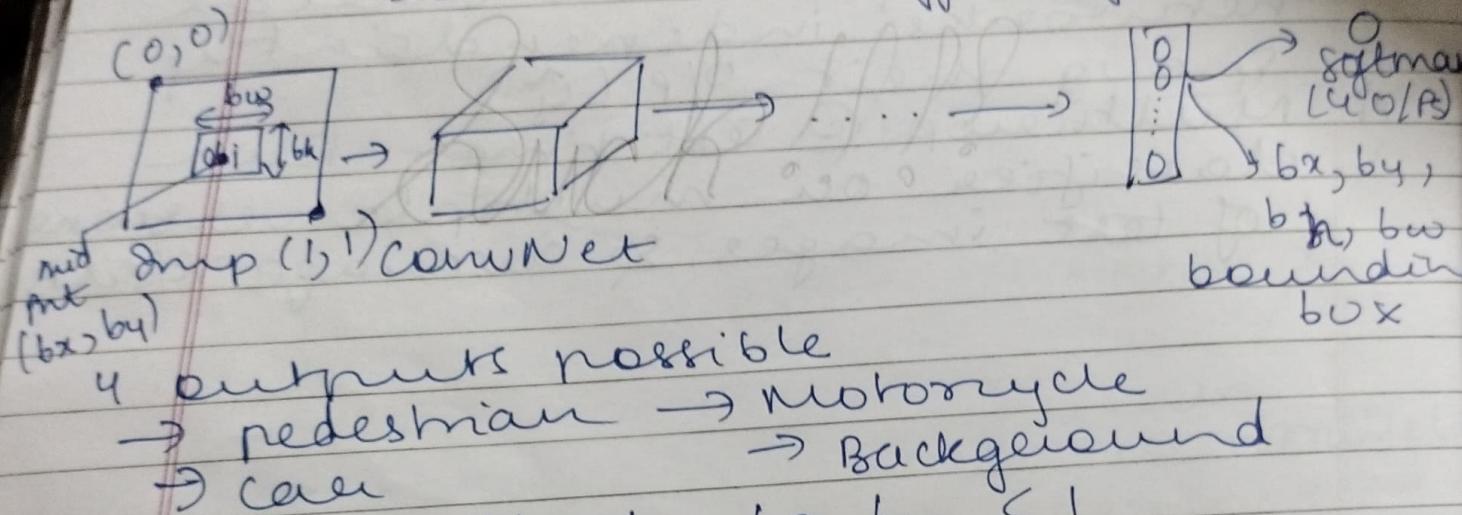


OBJECT LOCALIZATION

PAGE NO. / / /
DATE / / /

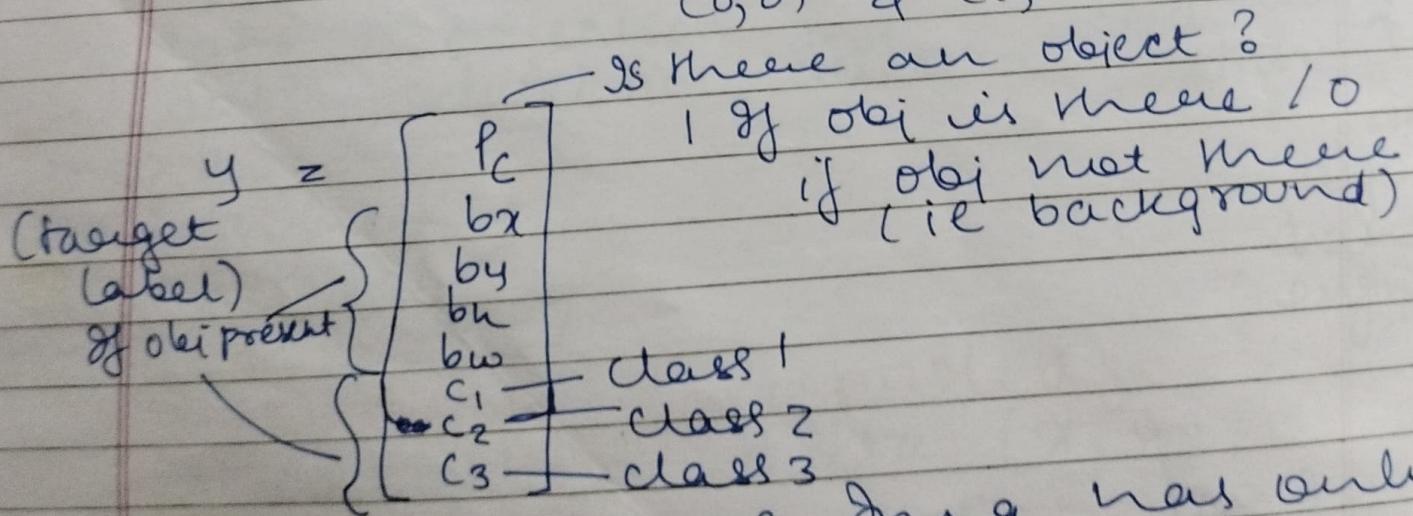
Image classification
classification with localization
Detection ?
Multiple objects of diff. types



- 4 outputs possible
 → pedestrian → motorcycle
 → car → Background

$$b_x, b_y, b_w \leq 1$$

(since opp. corners are $(0, 0)$ & $(1, 1)$)



We assume img has only 1 object

Eg if car in image

$$y = \begin{cases} 1 \\ b_x \\ b_y \\ b_w \\ 0 \end{cases}$$

Class 2 is car

if no obj

$$\hat{y} = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ \vdots \\ ? \end{bmatrix}$$

we don't care about the values as no obj found

$$\text{if } y_i = 1 \quad \delta(\hat{y}, y) = \sqrt{(\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2}$$

(y_1 to y_8 are the elements of y)

$$\text{if } y_i = 0 \quad \delta(\hat{y}, y) = (\hat{y}_i - y_i)^2 \quad \text{since we don't care of rest of the values}$$

- Here we squared to simplify but in practice we use a log function like softmax loss from c_1, c_2, c_3 to the softmax O/P
- We can use squared error loss something like squared error for bounding boxes
- If for P_c we can use logistic regression loss.

Landmark Detection

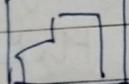
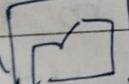
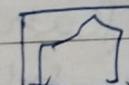
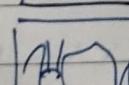
- If we want to find position for particular feature in an object we have the final layer output all the (x, y) coordinates along the

border of the landmarks

4) OBJECT DETECTION

→ using convnet to perform object detection using sliding windows detection algorithm

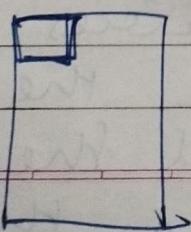
→ make a training set where the images containing the object to be detected are closely cropped ie the object takes up most of the img training set:

	x	
car		1
car		1
car		1
not car		0
not car		0

Used to train convnet to detect cars
Then we use it in sliding windows detection

→ Select a window of certain small size 

→ input it into the img

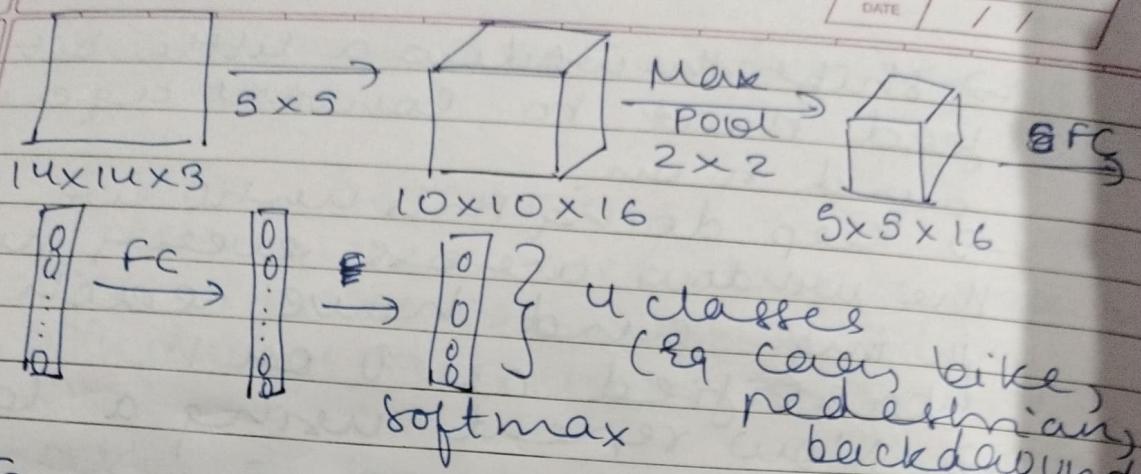


→ have the convnet make a prediction

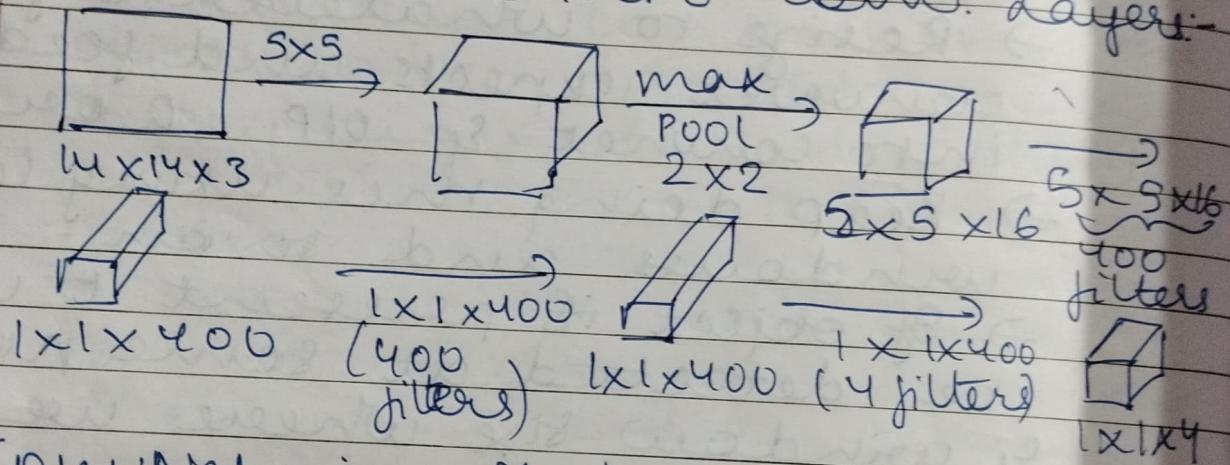
- Shift the window a little bit & feed that to convnet again and so on
- Keep doing this until we slide the window across every position in img. and have each position classified as 0 or 1
- Now repeat using a larger window
- Resize to whatever size convNet expects and feed it into convnet & O/P 0 or 1
- Keep doing this using larger windows and so on...
- If object is present it will be detected at some position & window size where we will get O/P as 1 ∞ 1
- Huge disadvantage of sliding windows detection is computational cost
- Using large slide hurts performance but small stride increases comp. cost

Convolutional implementation of Sliding Windows

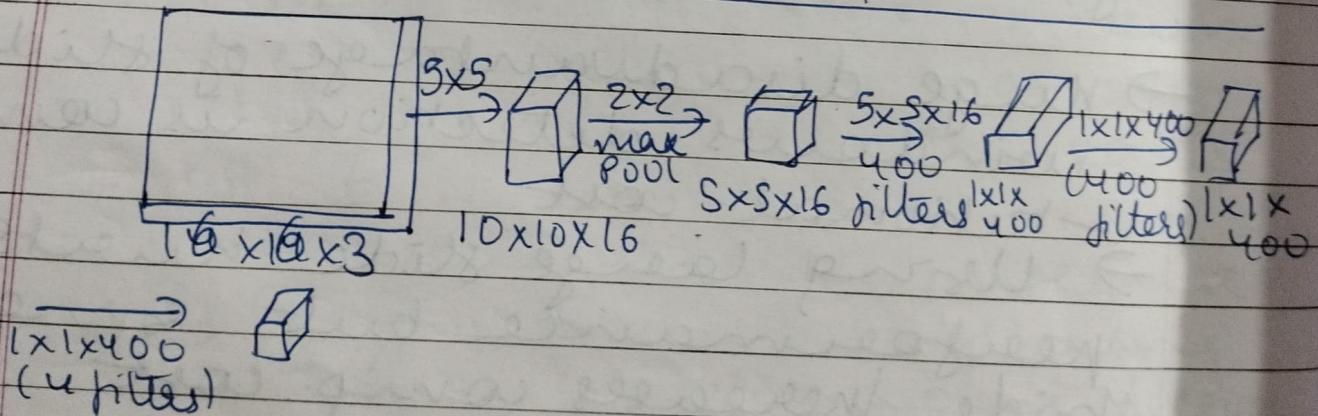
Training fully connected layers to convolutional layers



To make FC into conv. layers:

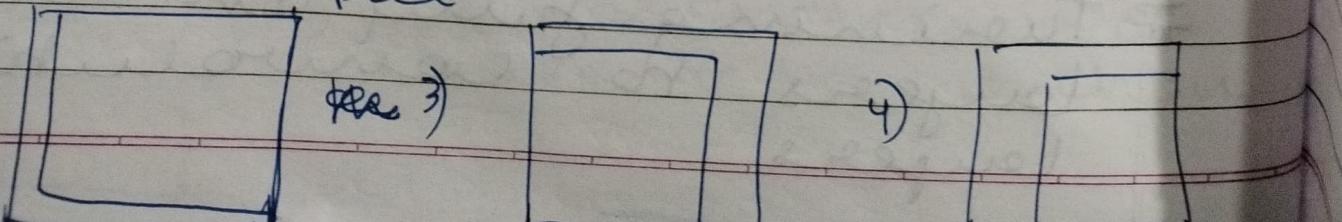


~~ConvNet in sliding window~~

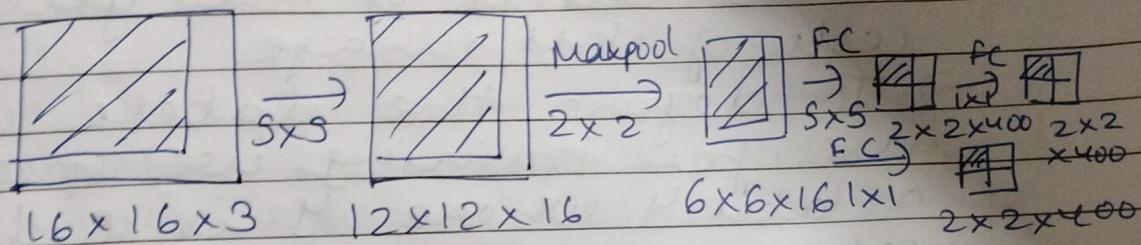


→ Suppose our image is $16 \times 16 \times 3$
Then make a $14 \times 14 \times 3$ window
& input that in convnet

New input



∴ We run 4 diff windows
 Thus the convnet
 → But lot of computation done
 by 4 convnets is highly duplicative
 → So convolutional implementation
 of sliding windows allows 4 results
 in the convnet to share computation



→ What convolutional implementation does is instead of forcing you to ~~combine~~^{run} all 4 propagation on 4 subsets it combines all 4 into 1 form of computation. It shares computation in regions that are common

Bounding Box predictions

→ How to get bounding box predictions to be more accurate

→ None of the windows form a perfect bound box
 Sometimes the shape of the perfect bounding is diff from the window we have taken

This can be solved by YOLO
- you only look once algo.

→ We place a grid on the input image we are going to take the img classification & localization also and apply it to each of the grid cells of this image

→ To define the labels we use for training we specify a label Y where Y is 8 dimensional vector as seen previously

→ first o/p is PC which is either 0/1 then b_x, b_y, b_w, b_h to specify bounding box if there is an image obj in the grid cell. and then c_1, c_2, c_3 ($c_4 = 0$ as obj present) → so we have a vector for each grid cell.

→ If no of grid cells is 7×7
→ Then target o/p volume is $3 \times 3 \times 8$

→ To train the NN input is Eg :- $100 \times 100 \times 3$ then we have usual convnet with conv, maxpool layers & so on. On that in the end we add more conv & max pool layers & so on that maps to a $3 \times 3 \times 8$ o/p vol (as grid is 3×3)

- Then we backpropagation
to train the network
- The advantage of this algo
is that N-N O/P possible
bounding boxes
 - The objects only assigned
to 1 grid cell ie no box
in which its mid point lie
 - b_x, b_y, b_h, b_w are specified
wrt the grid cell and
not the whole SPP frame
 b_x, b_y has to be b/w 0 & 1
as mid point is within
bounds of grid cell
 - b_w, b_h can be > 1

- ### # Intersection over Union (IoU)
- used to evaluate obj det algo
and used to make it worse
better
 - SPP O/P the intersection of
the actual bounding box
and the ideal bounding
box

Intersection = size of common area /
Over Union

Root for ideal box

size of common

two common

Obj

connect if
IOU > 0.5

→ we can use a more stringent threshold than 0.5 also

ii

Non Max Suppression

→ Algorithm may find multiple detections of same objects → Non max suppression is a way to make sure algo detects object only once
→ It looks at the probabilities associated with each of the detections. (actually P_c times $C_1 C_2 C_3 \rightarrow$ it takes the largest once as the most confident detection) → The non max suppression part looks at the $n \times n$ windows of all the ones with high overlap with high IOU get suppressed

→ STEP 1: Pick boxes with ~~large~~ P_c & ~~discard~~ all boxes with $P_c \leq$ some threshold value

→ STEP 2: Pick O/P box with largest P_c & O/P it as a prediction

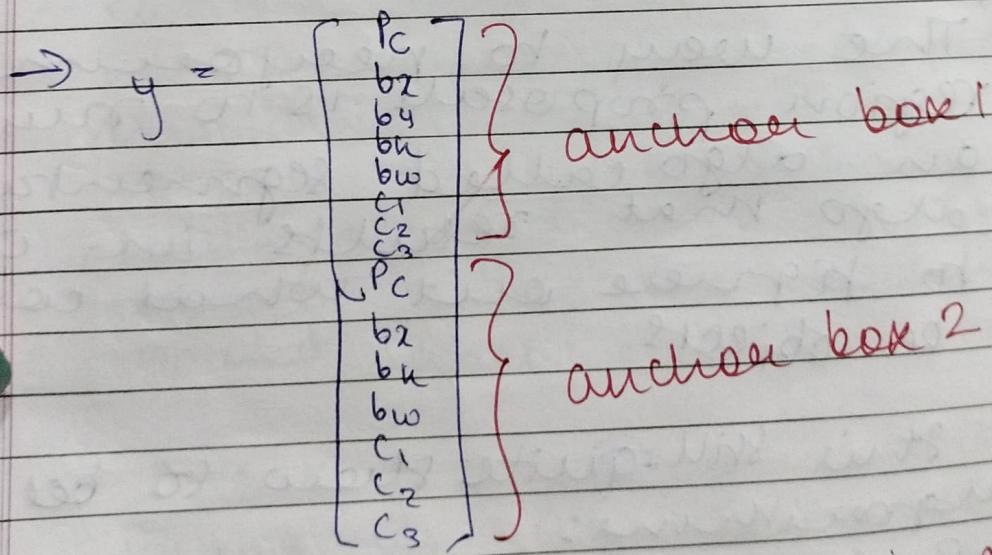
→ STEP 3: Discard any remaining boxes

→ If there are more than 1 type of obj then we have to independently carry non max suppression

on each o/p class :-

Anchor boxes

- If we have 2 objects with midpoints coinciding and lying on same grid cell then w/o the output y vector won't be able to detect both and will pick one of the 2.
- With the idea of anchor boxes is that we redefine 8 shapes called anchor boxes and associate 2 diff predictions with 2 diff boxes
- We can use more anchor boxes also (> 5)



Each object in training img is assigned to grid cell in at containing objects midpoint and anchor box for grid cell with highest IOU

Output = grid x grid x (87 no
of anchors
of boxes)

→ If we have 2 anchor boxes
for 3 obj then the algo doesn't
handle it well.

→ If 2 obj in a grid cell
have same anchor box
then also not handled well

Region proposals - RCNN

→ It picks just a few regions
that make sense to run
convolutional network rather
than running sliding window
on every single window
to select few windows & run
it on them

→ The way to implement
region proposal is to run
an algo called segmentation
algo that results in O/P
to figure out what could
be objects

It is still quite slow to be faster
algorithm:

→ R-CNN: Propose regions - Classify
proposed regions one at a
time. O/P label + bounding
box

→ fast R-CNN Propose regions
use convol. implementation

$\text{Output} = \text{grid} \times \text{grid} \times (8 \times \text{no of anchor boxes})$

- If we have 2 anchor boxes & 3 obj then the algo doesn't handle it well.
- If 2 obj in a grid cell have ~~strange~~ same anchor box shape then also not handled well

Region proposals - RCNN

- It picks just a few regions that make sense to run convolutional network rather than running sliding windows on every single window & select few windows & run it on them
- The way to perform region proposals is to run an algo called segmentation algo that results in O/P to figure out what could be objects

It is still quite slow so better algorithms:

- R-CNN: Propose regions. Classify proposed regions one at a time. O/P label + bounding box

- Fast R-CNN Propose regions. Use convol. implementation

of sliding windows to classify proposed regions

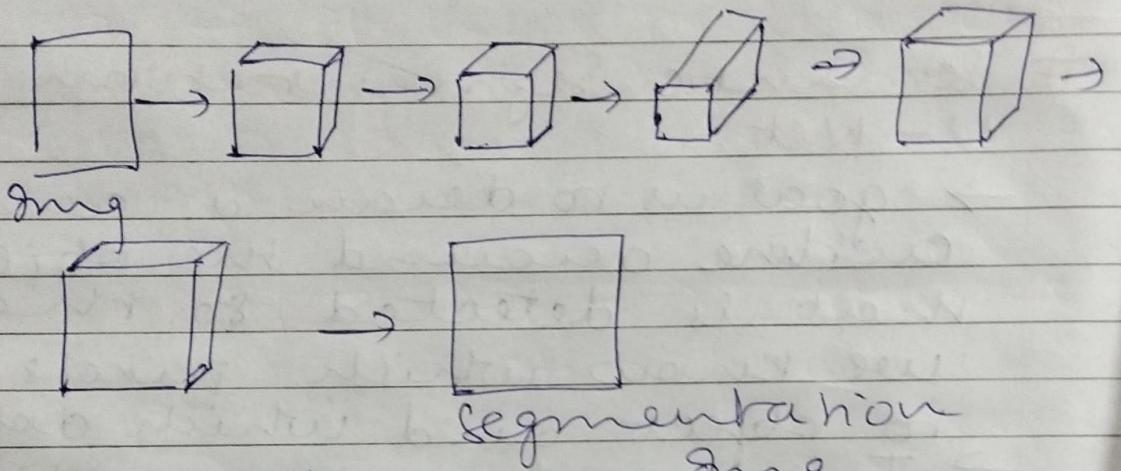
Faster R-CNN - convol. NN to propose regions

Semantic Segmentation with U-Net

- goal is to draw a careful outline around the object that is detected so that we know which pixels belong to object and which don't
- To segment out an obj from background we have 2 class labels 1 for Obj & 0 for not Obj. In this case job of segmentation algo is to off either 1/0 for every single pixel in this image → if we have one more obj we will have a 2nd class, class 2 and so on in which case the learning algo would label every pixel.

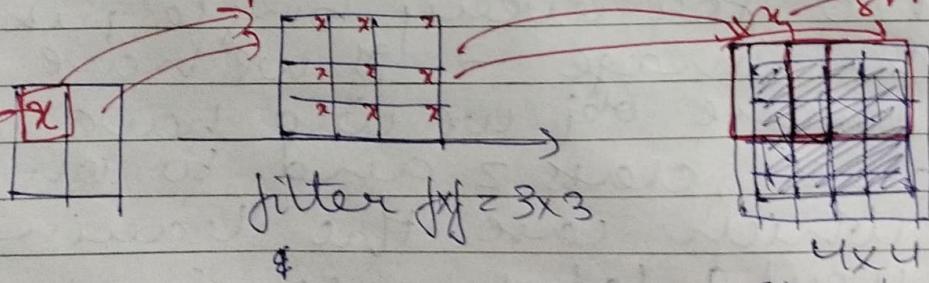
- One key step of semantic segmentation is as dimensions in CNN get smaller from left to right, it now needs to get bigger so they can gradually blow it up back to a full size image which

is a size we want for the O/P
 → as we go deeper height & width would go back up while no of channels will decrease



→ To image bigger again we need to implement transpose convolution.

Transpose convolutions



padding = 1

stride = 2

padding = 1

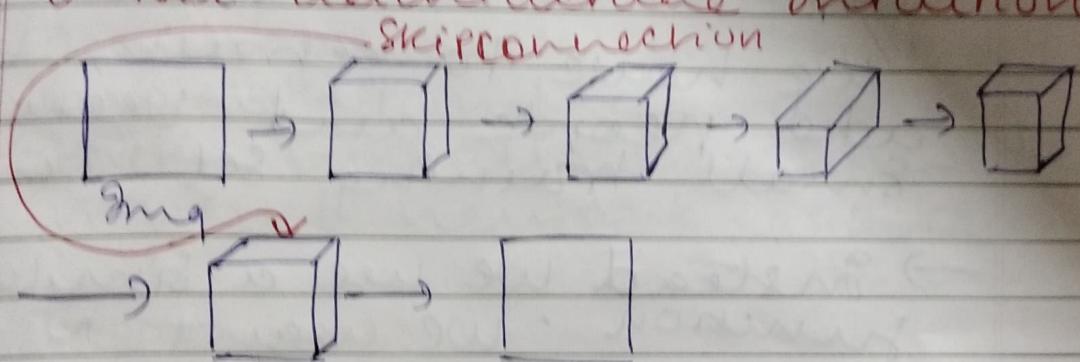
→ In transpose convolution we place filter on the output.

we take this value and multiply it to every element in filter and paste it on upper left

half of the filter

- stride and padding are kept output \neq not input
- The part of filter that overlaps no stride remains 0
- If more than 1 filters pass on a pixel of O/P we add them up.

U-Net architecture Intuition



→ ~~the final pixel O/P~~
 → From final layer to decide which region is a cat 2 types of info are useful. One is high level, spatial, high contextual which it gets from prior layer. But what is missing is a very fine grained spatial info. The skip connection ~~too~~ allows NN to take this very high res, low level feature info where it could capture from every pixel position

Face recognition

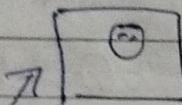
PAGE No.

DATE

11

one shot learning

→ most face recognition algo don't work well if we have only 1 training eg → it has to learn from just 1 training to recognize the person again

→  → CNN → softmax

we have to retrain to add people and very small training set to train

→ Instead we use a similarity function. We want $n \times n$ to learn a function which is going to denote d , which inputs 2 images & O/P. the degree of difference b/w 2 images

→ If same person O/P is small no

$$d(\text{img}_1, \text{img}_2) \leq \tau \text{ "same"} \\ \text{else} \quad > \tau \text{ "diff"}$$

for face recognition compare with all images in dataset and check if any has low d then they are same

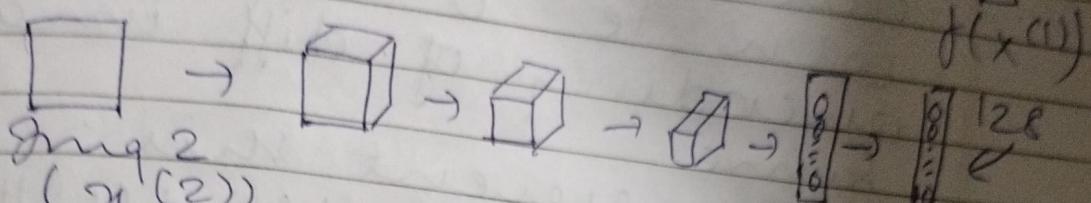
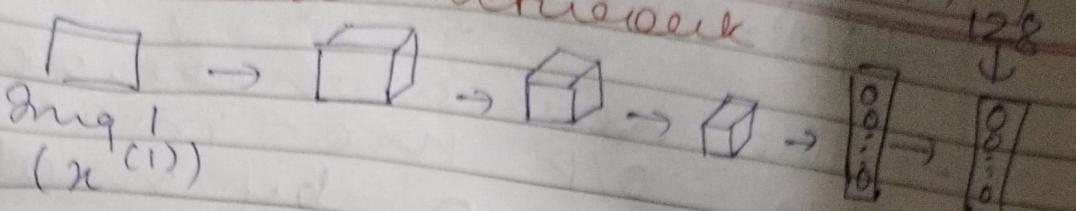
face
verification

#

Siamese Network

PAGE NO. _____
DATE _____

128



$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|$

→ Idea of running 2 identical CNN in two diff networks & and comparing them is called a siamese neural network architecture

Triplet loss

- One way to learn parameters of N.N so it gives us a good encoding for pairs of faces is to define & apply gradient descent on triplet loss function
- We will always be looking at 3 images an anchor img (A), a pos img (P) & -ve img (N)

$$\text{We want: } \|f(A) - f(P)\|^2 \leq \|f(A) - f(N)\|^2$$

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

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

→ The neural network can O/P the trivial soln ie 0 by making encoding of diff images same

→ To prevent this:

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

Pushes (anchor, +ve) & (anchor, -ve) hair away margin from each other. (hyperparameter)

LOSS FUNCTION

$$\begin{aligned} \mathcal{L}(A, P, N) = & \max(\|f(A) - f(P)\|^2 \\ & - \|f(A) - f(N)\|^2 + \alpha, 0) \end{aligned}$$

→ If loss is -ve then it O/P 0 if its +ve then it O/P 1

$$J_{\text{(cost func)}} = \sum_{i=1}^m \mathcal{L}(A^{(i)}, P^{(i)}, N^{(i)})$$

→ If u have 1 pic of each person you can't actually train this system. But if course, after having trained a system u can apply 1 shot learning problem.

PAGE NO.	11
DATE	

Choosing triplets APN

During training if $A, P \in N$ are chosen randomly
 $d(A, P) + \alpha \leq d(A, N)$ is easily satisfied

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

choose triplets that are hard to train on

$$d(A, P) + \alpha \leq d(A, N)$$

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

- face recognition verification & binary classification

→ we can take the 128 or even higher dimensional encoding
 If input is a logistic regression to make a prediction. Target O/P will be 1 if same or else 0
 This is binary classification
 → This is an alternative to triplet loss

→ O/P g will be a sigmoid function. Rather than feeding in these encodings we can take difference b/w encodings

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

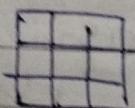
OR

$$\frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k}$$

Neural style transfer

What are deep convnets learning
 → if we want to visualize what hidden units in diff layers are computing:

- 1) start with a hidden unit in layer 1 and scan through training sets to find out what are the images / image patches that maximize that unit's activation ∵ we pause the training set then the NN figures out what in the image maximizes that particular unit activation



— Then we find 4 unit in layer 1 image patches that influence this layer (since 3×3)

→ If we do this for deeper layers a hidden unit will see a larger region of the image. Where at extreme end each pixel could hypothetically affect the O/P of these later layers of NN.

* Cost function

→ We define a cost function $J(G)$ that measures how good is a particular generated img and we use gradient descent to minimize $J(G)$ in order to generate this img.

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

This is a form of content img measuring & measures how similar contents of generated img and to contents of content img of or with s

$\alpha, \beta \rightarrow$ specify relative weighting b/w content cost & style cost

Find generated img G

- 1) Initiate G randomly $O_I : 100 \times 100 \times 3$ (eg)
- 2) Use gradient descent to minimize $J(G)$

$$G_t = G_{t-1} - \frac{\partial J(G)}{\partial G}$$

Content cost function

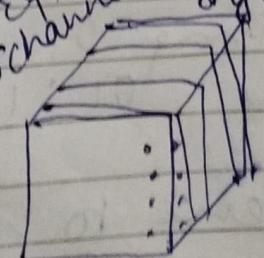
- For computing content cost layer l need shld be in between (if too deep or not very similar if too shallow or too similar to S)
- Use a pre-trained ConvNet and measure given a content img a^c & style img how similar are they
- Let $a^{CL(C)}$ & $a^{CL(S)}$ be activation of layer l on the image
- If $a^{CL(C)}$ & $a^{CL(S)}$ are similar both have similar content
- J content (C, S) = $\frac{1}{2} \| a^{CL(C)} - a^{CL(S)} \|^2$

Style cost function

- Say we use layer l's activation to measure style
- Define style as correlation b/w activations across channels

Layer l activation

(eg) style img



→ look at first 2 channels and see how correlated activations are

→ & get pairs of nos

at all positions (one in each channel)

- If 2 channels are correlated the features they detect usually occur together
- If they are uncorrelated the features won't occur together
- Then correlate in generated image also and this measures how similar its style is to Style Img

* Style matrix

Let $a_{ijk}^{[l]}$ activation at (i, j, k) in hidden layer L .
 $G_{kk'}^{[l]}$ is $n_c^{[l]} \times n_c^{[l]}$ - will measure correlation b/w channels

$k \in K'$

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

$$G_{kk'}^{[l](G)} = \sum_{i=1}^{n_h} \sum_{j=1}^{n_w} a_{ijk}^{(l)(S)} * a_{ijk'}^{(l)(S)}$$

$$G_{kk'}^{[l](G)} = \sum_{i=1}^{n_h} \sum_{j=1}^{n_w} a_{ijk}^{(l)(G)} * a_{ijk'}^{(l)(G)}$$

$$J_{\text{Style}}(S, G) = \| G^{[l](S)} - G^{[l](G)} \|^2$$

$$J_{\text{style}}(S, \alpha) = \frac{1}{(2n_H^{(l)} n_W^{(l)} n_C^{(l)})^2} \sum_K \sum_{K'} \left[G_{KK'} - G_{KK'}^{(\text{ref})} \right]$$

$$J_{\text{style}}(S, \alpha) = \sum_l \lambda^{(l)} J_{\text{style}}^{(l)}(S, \alpha)$$

allows us to use
diff layers in NN

generalizations of models