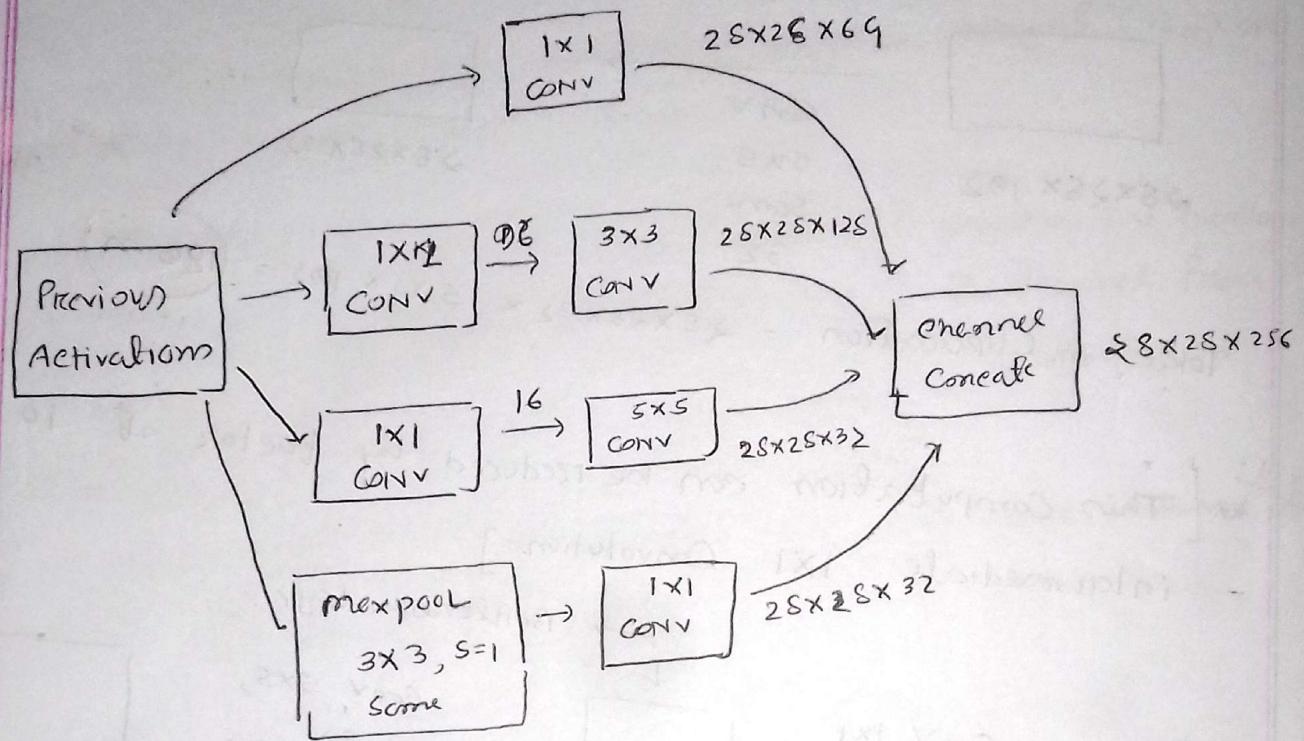


Fig:- one inception module.

* [Inception network is circled with a lot of inception block]

[slide 9 at 6:57 min for point number 27]

[inception network related papers etc]

P2: Practical advices for using ConvNets.

L-1 using open source implementation. [11:50 pm 27.01.19]

[open source implementation for fast use case].

→ open source implementation faster Project start out better.

L-2 Transfer Learning. [11:57 pm 27.01.19]

[** weights net take download one for fast network ↗ pre-train
one transfer learning faster than out.

** fast training net ↗ Data out weight pretrain weight 2nd
download one just train to softmax layer change ^{out} RNN

** There are several ways to freeze intermediate layers.

top-10% network one implement out 2nd.

→ if you have enough dataset you may not need to freeze
all the layers. (may freeze some of the layers).

[** Computer vision problem ↗ transfer learning faster or possibly =
=

L-3 Data Augmentation. [12:10 Am 28.01.19]

[data augmentation is often used to improve performance of a
Computer vision model]

Common Data augmentation method.

1. mirroring.

2. Random cropping.

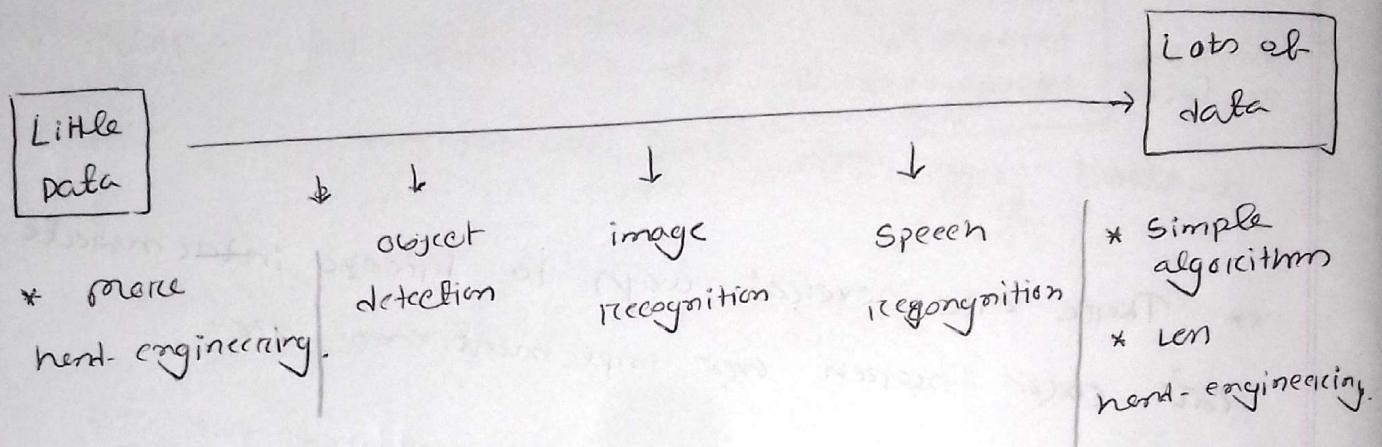
3. rotation.

4. shearing e.t.c.

** Color shifting is another data augmentation method.
Changing the (R, G, B) values of the pictures.

L-4 State of computer vision [12.30 pm 28.01.19]

Data vs. hand engineering.



Two sources of knowledge.

- Labeled data.
- Hand engineered features / network architecture / other components.

→ If we have little data transfer learning may help.
→ If we have little data transfer learning may help.
** Tip: For doing well on benchmarks / winning competitions.

1. Ensembling.

- Train several networks independently and average their outputs.

2. Multi-Crop at test time.

- Run classifiers on multiple versions of test images and average results.

3. Use open source code.

- Use architecture of networks published in the literature.
- Use open source implementation if possible.
- Use pretrained models and fine-tune on your dataset.

Programming Assignment :-

• will learn how to use keras.

→ keras is even higher level language than tensorflow.

** To train and test a model, there are 4 steps in keras,

1. Create the model.

2. Compile the model calling,

model.compile(optimizer="...", loss="...", metrics=["accuracy"])

model.compile(optimizer="...", loss="...", metrics=["accuracy"])

model.compile(optimizer="...", loss="...", metrics=["accuracy"])

3. Train the model on train data by calling,

model.fit(x="....", y="....", epochs=..., batch_size=...).

model.fit(x="....", y="....", epochs=..., batch_size=...).

model.fit(x="....", y="....", epochs=..., batch_size=...).

4. Test the model on test data by calling

model.evaluate(x=..., y=...)

model.evaluate(x=..., y=...)

[ResNet] implementation without a config) .

ResNet-50 model: (Computer vision & APPLY possible)

- Zero-padding pads the input with a pad of $(3, 3)$

Stage 1:-

- The 2D convolution has 64 filters of shape $(7, 7)$ and uses a stride of $(2, 2)$. name conv 1.
- BatchNorm is applied to channel axis of the input.
- Maxpooling with $(3, 3)$ window and $(2, 2)$ stride.

Stage 2:-

- The convolutional block uses three set of filters of size $[64, 64, 256]$, "B" is 3, "S" is 1 and block is "a".
- The 2 identity block uses three set of filters of size $[64, 64, 256]$, "B" is 3, and block size "b" and c.

Stage 3:-

- The convolutional block uses three set of filters of size $[128, 128, 512]$, "B" is 3, "S" is 2 and block is "a".
- The 3 identity block uses three set of filters of size $[128, 128, 512]$, "B" is 3 and block are "b", "c", "d"

Stage 4:-

- The convolutional block uses three set of filters of size [256, 256, 1024], "B" is 3, "S" is 2 and block is "a".
- The S identity block uses three set of filters [256, 256, 1024] "F" is 3, blocks are "b", "c", "d", "e", "f".

Stage 5:-

- The convolutional block uses three set of filters of size [512, 512, 2048], "F" is 3, "S" is 2 and block is "a".
- The 2 identity block uses three set of filters of size [512, 512, 2048], "F" is 3 and blocks are "b" and "c".
- The 2D Average Pooling uses a window of nhepc (2,2) and its name is avg-pool.
- . Then flatten the output.

- . The fully connected layer reduces its input to the number of classes using a softmax activation.

[Implementation details of ResNet-50]

v.1.9 v.1.9
** [deep-learning] model repository of github

1. VGG16
2. VGG19
3. ResNet50

4. Inception v3.
5. CRNN for music tagging.

own implement
own data

week-3 object detection.

L-1 object localization [11.00 pm 30.01.19].

[Localization is the process of finding any object within an image] → we have to put a bounding box around the object.

classification :- can classify an image based on some categories such as whether an image contains a car or not.

classification with localization :- Find an object within the image. There will be only one image somewhere middle of the image. and put a bounding box around the image.

Detection :- Localize multiple image object of multiple category within an image,

[slide 12 at 57:21 - 57:25 are the first]

Defining the target label J .

1. pedestrian	Need to output
2. car.	b_x, b_y, b_h, b_w ,
3. motorcycle.	class label (1-4)
4. background.	

$$Y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \rightarrow \begin{array}{l} \text{Is there any object?} \\ \text{boundary box parameters.} \\ \text{class parameters.} \end{array}$$

$$\mathcal{L}(\hat{J} - J) = \begin{cases} (\hat{j}_1 - j_1)^v + (\hat{j}_2 - j_2)^v + \dots + (\hat{j}_8 - j_8)^v & \text{if } J_1 = 1. \\ (\hat{j}_1 - j_1)^v & \text{if } J_1 = 0. \end{cases} \quad \begin{array}{l} \uparrow \text{number of parameters} \end{array}$$

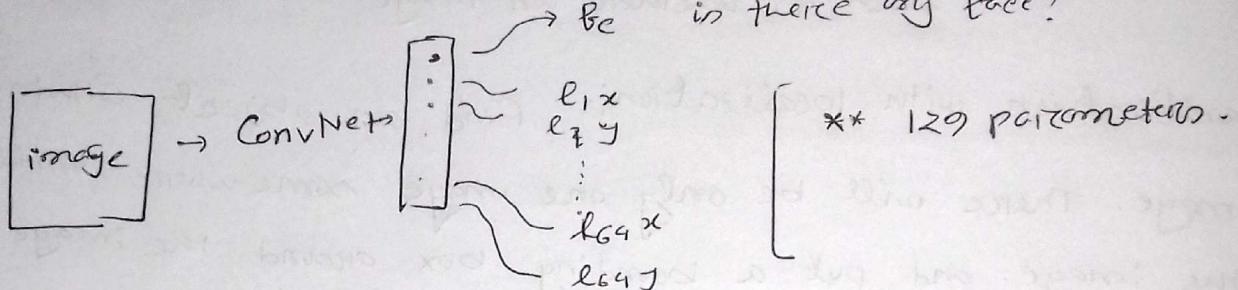
L-2 Landmark Detection

[11.25 pm 30.01.19]

[we can find the landmarks of the face within an image by using neural networks].

Suppose we have to find 64 landmark within an image

is there any face?



** people pose detection using Landmarks [recognish].
[Bind emotion from picture]

L-3 object detection

[11.35 pm 30.01.19]

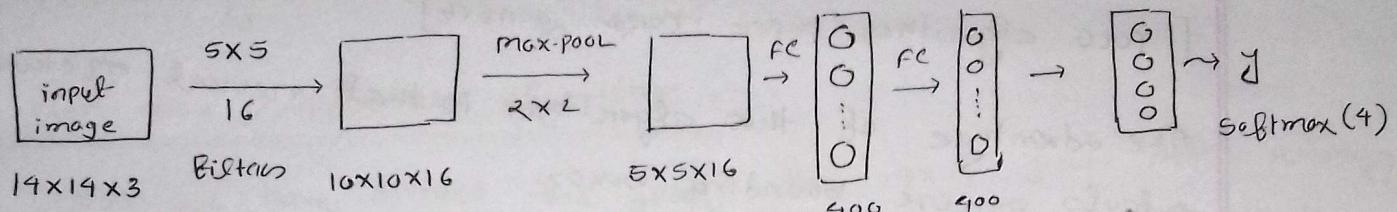
[object detection using sliding window]

- Sliding window is important & diff. than training image for each crop crop train crop.
- Take a window of fixed size ~~and~~ and slide it over the image with some stride for almost every part of the region.
- Using this process computational cost is really high. because we have to crop the image many times and pass it to the ConvNet.
- This problem can be solved.

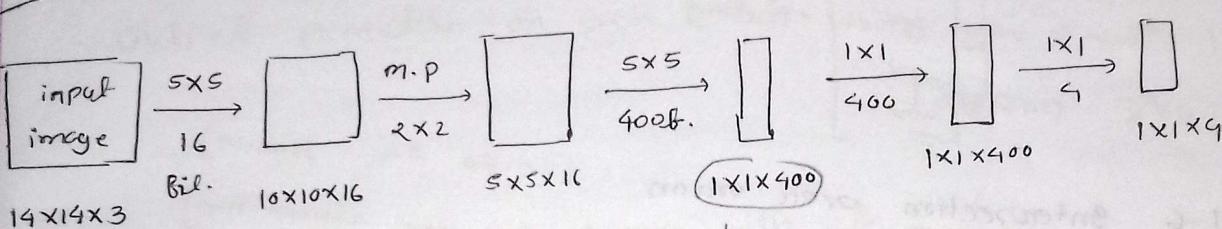
L-4 Convolutional implementation of sliding window.

[10:45 pm 30.01.17]

Turning FC layers into convolutional layers:-



→ Alternative



** Alternative implementation of FC layers using Convolution.

↳ This FC layer can be implemented using convolutions.

↳ mathematically it is equivalent to 400 fully connected unit

Convolutional implementation of Sliding window:-

[Sliding window shares a lot of computations between windows; convolutional implementation of sliding window do not recompute the shared computations as a result it become quite faster.]

[graph] slide \cup graph fnct \cup graph graph as graph or

⇒ [This algorithm has one problem that position of the bounding boxes will be not so accurate].

L-5 Bounding box prediction [1.00 Am 31.01.19]

[convNet does not quite accurately predict bounding box]

[YOLO algorithm paper taken from]

The advantage of this algorithm is that neural network outputs precise bounding boxes.

[image will be divided into grid cells.]

[This is very much efficient]

[IoU metric]

L-6 Intersection over union

[gt is used for evaluating object detection algorithm]

Intersection over union (IoU)

$$= \frac{\text{Size of intersection area}}{\text{Size of union area}}$$

"Correct" if $\text{IoU} \geq 0.5$.

gt is measure that shows whether an object is correctly localized or not.

Iou is a measure of the overlap between two bounding boxes.

[Non-max suppression ensures that an algorithm detects each object only once within an image.]

Non-max suppression algorithm :-

Let's we have divided the ^{image} grid into 19×19 grid. After running algorithm

Output prediction on each cell:

P_c
b_x
b_y
b_w
b_h

* output is considering that we have to detect only one object.

Now,

1. → discarded all the boxes with $P_c \leq 0.6$.

2. while there are any remaining boxes :-

→ Pick the box with largest P_c . output that as a prediction.

→ Discard any remaining box with $IoU \geq 0.5$ with the box output in the previous step.

[If any Predicted bounding box or after high overlap enter with another box before 2nd boundary of object it may detect twice, at other get one first of]

[As there are multiple object in an image we have to apply non-max suppression for each algorithm independently]

[Anchor boxes are predefined shapes of different objects].

Anchor box algorithm :-

Previously :- Each object in training is assigned to grid cell that contains that object's midpoint.

$$\text{So, output } J = 3 \times 3 \times 8$$

Consider 9 grid cells with image can have 3 objects

with two anchor boxes :-

Each object in training is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.

[object or anchor box is similar to the bounding box]

$$\text{output } J = 3 \times 3 \times 8 \times \text{number of anchor box}$$

** we can define any number of anchor box.

** [If two objects have same midpoint within same grid cell then anchor boxes help us to find the accurate bounding box for these objects].

~~** It is not work so well when number of anchor boxes are less than number of object in an image.~~

**

→ YOLO algorithm. [11.45 AM 31-01-19]

Training :- [Sliding ⌂ for finer object map output]

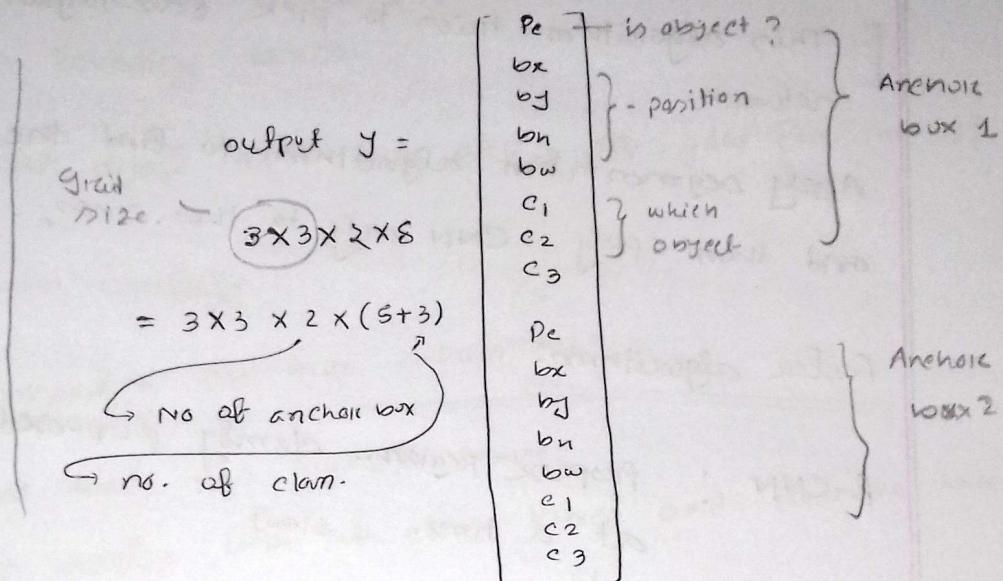
Suppose we have to detect 3 objects with 2 anchor boxes.

objects :-

1. Pedestrian

2. Car

3. motorcycle.



Making predictions :-

[stride is denoted earlier].

Output the non-max suppression outputs.

1. For each cell, get 2 (# no of anchor box) predicted boundary boxes.
2. Get rid of low probability predictions.
3. For each class (pedestrian, car, motorcycle) use non-max suppression to generate final prediction.

[That is YOLO object detection algorithm]

L-10 Region proposals. [11.55 AM 31.01.19]

Region proposal : R-CNN

[This algorithm tries to pick few regions to run convolutional network.

Apply segmentation algorithm to find the area of interest and then apply CNN only to this areas.

Faster algorithm.

R-CNN : propose regions. classify proposed regions one at a time. [slow]

Fast R-CNN : Propose regions. use convolution implementation of sliding windows to classify all proposed regions.
[time consume to propose regions].

Faster R-CNN : use convolutional network to propose regions.

[YOLO algorithm is much faster]. (very).

Programming Assignment:-

1. use object detection on a car detection dataset.
2. Deal with bounding boxes.
(notebook with detail about car recc. is good for reading)

→ calculate dimension carefully

steps of implementing non-max suppression :-

1. Select the box that has highest score.
2. Compute its overlap with all other boxes and remove boxes that overlap is more than iou-threshold.
3. Go back to step 1 and iterate until there is no more boxes with a lower score than the current selected box.

To implement this temporary have two built in function.

- tf.image.non_max_suppression()
- K.gather()

[Task 1] Implement model implementation [using chainer]

[Task 2] Implement model using tensorflow [using eager execution]

↳ for coding help refer to this video

** Important github repository link that one.

Summary of YOLO

1. Input image $(608, 608, 3)$.
2. The input image goes through a CNN, resulting in a $(19, 19, 5, 85)$ dimensional output.
3. After flattening the last two dimensions, the output is a volume of shape $(19, 19, 425)$.
 - Each cell in a 19×19 grid over the input image gives 425 numbers.
 - $425 = 5 \times 85$, because each cell contains prediction of 5 for 5 anchor boxes.
 - $85 = 5 + 80$ where 5 is because $(p_c, b_x, b_y, l_w, l_h)$ has 5 numbers and 80 is the number class we want to detect.
4. Select bounding boxes based on,
 - Score thresholding :- throw away boxes that have detected a class with a score less than the threshold.
 - Non-max Suppression :- Compute intersection over union and avoid selecting overlapping boxes.
[This is YOLO (You only look once) algorithm]

Week-4: Special applications: Face recognition & Neural style transfer.

P1: Face Recognition.

L-1 what is face recognition. [10.20 AM 01.02.19]

Face verification vs. Face recognition.

Verification:-

- Input image, name/ID.
- output whether the input image is that of the claimed person.

(Recognition is much harder than verification).

Recognition:-

- Has a database of K persons.
- Get an input image.
- output ID if the image is any of the K persons or "not recognized".

(For correct recognition accuracy of verification step must be higher).

L-2 one shot learning. [10.30 AM 01.02.19]

[One shot learning problem is that our system must be able to recognize a person given just one single image].

** It is difficult to do well on a deep learning problem with only one example of a class.

→ Learning from one example to recognize the person again.

** If a new employee joins the number of output will be changed and we have to retrain our whole system.

[It is not a good way].

Learning a "similarity" function.

$d(\text{img1}, \text{img2})$ = degree of difference between images

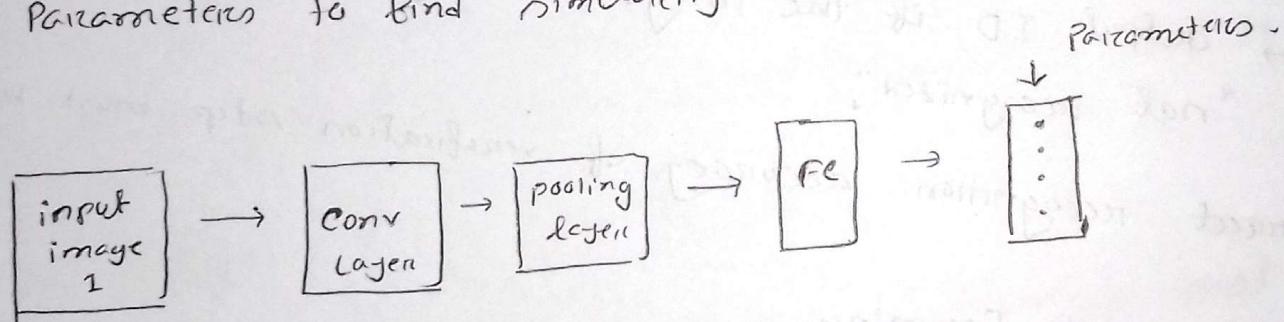
if $d(\text{img1}, \text{img2}) < \text{threshold}$ → "some"

$d(\text{img1}, \text{img2}) > \text{threshold}$ → "different"

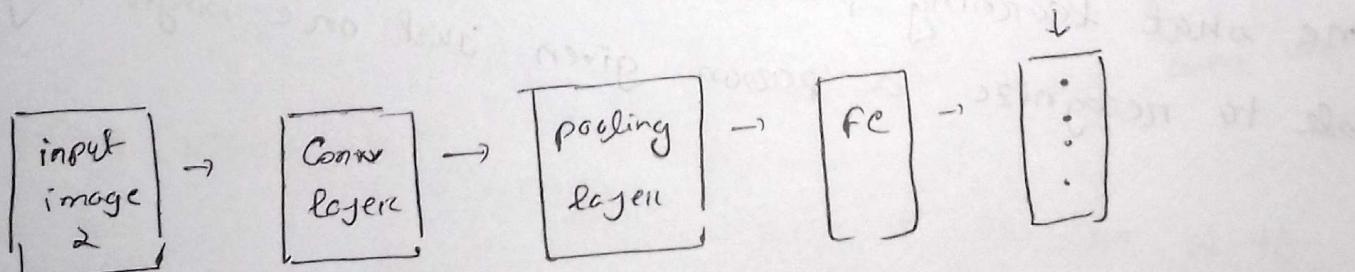
** This similarity function solves the problem of one-shot learning.

L-3 Siamese Network. [10:45 AM 01.02.17]

[This network is used to train similarity function]
Hence we have encode each image by definite number of parameters to bind similarity function.



"encoding of $x^{(1)}$ " = $B(x^{(1)})$.



"encoding of $x^{(2)}$ " = $B(x^{(2)})$

Fig: Siamese Network

Goal of Learning :-

Parameters of NN define an encoding $f(x^{(i)})$ this is encoding of $x^{(i)}$

Learn parameters so that :

→ if $x^{(i)}, x^{(j)}$ are the same person, $\|f(x^{(i)}) - f(x^{(j)})\|^2$ is small.

→ if $x^{(i)}, x^{(j)}$ are different persons, $\|f(x^{(i)}) - f(x^{(j)})\|^2$ is large.

L-4 Triplet loss. [11.00 AM 01.02.19]

[At each time we have to look at three images anchor, positive and negative. This is where triplet terminology comes from].

Objective Function :-

$$\|f(A) - f(P)\|^2 \leq \|f(A) - f(N)\|^2 \\ d(A, P)$$

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 \leq 0 \\ \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 \leq 0$$

if $f(A), f(P), f(N)$ are identical then $f(\text{img}) = \vec{0}$.
To avoid this we have to add a margin (α) so that positive image has certain degree of similarity greater than threshold.

so,

$$\|f(A) - f(P)\|^2 + \underline{\alpha} \leq \|f(A) - f(N)\|^2$$

$$\therefore \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \underline{\alpha} \leq 0.$$

Loss Function :-

Given 3 images A, P, N .

$$L(A, P, N) = \max \left(\|f(A) - f(P)\|_F^2 - \|f(A) - f(N)\|_F^2 + \alpha, 0 \right)$$

$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

Choosing the triplets A, P, N :
 $\begin{cases} (A, P) \rightarrow \text{same person} \\ (A, N) \rightarrow \text{different person} \end{cases}$

During training if (A, P, N) are chosen randomly

$d(A, P) + \alpha \leq d(A, N)$ is easily satisfied. As a result our network will not learn much.

so we have to choose triplet in such way that are hard to train on.

$$d(A, P) \approx d(A, N)$$

if we do that computational efficiency will increase.

popular way of naming algorithm in deep learning world.

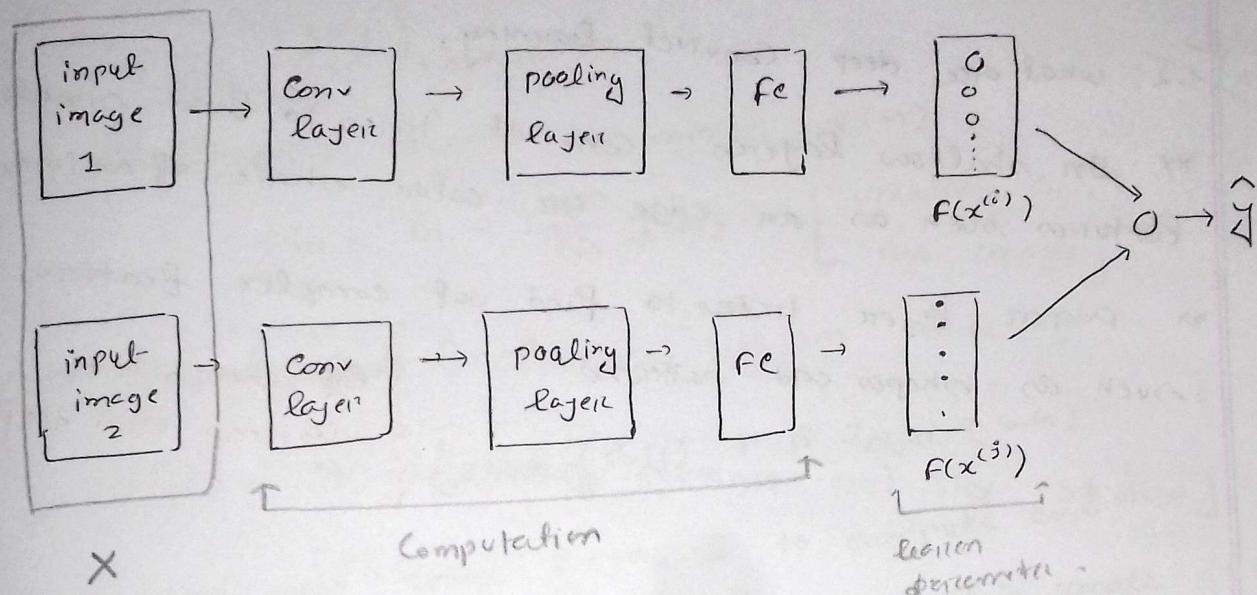
1. Net

2. Deep —

[online \rightarrow pretrained model weight w_{init}].

L-5 Face verification and Binary classification. [11.40 Am 01-02-19]

[Similarity function can also be learnt by using logistic Regression].



$$\hat{y} = \begin{cases} 1 & \text{if } x^{(i)} \text{ and } x^{(j)} \text{ are same person.} \\ 0 & \text{if } x^{(i)} \text{ and } x^{(j)} \text{ different person.} \end{cases}$$

$$\hat{y} = \sigma \left(\sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(j)})_k| + b_i \right)$$

→ number of parameters.

$$\text{This can also replaced by } = \frac{(f(x^{(i)})_k - f(x^{(j)})_k)}{f(x^{(i)})_k + f(x^{(j)})_k} \quad \boxed{\checkmark}$$

P2 : Neural style Transfer

L1 : what is neural style transfer. [12:30 pm 01-02-19]

[what Content image and style image do and how represent cat].

leads to important L2 : what are deep conv net learning. [12:40 pm 01-02-19]

- * In shallow layers conv net tries to find simple features such as an edge or color shade of an image.
- * Deeper layers tries to find out complex features such as shapes and patterns.

leads to ref from right hand side

face grayscale
radio接收器

L3 : cost function [12:50 pm 01-02-19].

for neural style transfer we have to formulate a cost function where from content image (c) and style image (s) we will generate image (g_1).

cost of generated image (g_1).

$$J(g_1) = \alpha J_{\text{content}}(c, g_1) + \beta J_{\text{style}}(s, g_1)$$

↓
How similar content image and generated image

↓
How similar content image and style image

α, β hyperparameters for relative weighting of content image and style image

find the generated image G_1 .

1. Initialize C_1 randomly.

$G_1 : 100 \times 100 \times 3$.

2. Use gradient descent to minimize $J(\theta)$

$$G_1 = G - \frac{\sigma}{\sigma_G} J(G). \quad \left[\begin{array}{l} \text{updating each pixel of} \\ \text{the image} \end{array} \right]$$

L-4 Content Cost function. [1.05 pm 01.02.10]

$$y(g_1) = \alpha \text{ Jcontent}(c, g_1) + \beta \text{ Jsingle}(s, g_1).$$

→ say you use hidden layer l to compute cost.
 → choose in middle if l is small

Say you use hidden layer l to compute
- [value of l is choose in middle if l is small it will
force generated image similar to content image. if l is large
generated image will not like content image].
(VGG network).

→ use pre-trained

Conv NCF (E.g.)

$\forall G_1 G_2$ network) .

→ Use pre-trained ConvNet
 Let $a^{[L](c)}$ and $a^{[L](g)}$ be the activation of layer L on

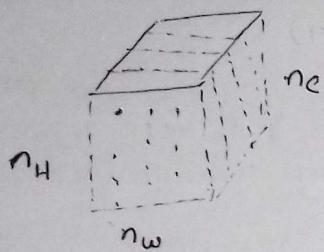
→ Let a & b be two images. Then both images must have similar features.

\rightarrow If $a^{[e]}(c)$ and $a^{[e]}(c_1)$ are given, then
 similar context. \rightarrow elementwise $a^{[e]}(c) \parallel a^{[e]}(c_1)$

$$[\text{Jointant } (C, G) =]$$

L-5 style cost function [7.45 PM 01.02.19]

Suppose you are using layer l 's activation to measure "style". Define style as correlation between activation across channels.



How correlated are the activations across different channels?

→ [Degree of correlation tells how often the high level feature occurs together and don't occur together in different parts of an image]

Style matrix :-

Let $a_{i,j,k}^{[l]}$ = activation at (i, j, k)

$n_H \quad n_w \quad n_c$
↓ ↓ ↓

(style matrix)

$G_l^{[l]}$ is $n_c^{[l]} \times n_c^{[l]}$ (style matrix)

$$G_{kk'}^{[l]} = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} a_{ijk'}^{[l]}$$

$[k, k' = 1, \dots, n_c^{[l]}]$

style image

Generalized
image.

$$G_{kk'}^{[l]} = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} a_{ijk'}^{[l]}$$

Style cost function :-

$$J_{\text{style}}^{[e]}(s, g) = \frac{1}{(2^{n_H} \cdot n_w \cdot n_c)} \sum_k \sum_{k'}^{\text{number of channel}} (g_{kk'}^{[e](s)} - g_{kk'}^{[e](g)})^2$$

We will get visually pleasant result if we take style cost function from different layers, \rightarrow some hyperparameters.

$$J_{\text{style}}(s, g) = \sum_e \lambda^{[e]} \cdot J_{\text{style}}^{[e]}(s, g)$$

So overall cost function:-

$$J(g) = \alpha J_{\text{content}}(s, g) + \beta J_{\text{style}}(s, g).$$

Now apply algorithm to minimize cost.

L-6 1D and 3D Generalization.

[2D \Rightarrow 25³ 1D & 3 3D data \Rightarrow 3 Convolutional network apply
cost 2725].

[1D data of any recurrent data neural network]

Programming Assignment:-

- Implement the neural style transfer algorithm.
- Implement novel artistic images using your algorithm.

The idea of using a network trained on a different task and applying it to a new task is called transfer learning.
Neural style transfer (NST) uses this technique.

Steps of a neural style transfer:-

1. Create an interactive session.
2. Load the content image.
3. Load the style image
4. Randomly initialize the image to be generated.
5. Load the VGG-16 model.
6. Build tensorflow graph.
 - Run the content image through the VGG-16 model and compute content cost.
 - Run the style image through the VGG-16 model and compute style cost.
 - Compute total cost
 - Define optimizer and learning rate.
7. Initialize the tensorflow graph and run it for a large number of iterations, updating generated image at every step.

** Previously train model use crop face recognition or ROI.

what you should remember:

1. Face verification solves an easier 1:1 matching problem.
Face recognition addresses a harder 1:k matching problem.
2. Triplet loss is an effective loss function for training a neural network to learn an encoding of a face image.

3.

view [FaceNet and DeepFace are two important papers for face verification and recognition.]

** Also look at official repository of FaceNet implementation.

[Completed
7.45 pm 02.02.19]