

# Week 1 Foundations

Thursday, December 14, 2017 8:45 PM

# CNNs

Thursday, December 14, 2017 8:46 PM

# Computer Vision

Classification, Object detection, Neural Style Xfer

Large images req very large weight matrices, making it difficult to avoid overfitting

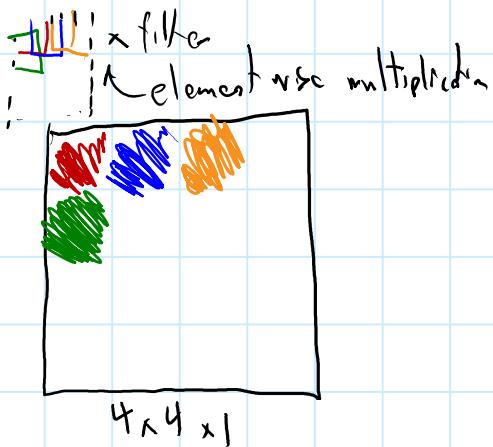
Sol'n? Convolutions!



## Edge Detection

Vertical edge detection,

$$\begin{array}{c} \text{convolve} \\ \curvearrowright \quad | \quad 0 \rightarrow \\ \star \quad | \quad 0 \rightarrow = \\ 6 \times 6 \times 1 \quad | \quad 0 \rightarrow \\ \text{image} \quad \quad \quad 3 \times 3 \end{array}$$



PyTorch: `:conv2d`    tf: `tf.nn.conv2d`    Keras: `Conv2D`

ex. image

$$\begin{bmatrix} 10 & 10 & 10 & 10 \\ 0 & 0 & 0 & 0 \end{bmatrix} \star \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} = \begin{bmatrix} 30 & 30 \\ 0 & 0 \end{bmatrix}$$



b7

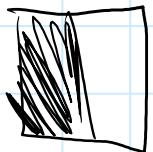
0 →

1 0 →

0 30 30

bright  
on  
leftdark  
on  
right

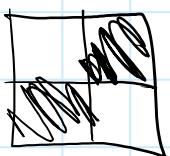
What if it's flipped?



→ Some convolution = same output \* -1

Horizontal

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



~~horizontal~~  
horizontal convolution =  $\begin{bmatrix} 0 & 0 & 0 & 0 \\ 30 & 10 & -10 & -30 \\ 30 & 10 & -10 & -30 \\ 0 & 0 & 0 & 0 \end{bmatrix}$

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

Sobel filter

(emphasizes center)

Can also use numbers themselves as parameters

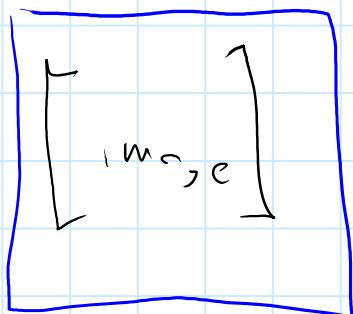
$$\begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{bmatrix}$$

# Padding

$$\begin{matrix} h \times n \\ \text{image} \end{matrix} \xrightarrow[\text{w/ filter}]{\text{Convolved}} f \times f = (h-f+1) \times (n-f+1)$$

Downsides: filters shrink images  
edge pixels are underutilized, center overemphasized

Soln: Padding



$$\begin{aligned} 6 \times 6 \text{ w/ } 1 \text{ px padding} &= 8 \times 8 \\ \text{put thru } 3 \times 3 \text{ filter} &\rightarrow 6 \times 6 \text{ output!} \\ \text{Output} &= (n+2p-f+1) \times (n+2p-f+1) \end{aligned}$$

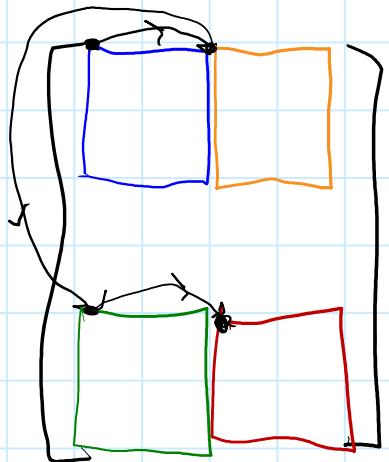
How much to pad? Term: valid & same

"Valid": no padding

"Same": Pad so output is  $n \times n$  ( $p = \frac{f-1}{2}$ )

# Strided Convolutions

$\text{Stride} = \text{pixel shift}$



Stride- $s$

1st

2nd

3rd

4th

$$\lfloor z \rfloor = \text{floor}(z)$$

$$\text{output} = \left\lceil \frac{n + 2p - f}{s} + 1 \right\rceil$$

(controls for stride, outside images)

Note in math texts convolution matrices are flipped before multiplication  
"cross-correlation"

$$\begin{matrix} 3 & 4 & 5 \\ 1 & 0 & 2 \\ -1 & 4 & 7 \end{matrix} \rightarrow \begin{matrix} 7 & 2 & 5 \\ 9 & 0 & 4 \\ -1 & 1 & 3 \end{matrix}$$

Why? preserves associativity

$$[(A * B) * C] = A * [(B * C)]$$

## Convolutions Over Volumes

Ex. RGB images

height width #channels  
 $n \times m \times 3$

6px

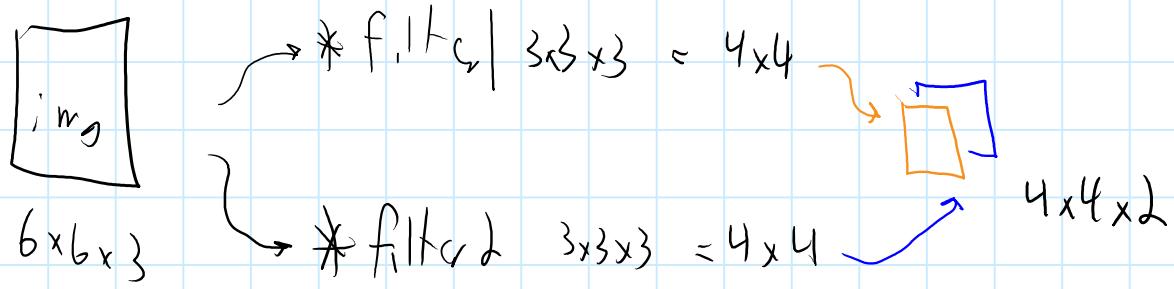
$$6 \times 6 \times 3 * \begin{array}{c} \square \\ \square \\ \square \end{array} = \begin{array}{c} \square \\ \square \\ \square \\ \square \end{array}$$



The channels of the filter scan over the respective channels of the image

## Multiple Filters

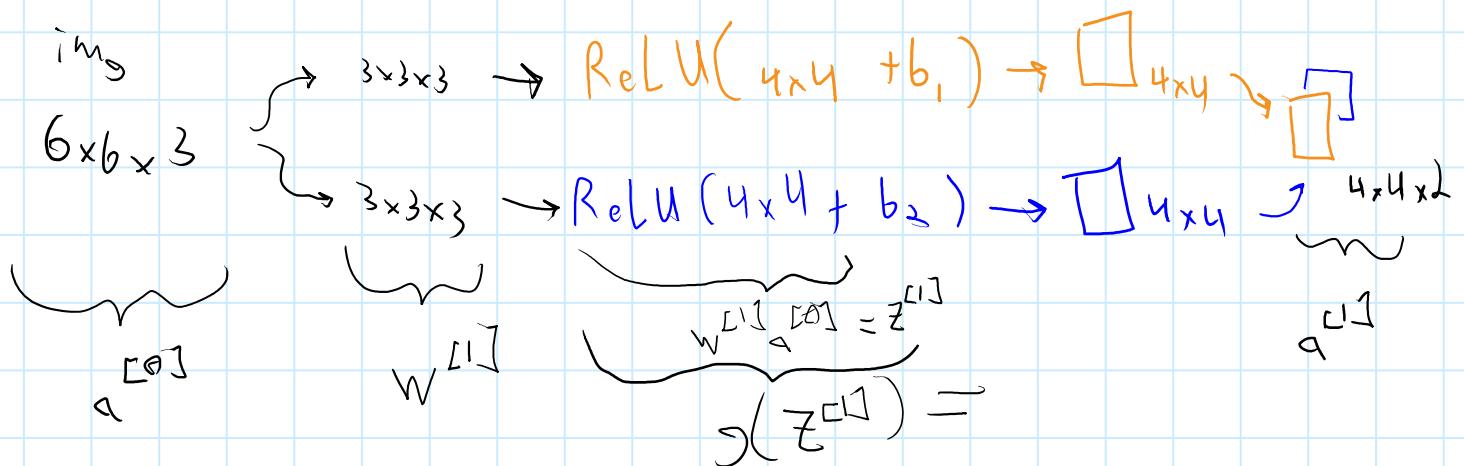
# Multple filters



$$h \times h \times n_c * f \times f \times n_c \rightarrow (h-f+1) \times (h-f+1) \times (n'_c)$$

↓  
 # filters  
 dependent on  
 stride as usual

One layer of a CNN



Notation

# Notation

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

$p^{[l]}$  = padding  
 $s^{[l]}$  = stride

$n_c^{[l]}$  = # filters

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

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

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

Each filter is  $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$

Activations:  $a^{[l]} \rightarrow n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

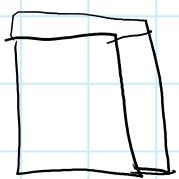
$A^{[l]} \rightarrow m \times n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

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

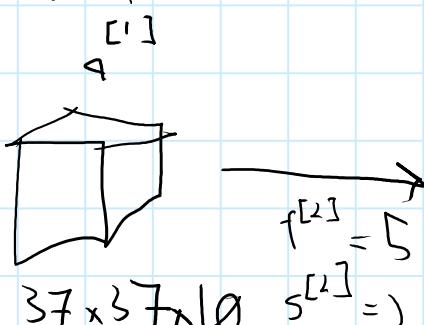
bias:  $n_c^{[l]} - (1, 1, \dots, n_c^{[l]})$



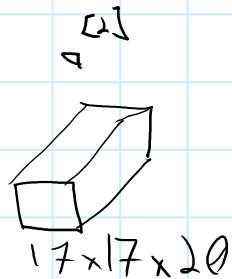
Simple CNN Example

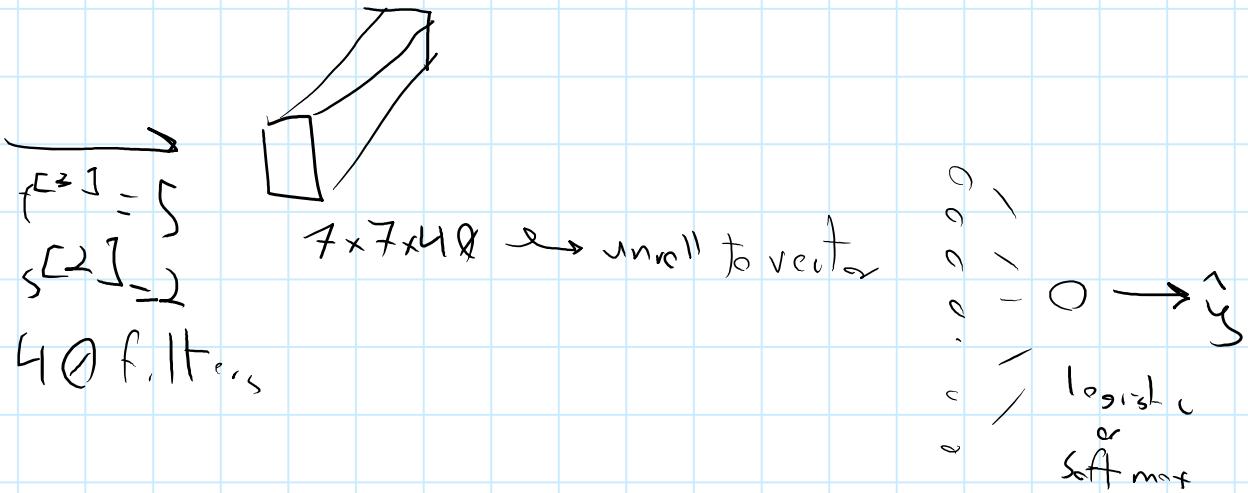


$39 \times 39 \times 3$   
 $f^{[1]} = 3$   
 $s^{[1]} = 1$   
 $p^{[1]} = 0$   
10 filters



$f^{[2]} = 5$   
 $s^{[2]} = 2$   
 $p^{[2]} = 0$   
20 filters





3 types of layers in a NN

Convolution (CONV)  
 Pooling (POOL)  
 Fully Connected (FC)

---

Pooling Layers

Max pooling

$$\begin{array}{cccc}
 1 & 3 & 2 & 1 \\
 2 & 4 & 1 & 1 \\
 1 & 3 & 2 & 3 \\
 5 & 6 & 1 & 2
 \end{array}
 \xrightarrow[f=2]{s=2} \begin{array}{cc}
 9 & 2 \\
 6 & 3
 \end{array}$$

if a feature is detected in an area, its presence is preserved

$n_c$  is preserved

$4 \times 4$

Average pooling

$$\begin{array}{cccc}
 1 & 3 & 2 & 1 \\
 . & . & . & .
 \end{array}$$

Not used as often as max pooling

$$\begin{array}{cccc}
 1 & 3 & 2 & 1 \\
 2 & 4 & 1 & 1 \\
 1 & 3 & 2 & 3 \\
 5 & 6 & 1 & 2
 \end{array} \rightarrow \begin{array}{cc}
 3.75 & 1.25 \\
 4 & 2
 \end{array}$$

Summary

$F$ : filter size

$s$ : stride

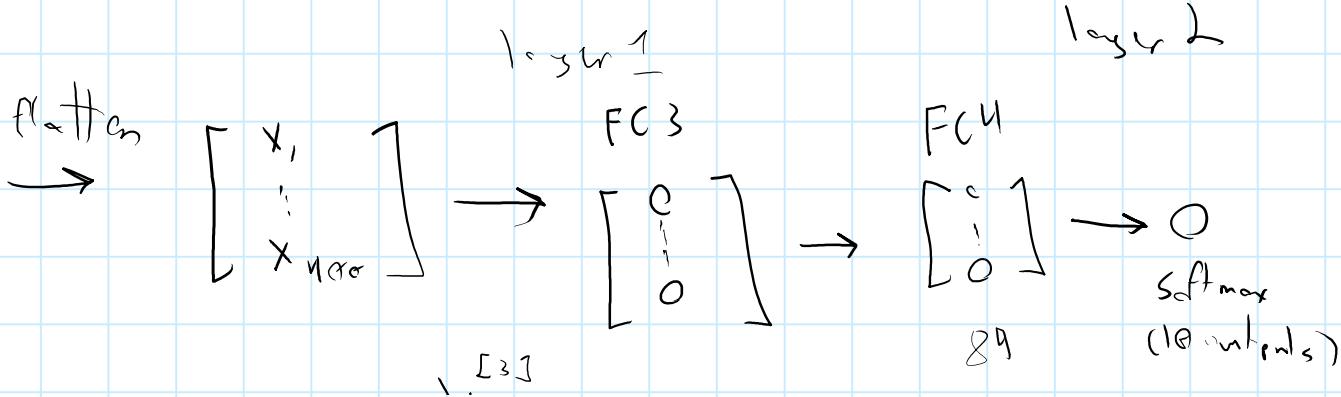
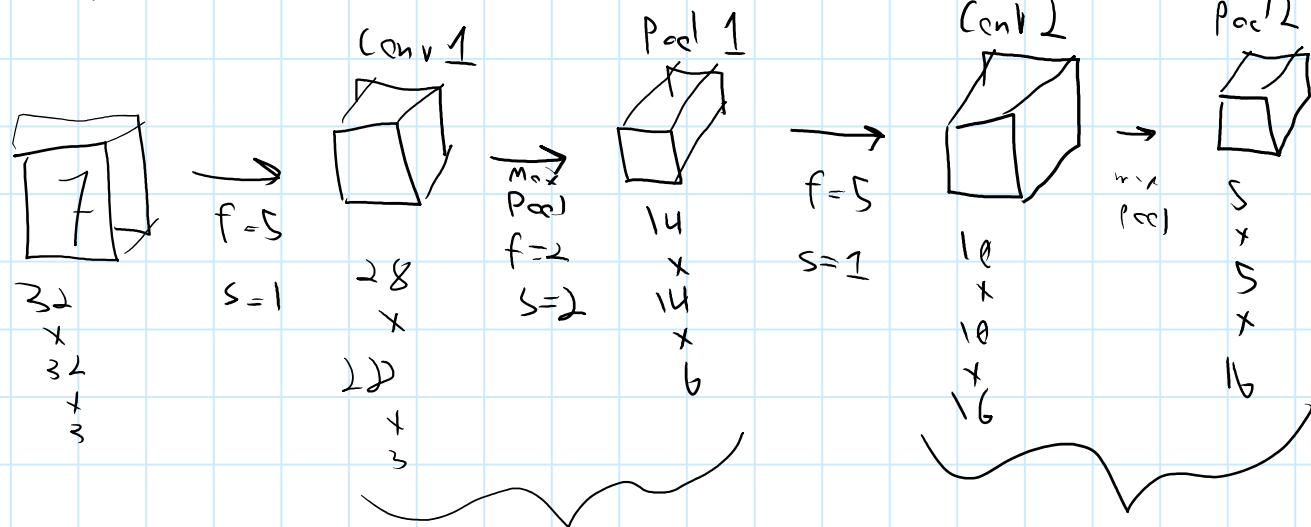
Max or avg

$p$ : padding (columns  $\geq 1$ )

$$h_h \times w_w \times n_c$$

$$\left\lfloor \frac{n_h - F}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_w - F}{s} + 1 \right\rfloor \times n_c$$

CNN Ex



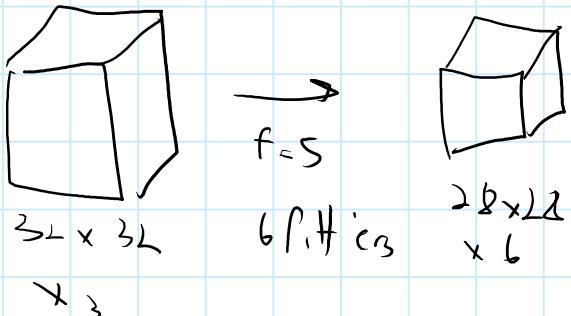
(120, 400)

## Why Convolutions?

Parameter sharing

each filter has  $5 \times 5 + 1$  params

$26 \times 6$  filters = 156 parameters



$\times 3$

$3 \otimes f_2$

4,704

$\xrightarrow{\text{SNN}}$   $\rightarrow 14M$  parameters!

Feature detectors are useful in different parts of an image (e.g., edge detection)

Connection sparsity:

Output values only depend on small # of outputs

# Week 2 Deep Conv Models - Case Studies

Thursday, February 1, 2018 11:12 PM

# Case Studies

Thursday, February 1, 2018

11:13 PM

Why look at case studies?

Intuition building

Outline:

Classic networks

- LeNet-5

- AlexNet

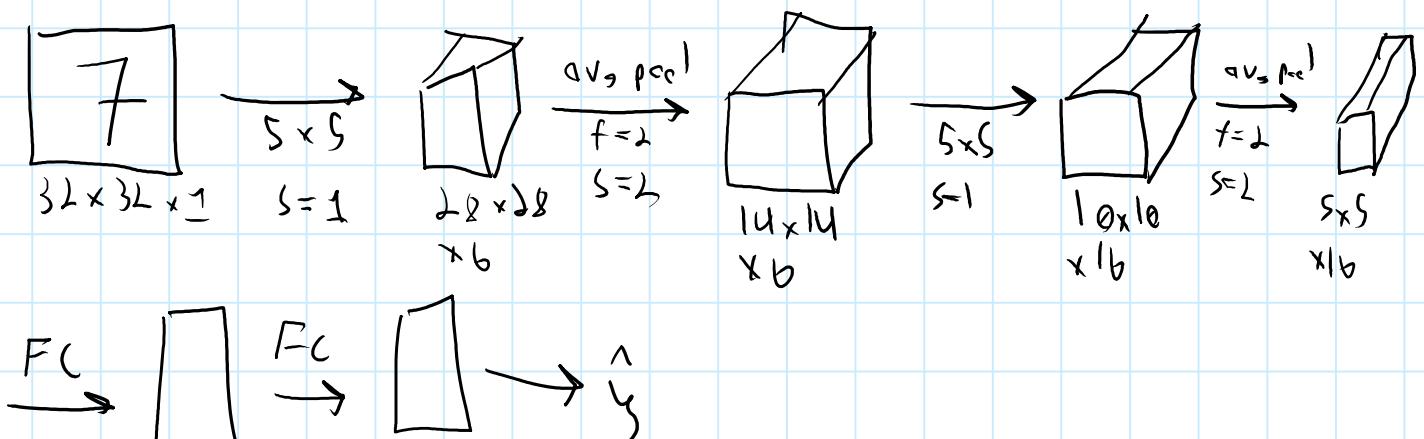
- VGG

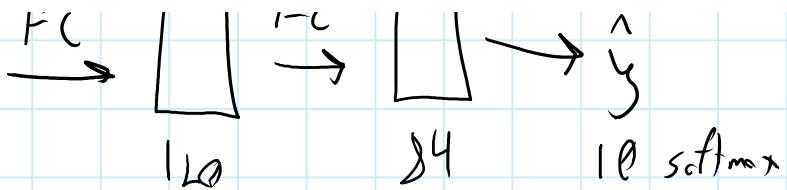
Residual Net (w, b)

Inception

Classic Networks

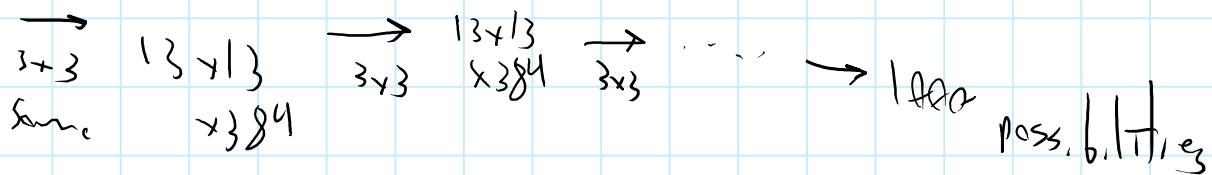
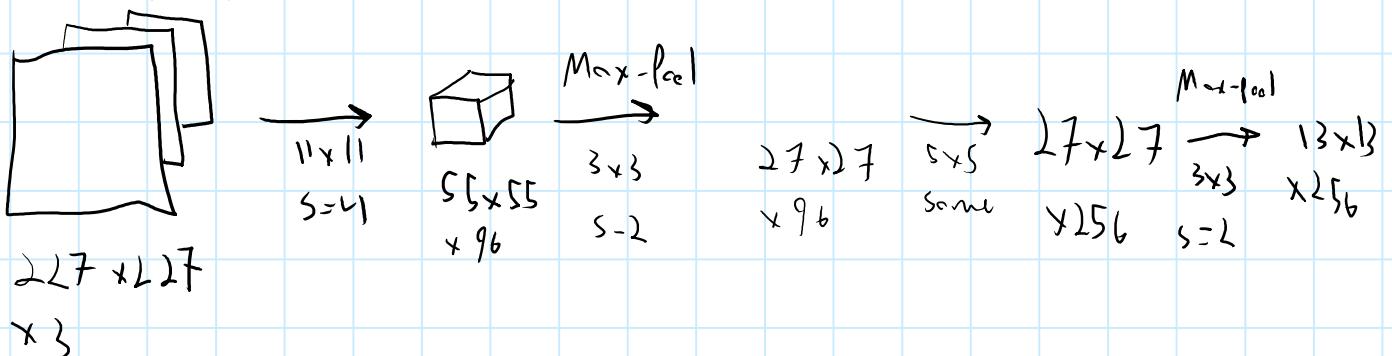
LeNet-5 Handwritten characters





$h_{H,W} \downarrow$     $n_c \uparrow$        $\text{Conv} \rightarrow \text{pool} \rightarrow \text{conv} \rightarrow \text{pool} \rightarrow \text{fc} \rightarrow \text{fc} \rightarrow \text{output}$

## Alex Net



Similar to LeNet but much bigger, used ReLU

Local Response Normalization (no longer used)

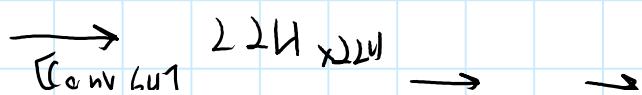
//

$\sqrt{G_1 G_2 - 16}$



$\text{Conv} = 3 \times 3$  filter,  $s=1$ , same

$\text{Max-Pool} = 2 \times 2$ ,  $s=3$

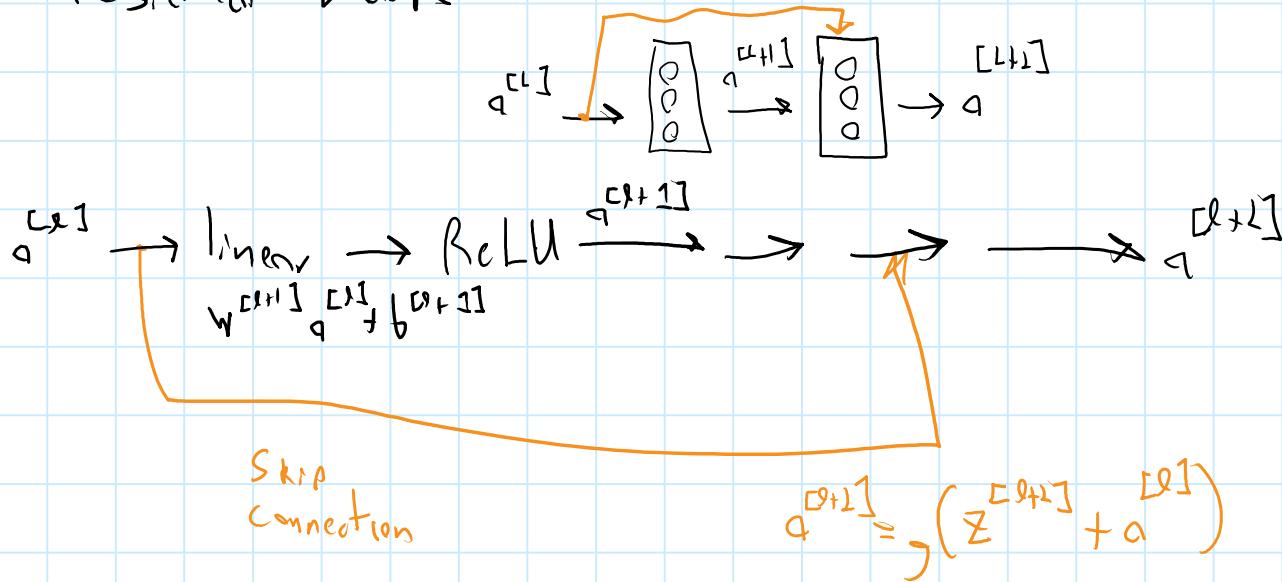


$224 \times 224$        $\times 2 \rightarrow \times 64$        $\text{pool}$        $\text{conv} \times 2 \dots$   
 $\times 3$

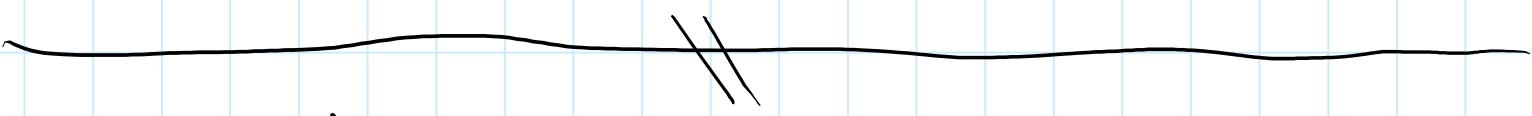
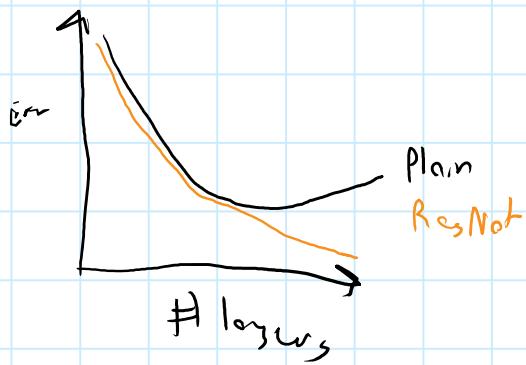


## Res Nets

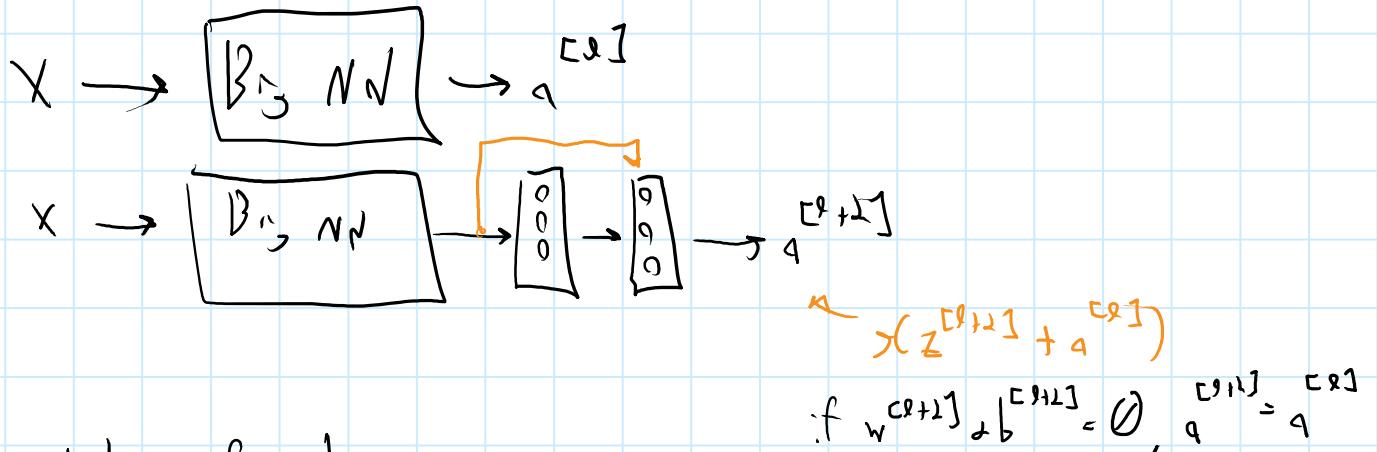
### Residual Block



Repeated every 2 layers, helps w/ vanishing/exploding gradients



## Why Res Nets Work



Identity function is easy to learn,  
still network has trouble w/ this

Mismatch size corrections

$$a^{[d+1]} \in \mathbb{R}^m, \text{ but } z^{[d]} \in \mathbb{R}^n$$

$$a^{[d+1]} = (z^{[d+1]} + W_s a^{[d]}) ; W_s \in \mathbb{R}^{m \times n}$$

Network in Network,  $1 \times 1$  convolution

$$\begin{array}{c} \text{Input} \\ 6 \times 6 \\ \times 32 \end{array} \star \begin{array}{c} \text{Filter} \\ 1 \times 1 \\ \times 32 \end{array} = \begin{array}{c} 6 \times 6 \times \\ \# \text{Filters} \end{array}$$

takes a slice through layers, like 32 neurons into 1 ReLU neuron

Can be used to reduce  $n_c$ , just like pooling can shrink  $n_H, n_W$   
Can also introduce nonlinearity

## Inception Network Key Ideas

Stack multiple filters into one output volume

Reducing computation cost:

Run  $1 \times 1$  before larger conv, get same volume but drastically

~ "bottleneck layer"

reduce cost



## Inception Network

# Week 3 Detection Algorithms

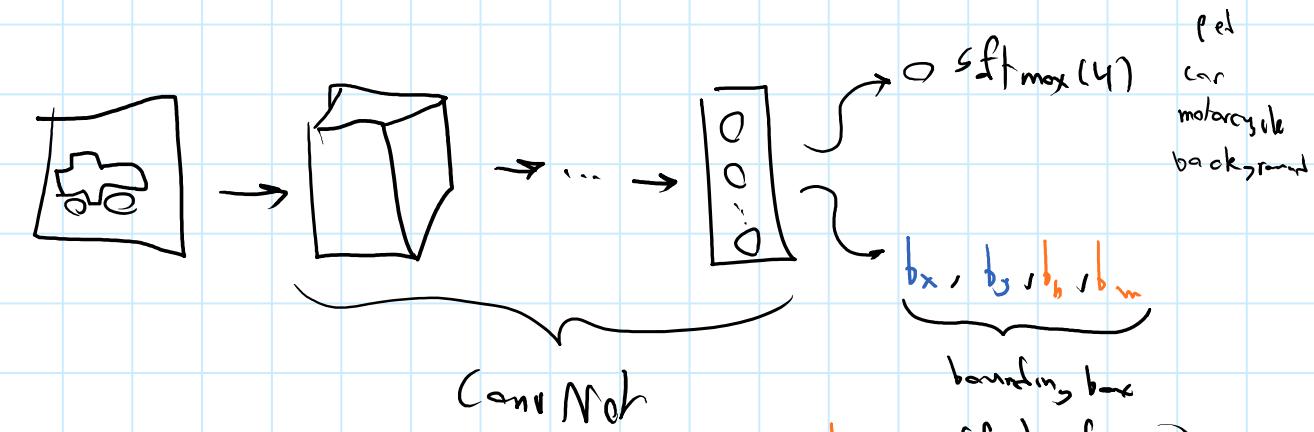
Monday, February 19, 2018 12:14 PM

# Detection Algorithms

Monday, February 19, 2018 12:15 PM

# Object Detection

Object Localization (bounding box around object in image)



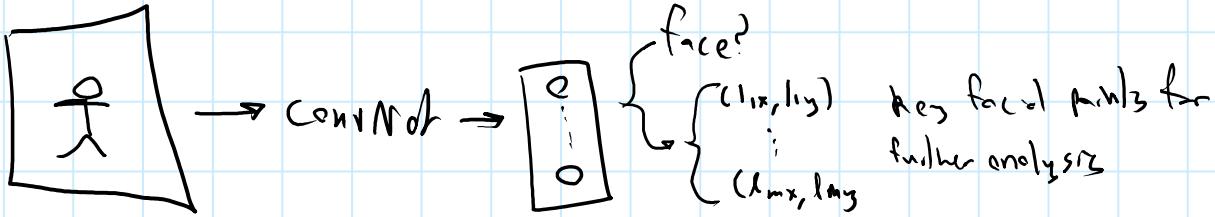
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_z \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

$p_c$  }  $P(\text{object})$   
 $b_x$  } bounding box  
 $b_y$  }  
 $c_1$  } which class,  
 $c_2$   
 $c_3$

$$\mathcal{L}(\hat{y}, y) \left\{ \begin{array}{l} y_i = 1 \rightarrow \sum_i (\hat{y}_i - y_i)^2 \\ y_i = 0 \rightarrow (\hat{y}_i - y_i)^2 \end{array} \right.$$

and mark Detection  
x & y coord of important features (e.g. corner of eye, body joints)

$\rightarrow$  rf



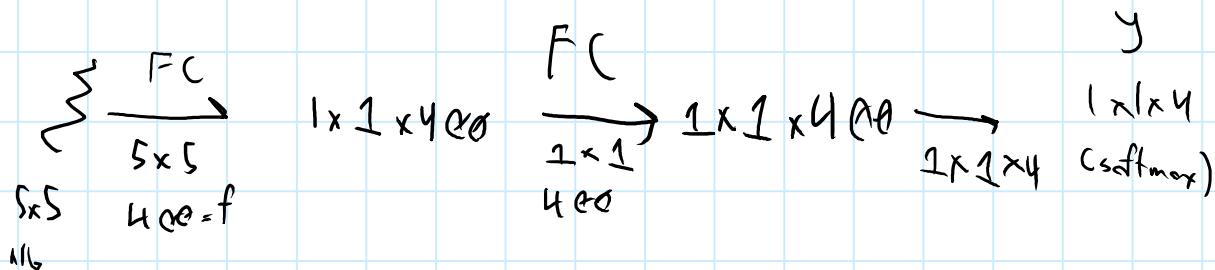
## Object Detection

Closely cropped labeled images, then sliding window detection  
Repeat w/ larger windows

High computational cost to be accurate, solved w/ convolutions

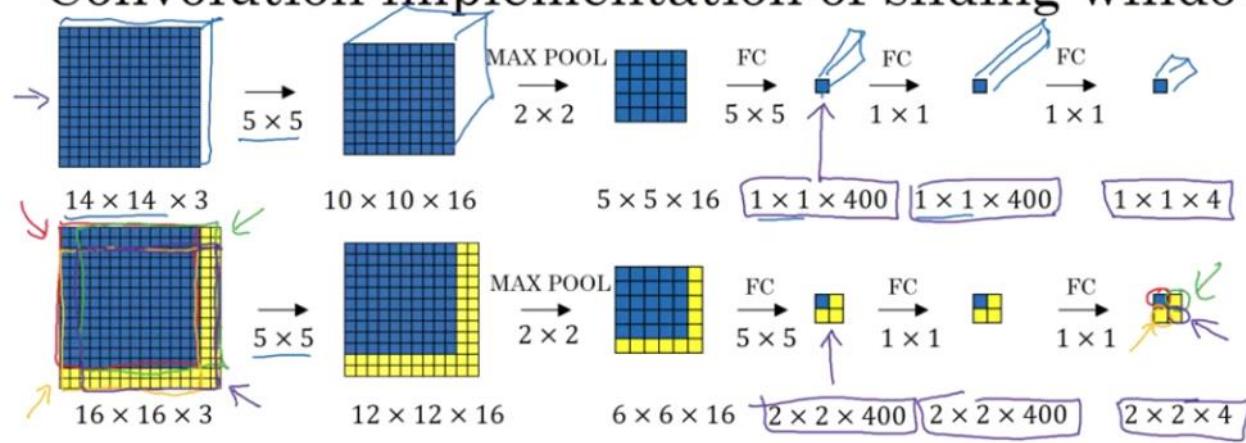
## CNN Implementation

Turning FC layers into CN



CNN implementation of sliding windows

# Convolution implementation of sliding windows



Screen clipping taken: 2/20/2018 12:42 PM

~~Bounding Box Predictions~~

YOLO Algorithm

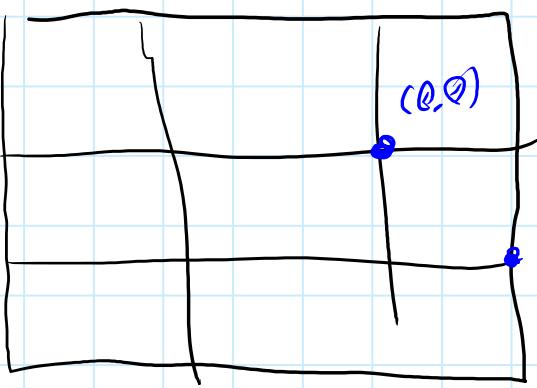
Divide image into grid,

For each one, each grid cell

( $n \times 3$  grid) will have a  $3 \times 3 \times 8$  output

$$y = \begin{bmatrix} p_1 \\ b \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Specifying bounding boxes



relative to grid cell

$b_h \text{ and } b_w$  may be  $> 1$



## Intersection Over Union (IoU)

Size of intersection

If  $\text{IoU} \geq 0.5$

Size of union

(of computed & ground truth bounding boxes)



Non-max

Suppression

Solution to issue of algorithm detecting single object multiple times

Take bounding box with  $\max p_c$ , negate ones which overlap  
repeat for all remaining,

$\text{IoU} \geq 0.5$

Best to carry out independently for each class of object ( $c_i$ )

---

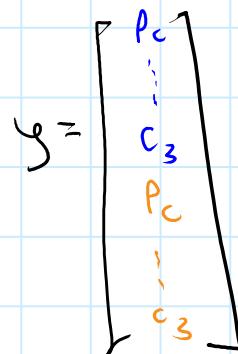
## Anchor Boxes

Predefined shapes for object classes

e.g.



Person



Each object in a training image is assigned to a grid cell that contains object's midpoint & anchor box w/ highest IoU

Has trouble marking 2 objects w/ similar anchor box overlapping

Can pick object shapes via k-means

---

## YOLO Algorithm

Ex. Detecting people & cars

Split into cells

For each cell set 2 bounding boxes

Eliminate low pc predictions

For each class, use non-max suppression to generate final predictions



Region Proposals      R-CNN

Run a segmentation algorithm to section out poss. objets

Place bounding boxes around "blobs", run classifier only on said blobs

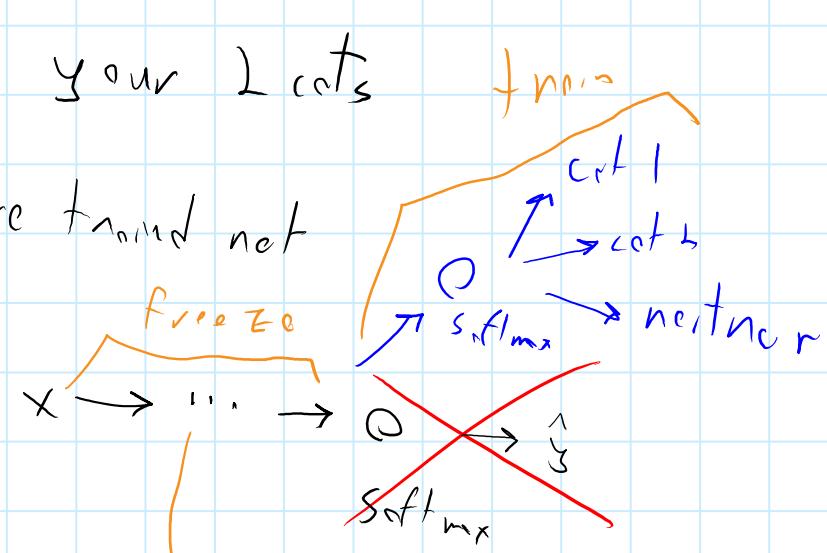
Speed up by using sliding windows over blobs

Mostly still slower than YOLO

# Practical Advice

Monday, March 19, 2018 9:41 PM

X for learning



Precompute activation to reduce compute time

Partial data? freeze fewer layers, train rest

Full set? Use preweights as "random" initialization

## Data Augmentation

Mirroring images

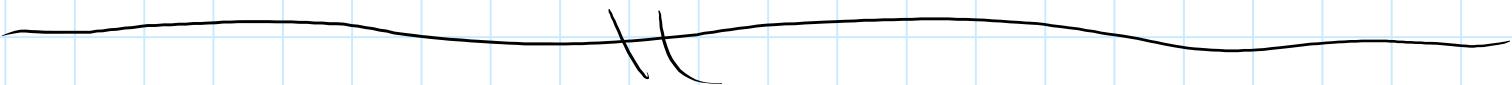
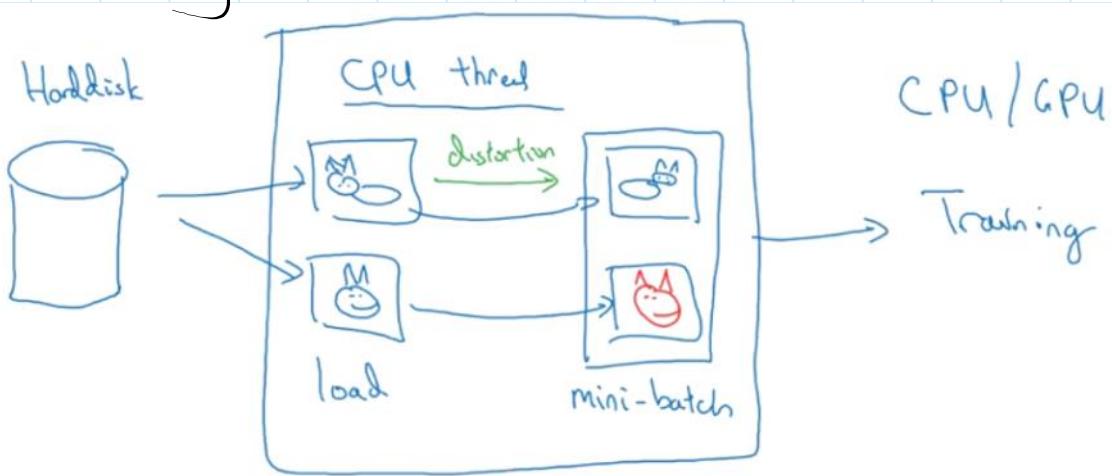
Random cropping

Rotation, shearing, local warping used less

Color shift (random permutation of RGB channels)

- can be done with PCA
- ex. image is heavily RB, little G  
will shift RB more than G

## Implementing



# Week 4 Special Applications

Monday, March 26, 2018 8:13 PM

# Face Recognition

Monday, March 26, 2018 8:14 PM

What is Facial Recognition?

Verifying v. Recognition [linearity check]

Verification

- Input: image, Name / ID
- Output: image matches credentials

Recognition

- Database of  $K$  persons
- Input image
- Output if image is any of  $K$  persons \*

\* Recall from Bayesian stats that even 99% is not accurate for large  $K$ ! 99% is not good enough for  $K=100$

---

One Shot Learning

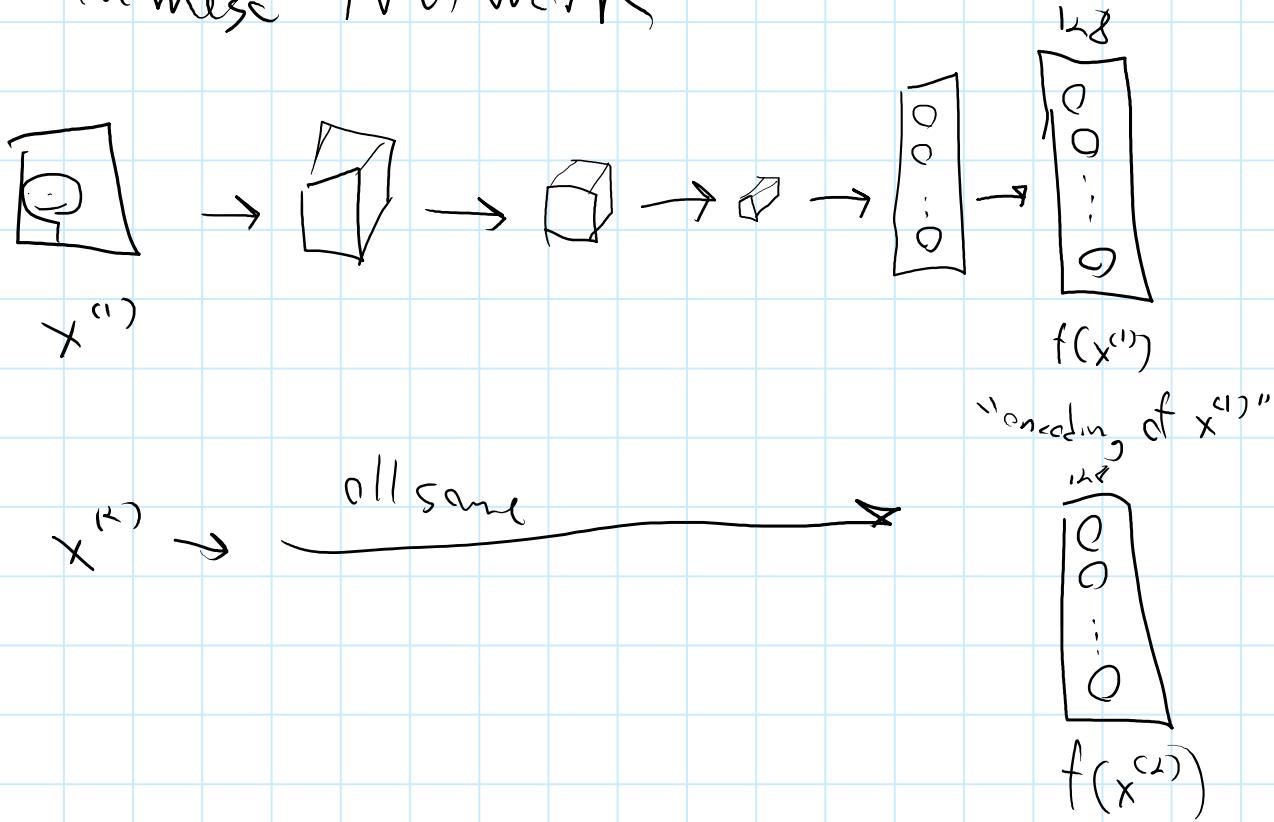
Learn from one example to recognize person again

Using a ConvNet is not efficient for this,  
use a similarity function

$d(\text{img}_1, \text{img}_2) = \text{° of diff between images}$

$$d(\text{img}_1, \text{img}_2) \begin{cases} \leq \tau & \text{"same"} \\ > \tau & \text{"different"} \end{cases}$$

## Siamese Network



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

Parameters of NN define an encoding,  $f(x^{(1)}) \leftarrow 128$   
 Learn parameters so that:

if  $x^{(1)}, x^{(2)}$  are the same person,  $\|f(x^{(1)}) - f(x^{(2)})\|_2^2$  is small  
 conversely, large if different

# Triplet Loss Function

[See Schaff et al., 2015, FaceNet]

A - anchor (reference)

$P - +$

$N -$

some  
different

$$\text{Want } \underbrace{\|f(A) - f(P)\|^2}_{d(A, P)} + \alpha \leq \underbrace{\|f(A) - f(N)\|^2}_{d(A, N)}$$

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

But a trivial sol'n is  $f(\text{img}) = \vec{0}$   
Prevent by  $\alpha$

## Loss Function

Given 3 images  $(A, P, N)$

$$\mathcal{L}(A, P, N) = \max(d(A, P) - d(A, N) + \alpha, 0)$$

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

Training set: ~10k pictures of 1k persons

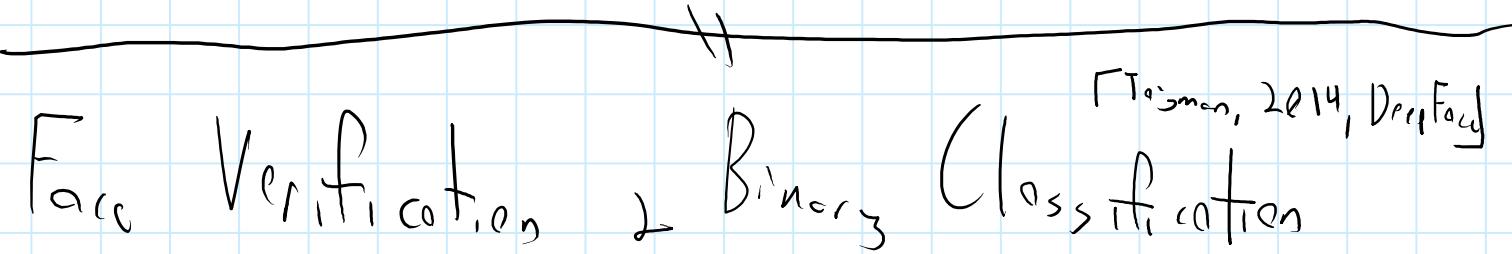
Issue: if  $A, P, N$  are chosen randomly,

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

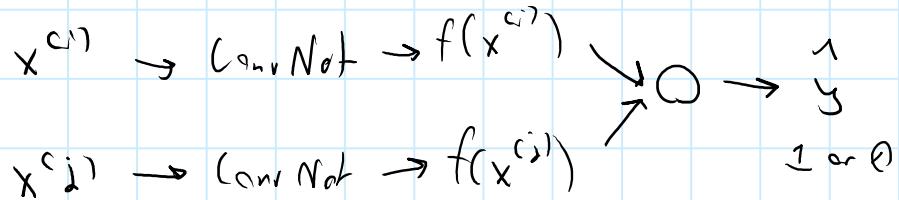
is easily satisfied

Choose triplets that are hard to train on

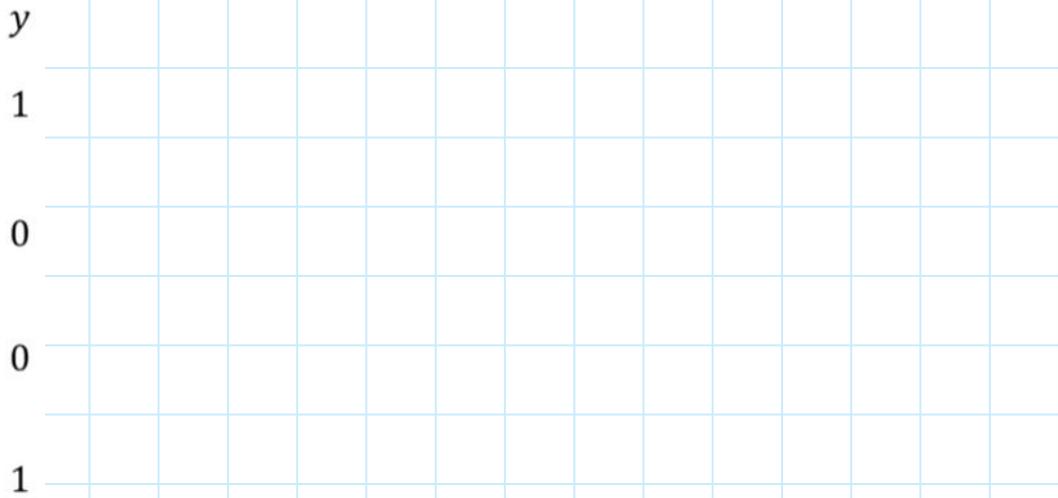
$$d(A, P) \approx d(A, N)$$



Alternative to triplet loss function



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



Screen clipping taken: 3/26/2018 11:09 PM

# Neural Style Transfer

Monday, March 26, 2018

8:14 PM

## What is Neural Style Transfer



Content (c)

Style (s)

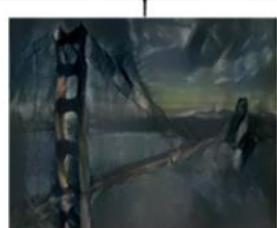


Generated image (g)



Content (c)

Style (s)



Generated image (g)

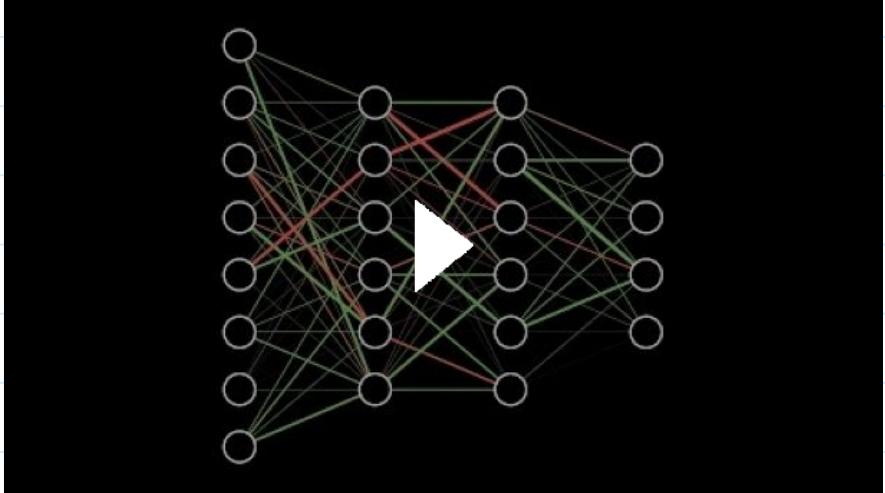
Screen clipping taken: 3/26/2018 11:13 PM

## What are Deep ConvNets learning?

Early layers detect features such as edges & colors  
As layers build features become more complex,  
Circles → eyes → faces, etc.

See 3Blue1Brown video for fantastic description,

[But what \\*is\\* a Neural Network? | Chapter 1, deep learning](#)



## Cost Function

[Gatys, 2015, A novel algorithm  
of artistic style]

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

Similarity to original content

Similarity to reference style

Rough Algorithm:

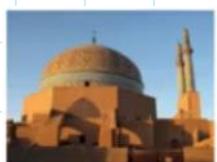
1) Initialize  $G$  randomly  $G : 100 \times 100 \times 3$  initially noisy

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

$$G := G - \frac{\partial}{\partial G} J(G)$$

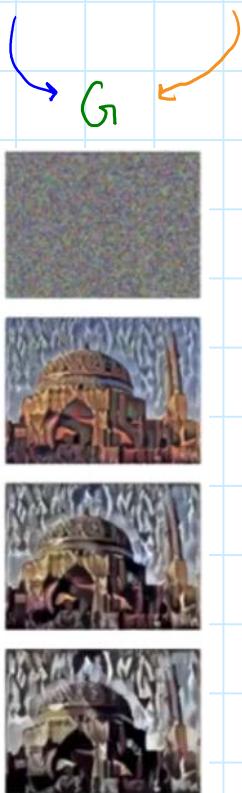
C

S



Screen clipping taken: 3/28/2018 9:01 PM

(. . .)




---

||

## Content Cost Function

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

- Must carefully select layer  $l$  for content cost
  - Too shallow  $\rightarrow$  too similar
  - Too deep  $\rightarrow$  feature placement will be arbitrary
- Use  $\text{protruding ConvNet}(L \xrightarrow{f_l(C)} V(G))$
- Let  $a \xrightarrow{f_l(C)}$  &  $a \xrightarrow{f_l(S)}$  be activation of layer  $l$  for each image
- If  $a \approx a$ , content is similar

$$J_{content}(C, G) = \| a \xrightarrow{f_l(C)} - a \xrightarrow{f_l(G)} \|_2^2$$

$$J_{content}(C, G) = \left\| \underbrace{a^{text} - a^{style}}_{L_2 \text{ norm of unaligned vectors}} \right\|^2$$

## Style Cost Function

What does "style" mean?

Correlation between activations across channels

e.g. similarity between R + Y channels at each voxel

Intuition  
correlation

Red channel detecting red lines (texture)

Yellow channel detecting orange haze (tint)

## Style Matrix

$$\begin{aligned} & [R] & H \times W \times C \\ & Q_{i,j,k} = \text{activation}_{\otimes}(i, j, k), \\ & G^{\text{feature}} \text{ is } n_c \times n_c, \\ & G_{k,k'}^{\text{feature}} = \sum_i \sum_j Q_{ijk} Q_{ijk'}, \quad G = \sum_k \sum_{k'} G_{k,k'} \\ & k=1, \dots, n_c, \\ & \text{"Unnormalized Cross Covariance"} \end{aligned}$$

Def for both  $S \leftarrow G$

$$\begin{aligned} J_{\text{style}}^{[\lambda]}(S, G) &= \|G^{[\lambda](S)} - G^{[\lambda](G)}\|_F^2 \\ &= \sum_k \sum_{k'} \left( G_{kk'}^{[\lambda](S)} - G_{kk'}^{[\lambda](G)} \right)^2 \quad \text{Frobenius Norm} \end{aligned}$$

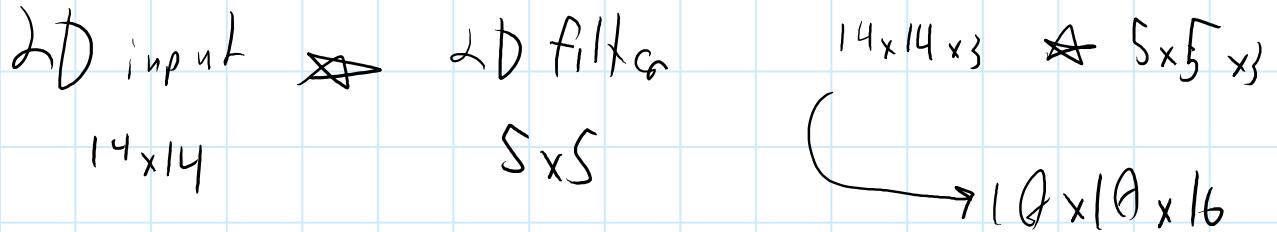
$$J_{\text{style}}(S, G) = \sum_k \lambda^{[k]} J_{\text{style}}^{[\lambda]}(S, G)$$

$\uparrow$   
hyperparameter

# 1D and 3D Generalizations

Wednesday, March 28, 2018 10:24 PM

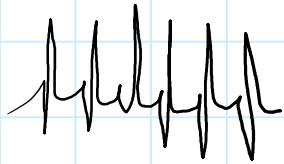
1 D



1D EKG Data

vector of y-values (pk heights)

14



$\star$  Gaussian filter

5



$14 \Rightarrow 5$

$10 \times 16$

$14 \times 16 \Rightarrow 5 \times 16$

$n_{channels}$   
 $6 \times 32$

3 D

e.g., CT Scan, channel is slice

$14 \times 14 \times 14 \times 1$

16 filters

$\star$

$5 \times 5 \times 5 \times 1$

$\rightarrow 10 \times 10 \times 10 \times 32$