

CONVOLUTIONAL

NEURAL NETWORKS

$n \times n$ image; $f \times f$ filter; padding = P ; stride = s

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

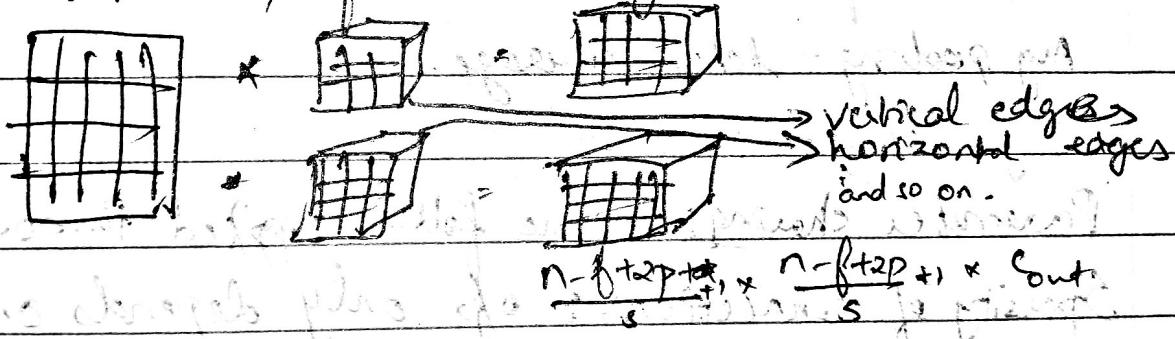
Convolution is better called technically as cross-correlation.

Convolutions and volumes

- * Use 3D filters corresponding to channels.

- * No. of channels in filters = # channels in images.

- * Count - specifies how many such channels are used.



ONE LAYER OF CNN:

No. of params is irrespective of the size of the ip image.

For ex: 3×3 filters of count 10 has

$$3 \times 3 = [27 + 1(\text{bias})] \times 10 = 280 \text{ params.}$$

Notations: $f[l]$: filter size | $p[l]$: padding | $s[l]$: stride

Input: $n_h^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]}$ $n_e^{[l]} = \# \text{filters}$

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

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

Activations:

$$a^{[l]} \rightarrow f^{[l]} \times n_w^{[l]} \times n_e^{[l]}$$

with batches $\rightarrow M \times a^{[l]}$

$$\text{bias: } n_c^{[l]} = (1, 1, 1, \dots, 1^{[l]})$$

CONVNET EXAMPLE \rightarrow Valid convolution

POOLING:

- max pooling: partitions
- large numbers - a feature is detected
- all features detected given by M.P. If no feature, low number is maintained.
- * applied for each channel. So 4 channels will have same $n_{in} = n_{out}$

Avg pooling: takes average.

Why CONVOLUTIONS?

Parameter sharing: \rightarrow same filter applied throughout the image.

Sparcity of connections: op only depends on a particular filter of inputs.

Translation invariance: A pic of cat shifted a couple of pixels is still a cat.

Learning op: $f^{[l+1]} = h^{[l]} \circ f^{[l]}$

Forward pass: $E^{[l+1]} = E^{[l]} \circ f^{[l]}$

$E^{[l+1]} = E^{[l]}_1 + E^{[l]}_2 + \dots + E^{[l]}_n$

DEEP CONVOLUTIONAL MODELS: Case Studies.

Architectures:

> LeNet-5 [60k params]

> AlexNet \Rightarrow 111rd to LeNet but ~~bigger~~ \Rightarrow ReLU

[60M] \Rightarrow Multiple GPU's \Rightarrow Local Response Normalization (LRN) - normalizing oval across channels to avoid high activations.

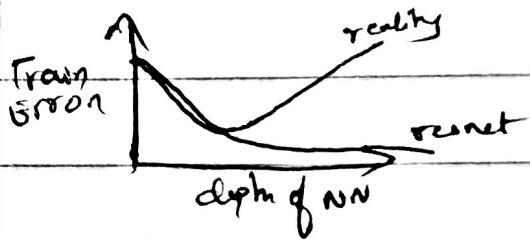
> VGG-16 net:

[130 M]

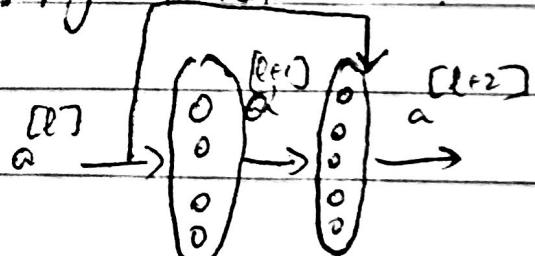
Resnets

makes use of skip connection.

In reality, train error \uparrow with deeper networks.



Why Resnets work?



$w^{(l+1)} \approx 0$, often regularized

$$a^{(l+1)} = g\left(w^{(l+1)} a^{(l+1)} + b^{(l+1)} + w_s a^{(l)}\right)$$

if diff dimension.

1x1 convolutions

- shrink the # channels, thus ↓ volume.

Inception Network (Motivation)

- uses diff filter sizes & max pool layers but maintains image dims.

1x1 convolution - bottleneck (reduces computation cost significantly - does not affect performance)

Inception network

- concatenation of inception blocks.

side branches: Take hidden layers & make a prediction

Transfer learning

- freeze the bottom layers of the trained Net by setting $lr = 0$ or $\text{trainable_params} = 0$.
- ↑ more examples & ↓ # layers freeze at top & ↑ layers trained at top
- Precompute one layer's activation for an image X to get a feature map & train softmax on it.

Data augmentation

1) Horizontally - horizontal flip.

2) Random cropping - should be reasonably large subsets.

3) Rotation / shearing / Local warping etc..

4) Color shifting : $R+20, G-20, B+20$; RGB are drawn from a probability distribution.

5) PCA color augmentation ; If img. X is more R, G but less B, PCA CA reduces more in R & not less in B.
(Refer AlexNet paper ↑)

* Have multiple CPU threads for loading & augmenting images.

> Hyperparameters - how much color shift, warp factor etc..?

State of Computer Vision

2 types of knowledge:

> labelled data > hand-engineering (Chacks)

* Do well on benchmark datasets : Tips:

- i) Ensembling (Train NN's and avg. their ops)
- ii) Multi-crop at test time - 10 crop of test image & average it.
- iii) QOL +

OBJECT DETECTION

- classification with localization
 - train or to output b_x, b_y - midpoints of bboxes & b_h, b_w & class labels

$$y = \begin{bmatrix} p_c & b_x & b_y & b_h & b_w & c_1 & c_2 & c_3 \end{bmatrix}^T$$

\downarrow
0/1 - obj act/not?

$$\text{obj detection loss } L(\hat{y}, y) = f((\hat{y}_1 - y_1)^2 + \dots + (\hat{y}_8 - y_8)^2) \text{ if } \hat{y}_1 = 1$$

$+ ((\hat{y}_1 - y_1)^2 - \text{if } \hat{y}_1 = 0)$

c_1, c_2, c_3 - softmax of

b_x, b_y, b_h, b_w - squared error loss.

p_c - logistic regression loss

Landmark detection

use consistent order of landmark points to get coordinates of landmarks.

$$y = \begin{bmatrix} f_x \\ f_y \\ l_{x1} \\ l_{y1} \\ l_{x2} \\ l_{y2} \end{bmatrix}$$

\rightarrow has face or no?

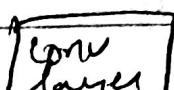
Object detection

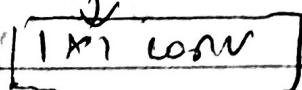
- Sliding windows detection, based on
 - start with small windows! d ~~now~~ slide across original image to get cropped images, that are fed into conv nets

CONVOLUTION IMPLEMENTATION OF SLIDING WINDOW:

- Turning FC into convolutional layers



conv layers \rightarrow  with # channels at M_{FC}



softmax

- Instead of using sliding windows, run the image across convolutions & in the final output, we get $n \times n \times 4$ instead of $1 \times 1 \times 4$, we can then check the no. o/p's to check if object is detected.

Problem: bounding box dims are not accurate.

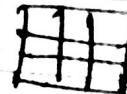
BOUNDING BOX PREDS:

YOLO - You Only Look Once

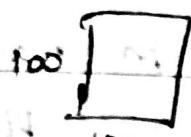
- * Split images as grid cells.
- * Perform localization & classification for each grid.
- * One $y = [...]$ for each grid cell.
- * An object belongs to the grid cell on which the object's midpoint falls.

If # grid cells = $O(1) - n \times n \times 8$ (approx)

($3 \times 3 \times 3$) on video.



NN +



\rightarrow conv \rightarrow max pool \rightarrow $3 \times 3 \times 8$ (y).

Limitations:

- Restricted to 1 object per grid
- * In practice, more objects do not occur in fixed grids

Efficiency + convolution implementation looks at the image only once
9 grids in parallel.

BBoxes: $b_x, b_y, b_w, b_h \rightarrow$ are 0/w. 0 and 1, as they are relative to the grid cell.

$b_x, b_y \rightarrow > 1$ as we can have a bigger car.

INTERSECTION OVER UNION (IoU)

- Evaluate CDA.
- Size of intersection \Rightarrow correct if $\text{IoU} \geq 0.5$
- \rightarrow IoU, Accuracy. (also measures similarity of 2 bboxes).

NON-MAX SUPPRESSION

- > Multiple detections of same objects. (ISSUE)

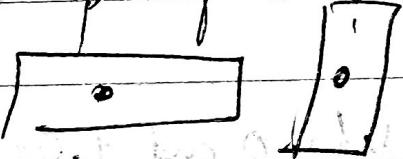
> In a 19×19 grid, multiple grids may detect an object, but the midpoint belongs to only one.

 - First take max (0.9), all other rectangles that have high overlap with max are suppressed (nms). Remove the rects that are suppressed. Finds max & suppresses others (non-max sup).

↳ Non Max Suppression

ANCHOR BOXES

> If there are > 1 object in grids, predefine shapes as anchor boxes.



> y-outputs are concatenated.

>

YOLO ALGORITHM

> puts all the above readings together.

REGION PROPOSAL

> RCNN's, picks regions to run CNN's on.

R-CNN: propose regions. Classify proposed regions & create a true-false label + bbox.