



→ Orthogonalization:

↳ "Things specific to a certain problem or a fun?".

↳ Chains of assumptions in ML =

1) Fit training set well on J (cost fun?).
↓

2) Fit dev set well on J : [Regularisation]
↓

3) Fit test set well on J
↳ Performs well in real world.

→ F_1 score → "Average" of Precision P & recall R

$$\left(\frac{2}{\frac{1}{P} + \frac{1}{R}} \right) \rightarrow \text{harmonic mean } \therefore !!$$

Dev set + Single number eval. matrix

better ≈
(speeds up cycle)

→ If there are multiple classifiers, average matrix is to be generated having lower avg. error).

→ Dev set and test set must be from distribution.

→ If we have a million examples, → 98% → training

$$\begin{cases} 1\% \rightarrow \text{dev} \\ 1\% \rightarrow \text{test} \end{cases}$$

→ Set your test set to be big enough to give high confidence in the overall performance of ur system.

↳ Sometimes train + Dev (no test)

↳ fine [test set becomes dev set]

→ If doing well on your metric + dev/test set doesn't correspond to doing well on your application, change your metric and/or dev/test set.

→ If your machine learning algorithm performs worse than a human, you can do:

- get labeled data from humans. (x, y)

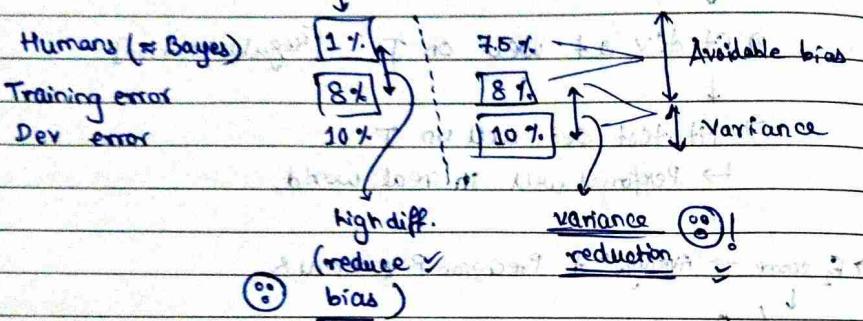
- Gain insight from manual error analysis:

Why did a person get this right?

- Better analysis of bias/variance.



Example:



→ Human performance gives an estimate for Bayes error

→ Surpassing human-level performance =

Team of humans 0.5% (Bayes error)

One human

$\times 1\% = 0.1\% = \text{Bias}$

Training error

$0.6\% + 0.1\% = 0.7\% = \text{Situation}$

Dev error

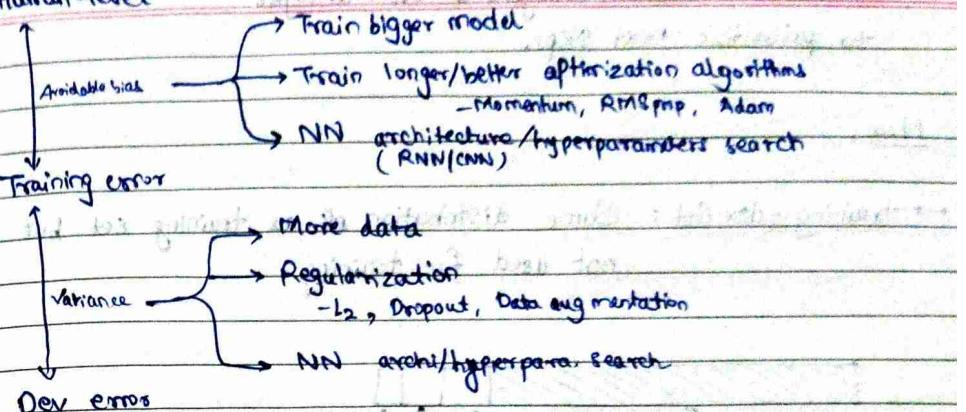
$0.8\% + 0.1\% = 0.9\% = \text{Variation}$

↔ Two fundamental assumptions of supervised learning =

1. You can fit the training set pretty well. [Low Available bias]
2. The training set performance generalizes pretty well to dev/test set. [Variance not too bad]

* Reducing (available) bias and variance →

Human-level



→ Incorrectly labelled examples →

↳ Deep learning algos are quite robust (insensitive) to random errors in training set.

↳ If mistlabelling in dev/test set → Perform error analysis with a column of Incorrectly labelled.

→ Ex. →

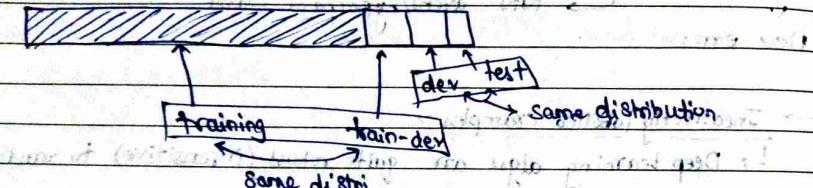
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100

- ↳ Apply some process to your dev and test sets to make sure they continue to come from the same distribution.
- ↳ Consider examining examples your algorithm got right as well as ones it got wrong.
- ↳ Train & dev/test data may now come from slightly diff distributions.

- 1) Setup dev/test set and metric.
- 2) Build initial system quickly.
- 3) Use Bias/Variance analysis & error analysis to prioritize next steps.

Now,

- training-dev set: Same distribution as the training set, but not used for training.



- No back propagation on train-dev, only on training.

eg.	Training error 1%	1%
	Training-dev error 9% variance	1.5%
	Dev error 10%	20% variance
	Human ... 0%	Human ... 0%
	J.E. 10% High bias	Variance 10% Avoidable bias (coz dev & train diff. distribution)
	T.D.E. 11%	11% variance
	DENE 12%	Data mismatch 20%

→ Ex:

Human level	"Human level" 4%.	6% ↑
Errors on examples trained on	"Training error" 4%.	6% ↑ avoidable bias
Errors on examples not trained on	"Training-dev" 10% "Dev/Test" 6% error	↓ variance data mismatched

Addressing Data Mismatch →

- Carry out manual error analysis to try to understand difference b/w training & dev/test sets.
- Make training data more similar; or collect more data similar to dev/test sets.

↳ How to?

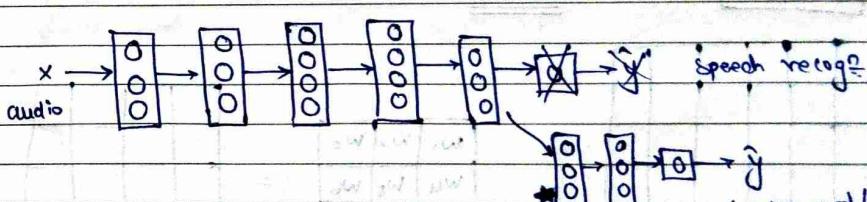
Artificial data synthesis →

eg.

$$\boxed{))} + \boxed{))} = \boxed{))}$$

"The quick brown fox jumps over the lazy dog" → Car noise ← Synthesized in-car audio

Transfer Learning →



When does transfer learning makes sense?

- Task A and B have the same input x .
- You have a lot more data for Task A & Task B.
- Low level features from A could be helpful for learning B.

When does multi-task learning makes sense = [Better] ↗ !!

- Training on a set of tasks that could benefit from having shared lower-level features.
- Usually: Amount of data you have for each task is quite similar.
- Can train a big enough neural network to do well on all tasks.

→ end-to-end learning -

↳ Pros:

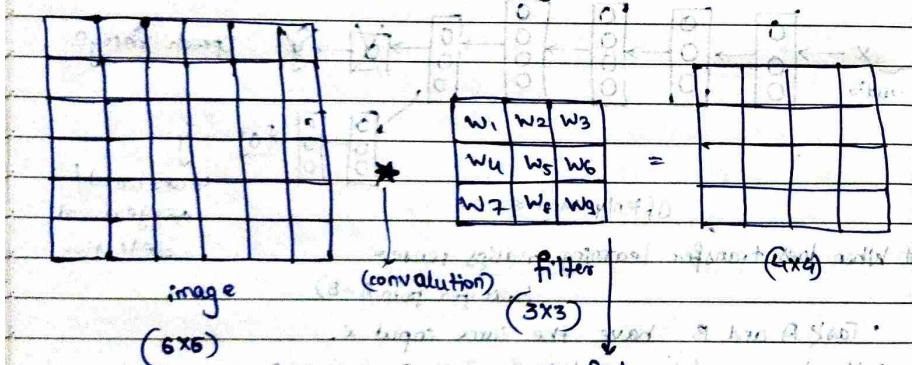
- 1) Let the data speak $x \rightarrow y \rightarrow$ "phonemes"
- 2) Less hand-designing of components needed.

↳ Cons:

- 1) May need large amount of data.
- 2) Excludes potentially useful hand-designed components.

↔ Key Question: Do you have sufficient data to learn a function of the complexity needed to map x to y ?

Course 04 → Week 1



e.g. of filters

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

NN properties
like contours/
edges.

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

sobel filter

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

Scharr filter

problems → 1) shrinking output
2) throwing away info. from edges.

① Padding →

↳ Adding a border to the image

$\text{f} = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

3×3

$f = 3$ 3×3

5×5 7×7

$n-f+1$ $7-3+1$

$n+2p-f+1$ $7+2 \cdot 1 - 3 + 1$

$= 28 = 5 \times 5 = (\text{preserved})$

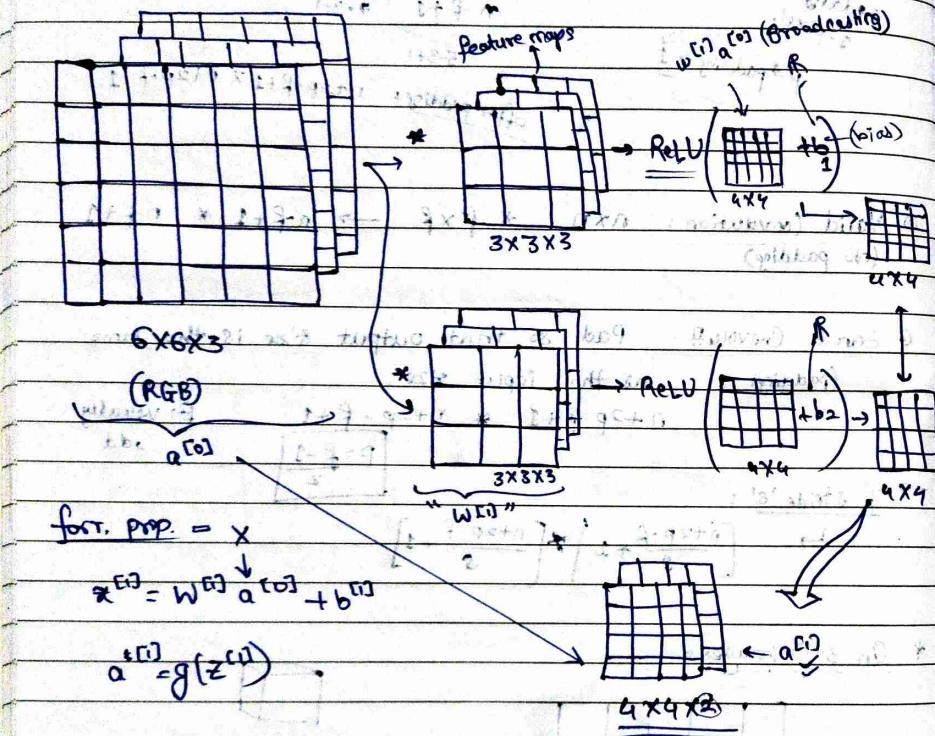
→ We can apply multiple filters to get multiple output images, we can stack them together for a combined effect.



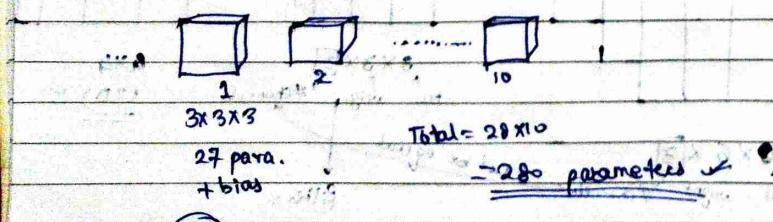
Summary: $n \times n \times n_c$ * $f \times f \times n_c$ → $(n-f+1) \times (n-f+1) \times n_c'$

$$\begin{matrix} 6 \times 6 \times 3 \\ 3 \times 3 \times 3 \end{matrix} \rightarrow \begin{matrix} 4 \times 4 \times 2 \\ \# \text{no. of filters} \end{matrix}$$

One Layer of CNN →



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



Notation →

→ If layer l is a convolutional layer:



$$f^{[l]} = \text{filter size}, \quad \text{Input: } n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$$

$$p^{[l]} = \text{padding}$$

$$s^{[l]} = \text{stride}$$

$$n_C^{[l]} = \text{no. of filters}$$

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

$$\text{Each filter's size: } f^{[l]} \times f^{[l]} \times n_C^{[l]}$$

$$\text{activations: } a^{[l]}$$

$$\rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

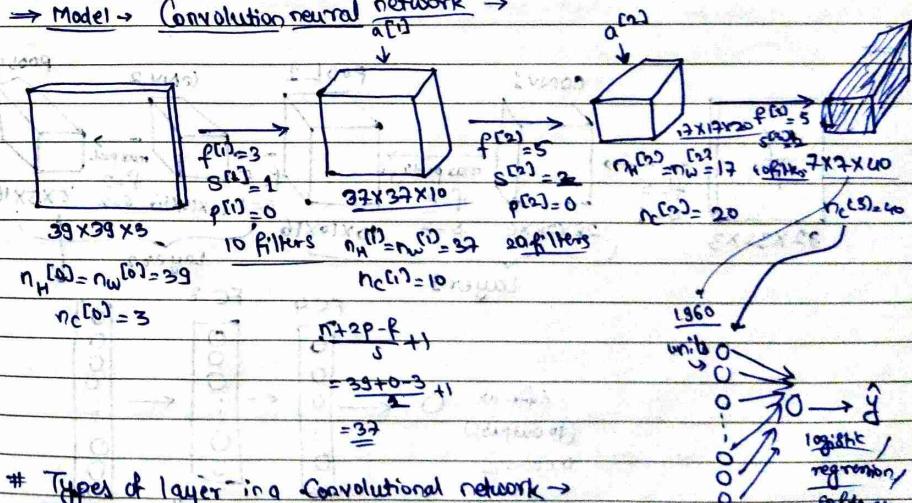
$$A^{[l]} \rightarrow m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]} \rightarrow \text{Batch / Gradi. Descent}$$

(set of m activ l)

$$\text{Weights: } f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]}$$

$$\text{bias: } n_C^{[l]} \cdot (1, 1, 1, n_C^{[l]})$$

→ Model → Convolution neural network →



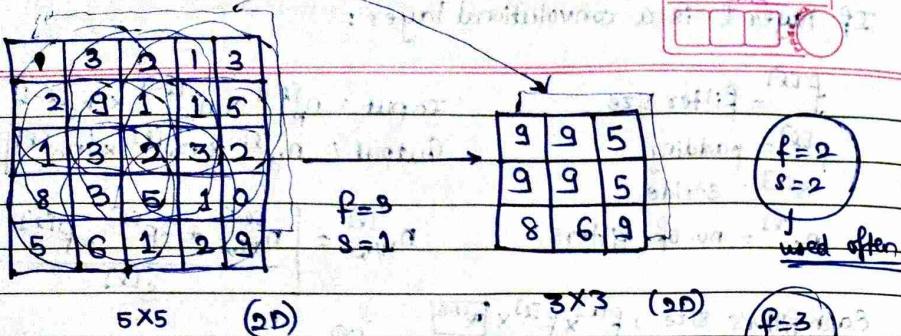
Types of layer in a convolutional network →

- Convolution (CONV)

- Pooling (POOL)

- Fully connected (FC)

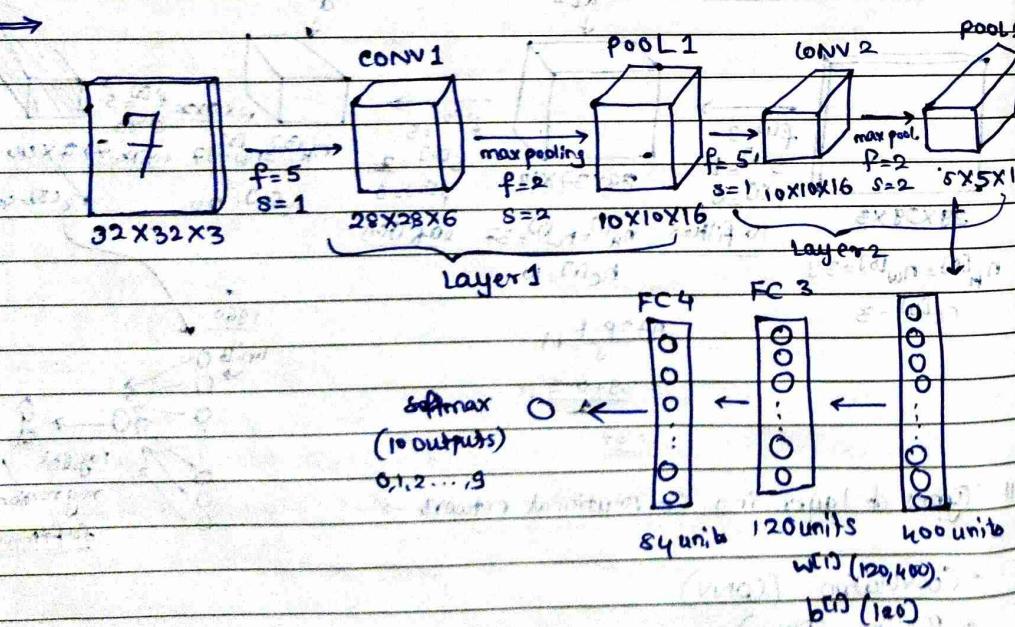
① Max Pooling : Used often



No parameters to learn!

$$\text{IF } (3D) \quad 5 \times 5 \times 2 \longrightarrow 3 \times 3 \times 2$$

② Average pooling : Instead of maximum, we take avg. of all nos.



→ Why Convolution?

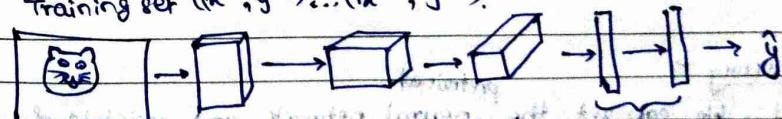
Relatively small no. of parameters (parameter sharing)

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

• Sparsity of connections: In each layer, each output value depends only on a small no. of inputs.

→ Putting it together →

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



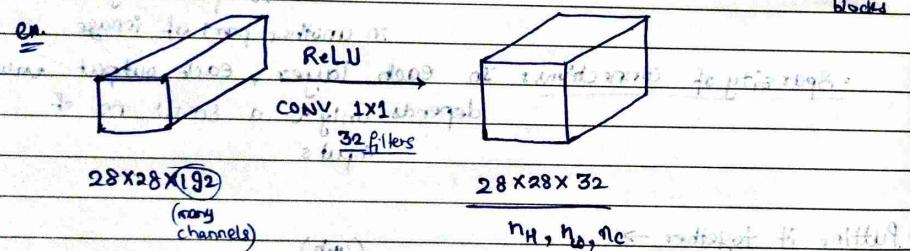
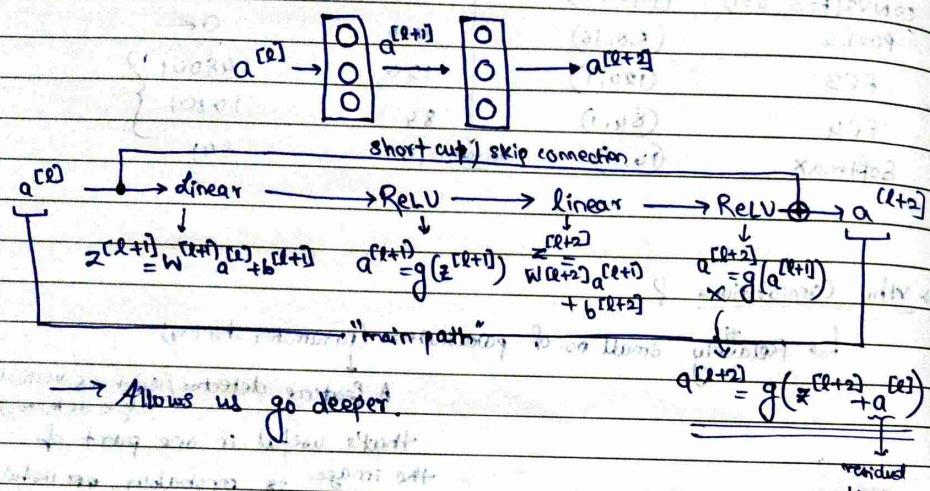
$$\rightarrow \text{cost } J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y^{(i)}, \hat{y}^{(i)})$$

→ Use gradient descent to optimize parameters to reduce J (or Adam).

→ Case Studies (Architectures of Neural Networks) =

3) LeNet-5 ✓

ResNets : Made up of 8 Residual blocks.

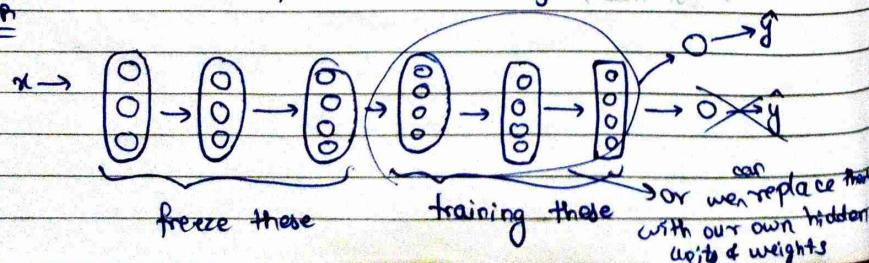


Transfer Learning →

④ !! We can get the neural network and weights of an open source model & change its (or replace it) softmax function with ours for our desired output.

We can freeze the hidden layers.

OP:



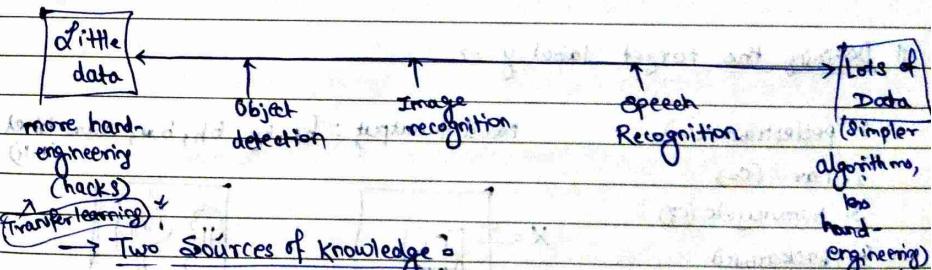
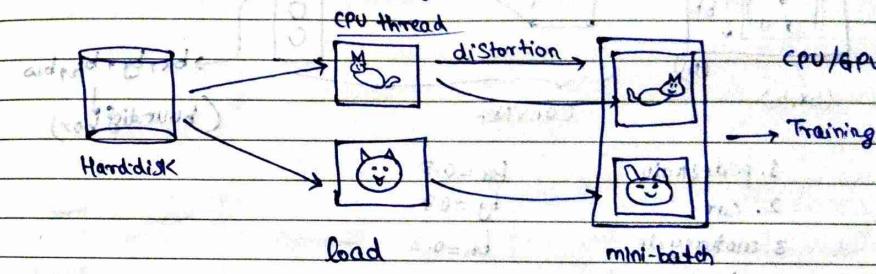
* Data Augmentation =

- Mirroring of image
- random cropping
- Shearing
- Local warping

* Color shifting →

→ Changing the RGB values of image such as adding more R, G than B (small changes).

4) Implementing distortions during training =

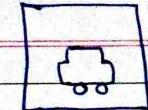


Tips for doing well on benchmarks →

- Train several networks independently & average their outputs. (Ensembling)
- Multi-crop at same time
 - Run classifier on multiple versions of test images & avg. results. (10-crop technique)
- Use pretrained models & fine-tune on your dataset

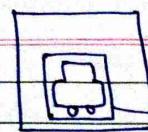
Week=3

Image Classification

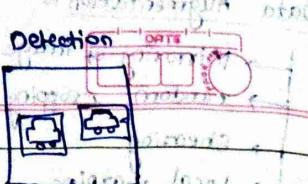


"Car"

Classify & localisation



"GT"

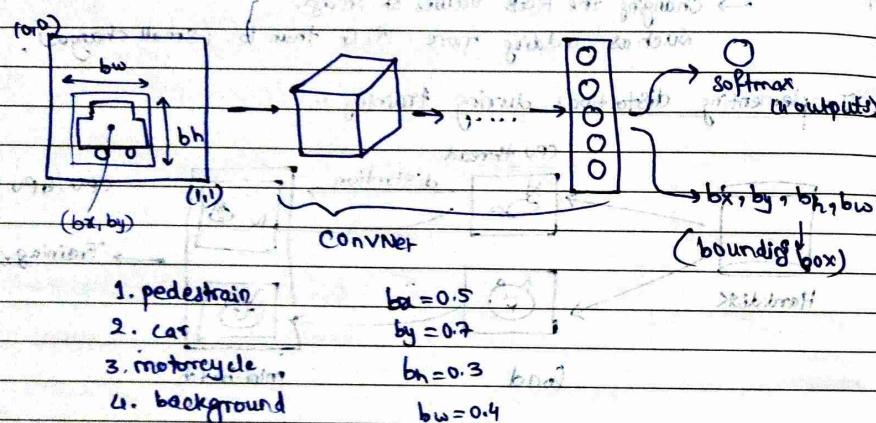


rectangle box

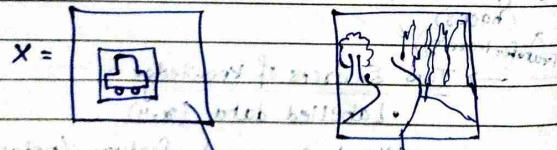
multiple

objects

1 Object

# Defining the target label $y \rightarrow$

1. pedestrian (c_1)
2. car (c_2)
3. motorcycle (c_3)
4. background



$$\mathcal{L}(\hat{y}, y) = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_n - y_n)^2$$

$$\text{if } y_i = 1 \rightarrow P_{c,i} = 1$$

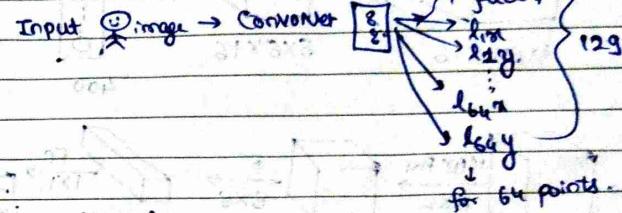
$$\text{if } y_i = 0 \rightarrow P_{c,i} = 0$$

$$y = \begin{cases} P_c & \text{Is there any obj?} \\ bx \\ by \\ bh \\ bw \\ c_1 \\ c_2 \\ c_3 \end{cases}$$

0	"don't care"
1	?
2	?
3	?
4	?
5	?
6	?

→ Landmark detection →

e.g. Key features of face like the jawline points, eye points, etc.

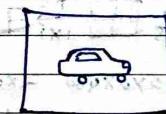


* Car Detection Example:

Training set: Output



1

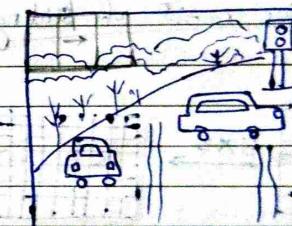


1



1

Sliding Window Algorithm →



we make this window & put it into ConvNet

→ we keep going (sliding) until we detect the car or all image.

→ we can now take a larger window & feed that to ConvNet

Disadv.: Very high computational cost.

Turning FC layer into Convolutional layer

→ Our algorithm can detect the object multiple times which is not good for us.

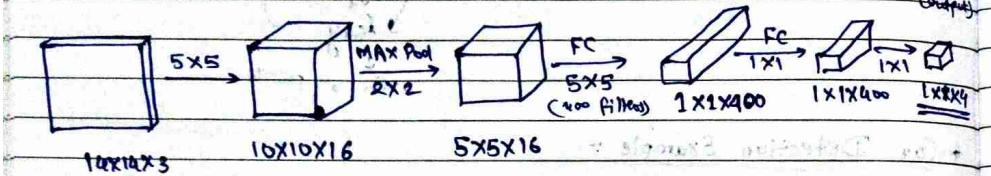
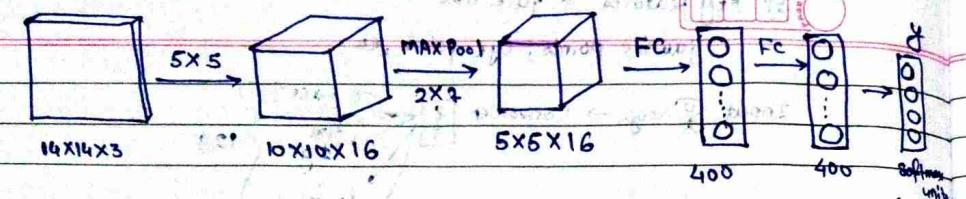
Therefore, ~~we need a step~~

Non-maximum suppression algorithm detects each object only one.

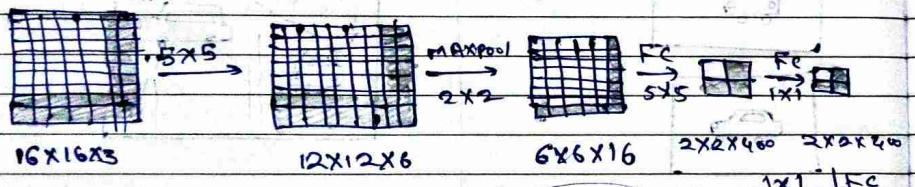
Non-max Suppression algorithm

→ First we look at the probability associated with each detection.

→ Maximum probability detection box suppresses other boxes with lower probability.



Convolution impl. of sliding windows

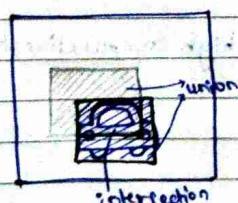


Less Computation

make all predictions at the same time

Instead of doing all the time

Evaluating localization of objects

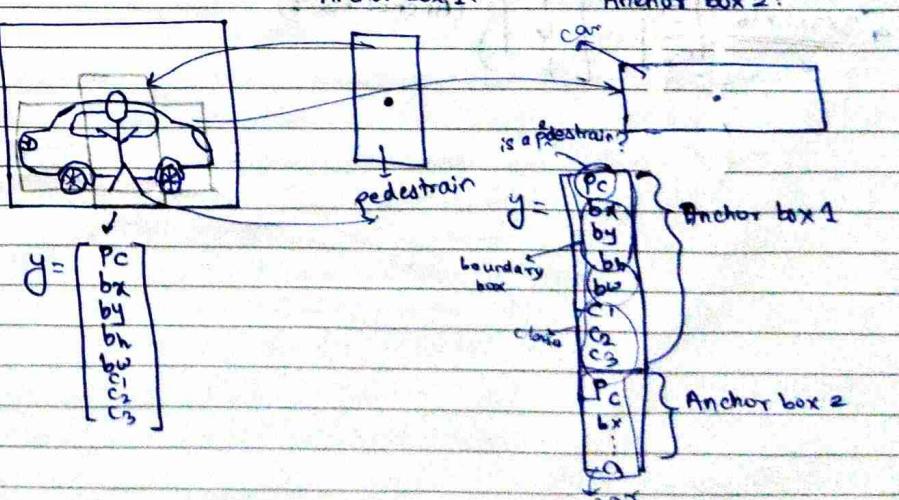


$$= \frac{\text{size of intersection}}{\text{size of union}}$$

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

predicted bounded box

Anchor Boxes



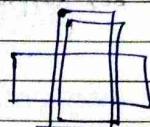
→ Previously: Each object in training image is assigned to grid cell that contains that object's midpoint.

→ With 2 anchor boxes: Each object in training image is assigned to grid cell that contains object's centre (midpt.) and anchor box for the grid cell with $\text{IoU} \rightarrow \text{highest}$

(grid cell, anchor box) pair

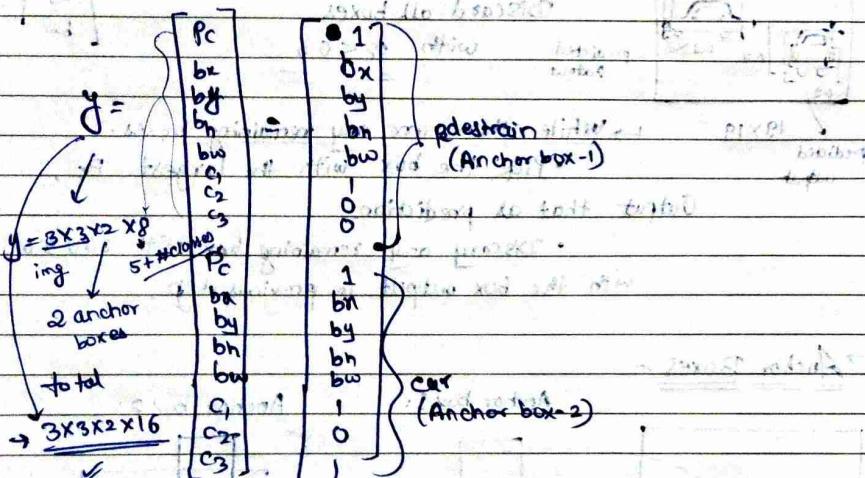
Output y:

$3 \times 3 \times 3$

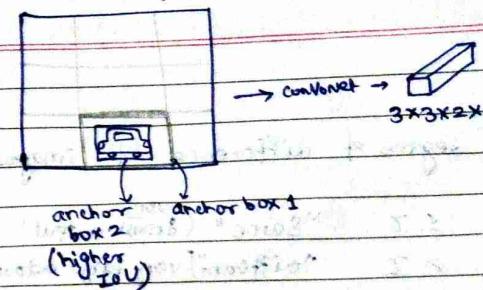


~~put~~:
3x3x16

previous ex. output →



YOLO Algorithm → Object Detection



DATE	
(anchor-1)	for this ex
Pc	0
bz	?
by	?
bh	:
bw	:
c1	:
c2	:
c3	:
Pc	1
bz	bz
by	by
bh	bh
bw	bw
c1	0
c2	1
c3	0
a.	

- For each grid cell, get 2 predicted bounding boxes.
 - Get rid of low probability predictions.
 - For each class (pedestrian, car, motorcycle) use non-max suppression to generate final predictions

- ⑥ R-CNN → The way we run the region proposal is called Region a segmentation algorithm. (Taking meaningful regions into consideration)

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

→ Fast R-CNN: Propose regions. Use convolution implementation of sliding window to classify all the proposed regions.

→ Feature - RNN: Use convolutional network to propose regions.

Week 4

99.9%

→ Face Verification: • Input image, name/ID

1:1 problem

- Output whether the input image is that of the claimed person.

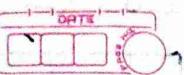
↳ Recognition :-

- Has a database of K persons.
- Get an input image.
- Output ID if the image is any of the K -persons (or "not recognised").

1:k problem

→ One-shot learning → Learning from one example to recognize the person again.

→ Triplet Loss function →

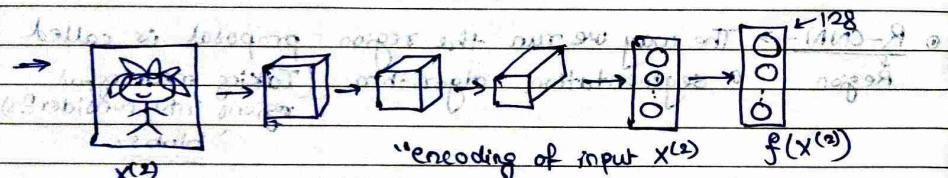
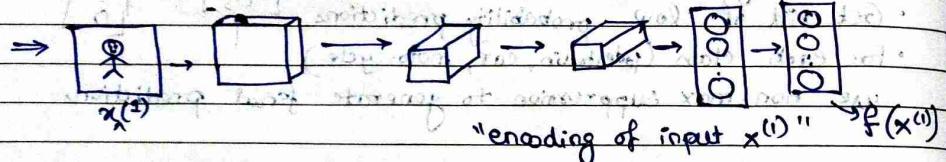


→ Similarity function:

$d(\text{img 1}, \text{img 2}) = \text{degree of difference b/w images}$

verif. If $d(\text{img 1}, \text{img 2}) \leq \tau$ "same" (small output)
ation $> \tau$ "different" (very large output)

Siamese Network →



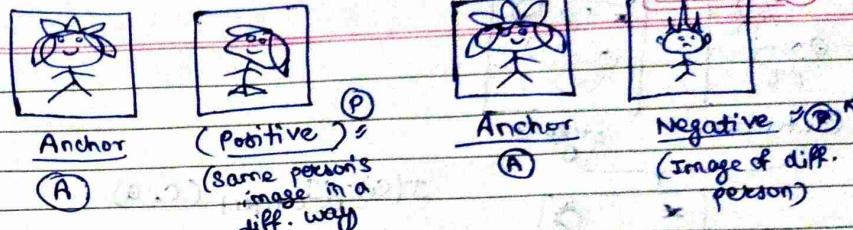
$$\rightarrow d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$$

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

Learn parameters so that:

• If $x^{(i)}, x^{(j)}$ are the same person, $\|f(x^{(i)}) - f(x^{(j)})\|_2^2$ is small 😊!

• If $x^{(i)}, x^{(j)}$ are the diff. person, $\|f(x^{(i)}) - f(x^{(j)})\|_2^2$ is large. 😭!



$$\text{Want/Objective: } \|f(A) - f(P)\|_2^2 \leq \|f(A) - f(N)\|_2^2$$

$\underbrace{\|f(A) - f(P)\|_2^2}_{d(A,P)} \quad \underbrace{\|f(A) - f(N)\|_2^2}_{d(A,N)}$

$$\Rightarrow \|f(A) - f(P)\|_2^2 - \|f(A) - f(N)\|_2^2 \leq \alpha (-\alpha)$$

$\downarrow \alpha$ another hyperparameter

→ Given 3 images A, P, N

$\overbrace{A, P, N}^{\substack{\text{diff.} \\ \text{diff.}}}$

$$L(A, P, N) = \max(\|f(A) - f(P)\|_2^2 - \|f(A) - f(N)\|_2^2 + \alpha, 0)$$

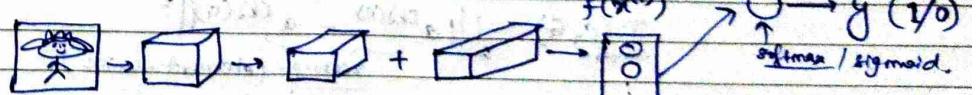
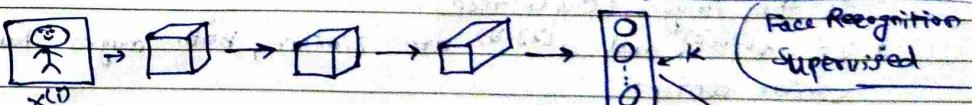
\hookrightarrow single triplet

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

• Choose triplets that're hard to train on,

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

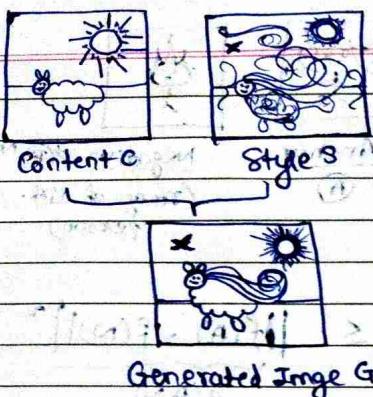
→ Use gradient descent to minimize the cost fun. ↴



$$\hat{y} = \sigma \left(\sum_{k=1}^{128} w_k (\|f(x^{(i)})_k - f(x^{(j)})_k \|_2^2) + b \right)$$

$$\approx \frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k} \quad X^2 \text{ Chi-square similarity}$$

→ Neural Style Transfer Cost function



$$J(G) = J_{\text{content}}(C, G)$$

$$+ \beta J_{\text{style}}(S, G)$$

Generated Image G

$\alpha, \beta \rightarrow$ hyperparameters.

→ Find the generated image G

1. Initiate G randomly

$G: 100 \times 100 \times 3$ → generates white noise image

2. Use gradient descent to minimize $J(G)$

$$G := G - \nabla J(G)$$

→ Content cost function:

$$\alpha J_{\text{content}}(C, G)$$

- Say you use hidden layer l (not too deep, not too shallow)
- Use pre-trained ConvNet. (Eg. VGG network)
- Let $a^{[l]}(C)$ and $a^{[l]}(G)$ be the activation of layer l on the images (C & G resp.)
- If $a^{[l]}(C)$ and $a^{[l]}(G)$ are similar, both images have similar content.

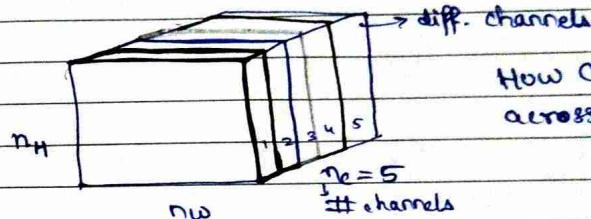
$$J_{\text{content}}(C, G) = \frac{1}{2} \| a^{[l]}(C) - a^{[l]}(G) \|^2$$

vectors. (element-wise)

→ Style cost function

→ Say you are using layer l 's activation to measure "Style".

Define style as correlation b/w activations across channels.



How correlated are the activations across diff. channels?

→ Using the degree of correlation between channels as a measure of style, which measures the degree to which the first channel is correlated with the second channel in the generated image.

→ Style Matrix:

Def $a^{[l]}_{i,j,k}$ = activation at (i, j, k) . $G^{[l]}(s)$ is $n_h^{[l]} \times n_w^{[l]}$ $\times n_c^{[l]}$

$$G_{kk'}^{[l]}(s) = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} \cdot a_{ijk'}^{[l]} \quad \leftarrow \text{style image}$$

$$G_{kk'}^{[l]}(G) = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]}(G) \cdot a_{ijk'}^{[l]} \quad \leftarrow \text{Generated image}$$

$\forall k, k' = 1, 2, \dots, n_c^{[l]}$

$$J_{\text{style}}^{[l]}(s, G) = \frac{1}{(2n_h^{[l]} n_w^{[l]} n_c^{[l]})^2} \sum_{k, k'} (G_{kk'}^{[l]}(s) - G_{kk'}^{[l]}(G))^2$$

$$\Rightarrow J_{\text{style}}(s, G) = \sum_l J_{\text{style}}^{[l]}(s, G)$$

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