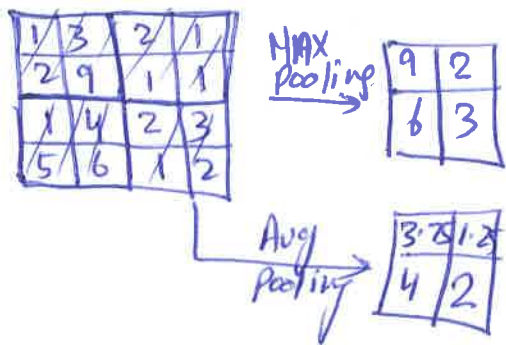


W1-29

MAX Pooling Example



W1-110

Neural Net Example

Very similar to LENET-5 (VANLECON)

	Input	Activation	# of Parameters
	(32, 32, 3)	$a^{[0]} = 3072$ (32x32x3)	0
Conv1 ($f=5, s=1$)	(28, 28, 8)	6272 (28x28x8)	208 (5x5 filter, 8 channels)
Pool1	(14, 14, 8)	1568	0
Conv2 ($f=5, s=2$)	(10, 10, 16)	1600	416 (5x5 filter, 16 channels)
Pool2	(5, 5, 16)	400	0
FC3	(120, 4)	120	4800 (120x400 = 48000)
FC4	(84, 1)	84	10,081 (40x84 = 3360 + 1 bias = 3361)
Softmax	(10, 2)	10	841 (10x84 = 840 + 1 bias = 841)

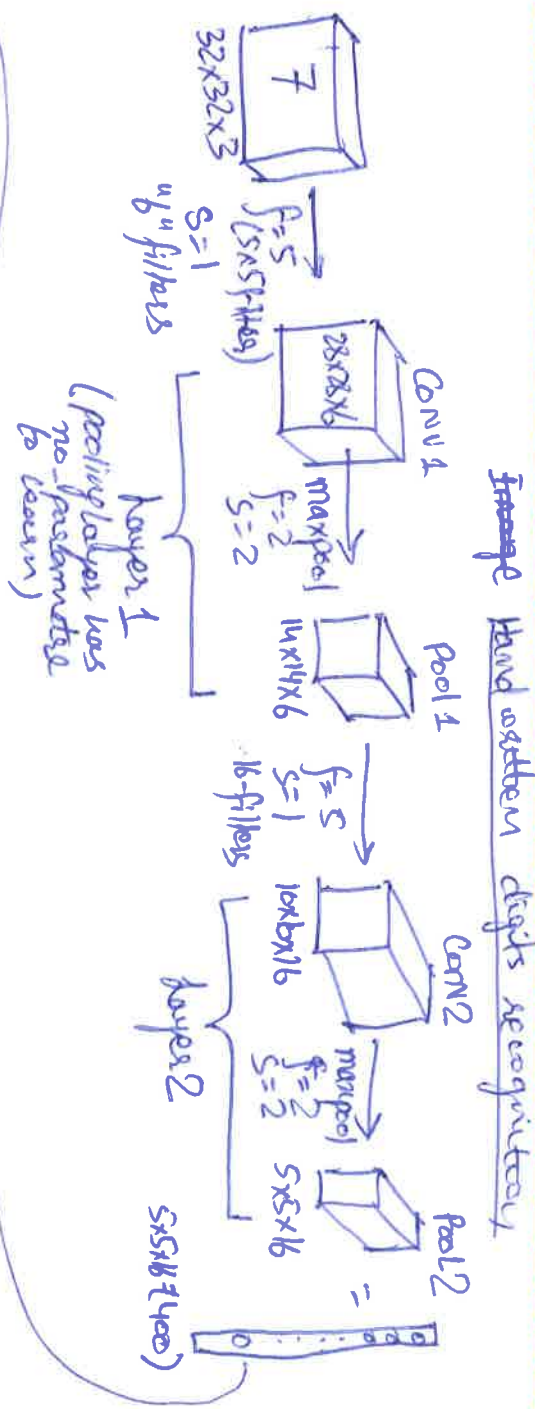


Softmax (10-outputs)

Generally $n_{in}, n_{out} \downarrow$ across layers

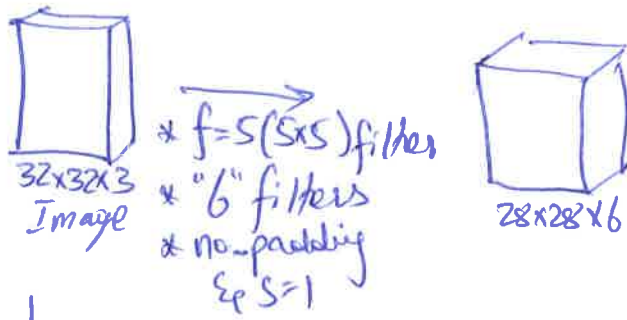
1) $n_{in}, n_{out} \downarrow$ across layers

2) This pattern is conv-pool-conv-pool



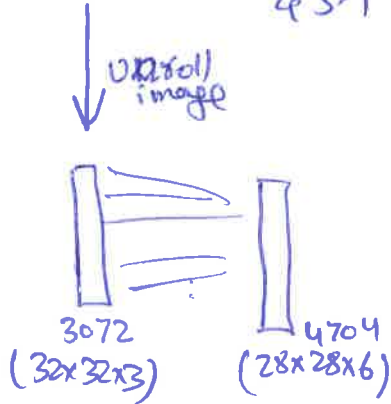
How Convolution reduces # of params?

CONV



of params here = $(5 \times 5 \times 3) \times 6 = 156$ params

FC Networks

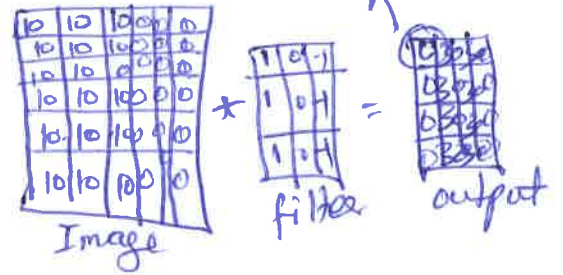


of params here = $(3072 \times 4704) \approx 14 \text{ M}$

Why Convolution

1) Parameter sharing

look at ~~the~~ image. same set of filter is applied to entire image. (9 parameters)



2) Sparsity of connections

In each layer, each output value (for example "A") depends on only a small number of inputs from image.

* Classic Networks

- 1) LeNet-5
- 2) AlexNet
- 3) VGG

* RESNet (152-layered N/w)

* Inception NN

} CNN

W2-L2

1) LeNet5 (60K params)

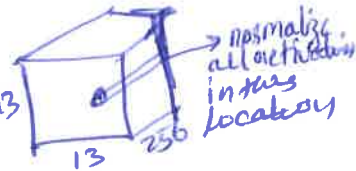
LeCun (98) → Gradient-based learning to document recognition

Advanced Notes

1) Sigmoid/tanh not Relu

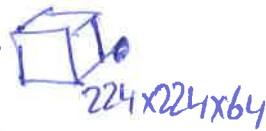
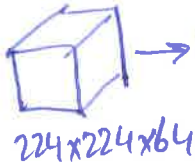
2) AlexNet

- (Krizhevsky et al 2012. ImageNet classification with deep CNN)
- Similarity with LeNet-5 (it has 60M params)
 - Used ReLU activations
 - Multiple GPUs
 - Local response Normalization layer



3) VGG-16 → 16-layers ~ 138M params
 (Simonyan & Zisserman 2015)
 Very deep CNN for large-scale image recognition

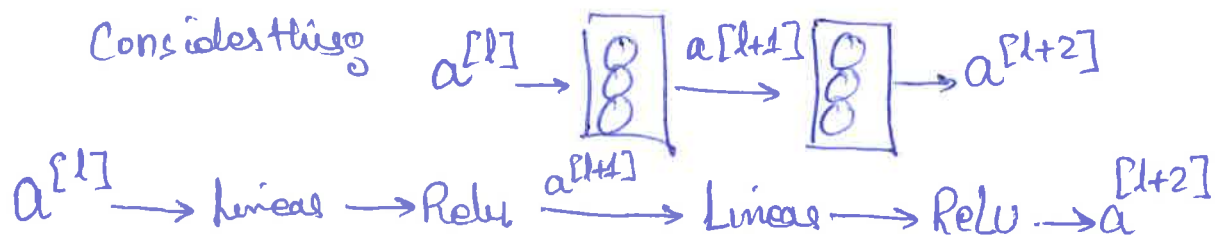
[conv64] x 2



easy paper
than VGG paper
than LeNet5

Residual block

Considering



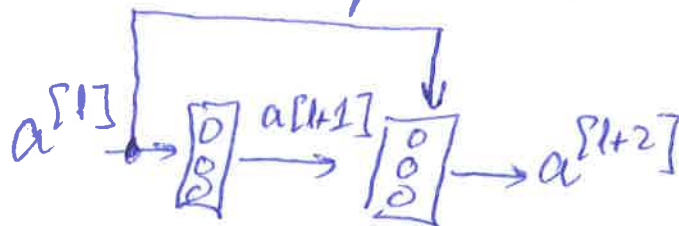
$$z^{[1+1]} = W^{[1+1]} a^{[1]} + b^{[1+1]}$$

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

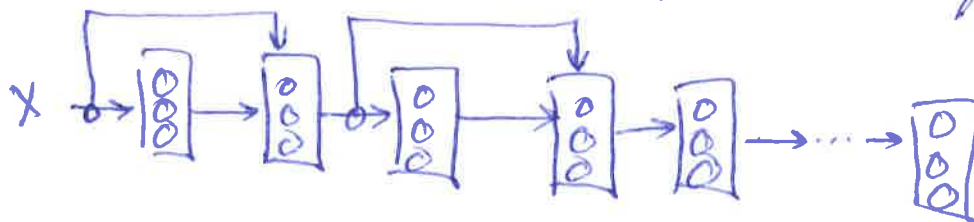
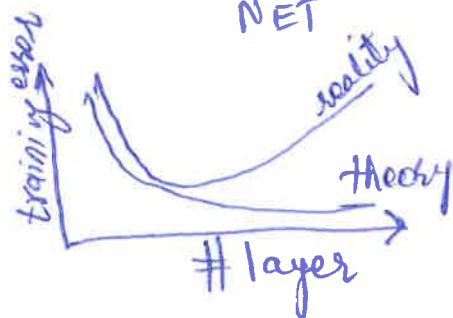
$$z^{[1+2]} = W^{[1+2]} a^{[1+1]} + b^{[1+2]}$$

$$a^{[1+2]} = g(z^{[1+2]})$$

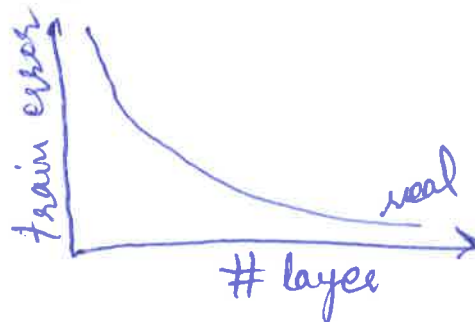
Now adding "shortcut or skip connection,"



$$\text{Here } a^{[1+2]} = g(z^{[1+2]} + a^{[1]})$$

Residual v/w (He et.al. 2015, Deep residual net for image recognition)Plain
NET

RESNET



W2-24

9

1) Why ResNet works

2) See the diagrams of ResNet in slides

[He et al. 2015 Deep residual networks for image recognition]

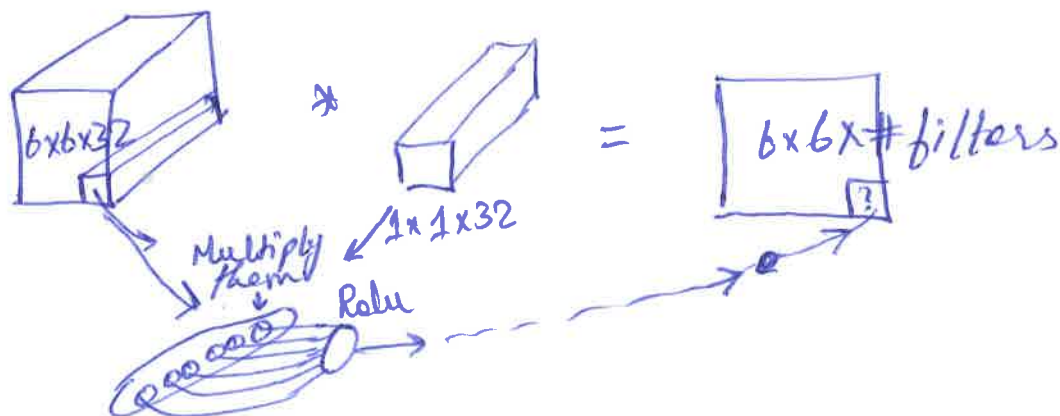
W2-25

1x1 Convolution (Lin et al. 2013 Network in Network)

$$\begin{bmatrix} 1 & 2 & 3 & 6 & 5 & 8 \\ 3 & 5 & 5 & 1 & 3 & 4 \\ \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} * \begin{bmatrix} 2 \end{bmatrix} = \begin{bmatrix} 2 & 4 & 6 & \dots \end{bmatrix}$$

1-channel (no effect)

It makes sense with more than one channel



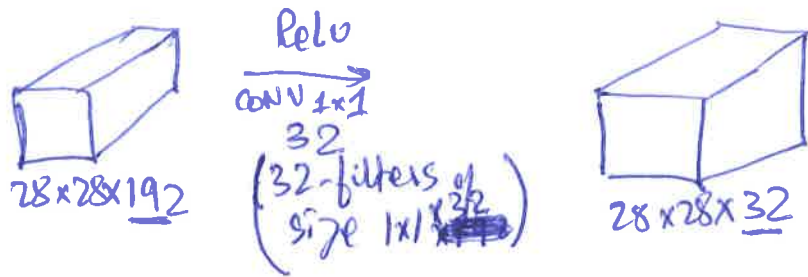
For more filters,

$$6 \times 6 \times 32 * \begin{matrix} 1 \times 1 \times 32 \text{ (1st filter)} \\ 1 \times 1 \times 32 \text{ (n-filter)} \end{matrix}$$

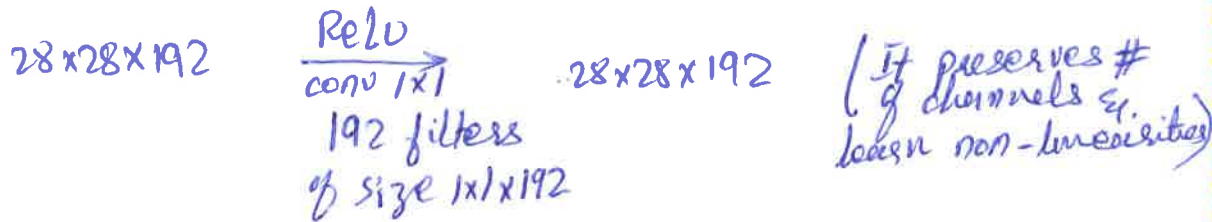
Here we will have "n" ReLU units.

Using 1x1 Convolution

(11)



* Pooling was used to decrease n_H, n_W but 1×1 convolution can be used to increase or decrease number of channels i.e. n_C .



W2-26

Motivation for inception network.
(Szegedy et al. 2014. Going deeper with convolutions)

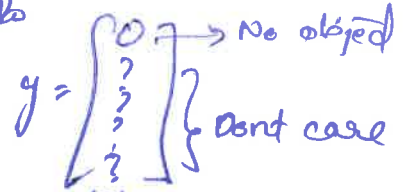
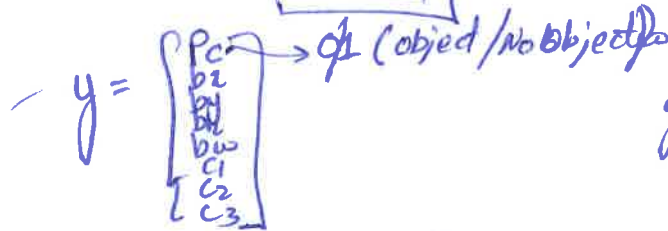
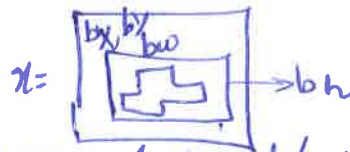
W3-11

13

Object detection with localization

- 1- C1 (Pedestrian)
- 2- C2 (car)
- 3- C3 (motorcycle)
- 4- background

Need to output b_x, b_y, b_h, b_w
class label (1-4)



$$d(\hat{y}, y) = (\hat{y}_1 - y_1)^2 + \dots + (\hat{y}_8 - y_8)^2 \quad (\text{Squared loss})$$

OR

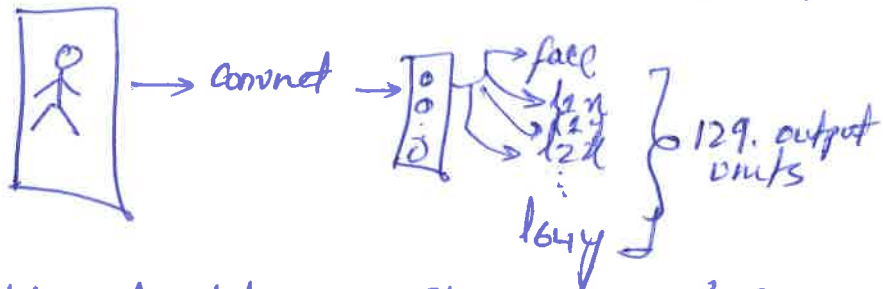
$$d(\hat{y}, y) = \text{logistic regression loss}(y_1) + \text{squared error loss}(b_x, b_y, b_h, b_w) + \text{log likelihood loss}(c_1, c_2, c_3)$$

W3-12

Landmark detection

It outputs (x, y) of image i.e. output keypoints of image

i.e.



Pose-detection

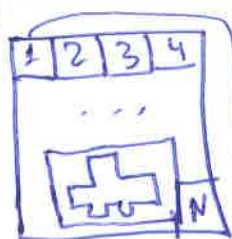
Annotate key positions of person's pose

W3-13

* Sliding windows detection (used for object detection)

W3-14

* Sliding window object detection using conv net is slow.



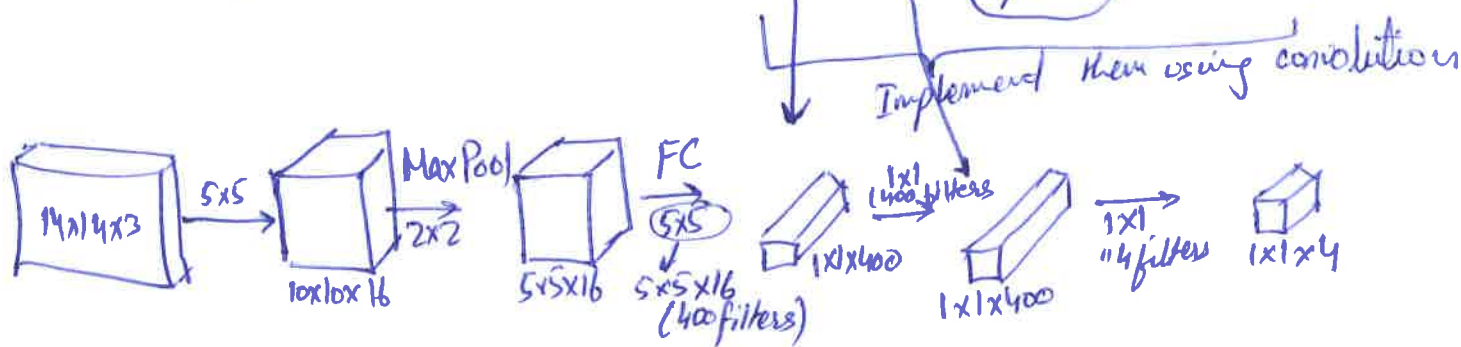
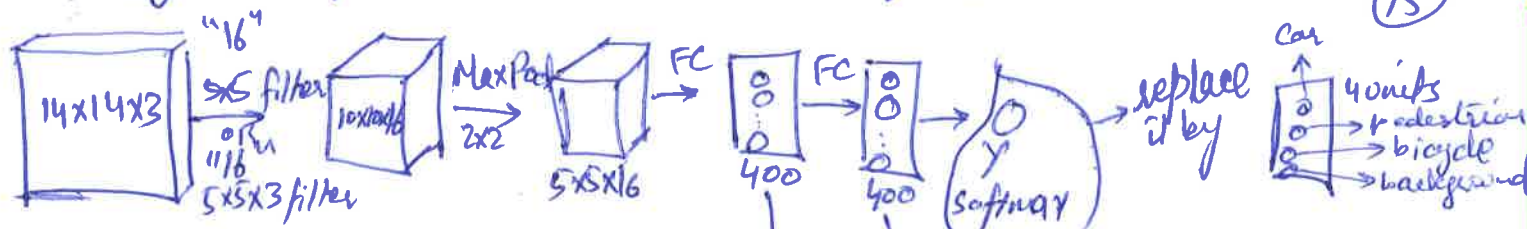
To detect car with its bounding box, we need to extract image patch as pass it to conv net. It will output



4 numbers. So, running sliding window over image is expensive as we have to pass 'N' patches to conv net to get the location of car. To fix this problem, let's see convolutional implementation of sliding window.

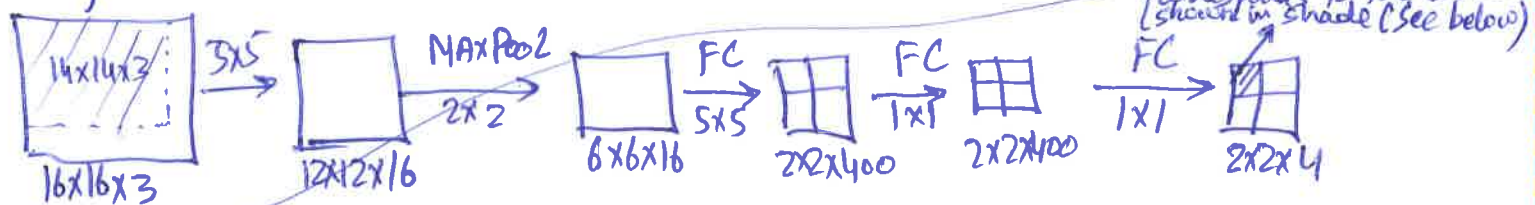
Turning FC layer into convolutional layers

(15)

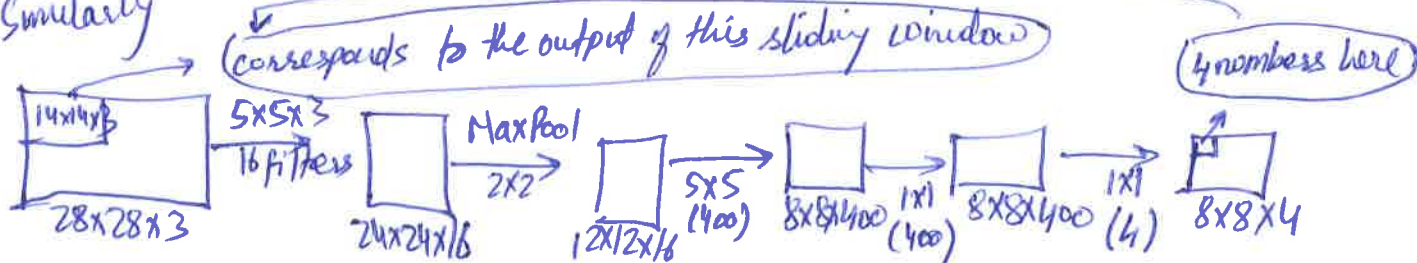


See Sermanet et al. 2014, Overfeat: Integrated recognition, localization & detection using convolutional networks

* Imagining you get an image of size 16x16x3. Show how can we do convolutional sliding window



Four numbers here are the result obtained by running 14x14x3 patch through convnet described above in the shaded region of 16x16x3 volume. Similarly



Thus through a single pass, we obtain the detection over entire 28x28x3 image, just by implementing kernel.

This type of object localization is not that much accurate. Now, we will have a look over Yolo algorithm that can make this accurate.

$$\arg\max_{h_k} \left[w_{em} f(x_k, h_k) + \sum_{j=1}^N w_{TRANS} f(h_k, h_j) + \sum_{\substack{i=1 \\ i \neq k}} w f(h_i, h_k) + w_{comp} f(h_k, y) \right]$$

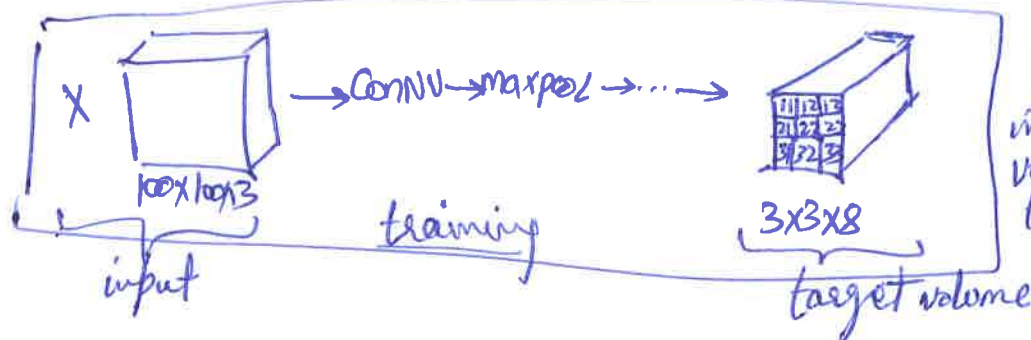
Diagram illustrating the process of placing objects into a 3x3 grid:

- A 100x100 bounding box is shown with two objects inside. The text "Bounding boxes are given" is written above it.
- An arrow labeled "Place" points to a 3x3 grid.
- The 3x3 grid is shown with the two objects placed in the bottom-left and bottom-right cells. The grid is labeled "3x3 grid" and "100" (referring to the original image size).
- The resulting feature vector y is shown as a column vector:

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$
- The feature vector y_{21} is shown as a column vector:

$$y_{21} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \end{bmatrix}$$
- The feature vector y_{23} is shown as a column vector:

$$y_{23} = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \end{bmatrix}$$
 The vectors y_{21} and y_{23} are grouped together with a brace and the number 8, indicating they are part of a set of 8 vectors.



This is the convolutional implementation, where volume "21" is similar to y_{21} ground truth vector

Evaluating object localization

Intersection over Union IoU = $\frac{\text{size of Intersection}}{\text{size of union}}$

"Correct" if $\text{IoU} \geq 0.5$ (In general)

- * IoU is the overlap of two bounding boxes.

Non-Max Suppression

Improves Yolo's accuracy of localizing bounding boxes.

Non-Max Suppression
Improves Yolo's accuracy of localizing bounding boxes. In object detection & localization, we have multiple detections of single object. Non-max suppression give only one location of detection by choosing the bounding box with highest probability.

Now have a look at non-max suppression algorithms.
It is simple



each output production is: $\begin{bmatrix} pc \\ bx \\ by \\ bh \\ bw \end{bmatrix}$ } C, C2, C3 are discarded right now

Discard all boxes with $P_c \leq 0.6$
while there are any remaining boxes:

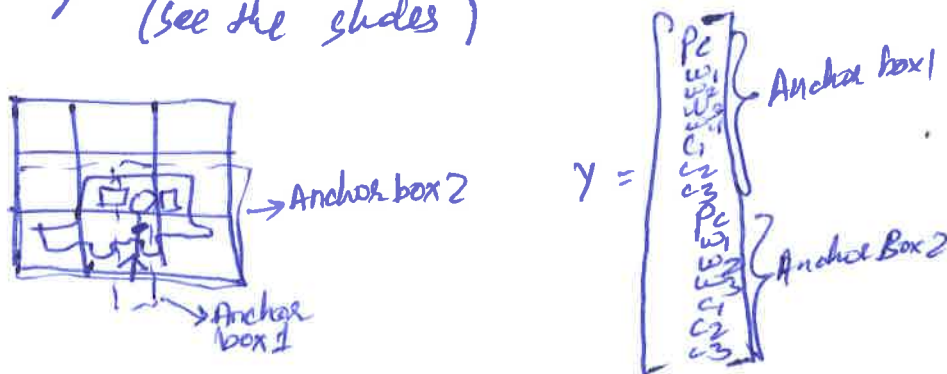
- 1) Pick the box with the largest pc. output that as prediction
- 2) Discard any remaining box with $IOU \geq 0.5$ with the box output in the previous step.

W3-28

Anchor Boxes

(19)

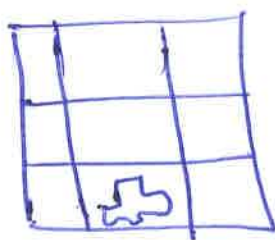
- * Anchor Boxes allow us to detect multiple objects in a single grid cell.
- (See the slides)



W3-29

YOLO Algorithm for object Detection (Redmon et al. 2015)

Training

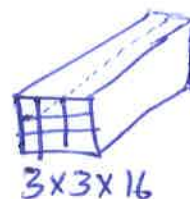


- * we have 2 anchor boxes
- * for each grid cell, $y = 2 \times 8$
 $\text{Anchor boxes } p, w_1, w_2, c_1, c_2, c_3$
- * Thus y for entire image
 $y = 3 \times 3 \times 2 \times 8$

Training:



→ ConvNet →



Prediction

See the slides

- * Get rid of low probability bounding box
- * Apply non-max suppression for each class

W4-216

Region Proposal : (R-CNN)

Girshik et al

2013, Rich feature hierarchies for accurate object detection & semantic segmentation

W4-L1

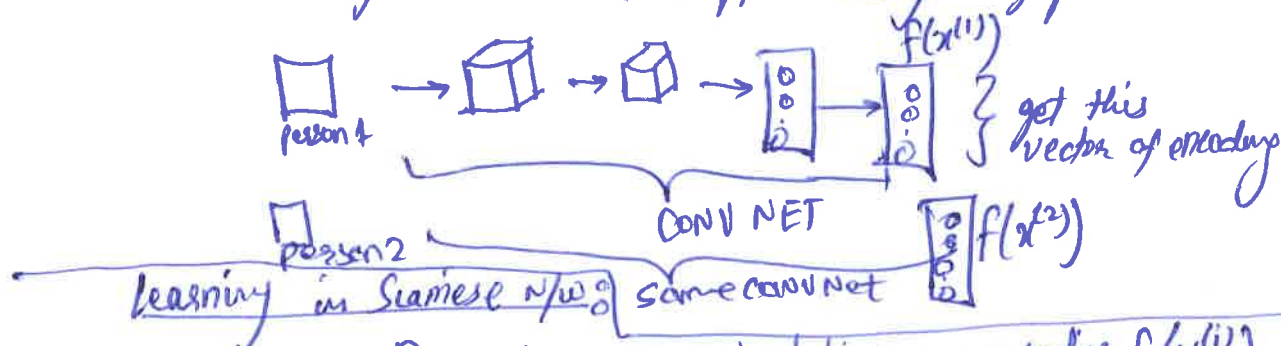
W4-L2

ONE-SHOT learning: learn from one example to recognize the person again.
It can be addressed by learning similarity function.

W4-L3

Siamese Network

Taigman et al. 2014. Deepface: closing the gap ...



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

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.

W4-L4

Triplet loss

- * How to define objective function for Siamese n/w.
- * Schroff et al. 2015, faceNet: A unified embedding for face recognition and clustering

Learning objective



we are dealing with three images, anchor positive & negative. That's why it is called triplet loss

we want:

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

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

margin

Loss function (Triplet loss function for Siamese n/w)

(23)

Given three image A (anchors), P (positive), N (negative)

$$d(A, P, N) = \max (\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

over all training set:

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

For training set

Choose triplets that are "hard" to train on.

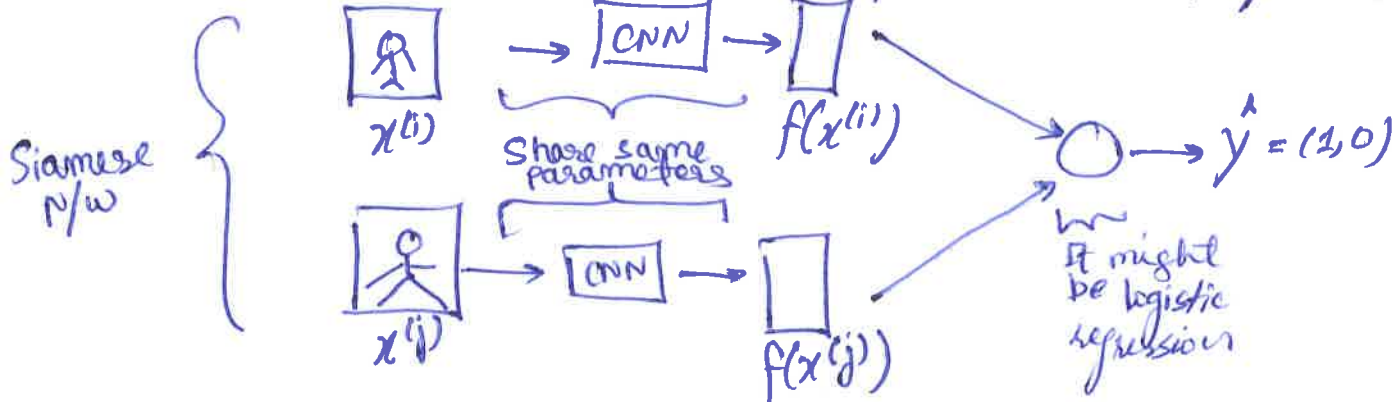
$$d(A, P) + \alpha \leq d(A, N)$$

such that $d(A, P) \neq d(A, N)$

W1-25

Face Verification & binary Classification

(Taigman et al. 2014 peepface closing the gap to human level performance)



Or

$$\sum_{k=1}^K |f(x^{(i)})_k - f(x^{(j)})_k|$$

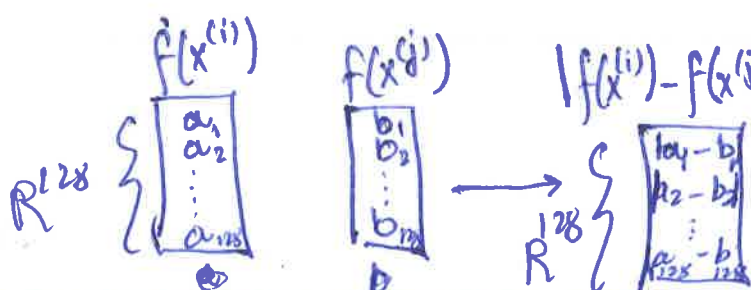
$$\text{If } f(x^{(i)}) \in \mathbb{R}^{128}$$

then

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

It could be Chi-Square similarity function as well

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



WS-26

NEURAL STYLE TRANSFER

(25)

WS-27

Visualizing CNN

Very interesting: Zeiler & Fergus 2013, Visualizing and understanding CNN

WS-28

Neural STYLE transfer cost function

Gatys et al. 2015. A neural algorithm of artistic style.

$$J(G) = \alpha \bar{J}_{\text{CONTENT}}(C, G) + \beta \bar{J}_{\text{STYLE}}(S, G)$$

WS-29

Content Cost function

WS-30

Style Cost function

WS-41

2D & 3D Convolution

