

## 4-1\_Convolution neural network

### 4.1 Computer vision

Computer vision is one of the areas that's been advancing rapidly thanks to deep learning. Deep learning computer vision is now helping self-driving cars figure out where the other cars and pedestrians around so as to avoid them.

Deep learning is even enabling new types of art to be created.

Two reasons excited about deep learning for computer vision :

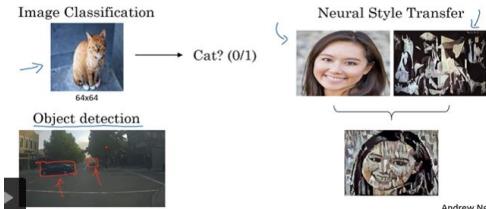
>1. First, rapid advances in computer vision are enabling brand new applications to view, though they just were impossible a few years ago.

And by learning these tools, perhaps you will be able to invent some of these new products and applications.

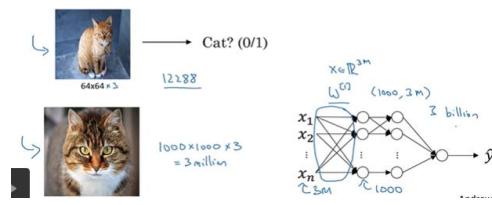
>2. Second, even if you don't end up building computer vision systems per se, because the computer vision research community has been so creative and so inventive in coming up with new neural network architectures and algorithms, is actually inspiring that creates a lot cross-fertilization into other areas as well.

e.g. while working on speech recognition, actually took inspiration from ideas from computer vision and borrowed them into the speech literature

#### Computer Vision Problems



#### Deep Learning on large images



#### 1. Examples

>1. Image classification

Input is like 64 x 64 image and try to figure out, is that a cat?

>2. object detection\_computer vision problem

e.g. while building a self-driving car, maybe don't just need to figure out that there are other cars in this image. But instead, need to figure out the position of the other cars in this picture, so that your car can avoid them.

In object detection, usually, not just figure out that these other objects but also draw boxes around them - they can be multiple objects/cars in the same picture.

>3. neural style transfer

Want a picture repainted in a different style.

So neural style transfer: have a content image, and a style image. And then have a neural network put them together to repaint the content image, but in the style of another image. So, algorithms like these are enabling new types of artwork to be created.

#### 2. Computer vision problem:

##### The inputs can get really big:

e.g. small image size/feature: 64 x 64 x 3=12288 --> weight matrix dimension in the first layers: (n1, 12288)

while big image size/feature might be: 1000 x 1000 x 3=3million-->first hidden layers, weight matrix dimension is (n1, 3million)

with so many parameters, it's difficult to get enough data to prevent neural network from overfitting, and the memory requirements to train in neural network with three billion parameters is a bit infeasible.

But for computer vision applications, we don't want to be stuck using only tiny little images. we want to use large images.

--> better implement the convolution operation, which is one of the fundamental building blocks of convolutional neural networks.

#### 4.1-2 edge detection example

The convolution operation is one of the fundamental building blocks of a convolutional neural network

##### 1. Compute vision problem

In previous videos, get intuition how the early layers of the neural network might detect edges and then the some later layers might detect cause of objects and then even later layers may detect cause of complete objects like people's faces.

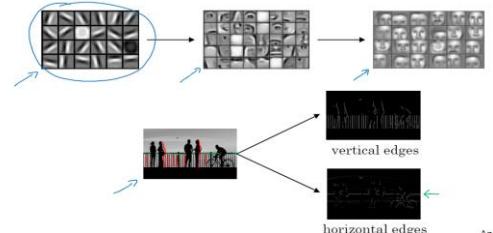
this item see how you can detect edges in an image.

##### 1.1 Edge detect

e.g. detect object in one image, first thing might detect vertical edges in this image: vertical lines, where the buildings are, as well as kind of vertical lines idea all lines of these pedestrians and so those get detected in this vertical edge detector output.

And might also want to detect horizontal edges so for example, there is a very strong horizontal line where this railing is and that also gets detected sort of roughly here.

#### Computer Vision Problem



##### 2. Vertical edge detection

E.g: grey image 6x6 -->detect vertical edge

>1.method: using a filter/kernel: to convolt the image: kernal : input \* kernal

>2.Convolution calculation:

>>>1. do element-wise products between kernal and the areas covered in image

>>>2. Add up all the result,

>>>3. then move kernal horizontally step by step to get different horizontal result

>>>4. and then move vertically step by step to get vertical result.

note:

1. '\*' here means convolution, while in python only means element-wise multiply.

2. Image , filter, and result matrix have different dimensions

##### Vertical edge detection

$$3 \times 1 + 1 \times 1 - 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + 1 \times 1 + 8 \times -1 + 2 \times -1 = -5$$

3	0	1	2	7	4
1	5	6	3	0	1
2	4	5	6	7	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

$6 \times 6$

1	0	-1
0	1	0
0	0	1

$3 \times 3$

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

$4 \times 4$

e.g. in the last: get to know the edge between light and deep dark area.

### 3. Vertical edge detection example

filter is defined to detect vertical edge: three columns 1, 0, -1 , could visualize as color: bright, middle-grey, dark.

Input image: visualize as bright , middle-grey.

Conv. Result: grey-bright-grey: the light region corresponds to having detected the vertical edge in middle of input image. bright region in the middle of result visualization means there is a strong vertical edge right down the middle of the image.

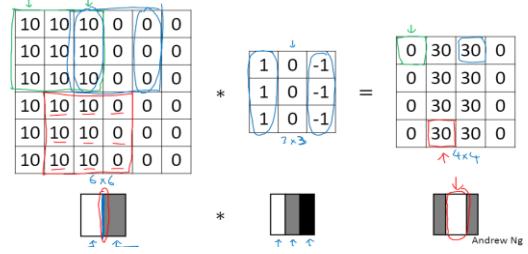
Note: detected edge seems really thick 30, due to working on a very small image . when using a larger image e.g 1000 x 1000, the conv. will do a pretty good job in detecting vertical edge.

one intuition to take away from vertical edge detection:

a vertical edge is a three by three region, since using a  $3 \times 3$  filter- where bright pixels on the left, do not care that much what is in the middle, and dark pixels on the right. The middle in this  $6 \times 6$  image is really where there could be bright pixels on the left and darker pixels on the right and that is why it thinks its a vertical edge over there.

The convolution operation gives you a convenient way to specify how to find these vertical edges in an image.

### Vertical edge detection



### 4.1-3 more edge detection examples

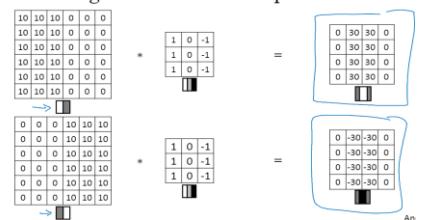
Item purpose: learn the difference between positive and negative edges, that is, the difference between light to dark versus dark to light edge transitions. And also see other types of edge detectors, as well as how to have an algorithm learn.

#### 1. Vertical edge detection examples

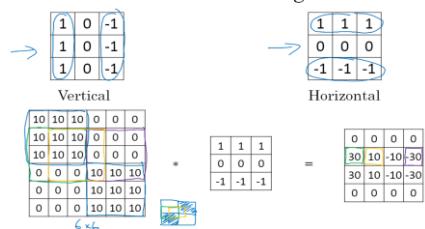
eg.: image flipped: dark -->bright,  
filter no change -->result image: get -30 (shows dark to bright rather than 30-bright to dark)in the middle:  
grey-dark-grey

if don't care which of these two cases it is, could take absolute values of this output matrix. But this particular filter does make a difference between the light to dark versus the dark to light edges.

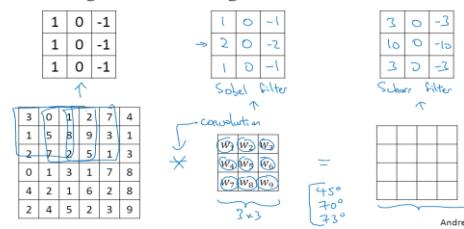
#### Vertical edge detection examples



#### Vertical and Horizontal Edge Detection



#### Learning to detect edges



#### 2. Horizontal Edge detection:

a horizontal edge : is a  $3 \times 3$  region where the pixels are relatively bright on top and relatively dark in the bottom row.

Input image: right top and left down part are bright , others dark.

Filter: horizontal filter: 1; 0; -1: bright-edge - dark.

Result image:

30: corresponds to input image right top area-bright;

-30: correspond to input image left down area-dark;

10, -10: corresponds to transmission area in the input image. when input image size is bigger, the result value for transmission will be smaller.

Summary, different filters allow you to find vertical and horizontal edges.

#### 3. Detect filter

**Sobel filter:** 1, 2, 1; put a little bit more weight to central row (central pixl), make it a little bit more robust

**Sharr filter:** different properties . is vertical kernal, if turn 90 degree, get horizontal kernal.

complicated pic edge: kernel /detect filter number could not be hand-picked, treat these parameter instead learn by using back propagation.

may could learn not only vertical and horizontal edge, but also 45, 30, 71 degree edge, etc.

### 4.1-4 padding

#### 1. Conv. cons

pic dimension: mxn, kernel dimension axb

-->Conv result image dimension:  $(m-a+1) \times (n-b+1)$ , shrink than original image

cons:

1. result shrink: may loose original feature

esp. in deep neural network, if shrink every layer, then end up with a very smaller image.

2. image corner used only once while middle pixel in image convoluted many times-->lose/throwing away many info. near the edge of image.

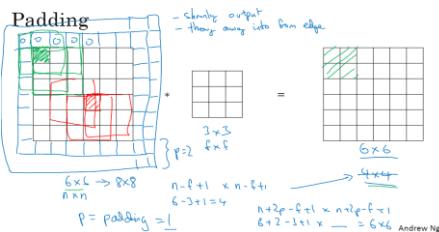
#### 2. Padding

addressing:for convolution cons, padding image

##### >1. Padding effect:

padding pixels around edges (padding 0)-->increase result dimension: padding influence around pixels in result, which influence by input image corner.

-->this effective maybe not quite throwing away but counting less the information of the edge of a corner of the image is reduce.



#### Valid and Same convolutions

$$\text{"Valid": } n \times n \times f \times f \rightarrow \frac{n-f+1}{2} \times \frac{n-f+1}{2}$$

"Same": Pad so that output size is the same as the input size.

$$n \times n \times f \times f \rightarrow \frac{n+2f-f-1}{2} \times \frac{n+2f-f-1}{2} \Rightarrow P = \frac{f-1}{2}$$

f is usually odd

$P = \frac{f-1}{2}$

$P = 2$

Andrew Ng

### 3.Padding method:

>1. 'valid' convolution: no padding

>2. Same convolution: add padding so output size is same as input size

-->padding size  $p=(f-1)/2$

#### 3.1 kernel size always odd, rarely is even number in computer vision, as:-

>1. if kernel size is even, then need asymmetric padding: padding more on left side/right side.

when kernel size is odd, padding size will be symmetric and padd more natural: padding same size on left and right.

>2 odd size kernel has a central position/pixel, could talk about the position of filter.

sometimes in computer vision it's nice to have a distinguisher, it's nice to have a pixel, you can call the central pixel so you can talk about the position of the filter.

Probably get just fine performance even if use an even size kernel, but if stick to the common computer vision convention, usually just use odd number  $f$  for kernel size.

### 4.1-5 strided convolutions

Strided convolutions is another piece of the basic building block of convolutions as used in Convolutional Neural Networks.

#### 1. Strided conv.

filter/kernel moving in the input image, instead move one step, move 's' step every time.

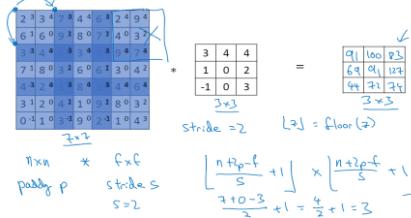
convolution result size:  $(\text{image size} + 2 * \text{padding size} - \text{kern size}) / \text{strided step} + 1$

Note:

if result size calculation  $(n+2p-f)/s+1$  is not integer: round down calculation (floor of result size), means:

Implement principle: do conv computation only if kernel is fully contained within the image/image + padding. if any part of kernel hangs outside of image+padding, do not convolution computation.

#### Strided convolution



#### Summary of convolutions

$n \times n$  image     $f \times f$  filter

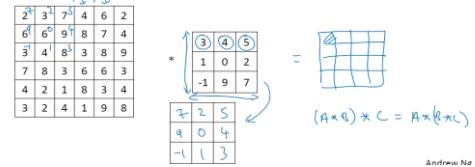
padding  $p$     stride  $s$

Output size:

$$\left[ \frac{n+2p-f}{s} + 1 \right] \times \left[ \frac{n+2p-f}{s} + 1 \right]$$

Technical note on cross-correlation vs. convolution

Convolution in math textbook:



### 2. Summary of conv

Result image size/dimension:  $(n+2p-f)/s+1$

Note: we can choose all of these numbers ( $n, p, f, s$ ) so that there is an integer result, yet sometimes don't have to do that and rounding down is just fine as well.

### 3. Technical note:

convolution in math textbook:

before doing the element-wise product and summing, flip filter/kernel on the horizontal as well as the vertical axis.

official convolution is: image \* flipped kernel (flipped vertically, and then horizontally-mirrored vertical, then mirror horizontal).

double flipping enable convolution operator to enjoy associativity property:  $(A+B)*C=A*(B+C)$

cross-correlation : image \* kernel

note: in deep learning will not flipped kernel, so actually is doing cross-correlation convolution, but normally just call it convolution operator.

### 4.1-6 convolution over volumes

#### 1. Conv. On RGB images/volume

E.g.: Detect features instead of on grey image but RGB image: a stack of three 2D images:  $n_h, n_w, n_c$  (channel number)

Filter: 3D filter: also have 3 layers/stack number, corresponding to input image:  $f \times f \times n_c$  (input image channel number): have same 3rd dimension/channel number as input image channel.

Result image: 2D:  $[(n_h + 2p - f)/s + 1, (n_w + 2p - f)/s + 1, 1]$  (only one layer-2D)

Volume conv computation: element-wise product (over all channels) and sum; then move filter same as 2dimension, move kernel to right/down by stride step (3rd dimension direction no move as cover all 3rd dimension data by kernel (same 3rd dimension size), in each step calculation

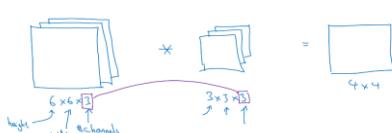
#### Note:

Filter design: could have different design/parameter for each channel.

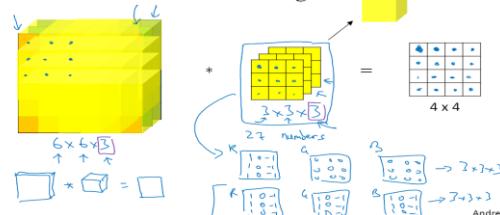
e.g: if want filter to capture edge only in the red channel of input image, could set edge detect parameter in the filter layer/channel that corresponds to red layer/channel of input image, and set to filter layers that corresponds to input green, blue channel.

or if want to capture edge while do not care color, then could set same filter parameter in all filter layers/channels.-detect edges in any color of input image.

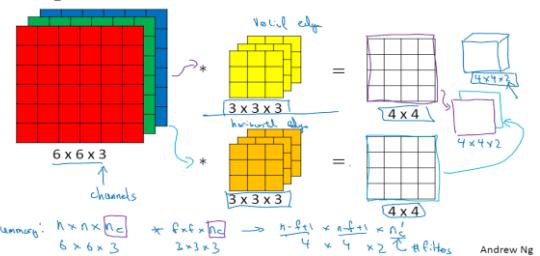
#### Convolutions on RGB images



#### Convolutions on RGB image



#### Multiple filters



#### 2. Multiple filters

one filter could capture one feature of input image, if want to capture other feature, use another filter-->use multiple filters at same time

#### >2.1 Process:

>>1. use multiple feature detectors/filters at same time;

>>2. Stack each filter conv. Result (2D) together

#### 2.2 Volume Result: 3D, 3rd dimension is filter/detector number.

same input with multi kernal-->stack output together to get volume output  
output channels number = feature numbers that are detecting.

#### 2.3 Dimension:

size: input:  $n \times n \times n[c]$ , kernel:  $f \times f \times n[c]$ , padding:  $p$ , stride step:  $s$   
-->output size:  $[(n+2p-f)/s+1] \times [(n+2p-f)/s+1] \times n[k]$  (kernel number)

Question: could use different size filter? while stacking together , use padding?

## 4.1-7 one layer of convolutional network

### 1. One layer of convolutional neural network

>1. for each filter cov. result: take it as  $w \times h$  result, add bias  $b$  (singal real number-broadcasting while adding to conv. result)-get  $Z[1]$ , and then apply activation function--->final output for each filter  $A[1]$   
for i-th filter in neural layer 1:

$Z[1]_i = \text{conv result} + b[1]_i$

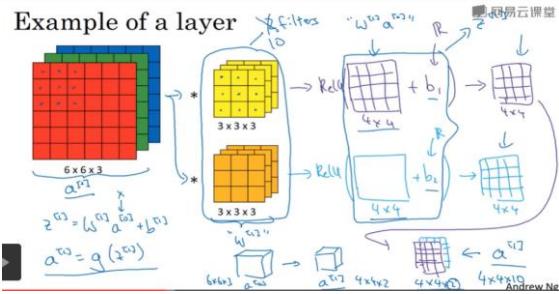
$A[1]_i = g(Z[1]_i)$  g function here: just stack  $Z[1]_i$  together.

(each filter as one hidden units in 1 layer: compute  $Z$ : conv result + b-bias;  
all filters: -->all hidden units in 1 layer;-->activation function stakcking all hidden units computation result together).

note: bias  $b[1]_i$  added for each filter could be different.

>2. Stack up all filter final output  $A[1]_i$ , get convolutional neural network output  $A[1]$ .

Note: cov. function of input image and filters, plays a role similar to  $W[1] \cdot A[0]$ -linear operation  
then add bias, apply activation function and output. -->this is one layer of a convolutional net.



### 2. Parameter number in one layer

One filter parameter number:  $(f \times f \times n_c) + 1$  (bias-single number)  
multi filter parameter: filter number x one filter parameter number

Note: parameter number not influenced by input image size. So once learned certain number feature detectors that work, you could apply this even to large images. And the number of parameters still is fixed and relatively small.

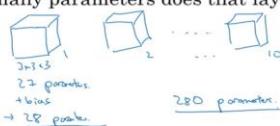
e.g: neural network with 10 filters of  $3 \times 3 \times 3$ , paramters in this neural network layer is  $[3 \times 3 \times 3 \text{ (filter parameter)} + 1 \text{ (bias)}] \times 10 = 280$ , parameter number will not be influenced by later , even much larger image size like  $1000 \times 1000 \times 3$ .

**Convolutional neural nets property: less prone to overfitting.**

-->could use filter with small parameters to detect large size image feature--->make algorithm less prone to over fitting.

### Number of parameters in one layer

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?



### 3. Summary of notation:

if layer l is a convolution layer: dimension /note for parameters  
input:  $n_h[l-1] \times n_w[l-1] \times n_c[l-1]$  ( $n_c[l-1]$ : number of channel of last layer output)

filter size:  $f[l]$ :  $f_l \times f_l \times n_c[l-1]$  (filter channel same as input channel)

padding:  $p[l]$

stride:  $s[l]$

filter number:  $n_c[l]$

weights:  $f[l] \times f[l] \times n_c[l-1] + 1 \times n_c[l]$

bias:  $b[l]_j$  (dimension =  $n_c[l]$ ): one bias-real number for each filter; in coding for convinence :  $1 \times 1 \times n_c[l]$ (four dimension)

output:  $n_h[l] \times n_w[l] \times n_c[l]$  (output channel same as filter number-feature number)

$$n_h[l] = [n_h[l-1] + 2p[l] - f[l]] / s[l] + 1$$

$n_c[l]$ : filter number in layer l.

Activation:  $A[l]$ : same dimension:  $n_h[l] \times n_w[l] \times n_c[l]$

m: ordering/index of training example

could search /index in this way, first look for sample ordanance, then check it's output.

## 4.1-8 a simple convolutional network example

### 1. Convnet Example

e.g: image classificaiton use Convnet  
>1. input: 39x39x3;

>2. conv in layer 1:  
filter in layer 1:  $3 \times 3 \times 3$ ; ten filters;  $s=1$ ,  $p=0$   
--> layer1 output:  $37 \times 37 \times 10$

>3. Conv in layer 2:  
filter in layer 2:  $5 \times 5 \times 10$ ,  $s=2$ ,  $p=0$ , 20 filters  
-->layer2 output:  $17 \times 17 \times 20$

>4conv in layer3:  
filter :  $5 \times 5 \times 20$ ,  $s=2$ ,  $p=0$ , 40 filters  
-->layer 3 output:  $7 \times 7 \times 40$   
compute  $7 \times 7 \times 40 = 1960$  features of image 39x39x3

>5.unroll/flatten output into one vector (1960,1)--->feed to activation function

activation function: logistic or softmax function (depending on trying to recognize car or not car, or recognize any one of k different objects)

>6. Get final output  $y^A$  from activation function.

### Note:

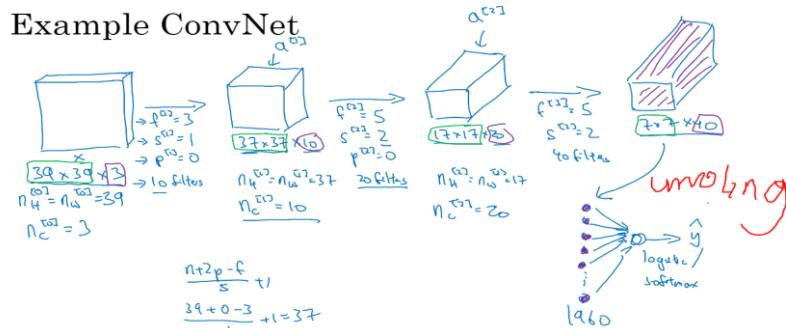
1. A lot of the work in designing convolutional neural net is selecting hyperparameters like: deciding filter size, stride step, padding size and how many filters to use in each layer.

2. General trend in a lot of convolutional neural networks:

Conv output size decreasing with layer become deeper , while channel increase bigger:

As go deeper in a neural network, the conv result image/visulization's height and width will stay the same for a while and gradually trend down as go deeper in the neural network.

Whereas the number of channels will generally increase.



## 2. types of layer in convolutional network:

In a typical ConvNet, there are usually three types of layers:

- >1. Conv layer
- >2. Pooling layer
- >3. Fully connected layer- FC.

And although it's possible to design a pretty good neural network using just convolutional layers, most neural network architectures will also have a few pooling layers and a few fully connected layers.  
which are a bit simpler than convolutional layers to define.

## Types of layer in a convolutional network:

- Convolution (conv) ← }
- Pooling (pool) ← }
- Fully connected (FC) ← }

## 4.1-9 pooling

Other than convolutional layers, ConvNets often also use pooling layers to: 1. reduce the size of the representation, to speed the computation; 2. as well as make some of the features that detects a bit more robust.

### 1. Pooling layer: Max pooling:

Max pooling: each of the outputs will just be the max from the corresponding reshaded region in the input image.

E.g. right : input 4x4, filter 2x2, s=2-->result: 2x2  
each result corresponds the max of reshaped region-covered by filter, in the input image.

#### 1.1 Intuition of max pooling:

If take each pixel of input image as some set of features or the activations in some layer of the neural network, then a large number in the pixel means that it's maybe detected a particular feature.  
e.g: the upper left-hand quadrant (象限) has this particular feature-max value 9. It maybe a vertical edge or maybe a eye, etc. . Clearly, that feature exists in the upper left-hand quadrant.  
Whereas this feature in upper right quadrant, maybe is cat eye detector-max value 2, it doesn't really exist in the upper right-hand quadrant.

(note: normally conv followed by pooling:  
conv result: each layer result is one filter conv result, capturing one feature of image-->followed by max pooling, to filter out the particular feature captured by this one filter in image.  
max pooling works on each layer separately, filter out particular feature captured by each filter in image.)

#### 1.1 Intuition of max pooling:

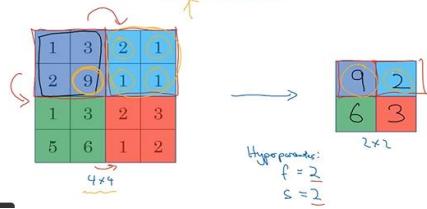
Max operation function : as long as features are detected anywhere in one of these quadrants , it then remains preserved in the output of max pooling.  
if these features detected anywhere in this filter, then keep a high number. But if this feature not detected, so maybe this feature doesn't exist in the upper right-hand quadrant. Then the max of all those numbers is still itself quite small.  
So maybe that's the intuition behind max pooling.

#### note:

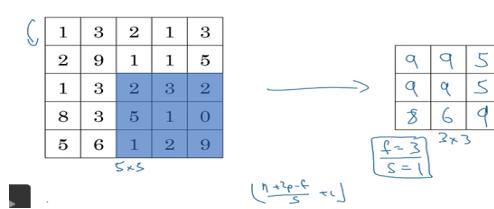
max value in one pixel of image could be many pixel's value in pooling result: as this input pixel will be involved in max pooling many times.-below eg. value 9 occur many times in different pixels of result.

Main reason people use max pooling is because it's been found in a lot of experiments to work well, and the intuition above is the fancy real underlying reason that max pooling works well in confidence.

### Pooling layer: Max pooling



### Pooling layer: Max pooling



### 1.2 . Max pooling property:

#### > Parameters:

Only Have hyperparameters, size f, stride step s. yet no parameters to learn-nothing for gradient descent to learn.  
once fix f, & stride step, it is a fixed computation and gradient descent change nothing.

max pooling result size: same computatio as conv. :  $(n+2p-f)/s + 1$ .

### 1.3 Max pooling on volume image:

> Max pooling out has same dimension on the 3rd direction, as .input image , when it is 3D dimension  
pooling output channe number: =image channel number

>Computation: comput max pooling on each input channel independently.

#### 1.4. Summarize:

max pooling filter have no channel: calculate image of each channel by same filter, and tack output together.

while in convolution neural network output has no channel for each filter (final output tack all filters' output together, final output channel number =filter number)

#### > Pooling :

image with channels C--> filter without parameter (once f size and stride step fix), works on each channel separately-->tack isolate channel output -->final output channel C  
(filter out feature captured in each layer-if this layer result is conv result, that means to filter out /find max feature captured by each filter)

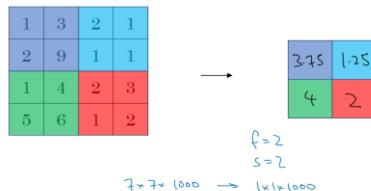
#### >convolution:

image with channels C-->filter with same channel-->output only 1 channel.

(each filter works on all input layers, to capture one input feature: by computing on filter covered input region through all channels).

note: while computing image style matrix: element wise same position of two channels and then sum-squared matrix: channel number x channel number)

### Pooling layer: Average pooling



### Summary of pooling

#### Hyperparameters:

- |                        |            |
|------------------------|------------|
| f : filter size        | $f=2, s=2$ |
| s : stride             | $f=3, s=2$ |
| Max or average pooling |            |

→ p: padding

No parameters to learn!

$$\frac{n_h \times n_w \times n_c}{\lfloor \frac{n_h-f+1}{s} \rfloor \times \lfloor \frac{n_w-f+1}{s} \rfloor \times n_c}$$

## 2. Pooling layer: Average pooling

instead of taking the maxes within each filter, take the average.

max pooling is used much more often than average pooling with one exception, which is sometimes very deep in a neural network: might use average pooling to collapse your representation from e.g.  $7 \times 7 \times 1,000$ . An average over all the the spacial extents to get  $1 \times 1 \times 1,000$ .

## 3. Summary of pooling

3.1 hyperparameters for pooling are:

1> f: the filter size

2> s: stride

Common choices of parameters might be  $f=2, s=2$ .

This is used quite often and this has the effect of roughly shrinking the height and width by a factor of above two.

3> the other hyperparameter is: using max pooling or average pooling

4> padding

Very rarely used when you do max pooling, usually, you do not use any padding, although there is one exception that we'll see next week as well.

### Size:

Input volume size:  $N_H \times N_W \times N_C$

Output volume size:  $((N_H-f/s) + 1) \times ((N_W-f/s) + 1) \times N_C$ .

### Note:

1. input channels number = output channels number -because pooling applies to each of your channels independently.

2. In pooling there are no parameters to learn.

So when we implement that crop, you find that there are no parameters that backdrop will adapt through max pooling.

3. In pooling there are just above hyperparameters, that set once, maybe set by hand or set using cross-validation

And then beyond that, you are done. Its just a fixed function that the neural network computes in one of the layers, and there is actually nothing to learn.

## 4.1-10 convolution neural network example

### 1. Neural network example\_(Inspired by LeNet-5)

e.g: recognize image number

#### 1.1. Layers

>1. input image:  $32 \times 32 \times 3$

>2. layer 1

> Conv\_1:  $f=5, s=1, 6$  filter-->output  $28 \times 28 \times 6$

>. Max pooling\_1:  $f=2, s=2, 6$  filter--->output  $14 \times 14 \times 6$

>3. Layer2

>Conv\_2:  $f=5, s=1$  filter-->output  $10 \times 10 \times 1$

>Max Pool\_2:  $f=2, s=2, 16$  filter--->output  $5 \times 5 \times 16$

>4. Flatten result into vector (400, 1)

>5. FC3: fully connected-->output vector (120,1)

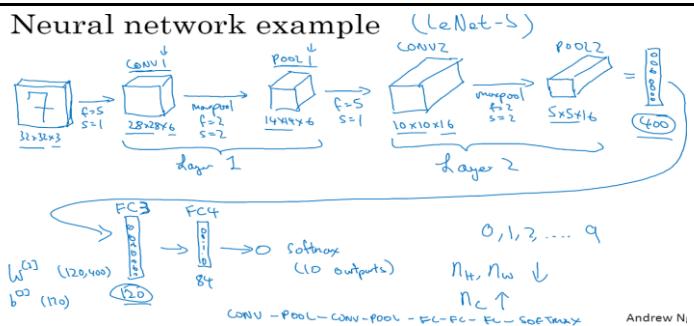
Fully connected layer just like the single standard neural network:

flattened vector (400,1) , 400 units densely connected to 120 units, by weight matrix -dimension (120, 400)

>6. FC4: -->output (84, 1)

feed to activation function.

>7activation function with softmax:



### 1.2 Note:

1. two convention of defining one layer:

> One convention: convn + pool together is one layer.

> 2nd convention: conv itself is one layer, pooling layer as one layer.

When people report the number of layers in a neural network usually people just record the number of layers that have weight/parameters. And because the pooling layer has no weights/ parameters, only a few hyper parameters, here use a convention that Conv 1 and Pool 1 shared together, treat that as Layer 1. (when I count layers, I'm just going to count layers that have weights)

2. Guideline for Hyper parameters:

seems like there a lot of hyper parameters, one common guideline is to actually not try to invent your own settings of hyper parameters, but to look in the literature to see what hyper parameters work for others. And to just choose an architecture that has worked well for someone else, and there's a chance that will work for your application as well.

### 1.3 Activation shape, size and parameters

List in the table , for each layer, the activation/input shape, size, and parameters in neural network

>Activation, a0 has dimension 3072. Well it's really  $32 \times 32 \times 3$ . And there are no parameters

>Max pooling layers don't have any parameters;

>conv layers tend to have relatively few parameters

And in fact, a lot of the parameters tend to be in the fully connected layers of the neural network.

> activation size tends to maybe go down gradually as you go deeper in the neural network.

If it drops too quickly, that's usually not great for performance as well.

e.g.: starts first there with 6,000 and 1,600, and then slowly falls into 84 until finally you have your Softmax output.

A lot of ConvNets will have properties will have patterns similar to above.

### Summary:

So basic building blocks of neural networks is , convolutional neural networks, the conv layer, the pooling layer, and the fully connected layer.

A lot of computer vision research has gone into figuring out how to put together these basic building blocks to build effective neural networks. And putting these things together actually requires quite a bit of insight. I think that one of the best ways for you to gain intuition is about how to put these things together is see a number of concrete examples of how others have done it. So what I want to do next week is show you a few concrete examples even beyond this first one that you just saw on how people have successfully put these things together to build very effective neural networks.

And through those videos next week I hope you hold your own intuitions about how these things are built. And as we are given concrete examples that architectures that maybe you can just use here exactly as developed by someone else or your own application.

## Neural network example

	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	~ 3,072 $a^{(0)}$	0
CONV1 ( $f=5, s=1$ )	(28,28,8)	6,272	208 ↙
POOL1	(14,14,8)	1,568	0 ↙
CONV2 ( $f=5, s=1$ )	(10,10,16)	1,600	416 ↙
POOL2	(5,5,16)	400	0 ↙
FC3	(120,1)	120	48,001 ↙
FC4	(84,1)	84	10,081 ↙
Softmax	(10,1)	10	841

## 4.1-11: why convolutions?

### 1. Conv Advantage:

e.g.: image 32x32x3 (3072)----> conv: f=5, 6 filters----->conv result 28x28x6 (4704)

if use standard neural network/fully connection-->weight matrix is (4704, 3072)-no problem today to train, but this is a quite small image and yet with a lot of parameters. if use larger image, then parameters will be infeasible large.

yet with conv layer: each filter have 5x5 +1-bias=26 paramters, 6 filter have 156 parameters. these small parameters due to:

#### >1. parameter sharing in conv:

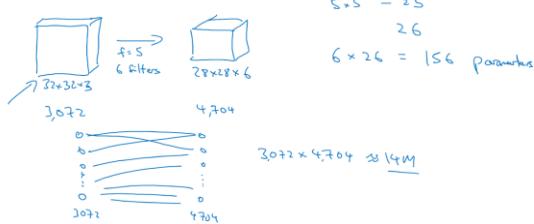
a feature detector (like vertical edge detector) that's useful in one part of image is probably useful in another part of the image.

e.g.: if you've figured out 3x3 filter for detecting vertical edges, you can then apply the same 3x3 filter over here, and then the next position over, and the next position over, and so on. And so, each of these feature detectors, each of these outputs can use the same parameters in lots of different positions in your input image , in order to detect say a vertical edge or some other feature.

this should be true in low level feature, as well as high level feature -e.g maybe detecting the eye that indicates a face or a cat or something is there)

But being with a share in this case the same nine parameters of a filter to compute all 16 of these outputs, is one of the ways the number of parameters is reduced.

### Why convolutions



The second way that convolutions get away with having relatively few parameters is by having sparse connections

#### >2. sparsity connection:

In each layer, each output value depends only on a small number of input.  
(standard neural network: each input unit connected to layer output)

e.g: right pic, conv result first value 0 = filter \* right corner region-9 values of input image, decided only by filter and part region of input -not all input image.

each output value in conv result, depend only on a small number of input( filter size)=element-wise and sum of filter and it's specified region in image. while standard NN, each output value/unit  $z[i][j]$  value depending on all input=  $w[i][j]X$

so, through these two mechanisms-sharing parameter and sparsity connection, a neural network has a lot fewer parameters which allows it to be trained with smaller training cells and is less prone to be over-fitting (over fitting?)

#### 3. good at capturing translation invariance.

e.g: picture of cat shift couple of pixels to right, still pretty clear a cat.  
convolutional structure helps neural network encode the fact that a picture shifted a few pixels should result in pretty similar feature, and should probably be assigned the same output label.

-->the fact that applying the same filter in all the position of image, both in earlier layers and later layers, hopes network automatically learn to be more robust, to better capture the desirable property of translation invariance.

### 2. Putting it all together and see how to train one of these networks

e.g: build a cat detector  
have a labeled training sets: X is an image. y can be binary labels, or one of K causes.

And choose a convolutional neural network structure:

inserted the image and then having neural convolutional and pooling layers and then some fully connected layers followed by a software output that then operates Y hat.  
The conv layers and the fully connected layers will have various parameters, W, as well as bias's B.

Define a cost function similar to previous courses:

>>> Randomly initialized parameters W and B,  
>>> Compute the cost J, as the sum of losses of the neural networks predictions on the entire training set.  
>>> Train this neural network: use gradient descent or some of the algorithm like, gradient descent momentum, or RMSProp or Adam, or something else, in order to optimize all the parameters of the neural network to try to reduce the cost function J.

--> If do this, can build a very effective cat detector or some other detector

#### Summary:

This week have seen all the basic building blocks of a convolutional neural network, and how to put them together into an effective image recognition system.

There're just a lot of the hyperparameters in convolution neural networks, next week will show a few concrete examples of some of the most effective convolutional neural networks, so you can start to recognize the patterns of what types of network architectures are effective.

And one thing that people often do is just take the architecture that someone else has found and published in a research paper and just use that for your application. And so, by seeing some more concrete examples next week, you also learn how to do that better.

Note: 空洞卷积 (dilated convolution)

#### >1. parameter sharing in conv:

**Intuition:** The same feature seems like it will probably be useful in one position of input image, has a good chance of being useful for other positions of the image. -->So, maybe don't need to learn separate feature detectors for different position of the image.

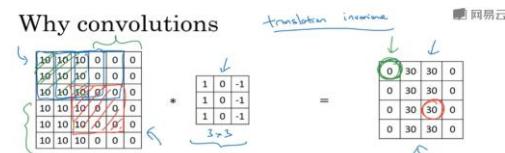
And maybe you do have a data set where you have the upper left-hand corner and lower right-hand corner have different distributions, so, they maybe look a little bit different but they might be similar enough, they're sharing feature detectors all across the image, works just fine.

#### Summarize:

1. Do not need different feature detector /filter /kernel in different position of the image . different position of image data distribution maybe different, but they might similar enough sharing feature detector all across image and works fine.

2. standard neural network: each input unit connected to layer output (by unique weight for each input unit), weight is huge  $n[l] \times n[l-1]$ ; while conv parameter is only (filter h size x w size x filter channel+bias) x filter number, quite small quant (shared weight parameters for input units).

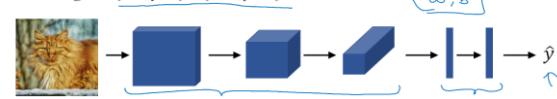
3. conv could use small parameter, as: same parameters can use in different image positions



**Sparsity of connections:** In each layer, each output value depends only on a small number of inputs.

### Putting it together

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



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

Use gradient descent to optimize parameters to reduce  $J$

<https://www.cs.toronto.edu/~guerouj/courses/24310/2017F/07%20-%20Convolutional%20Neural%20Networks%20-%20Part%201%20-%20Introduction%20and%20Forward%20Pass.html>

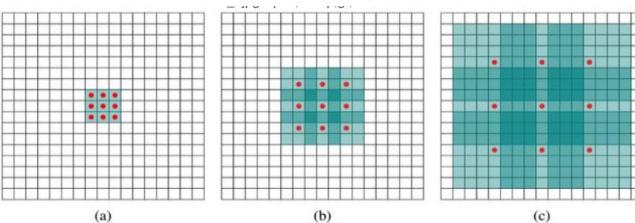


Image size: h x w;  
filter size: f x f;  
padding: p;  
stride: s  
dilate conv size: d (空洞数)

-->result size: { h + 2p - (f+d\*(f+1)) } / s + 1  
(standard conv filter + dilate conv-->new filter size: f\*(1+d) + d

(a)图对应3x3的1-dilated conv，和普通的卷积操作一样，(b)图对应3x3的2-dilated conv，实际的卷积kernel size还是3x3，但是空洞是1，也就是对于一个7x7的图像patch，只有9个红色的点和3x3的kernel发生卷积操作，其余的点略过。也可以理解为kernel的size为7x7，但是只有图中的9个点的权重不为0，其余都为0。可以看到虽然kernel size只有3x3，但是这个卷积的感受野已经增大到了7x7（如果考虑到这个2-dilated conv的前一层是一个1-dilated conv的话，那么每个红点就是1-dilated的卷积输出，所以感受野为3x3，所以1-dilated和2-dilated合起来就能达到7x7的感受野）。(c)图是4-dilated conv操作，同理跟在两个1-dilated和2-dilated conv的后面，能达到15x15的感受野，对比传统的conv操作，3层3x3的卷积加起来，stride为1的话，只能达到(kernel-1)\*layer+1=7的感受野，也就是和层数layer成线性关系，而dilated conv的感受野是指数级的增长。

dilated的好处是不做pooling损失信息的情况下，加大了感受野，让每个卷积输出都包含较大范围的信息。在图像需要全局信息或者语音文本需要较长的sequence信息依赖的时候，都能很好的应用dilated conv，比如图像分割[3]、语音合成WaveNet[2]、机器翻译ByteNet[1]中，简单贴下ByteNet和WaveNet用到的dilated conv结构，可以更形象的了解dilated conv本身。

## 转置卷积

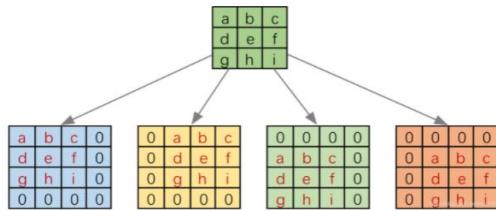
### 普通卷积（直接卷积）

普通的卷积过程大家都很熟悉了，可以直观的理解为一个带颜色小窗户（卷积核）在原始的输入图像一步一步的滑动，来通过加权计算得出输出特征。如下图。

$$\begin{array}{c} \begin{array}{|c|c|c|c|} \hline 3 & 3 & 2 & 1 \\ \hline 0 & 0 & 1 & 3 \\ \hline 3 & 1 & 2 & 2 \\ \hline 2 & 0 & 0 & 2 \\ \hline \end{array} \end{array} * \begin{array}{c} \begin{array}{|c|c|c|} \hline 0 & 1 & 2 \\ \hline 2 & 2 & 0 \\ \hline 0 & 1 & 2 \\ \hline \end{array} \end{array} = \begin{array}{c} \begin{array}{|c|c|} \hline 12 & 12 \\ \hline 10 & 17 \\ \hline \end{array} \end{array}$$

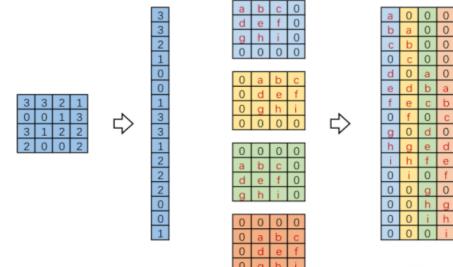
输入              卷积核              输出

但是实际在计算机中计算的时候，并不是像这样一个位置一个位置的进行滑动计算，因为这样的效率太低了。计算机会将卷积核转换成等效的矩阵，将输入转换为向量。通过输入向量和卷积核矩阵的相乘获得输出向量。输出的向量经过整形便得到我们的二维输出特征。具体的操作如下图所示。由于我们的3x3卷积核要在输入上不同的位置卷积4次，所以通过补零的方法将卷积核分别置于一个4x4矩阵的四个角落。这样我们的输入可以直接和这四个4x4的矩阵进行卷积，而舍去了滑动这一操作步骤。



[https://blog.csdn.net/tsyccnh/article/details/87357447?ops\\_request\\_misc=%252B%2522request%2525Fid%2522%253A%2522162865030316780357272834%2522%252C%2522cm%2522%253A%2522%2522040713.130102334.%2522%2527D&request\\_id=162865030316780357272834&biz\\_id=0&utm\\_medium=distribute.pc\\_search\\_result.none-task-blog-2-all-top\\_positive-default-1-87357447.pc\\_search\\_download\\_positive&utm\\_term=%E8%BD%AC%E7%BD%AE%E5%8D%B7%E7%A7%AF&spm=1018.2226.3001.4449](https://blog.csdn.net/tsyccnh/article/details/87357447?ops_request_misc=%252B%2522request%2525Fid%2522%253A%2522162865030316780357272834%2522%252C%2522cm%2522%253A%2522%2522040713.130102334.%2522%2527D&request_id=162865030316780357272834&biz_id=0&utm_medium=distribute.pc_search_result.none-task-blog-2-all-top_positive-default-1-87357447.pc_search_download_positive&utm_term=%E8%BD%AC%E7%BD%AE%E5%8D%B7%E7%A7%AF&spm=1018.2226.3001.4449)

进一步的，我们将输入拉成长向量，四个4x4卷积核也拉成长向量并进行拼接，如下图。



我们记向量化的图像为  $I$ ，向量化的卷积矩阵为  $C$ ，输出特征向量为  $O$  则有：

$$I^T * C = O^T$$

如下图所示。

$$\begin{array}{c} \begin{array}{|c|c|c|c|} \hline a & b & c & 0 \\ \hline b & c & 0 & 0 \\ \hline c & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array} \end{array} \times \begin{array}{c} \begin{array}{|c|c|c|c|} \hline a & b & c & 0 \\ \hline d & e & f & 0 \\ \hline g & h & i & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array} \end{array} = \begin{array}{c} \begin{array}{|c|c|c|c|} \hline 12 & 12 & 10 & 17 \\ \hline \end{array} \end{array}$$

我们将一个1x16的行向量乘以16x4的矩阵，得到了1x4的行向量。那么反过来将一个1x4的向量乘以一个4x16的矩阵是不是就能得到一个1x16的行向量呢？没错，这便是转置卷积的思想。

$$\begin{array}{c} \begin{array}{|c|c|c|c|} \hline a & 0 & 0 & 0 \\ \hline b & 0 & 0 & 0 \\ \hline c & 0 & 0 & 0 \\ \hline d & 0 & 0 & 0 \\ \hline e & 0 & 0 & 0 \\ \hline f & 0 & 0 & 0 \\ \hline g & 0 & 0 & 0 \\ \hline h & 0 & 0 & 0 \\ \hline i & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array} \end{array} \times \begin{array}{c} \begin{array}{|c|c|c|c|} \hline 1 & 2 & 1 & 1 \\ \hline 0 & 1 & 1 & 1 \\ \hline 1 & 1 & 2 & 2 \\ \hline 2 & 2 & 2 & 2 \\ \hline 1 & 1 & 2 & 2 \\ \hline 0 & 0 & 0 & 2 \\ \hline \end{array} \end{array} = \begin{array}{c} \begin{array}{|c|c|c|c|} \hline 12 & 12 & 10 & 17 \\ \hline \end{array} \end{array}$$

## 转置卷积

一般的卷积操作（我们这里只考虑最简单的无padding, stride=1的情况），都将输入的数据越卷越小。根据卷积核大小的不同，和步长的不同，输出的尺寸变化也很大。但是有的时候我们需要输入一个小的特征，输出更大尺寸的特征该怎么办呢？比如图像语义分割中往往要求最终输出的特征尺寸和原始输入尺寸相同。但在网络卷积池化化的过程中特征图的尺寸却逐渐变小。在这里转置卷积便派上了用场。在数学上，转置卷积的操作也非常简单，把正常的卷积的操作反过来即可。

对应上面公式，我们有转置卷积的公式：

$$O^T * C^T = I^T$$

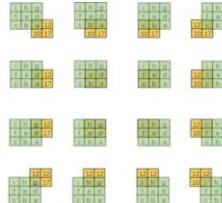
如下图所示：

$$\begin{array}{c} \begin{array}{|c|c|c|c|} \hline a & b & c & 0 \\ \hline d & e & f & 0 \\ \hline g & h & i & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array} \end{array} \times \begin{array}{c} \begin{array}{|c|c|c|c|} \hline 1 & 2 & 1 & 1 \\ \hline 0 & 1 & 1 & 1 \\ \hline 1 & 1 & 2 & 2 \\ \hline 2 & 2 & 2 & 2 \\ \hline 1 & 1 & 2 & 2 \\ \hline 0 & 0 & 0 & 2 \\ \hline \end{array} \end{array} = \begin{array}{c} \begin{array}{|c|c|c|c|} \hline 12 & 12 & 10 & 17 \\ \hline \end{array} \end{array}$$

这里需要注意的是这两个操作并不是可逆的，对于同一个卷积核，经过转置卷积操作之后并不能恢复到原始的数值，保留的只有原始的形状。

所以转置卷积的名字就由此而来，而并不是“反卷积”或者是“逆卷积”，不好的名称容易给人以误解。

这是一个很有趣的场景。结合整体来看，仿佛有一个更大的卷积核在2x2大小的输入滑动。但是输入大小，每一次卷积只对对应的核的一部分。我们来把更大的卷积核补全，如下图：



这里和直接卷积有很大的区别，直接卷积我们是用一个“小世界”去看一个“大世界”，而转置卷积是用一个“大世界”的一部分去看“小世界”。这里有一点需要注意，我们定义的卷积核是左上角为a，右下角为i，但在可视化转置卷积中，需要将卷积核旋转180°后再进行卷积。由于输入图像大小，我们按照卷积核大小来进行补零操作，两边的补零数量固定，即3-1，这样我们就能将一个转置卷积操作转换为对应的直接卷积，如下图：

$$\begin{array}{c} \begin{array}{|c|c|c|c|} \hline a & 0 & 0 & 0 \\ \hline b & 0 & 0 & 0 \\ \hline c & 0 & 0 & 0 \\ \hline d & 0 & 0 & 0 \\ \hline e & 0 & 0 & 0 \\ \hline f & 0 & 0 & 0 \\ \hline g & 0 & 0 & 0 \\ \hline h & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array} \end{array} \times \begin{array}{c} \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 0 & 1 \\ \hline \end{array} \end{array} = \begin{array}{c} \begin{array}{|c|c|} \hline 12 & 12 \\ \hline \end{array} \end{array}$$

总结一下将转置卷积转换为直接卷积的步骤：(这里只考虑stride=1, padding=0的情况)

1. 对输入进行四边补零，单边补零的数量为k-1

2. 将卷积核旋转180°，在新的输入上进行直接卷积

### 形象化的转置卷积

但是仅仅按照矩阵形式来理解转置卷积似乎有些抽象，不像直接卷积那样理解的直观。所以我们也来尝试一下可视化转置卷积。前面说了在将直接卷积向量化的时候是将卷积核补零然后拉成列向量，现在我们有了一个新的转置卷积矩阵，可以将这个过程反过来，把16个列向量再转换成卷积核。以第一列向量为例，如下图：

$$\begin{array}{c} \begin{array}{|c|c|c|c|} \hline 12 & 12 & 10 & 17 \\ \hline \end{array} \end{array} \times \begin{array}{c} \begin{array}{|c|c|} \hline a & 0 \\ \hline 0 & 0 \\ \hline \end{array} \end{array} \Rightarrow \begin{array}{c} \begin{array}{|c|c|} \hline a & 12 \\ \hline \end{array} \end{array}$$

**10|17**      **0|0**      ✓      **10|17**

这里将输入还原为一个 $2 \times 2$ 的张量，新的卷积核由于只有左上角有非零值直接简化为右侧的形式。对每一个列向量都做这样的变换可以得到：

**12|17**      **12|17**      **12|17**      **12|17**

**12|17**      **12|17**      **12|17**      **12|17**

**12|17**      **12|17**      **12|17**      **12|17**

**12|17**      **12|17**      **12|17**      **12|17**

## 4-2\_Examples Convolution neural network

### 4.2-1 Why look at case studies

#### 1. Why look at case studies

Last week we learned about the basic building blocks, such as convolutional layers, pooling layers, and fully connected layers of confidence. It turns out the past few years of computer vision research has been working on how to put together these basic building blocks to form effective convolutional neural networks. One of the best ways for user gain intuition is to see some of these examples.

It turns out that a neural network architecture that works well on one computer vision tasks often works well on other tasks as well, e.g: If someone else's train a neural network or has figured out a neural network architecture, there's very good at recognizing cats and dogs and people. But you have a different computer vision tasks like maybe you're trying to build a self-driving car, you might well be able to take someone else's neural network architecture and apply that to your problem.

After the next few videos, you'll be able to read some of the research papers from the field of computer vision, and I hope that you might find it satisfying as well.

### 2. Outline

#### 2.1 Classic networks:

>LeNet-5 network : Came from 1980s

> AlexNet

which is often cited in the VGG network.

>VGG

These are examples of pretty effective neural networks, and you see ideas from these papers that will probably be useful for your own work as well.

#### 2.2 ResNet or called residual network

ResNet neural network trained a very deep 152 layer neural network. It has some very interesting tricks, interesting ideas how they do that effectively.

#### 2.3 Inception neural network

After seeing these neural networks, I think you have much better intuition about how to build effective convolutional neural networks. Even if you end up not working computer vision yourself, you find a lot of the ideas from some of these examples, such as ResNet, Inception network, many of these ideas are cross fertilizing, are making their way into other disciplines. you'll find some of these ideas very interesting and helpful for your work.

### 4.2-2classic neural network:

#### 1. LeNet-5

##### 1.2 Summarize LeNet-5:

1. Lenet5-trained on grey image, image has no channel.

2. at earlier time:

- >1. avg pooling use more often , while now max pooling more often used
- >2. Padding not used-->result shrink after conv

3. not use softmax as activation, but another one which not in use today.

4. Lenet 5 has 60,000 parameters.

5. Reuse's height and width goes down while channel number goes up.

6. Pattern still quite common that: conv/several convs +pool+conv+pool+...+fc+fc

Note: only for paper reading:

1. Activation function used is sigmoid/tanh. Not Relu/softmax which use more often today.

2. Due to computation problem, did not use same filter for all input channels( conv filter channel /= input channel); different filters look at different channels of the input block.

3. Non-linearity (sigmoid) was used after pooling

this isn't really done right now

## Outline

### Classic networks:

- LeNet-5 ↩
- AlexNet ↩
- VGG ↩

ResNet (152)

### Inception

#### 1. LeNet-5

##### 1.1 Architecture\_LeNet 5

e.g: Recognize hand-written digits. Image: 32 x 32 x1

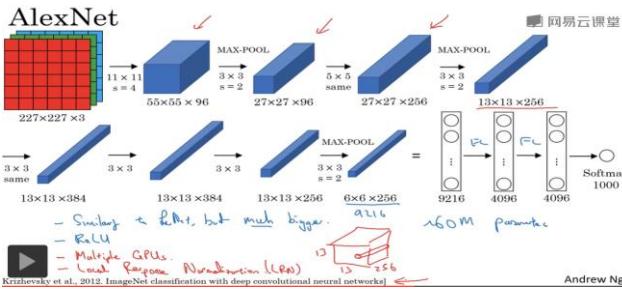
>Layer 1: conv + pooling  
>> conv1: filter 5x5, s=1, p=0, 6 filters---->result 28x28x6  
>> Avg pool: f=2, s=2 --->result 14x14x6

>Layer 2: conv +pooling  
>>conv2: filter 5x5,s=1, 16 filters---->result 10x10x16  
>>avg pool: f=2,s=2--->result 5 x 5 x 16.

>Flatten result to vector (400,1)

>Layer 3: FC: weight (400, 120)-->result vector (120,1)

>Softmax-->prediction  $y^*$ , 10 outputs probability



### 2. AlexNet

#### 2.2 Compare vs LeNet:

Much better performance than LeNet , due to:

>1. similar to Lenet, but much bigger

AlexNet has about 60,000 parameters., while AlexNet has 60 Million parameters;

And the fact that they could take pretty similar basic building blocks that have a lot more hidden units and training on a lot more data, they trained on the image that dataset that allowed it to have a just remarkable performance.

#### >2. Relu activation function used

>>Had a complicated way of training on two GPUs

when this paper was written, GPUs was still a little bit slower, so it had a complicated way of training on two GPUs. And the basic idea: was that, a lot of these layers (e.g layer 1, layer2, layer 3 )was actually split across two different GPUs and there was a thoughtful way for when the two GPUs would communicate with each other.

>>Had Local Response Normalization layer. And this type of layer isn't really used much, which is why I didn't talk about it.

LRN basic idea: e.g take one blocks/layer result.e.g 13 x 13 x256, look at one position on the 2D 13 x 13, So one position height and width, and look down this across all the channels, look at all 256 numbers and normalize them.

LRN motivation: for each position in this 13 x 13 image, maybe you don't want too many neurons with a very high activation.

note: many researchers have found LRN doesn't help that much so this is one of those ideas I guess I'm drawing in red because it's less important for you to understand this one. And in practice, I don't really use local response normalizations really in the networks language trained today.

Before AlexNet, deep learning was starting to gain traction in speech recognition and a few other areas, but it was really just paper that convinced a lot of the computer vision community to take a serious look at deep learning to convince them that deep learning really works in computer vision. And then it grew on to have a huge impact not just in computer vision but beyond computer vision as well

AlexNet had a relatively complicated architecture, there's just a lot of hyperparameters.

--take pretty similar basic building block as LeNet and just have a lot more hidden units and trained on a lot more data, which trained on the image and data-->remarkable performance

#### 2. AlexNet

##### 2.1 Architecture

input image:color 227 x227 x3

>Layer 1: conv + pooling  
>> conv1: filter 11x11 s=4, p=0, 96 filters---->result 55 x 55 x 96  
s=4, shrink image size 4 times  
>> Max pool: f=3, s=2 --->result 14x14x6

>Layer 2: conv +pooling  
>>conv2: filter 5x5 same padding, 256 filters---->result 27 x27 x256  
>>Max pool: f=3,s=2--->result 13 x 13 x 256.

>Layer3: Conv3: same padding 384 filters-->13x13x384

>Layer4: conv only : same padding, 384 filters-->13 x13 x 384

>Layer5: Conv only: same padding: 256 filters, -->13 x13 x256

>Layer 6: max pool: f 3x3, s=2--->6 x 6 x 256

>Flatten result to vector (9216,1)

>Layer 7: FC: weight (4096, 9216)-->result vector (120,1)

>Layer 8: FC: weight (4096, 4096)-->result vector (4096,1)

>Softmax-->prediction  $y^*$ , 1000 outputs probability

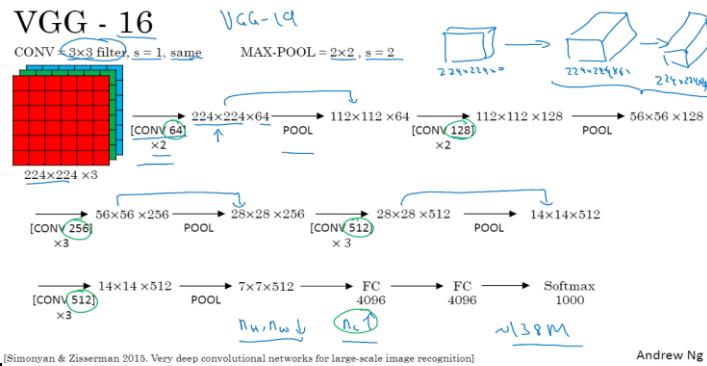
### 3. VGGNet-16

#### 3.0 Intuition

A remarkable thing about the VGG-16 net is instead of having so many hyperparameters, use a much simpler network:  
 Focus on just having:  
 Conv-layers = 3x3 filters  $s=1$ , same padding  
 max pooling layers :  $2 \times 2$ ,  $s=2$

One very nice thing about the VGG network was it really simplified this neural network architectures.

VGGNet-16: has 16 layers that have weights: 13 conv layers, 2 FC layer, 1 softmax layer.



[Simonyan & Zisserman 2015. Very deep convolutional networks for large-scale image recognition]

Andrew Ng

### 3. VGGNet-16

#### 3.2 Summary: VGG-16

1. Large network: 138 million parameters; (Le-Net 60,000 , Alex Net-60 million)

2. Architecture's simplicity made it quite appealing: quite uniforme

>>1. Few conv-layers followed by a pooling layer

which reduces the height and width,

>>2. Conv layer filter number is double: from 64 to 128, to 256, 512

Maybe authors thought 512 was big enough on the game here.

This sort of roughly doubling on every step, or doubling through every stack of conv-layers was another simple principle used to design the architecture of this network.

(channel double by conv-size no change, size half by pooling-channel number no change)

The relative uniformity of this architecture made it quite attractive to researchers.

Disadvantage of VGG-16:

it is a pretty large network in terms of the number of parameters had to train.

Recommend reading paper: starting with the AlexNet paper followed by the VGG net paper and then the LeNet paper is a bit harder to read but it is a good classic once you go over that.

### 4.2-3 Residual network.

Very, very deep neural networks are difficult to train because of vanishing and exploding gradient types of problems.

Skip connections which allows you to take the activation from one layer and suddenly feed it to another layer even much deeper in the neural network. And using that, you'll build ResNet which enables you to train very, very deep networks. Sometimes even networks of over 100 layers.

#### 1. Residual block

e.g: Two layers :

>1. Main path is as below: from  $a[l]$  to  $a[l+2]$ , need to go through all the steps.

$$a[l] \xrightarrow{\text{Linear function: } Z[l+1] = W[l+1]a[l] + b[l+1]} \xrightarrow{\text{activation function ReLU: } a[l+1] = g(Z[l+1])} \\ \xrightarrow{\text{Linear function: } Z[l+2] = W[l+2]a[l+1] + b[l+2]} \xrightarrow{\text{activation function ReLU: } a[l+2] = g(Z[l+2])}$$

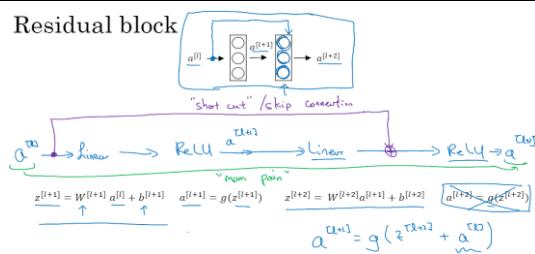
>2. Residu net: take  $a[l]$  and just forward it, copy it, much further into the neural network

$$\text{Add } a[l], \text{ before applying activation function (e.g. the ReLU non-linearity): } a[l+2] = g(Z[l+2] + a[l])$$

So rather than needing to follow the main path, the information from  $a[l]$  can now follow a shortcut to go much deeper into the neural network.

Note:  $a[l]$  is being injected in the a further layer: after linear part, but before activation part.

Also called: short cut/ skip connection: refers to  $a[l]$  just skipping over a layer or kind of skipping over almost two layers in order to process information deeper into the neural network.



#### 2. Residual Network:

using residual blocks allows you to train much deeper neural networks. And the way you build a ResNet is by taking many of these residual blocks, blocks like these, and stacking them together to form a deep network.

**Residual block:** one layer output  $A[l]$  skip few layers to add directly before activation function.

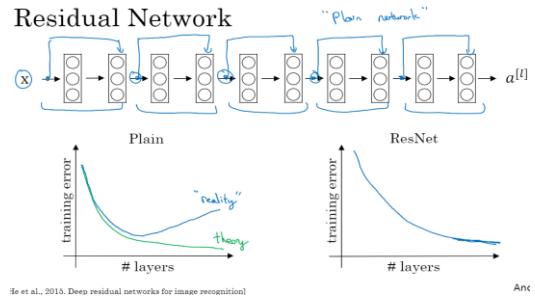
add in the a few further layer: between linear part and activation part: Influence layer output =  $g(Z[L\text{-influenced}] + A[l])$

$Z'$ : now actually have two patch: main patch (normal one) + short patch (skipped patch).

**Residual network motivation:** train deep network.

**Residual network:** tacks many residual block together to form deep network.

e.g.: To turn this into a ResNet, add all those skip connections although those short like a connections like so: every two layers ends up with that additional change that we saw on the above example to turn each of these into residual block. So above this picture shows five residual blocks stacked together, and this is a residual network.



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

Andrew Ng

Anc

### 3. Training error Residual Network

If use standard optimization algorithm such as a gradient descent or one of the fancier optimization algorithms to the train algorithm:

>1. Training error for network without residual block:

for a plain (no residual connection) network, without all the extra residual/ short cuts or skip connections:  
 >> Empirically, as increase the number of layers, the training error will tend to decrease after a while but then they'll tend to go back up.

>>>While in theory as make a neural network deeper, it should only do better and better on the training set.

In theory, having a deeper network should only help, but in practice having a plain network, so no ResNet, very deep means that all your optimization algorithm just has a much harder time training.

so, in reality, your training error gets worse if you pick a network that's too deep.

>2. Training error for network with ResNet:

Even as the number of layers gets deeper, training error kind of keep on going down. Even if we train a network with over a hundred layers/over a thousand layers.

#### Summary:

By taking these activations be at X or these intermediate activations, and allowing it to go much deeper in the neural network, this really helps with the vanishing and exploding gradient problems and allows you to train much deeper neural networks without really appreciable loss in performance (maybe at some point, this will plateau / flatten out, and it doesn't help that much deeper and deeper networks).

But ResNet has certainly been effective at helping train very deep networks.

### 4.2-4 why Residual network works

Intuition: how you can make algorithm deeper and deeper without really hurting your ability to at least get them to do well on the training set.

Doing well on the training set is usually a prerequisite to doing well on your hold up or on dev / test sets. So, being able to at least train ResNet to do well on the training set is a good first step toward that.

#### 1. Why do residual networks work

If make a network deeper, it can hurt the ability to train the network to do well on the training set. And that's why sometimes don't want a network that is too deep. But this is not true or at least is much less true when you training a ResNet.

e.g:

Algorithm 1-without residual block: X---->Bin NN----->a[i]

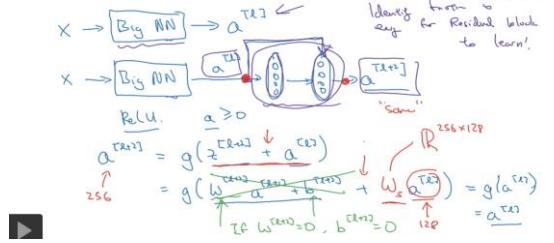
Algorithm 2-with residual block: X---->Bin NN----->a[i] (also skip to a[i+2]-----a[i+1]-----a[i+2]

Algorithm 2 is modified based on algorithm 1 to a little bit deeper with a ResNet block, a residual block with that extra short cut.

For the sake our argument, using Relu activation functions in the neural network. So, all the activations  $a[i] \geq 0$ , from algorithm 2,  $a[i+2] = g(w[i+2] \cdot a[i+1] + b[i+2] + a[i])$ .

if using regularization,  $w[i+2], b[i+2]$  close to 0, then  $a[i+2] \approx a[i]$ . so just get back,  $a[i]$ .

#### Why do residual networks work?



#### Note: skip dimension issue:

>1. In the e.g., assuming that  $a[i+2]$  and  $a[i]$  have the same dimension.

While in ResNet is a lot of use of same convolutions so the dimension of  $a[i]$  is equal to the  $a[i+2]$  dimension, then we can actually do this short circle connection, because the same convolution preserve dimensions, and so makes that easier for you to carry out this short circle and then carry out this addition of two equal dimension vectors.

>2. In case the input  $a[i]$  and output  $a[i+2]$  has different dimensions , add an extra matrix  $W_s$ ,  $W_s \cdot a[i]$ , have same dimension as  $a[i+2]$ .

$W_s$  could be a matrix of parameters we learned, or a fixed matrix that just implements zero paddings to make  $a[i]$  have same dimension as  $a[i+2]$ -if  $a[i+2]$  dimension is bigger

#### 1.1 Reason: Residual network works , due to:

>1. identity function  $a[i+2] \approx a[i]$ , is easy for residual block to learn.

It's easy to get  $a[i+2] \approx a[i]$  because of this skip connection. Means that adding these two layers in neural network, doesn't really hurt neural network's ability to do as well as this simpler network without these two extra layers, because it's quite easy for it to learn the identity function to just copy  $a[i]$  to  $a[i+2]$  using despite the addition of these two layers.

And this is why adding two extra layers, adding this residual block to somewhere in the middle or the end of this big neural network it doesn't hurt performance.

>2. If all of these hidden units actually learned something useful then maybe you can do even better than learning the identity function.

>3. For very deep plain nets in very deep network without this residual of the skip connections , when network is deeper and deeper, it's actually very difficult for it to choose parameters that learn even the identity function, which is why a lot of layers end up making your result worse rather than making your result better.

Main reason the residual network works is that it's so easy for these extra layers to learn the identity function, kind of guaranteed that it doesn't hurt performance and then a lot the time you maybe get lucky and then even helps performance. At least is easier to go from a decent baseline of not hurting performance and then great in decent can only improve the solution from there.

#### 2. ResNet example on images

e.g on images

a plain network: an image ----> a number of conv layers ---> softmax output at the end.  
 Turn this into a ResNet: add those extra skip connections in plain network

There are a lot of 3x3, same convolutions and that's why you're adding equal dimension feature vectors. So rather than a fully connected layer, these are actually convolutional layers but because the same convolutions, the dimensions are preserved and so the  $z[i+2]$  plus  $a[i]$  by addition makes sense.

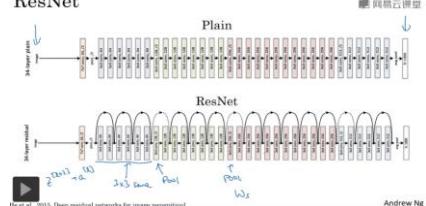
if skip to the pooling layers,

1. make an adjustment to the dimension with matrix  $W_s$ :  $a[i+2]$  dimension smaller than  $a[i]$ , matrix  $W_s$ , parameters learned to learn?

2. or could 0 padding  $z[i+2]$  to  $a[i]$  dimension?

if  $a[i+2]$  next layer is conv, and use same padding, then  $a[i+2]$  need 0 padding for same padding-->possible to 0 padding  $z[i+2]$  first, adding  $a[i]$ ?

#### ResNet



### 4.2-5 network in network and 1x1 convolution

#### 1. for input 2D data:

Input just multiply filter value, not much useful.

e.g: right : filter value=2, -->each pixel of input just multiply 2.

#### 2. 1x1 convolution for input volume data:

Filter size:  $1 \times 1 \times$  channel number (= input channel)

Computation: take one position in input, and look at it through all channels: do the element-wise product and then sum, -->then use Relu.

(why use Relu-non-linear activation for 1x1 conv?

standard conv fxf size: activation function is just stacking conv result of different filters together.)

#### 3. 1x1 conv. intuition:

take right e.g: input image  $6 \times 6 \times 32$ , filter  $1 \times 1 \times 32$

Intuition:

basically having a fully connected neural network that applies to each of the 32 different positions (channels). Fully connected neural network inputs 32 numbers and outputs number of filters, outputs. By doing those at each of the 36 positions, each of the six by six positions, end up with an output that is  $6 \times 6$  number of filters. **This can carry out a pretty non-trivial computation on your input volume.**

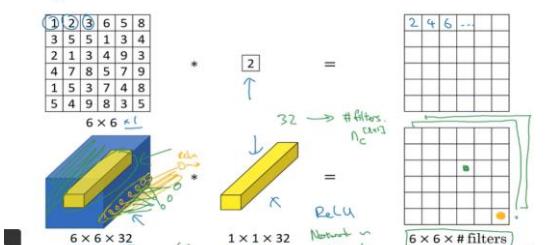
(same padding, with standard conv  $f \times f \times n_c$ : trivial computation on input volume-each output compute on input region covered by filter size, while  $1 \times 1$  conv only compute one pixel and not repeated.)

Take each position of input (through all channels):  $1 \times 1 \times 32$ , as an input, and one filter ( $1 \times 1 \times 32$ ) as the weight of one unit/neuron in the hidden layer, different filters as different units in the hidden layer, with different weight parameter ( $1 \times 1 \times 32$ ), to capture input different feature.

As if have one neuron that is taking input 32 numbers. Multiplying each of these 32 numbers in one slice in the same position, height, and width, but these 32 different channels, multiplying them by 32 weights. **In applying a ReLU, nonlinearity to it and then outputting the corresponding value in the result image.**

If have multiple filters, then it's as if you have multiple units to taking as input all the numbers in one slice and then building them up into an output,  $6 \times 6 \times$  number of filters.

#### Why does a $1 \times 1$ convolution do?

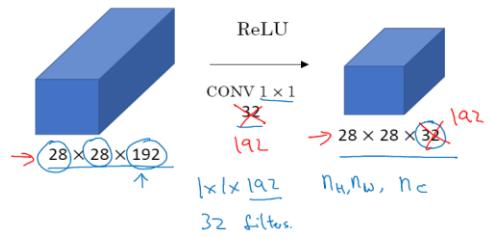


#### Summarize: 1x1 conv

Fully connection in channel direction, applied to each position of image, output dimension:  $H(\text{input}) \times W(\text{input}) \times \text{filter numbers}$

input channel number : could only shrink by convolution (filter numbers)

## Using $1 \times 1$ convolutions



### 2. Using $1 \times 1$ CONV.

E.g: input  $18 \times 28 \times 192$  -->filter cov  $1 \times 1 \times 32$  (with Relu) ---->result  $28 \times 28 \times 32$

if shrink input height & width -->pooling layer

if shrink channel -->conv layer: standard conv layer or  $1 \times 1$  CONV.

**The effect of a one-by-one convolution is it just has nonlinearity:** allows to learn a more complex function of network by adding another conv layer, and allows to shrink the number of input channels or keep it the same or even increase it if you want.  
(linear computation on each position through all channels and then apply Relu-->learn complex function;  
shrink input channels by applying multi filters)

note: when to use  $1 \times 1$  conv: region considered is narrower than standard conv, only complex function learned by added Relu??

### 4.2-6 Inception network motivation

When designing a layer for a ConvNet, have to pick, filter size, add pooling layer or not, etc.

Inception network make the decision by its own, this makes the network architecture more complicated, but it also works remarkably well.

#### 1. Motivation for inception network

choosing filter size in conv, or whether or not adding conv/pooling layer-->inception do it by itself.

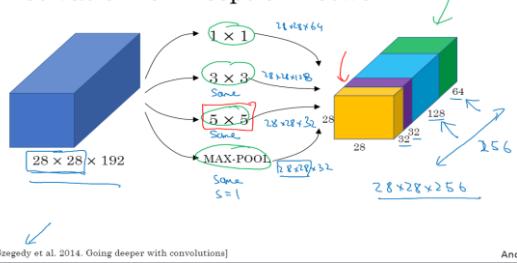
E.g: input  $28 \times 28 \times 192$ :  
use conv: 3 filters size:  $1 \times 1 \times 64$ ,  $3 \times 3$ ,  $5 \times 5$  filter--same padding,  $5 \times 5$  filter--same padding-->output H, W same, only different channel.  
use one max-pool: and use same padding for pooling result having same dimension as conv result-->this is the only case when to use same padding for max-pool.

stacks result together--> $28 \times 28 \times 256$

Above this is the heart of the inception network : use different size filter, max pool.

The basic idea: instead to pick one of these filter sizes or pooling, do them all and just concatenate all the outputs, and let the network learn whatever parameters it wants to use, whatever the combinations of these filter sizes it wants.  
(fixed hyperparameters: filter size, step, padding as right e.g?)

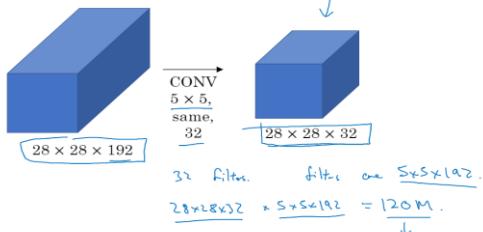
#### Motivation for inception network



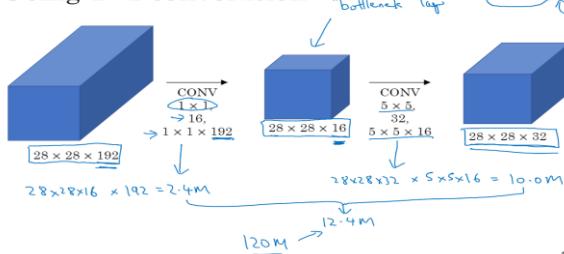
[Szegedy et al. 2014. Going deeper with convolutions]

And

#### The problem of computational cost



#### Using $1 \times 1$ convolution



Andrew Ng

#### Summary:

When building a layer of a neural network and you don't want to have to decide filter size, or pooling layer, the inception module does them all, and concatenates the results.

And then we run into the problem of computational cost, while using a  $1 \times 1$  convolution can create this bottleneck layer thereby reducing the computational cost significantly.  
And it turns out that so long as you implement this bottleneck layer within reason, you can shrink down the representation size significantly, and it doesn't seem to hurt the performance, but saves you a lot of computation.

So these are the key ideas of the inception module.

(Note: computation cost saving by  $1 \times 1$  conv:  
1x1 conv: shrink input channel very much (e.g. from 192 to 16); this bottleneck is caused by shrinking, much smaller channel--while size no change and only compute one time for each position.  
could take it as capturing input feature in each position (through all channels-combining all original features), use linear function/fully connection + Relu-->get much smaller , yet complex features for each position by small number of filters.  
then use standard conv (same padding): to capture more features (channel increase from 16 to 32):  
but combining multi-positions of input (filter size region), and all of these positions' features/channel.)

motivation for  $1 \times 1$  conv: only for computation saving? When to intro. ?by shrinking channels/features and then increase again, will original features/info. lost? how to choose  $1 \times 1$  conv filter number?

### 2. Computational cost:

above e.g:

>1. conv filter  $5 \times 5 \times 192$ , 32 filter

input  $28 \times 28 \times 192$ --> $5 \times 5$  filter conv,  $28 \times 28 \times 32$ .

computation cost:

multiplication: out value number x computation for each value:  $28 \times 28 \times 32 \times (5 \times 5 \times 192) = 120$  million (expansion computation)

Note:  $1 \times 1$  CONV will reduce computation by 10 times.

>2.  $1 \times 1$  convolution: filter size:  $1 \times 1 \times 192$ , 16 filter, +  $5 \times 5 \times 16$ , 32 filter

input  $28 \times 28 \times 192$ --> $1 \times 1$  conv, output  $28 \times 28 \times 16$ --> $5 \times 5$  filter conv,  $28 \times 28 \times 32$ .

input dimension and output dimension same as above con filter  $5 \times 5 \times 192$ , but only add a much smaller intermediate volume  $28 \times 28 \times 16$ —also called bottleneck layer, only 16 channel instead of 192 channel of input.

computation cost:

bottleneck: output value number x each number computation =  $28 \times 28 \times 16 \times 192 = 2.4$  million

conv  $5 \times 5$  layer computation:  $28 \times 28 \times 32 \times (5 \times 5 \times 16) = 10$  million

Total computation: 12 million

-->reduce computation from 120 million multiplication to 12 million multiplication.

here only consider multiplication number, while adding function is same number

### 4.2-7 Inception network

#### 1. Inception module:

take input from some previous layer, e.g  $28 \times 28 \times 192$ , using filters as below:

>1.  $1 \times 1$  conv + standard conv: for save computation, add  $1 \times 1$  conv before standard conv filter

>2. only  $1 \times 1$  conv:

>3. Pooling: +  $1 \times 1$  conv

>>>pooling use same padding here , for keep result same dimension with above conv result dimension, for stack result together

Note: channel number do change by pooling, still 192.

>>>  $1 \times 1$  conv, used to shrink channels, as pooling do not shrink channel.

note: maxpooling +  $1 \times 1$  conv: -->get 'max' value for each feature/channel, then compute more complex, new feature(by  $1 \times 1$  conv filter)

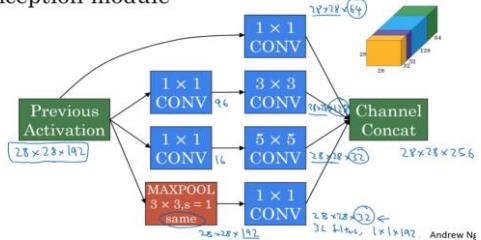
better than  $1 \times 1$  conv + maxpool??, (or  $1 \times 1$  conv only for saving computation , add after maxpool is for future computation saving in next layer?)

>4. Channel concat: Stack above result together

So this is one inception module, and what the inception network does, is, more or less, put a lot of these modules together.

(inception module output: size no change, only channel number changed)

#### Inception module



Andrew Ng

## 2. Inception network:

Put inception modules together.

Right example: looks really complicated. But for the one block there, is basically the inception module that described above.

The inception network is just a lot of these blocks repeated to different positions of the network. (note: there is one position, adding max pooling between two blocks for shrinking height and width).

Note:

There are side branches added in the inception network: Regularizing effect and prevent overfitting (how?).

takes some hidden layer and it tries to use that to make a prediction- softmax output a: take hidden layer out and passes through a few layers like a fully connected layers. And then has the softmax try to predict what's the output label. hidden layer output: ->full connection->activation function-->output.

This is to ensure that features computed even in the hidden units, even in intermediate layers , are not too bad to predicting output class of a image.

note: side branches: if this result ok, then could just take these result as final result-ignore future layers output (kind of reduce layer number, prevent overfitting?)?

Note:

Inception 'name': also called GoogLeNet, developed by authors at Google.

there is a movie 'inception', cite this name as motivation for needing to build deeper new networks

## Summary:

Inception network, is largely the inception module repeated a bunch of times throughout the network, now these later versions like inception v2, inception v3, inception v4, combined with the residual idea of having skipped connections, and that sometimes works even better.

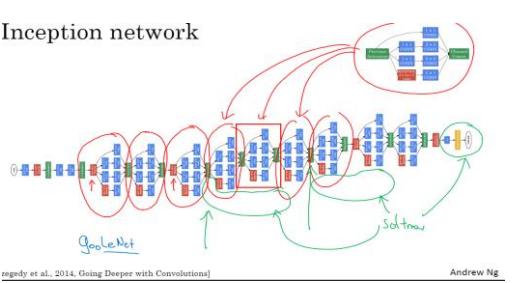
But all of these variations are built on the basic idea of coming up with the inception module and then stacking up a bunch of them together.

## 4.2-8 open source implements

It turns out that a lot of these neural networks are difficult or finicky to replicate because a lot of details about tuning of the hyperparameters such as learning decay and other things that make some difference to the performance. And it's sometimes difficult even at the top universities to replicate someone else's published work just from reading their paper. Fortunately, a lot of deep learning researchers routinely open source their work on the Internet, such as on GitHub. Recommend to use open source implementation/architecture(open source internet: GitHub) as a starter.

### Search on Github:

> 1.1 search in Google: e.g 'Resnet GitHub'-->open the one gonna to copy and click 'download'-->copy 'https', get URL and copy it.  
>1.2: in cp command type: type 'git clone' +copied address, enter,->source downloaded to local dic.  
>1.3: Open the download file to look into detail for the coding.  
use transfer appliton to use architecture others have already training with large data/ multi GPUs...



## Summary:

When developing a computer vision application, a very common workflow would be:

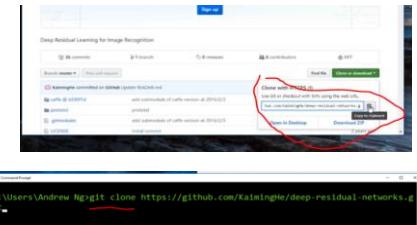
>1. Pick an architecture that you like,( maybe one of the ones you learned about in this course. Or maybe one that you heard about from a friend or from some literature. )

>2. Look for an open source implementation and download it from GitHub to start building from there.

One of the advantages of doing so also is that sometimes these networks take a long time to train, and someone else might have used multiple GPUs and a very large dataset to pretrain some of these networks. It's. But because so many vision researchers have done so much work implementing these architectures, I found that often starting with open-source implementations is a better way, or certainly a faster way to get started on a new project.

And that allows to do transfer learning using these networks .

but if implementing computer vision from scratch, then your workflow will be different.



## 4.2-9 practical advice of convnets \_transfer learning

If you're building a computer vision application rather than training the ways from scratch, from random initialization, often make much faster progress if you download weights that someone else has already trained on the network architecture and use that as pre-training and transfer that to a new task that you might be interested in.

The computer vision research community has been pretty good at posting lots of data sets on the Internet like Image Net, or MS COCO, or Pascal types of data sets, these are the names of different data sets that people post online and a lot of computer researchers have trained their algorithms on.

And use transfer learning to sort of transfer knowledge from some of these very large public data sets to your own problem.

### 1. Transfer learning

download source (code and source) others have trained , use it as pre-training then transfer to own task.  
E.g: cat classification: is A, B or neither ..small training set.

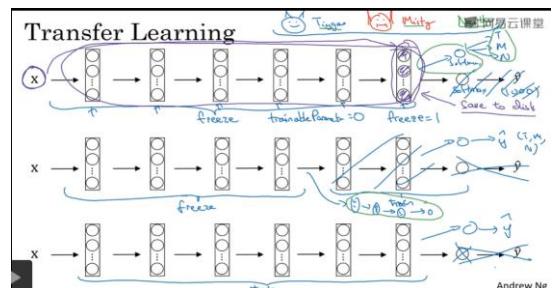
#### 1.1 small training set for new application/task

> 1. download Imagenet (both code and weight), which have 1000 classification (softmax) ,output; softmax-->1000output  
> 2. Remvoe softmax layer and outp, create new softmax layer and output  
>3. Frozen earlier layers' parameter, train parameters associated with your own softmax layer with three possible outputs, on the small data real cared.

Note:

1. frozen open source parameter method: different framework have different method to specify whether or not to train parameter associated with a particular layer. like: put trainableParameter=0, means do not train weight in this layer. or freeze =1

2. Trick for fast computing /training on softmax layer: could save the last frozen weight layers output ( based on all real cared data), take it as 'input' to train shallow neural netork-softmax layer.: do not need to computated every time on previous layer activation of every epoch or taking pass through training set.



## 1.2 Large training set for new application/task

Freeze less low layers, and train more top later layers, and the later layer and softmax new created.

Having your own output unit: could take the last few layers away and just use that as initialization and do gradient descent from there or you can also blow away these last few layers and just use your own new hidden units and in your own final softmax outputs.

Either method worth trying, maybe one pattern is if you have more data, the number of layers you've freeze could be smaller and then the number of layers you train on top could be greater.

If have bigger data set, maybe of enough data not just to train a single softmax unit, but to train some other modest-sized neural network that comprises the last few layers of this final network that you end up using.

#### 1.3 Have a lot of data for new application:

Could: take this open source network and weights and use the whole thing just as initialization and train the whole network, with your own softmax output.

-->Use the weights download just as initialization, replace random initialization and then could do gradient descent, training updating all the ways and all the layers of the network.

## Summary:

This's transfer learning for the training of ConvNets. In practice, because the open data sets on the internet are so big and the ways you can download that someone else has spent weeks training has learned from so much data, you find that for a lot of computer vision applications, you just do much better if you download someone else's open source ways and use that as initialization for your problem.

In all the different disciplines, different applications of deep learning, computer vision is one where transfer learning is something that should almost always do unless, you have an exceptionally large data set to train everything else from scratch .  
But transfer learning is just very worth seriously considering unless you have an exceptionally large data set and a very large computation budget to train everything from scratch by yourself.

## 4.2-10 practical advice of convnets \_data augmentation

Most computer vision task could use more data. And so data augmentation is one of the techniques that is often used to improve the performance of computer vision systems. Computer vision is a pretty complicated task: input this image, all these pixels and then figure out what is in this picture. And it seems like you need to learn the decently complicated function to do that. In practice, there almost all competing visions task having more data will help. This is unlike some other domains where can get enough data.

Today, state of computer vision is : for the majority of computer vision problems, we feel like we just can't get enough data. And this is not true for all applications of machine learning, but it does feel like it's true for computer vision.

That means when you're training in computer vision model, often data augmentation will help. And this is true whether you're using transfer learning or using someone else's pre-trained ways to start, or whether you're trying to train something yourself from scratch.

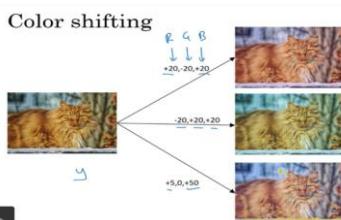
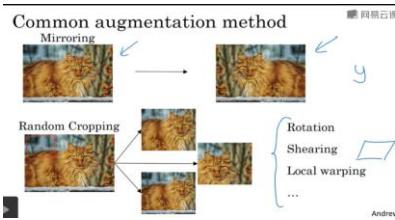
### 1. Common Data augmentation:

>1. Mirroring on the vertical  
good to use if mirrored pic still keeps the feature you want

>2. Random cropping  
not good option , as have chance to crop content /features mattered.  
Yet in practice works well as long as the random crops are **reasonable large subset** of the actual image

Mirroring and random cropping are frequently used.

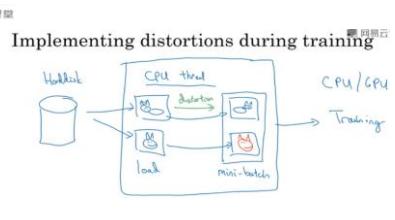
3. Rotate, shearing, local warping.. : a little bit less used may due to complexity.



### 2. color shifting: Special Data augmentation: commonly used

>1. color shifting: take different values of color channel (R, G, B), and use them to distort color channels.  
color shifting motion: image color could be changed by environment (sunlight, indoor illumination) --> algorithm more robust to image color change.

Note: there are different ways to sample color R, G, B, one is PCA  
color distortion: PCA (principal component analysis)  
color augmentation: keep overall color of the tint same (使总体的颜色保持一致) : weight of color RGB same: if image too purple: R&G weight high, PCA will add a lot subtract to R, G, but little to B.



### 3. Implementing distortions during training:

#### 1. Harddisk

Have training data stored in a hard disk . If I have a small training set, you can do almost anything and it is okay. But the very last training set and below is how people will often implement it.

#### 2. CPU thread

A CPU thread is constantly loading images of hard disk as well as implementing whatever distortions are needed to form a batch or really many batches of data. And this data is then constantly passed to some other thread or some other process for implementing training and this could be done on the CPU or really increasingly on the GPU if you have a large neural network to train.

#### 3. CPU/GPU

So, a pretty common way of implementing data augmentation is to really have one thread /multiple threads that is responsible for loading the data and implementing distortions, and then passing that to some other thread or some other process that then does the training.

store original data in hard disk-->1 /more thread : resp. for loading data and add distortion (data augmentation)-->pass to other thread/process for training

note: CPU thread for loading+ adding distortion, and training on CPU/GPU, could run in parallel.

#### Summary:

Similar to other parts of training a deep neural network, the **data augmentation** process also has a few hyperparameters: such as **how much color shifting** to implement and exactly what parameters to use for random cropping.

So, similar to elsewhere in computer vision, a good place to get started might be to use **someone else's open source implementation** for how they use data augmentation.

But if you want to capture more in variances, while someone else's open source implementation isn't, it might be reasonable also to tune these hyperparameters yourself.

## 4.2-11 State of computer vision

Deep learning has been successfully applied to **computer vision, natural language processing, speech recognition, online advertising, logistics, many problems**. There are a few things that are unique about the application of deep learning to computer vision, about the status of computer vision.

### 1. Data vs. hand-engineering

Can think of most machine learning problems as falling somewhere on the spectrum between having relatively little data to have lots of data:

#### 1.1: Data spectra

>1. Speech recognition:  
today we **have a decent amount of data for speech recognition** and it's relative to the complexity of the problem.

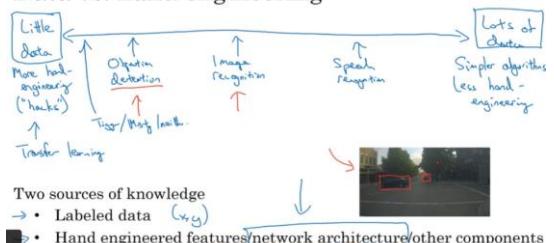
#### 2. Image recognition:

even though there are reasonably large data sets today for image recognition or image classification, because image recognition is just a complicated problem to look at all those pixels and figure out what it is. It feels like even though the online data sets are quite big like over a million images, feels like we **still wish we had more data**.

#### 3. Object detection

Object detection: look in the picture and actually putting the bounding boxes, telling where in the picture the objects are. As the cost of getting the bounding boxes is just **more expensive** to label the objects and the bounding boxes, we tend to have less data for object detection than for image recognition.

### Data vs. hand-engineering



### 1.3: Algorithm Two source of knowledge:

1. labeled data (x, y) , use for supervised learning

2. hand engineering

There are lots of ways to hand-engineer a system: carefully hand **designing features**, carefully hand designing **network architectures** or maybe other components of system.

So when don't have much labeled data, have to call more on hand-engineering .

As computer vision is trying to learn a **really complex function**, and it often feels like we don't have enough data for computer vision. Even though data sets are getting bigger and bigger, often we just don't have as much data as we need.

This is reason why:

1. computer vision historically and even today has relied more on hand-engineering.

### 2. Computer vision has developed rather complex network architectures.

In the absence of more data, the way to get good performance is to spend more time architecting. when you don't have enough data hand-engineering is a very difficult, very skilful task that requires a lot of insight. And someone that is insightful with hand-engineering will get better performance, and is a great contribution to a project.

#### 1.2 Note on average:

>1. when have a lot of data: tend to use simpler algorithms as well as less hand-engineering.

So, there's just less need to carefully design features for the problem, but instead you can have a **giant neural network, even a simpler architecture**, and have a neural network. Just learn whatever it wants to learn when you have a lot of data.

>2. Don't have that much data : tend to engage in more hand-engineering. And there are more hacks(黑客攻击). **how???** when don't have much data then hand-engineering is actually the best way to get good performance.

## 1.4 Computer vision status

>1. Historically due to very small data sets, computer vision literature has relied on a lot of hand-engineering.

>2. Latest: still a lot of hand-engineering of network architectures in computer vision.

In the last few years, the amount of data with computer vision task has increased dramatically, resulted in a significant reduction in the amount of hand-engineering that's being done. yet the data is not big enough, so still a lot of hand-engineering architecture.

-->very complicated hyper parameters choices in computer vision, more complex than a lot of other disciplines.

### Object detection:

And usually have smaller object detection data sets than image recognition data sets, and the algorithms become even more complex and has even more specialized components.

### Transfer learning:

Fortunately, one thing that helps a lot when you have little data is transfer learning. used a lot for when you have relatively little data.

## 2. Tips for doing well on benchmarks/winning competitions

Computer vision researcher are really into doing well on standardized benchmark data sets and on winning competitions: easier to get the paper published. So, there's just a lot of attention on doing well on these benchmarks.

>Positive side : helps the whole community figure out the most effective algorithms.

> yet, in the papers people do things that allow you to do well on a benchmark, but that wouldn't really use in a production or a system .

### Tips for doing well on benchmarks:

#### >1. Ensembling

Training several neural networks (different weight, layers) independently and average their output ( $y^*$ ) (not average weights of different network): maybe 1% or 2% better performance on benchmark.

Wasting computation and time, not use in real customer:

As ensembling means that to test on each image, might need to run an image through maybe 3 to 15 different networks quite typical. This slows down your running time by a factor of 3 to 15, or sometimes even more.

Ensembling is one of those tips that people use doing well in benchmarks and for winning competitions. But almost never use in production to serve actual customers. I guess unless you have huge computational budget and don't mind burning a lot more of it per customer image

## Tips for doing well on benchmarks/winning competitions

### Ensembling

- Train several networks independently and average their outputs

### Multi-crop at test time

- Run classifier on multiple versions of test images and average results



## 2. Tips for doing well on benchmarks:

### Summarize:

Ensembling and multi-crop are techniques that are used much more for doing well on benchmarks than in actual production systems.

And one of ensembling problems is that you need to keep all these different networks around, this takes up a lot more computer memory.

Multi-crop at least keep one network around. So it doesn't suck up as much memory, but it still slows down run time quite a bit.

## 3. Use open source code:

### >1. Use architectures of networks published in the literature

In fact, do not tend to use these two techniques when building production systems even though they are great for doing better on benchmarks and on winning competitions. Because a lot of the computer vision problems are in the small data regime, other people have done a lot of hand-engineering of the network architectures.

And a neural network that works well on one vision problem often may be surprisingly, work on other vision problems as well.

So, to build a practical system often you do well starting off with someone else's neural network architecture. And

### >2. Use an open source implementation if possible

Because the open source implementation might have figured out all the finicky details like the learning rate, case scheduler, and other hyper parameters.

### >3. use pretrained model and fine-tune on your data set

Someone else may have spent weeks training a model on half a dozen GPU and on over a million images. By using someone else's pretrained model and fine tuning on your data set, you can often get going much faster on an application.

But have the compute resources and the inclination, or want to invent your own computer vision algorithm, could train networks from scratch.

## Use open source code

网易云课堂

- Use architectures of networks published in the literature
- Use open source implementations if possible
- Use pretrained models and fine-tune on your dataset

## 4.2-13 Mobile

Have learned about the ResNet architecture, inception net. while MobileNets is also another foundational convolutional neural network architecture used for computer vision.

Using MobileNets will allow you to build and deploy new networks that work even in low compute environment, such as a mobile phone.

### 1. Why need another neural network architecture

Other neural networks have learned about so far are quite computationally expensive. If you want your neural network to run on a device with less powerful CPU or a GPU at deployment, then there's another neural network architecture called the MobileNet that could perform much better.

### >1. Normal convolution function:

input image ( $n \times n \times n_c$ -number of channels) -->filter ( $f \times f \times n_c$ )

e.g.  $6 \times 6 \times 3$ , convolve it with a filter  $3 \times 3 \times 3$ . The way you do this is you would take this filter which I'm going to draw as a three-dimensional yellow block and put the yellow filter over there. There are 27 multiplications you have to do, sum it up, and then that gives you this value. Then shift a filter over, multiply the 27 pairs of numbers, add that up that gives you this number, and you keep going to get this, to get this, and then this, and so on, until you have computed all four by four output values. We didn't use padding in this case, and we use a stride of one, which is why the output size a bout by a bout is a bit smaller than the input size. That is this four by four instead of six by six. Rather than just having one of these three by three by three filters, you may have some  $n_c$  prime filters and if you have five of them, then the output will be four by four by five, or a bout by a bout by  $n_c$  prime. Let's figure out what is the computational cost of what we just did. It turns out the total number of computations needed to compute this output is given by the number of filter parameters which is three by three by three in this case, multiplied by the number of filter positions. That is a number of places where we place this big yellow block, which is four by four, and then multiplied by the number of filters, which is five in this case. You can check for yourself that this is the total number of multiplications we need to do, because at each of the locations we plopped down the filter, we need to do this many multiplications, and then we have this many filters. If you multiply these numbers out, this turns out to be 2,160. We'll come back. You see this number again later in this video, when we come up with the depthwise separable convolution that will be able to take as input a six by six by three image and outputs four by four by five set of activations that were fewer computations than 2,160. Let us see how the depthwise separable convolution does that. In contrast to the normal convolution which you just saw, the depthwise separable convolution has two steps. You're going to first use a depthwise convolution followed by a pointwise convolution. It is

## 4-3 Object Detection

### 4.3-1 Object localization

Object detection is one of the areas of computer vision that's just exploding and is working so much better than just a couple of years ago. In order to build up to object detection, you first learn about object localization.

#### 1. What are localization and detection

##### 1.1 Definition

>1. **Image classification task:** is an algorithm looks at this picture and might be responsible for saying this is a car. So that was classification.

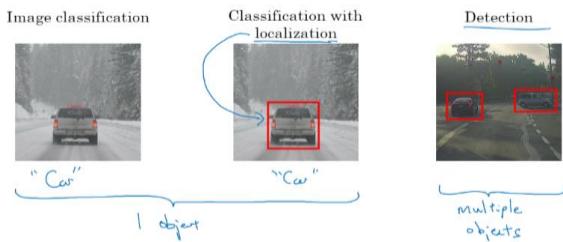
The problem you learn to build in your network to address later on this video is classification with localization.  
>2. **classification with localization:** Which means not only do you have to label this as say a car but the algorithm also is responsible for putting a bounding box, or drawing a red rectangle around the position of the car in the image.

The term localization refers to figuring out where in the picture is the car you've detected.  
Later this week, you then learn about the detection problem.

>3. **Detection problem:** there might be **multiple objects** in the picture and you have to detect them all and and localized them all.

If doing this for an autonomous driving application, then you might need to detect not just other cars, but maybe other pedestrians and motorcycles and maybe even other objects.

### What are localization and detection?



#### 2. Classification with localization:

##### 2.1 Standard classification pipeline:

Input a picture into a ConvNet with multiple layers , and this results in a vector features that is fed to **maybe a softmax unit** that outputs the predicted class.

e.g. Building a self driving car, maybe object categories are a pedestrian, or a car, or a motorcycle, or a background. So these are algorithm classes, they have a softmax with four possible outputs.

##### 2.2 Localize object in the image: modify on the standard classification pipeline

To recognize and also localize object in the image, change neural network to have a few more output units that output a bounding box.

e.g. can have the neural network output four more numbers, (bx, by, bh, and bw), and these four numbers parameterized the bounding box of the detected object.

Notational convention: upper left of the image is (0,0), right down corner (1,1), object central positon and size info. labeled based on image size.  
Specifying the bounding box by: box midpoint (bx, by) + box size (bh, bw)

##### 2.3 Define target label:

**Training set: (x, y)**

x: image;

y: labeled: object classification + box location (bx, by, bh, bw)

Training set contains not just the object cross label, which a neural network is trying to predict up here, but it also contains four additional numbers. Giving the bounding box then you can use supervised learning to make your algorithm outputs not just a class label but also the four parameters to tell you where is the bounding box of the object you detected.

##### Target label y definition:

$y^{\wedge} = [P_c, bx, by, bh, bw]^T$

>1. **Pc:** is there a object or not (1: yes, 0 :no)

(there is a object probability that one of classes trying to detect is there.). if P<sub>c</sub> is 0, then do not care rest output b,c,?, could be whatever.

>2. **bx, by, bh, bw:** object position and size. If P<sub>c</sub> is 0, then bx, by, bh, bw is 0

>3. **C1, c2, c3:** object category , detected object is which category, if c1,then 1,0,0; at most one of these class category should be 1.

(only one project/class in image, softmax output is probability of each class-probability sum =1, take the max probability as prediction clas, and loss =  $y(k) \cdot \log(h(k))$ )

training set: image, + labeled target (only have two case: P<sub>c</sub> is 1, P<sub>c</sub> is 0)

label training set need bounding boxes (bx, by, bh, bw) in the labels.

e.g1: left pic,\_there is one car in image ->label y = [1; bx; by; bh; bw; 0; 1; 0]

e.g2: right pic\_no car/object in image-->label y=[0; ?; ?; ?; ?; ?], pc set to 0, and rest part do not care. (make random value?? in label data??)

##### Summary:

So that's how you get a neural network to not just classify an object but also to localize it. The idea of having a neural network output a bunch of real numbers to tell you where things are in a picture turns out to be a very powerful idea.

### 4.3-2 Landmark

neural network output four numbers of bx, by, bh, and bw to specify the bounding box of an object you want a neural network to localize. In more general cases, you can have a neural network just **output X and Y coordinates of important points and image, sometimes called landmarks**, that you want the neural networks to recognize.

##### 1.2 objects number in image:

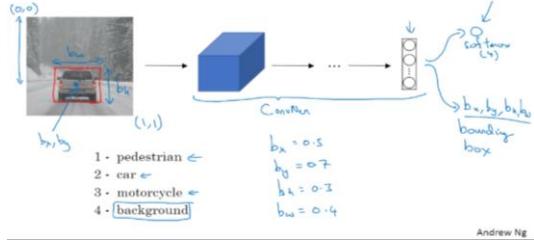
1. The classification and the classification of localization problems usually have **one object**: one big object in the middle of the image that you're trying to recognize or recognize and localize.

##### 2. In detection problem there can be multiple objects.

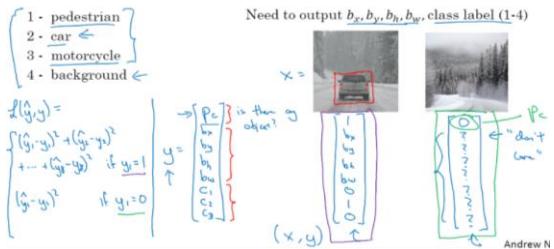
And in fact, maybe even **multiple objects of different categories** within a single image.

So the ideas learned about for image classification will be useful for classification with localization. And that the ideas learned for localization will then turn out to be useful for detection.

### Classification with localization



### Defining the target label y



### 3. Loss function (one input/image)

output  $y^{\wedge} = [y^{\wedge}1, y^{\wedge}2, y^{\wedge}3, y^{\wedge}4, \dots, y^{\wedge}8]^T$  (8 dimension)

labeled target  $y = [y_1, y_2, y_3, y_4, \dots, y_8]^T$  (8 dimension)

##### Loss for one signal example:

> . if  $y_1=1$  (labeled  $p_c=1$ ) --> $L(y^{\wedge}, y) = \sum_{i=1}^8 (y_i - y_i^{\wedge})^2$ : sum square error over all dimension

note: loss function here is simplified by using squared error to all output unit.  
in reality could use different method for different  $y^{\wedge}$ :  $P_c$  could use logistic regression, object position and size could use squared error, classification ( $c_1, c_2, c_3$ )could use log loss (output one value)

>2. if  $y_1=0$  (labeled  $p_c=0$ ) --> $L(y^{\wedge}, y) = (y^{\wedge}1 - y_1)^2$ , as when there is no object, all cared is how accurately the neural network output  $P_c$ .

## 1. Landmark detection:

### 1.1 Example

e.g1. Face detection :  
Right pic, output coordinance of 4 corners of eyes.

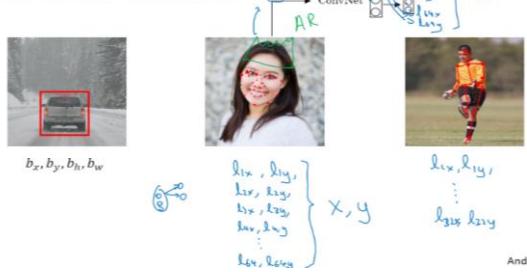
-->label output add 8 number for the corner coordinance: l1x, l1y, ...l4x, l4y., train neural network.  
Take 64 points as landmark in image , -->Neural network output 128+1 value, tell if there is a face as well as all the landmarks on the face.

E.g 2: pose detection:

Could define a few key positions like the midpoint of the chest, the left shoulder, left elbow, the wrist, and so on, and just have a neural network to annotate key positions in the person's pose as well.  
By having a neural network output, all of those points I'm annotating, you could also have the neural network output the pose of the person.  
And of course, to do that you also need to specify on these key landmarks and use their coordinates to specify the pose of the person.

Being able to **detect these landmarks** on the face, this is also a key building block for the computer graphics effects that warp the face or drawing various special effects like putting a crown or a hat on the person.

## Landmark detection



## 1.2: Training set

for training set (X, Y): . images X , labels Y -People have to go through all the images X and **laboriously annotate** all of these landmarks.

> Landmark: Output x, y coordinance of important points in image

>Method:

Modify output to add more output in  $y^A$  for landmark coordinance.: adding a bunch of output units to output the x, y coordinantes of different landmarks want to recognize.

Note:

1. Label training set need also contains all of these landmarks' coordinance

>2. the identity of landmark 1 must be consistent across different images, labels of landmark have to consistent across differnt images (landmarks have same postion meaning in all input image, if landmark1 means left eye right corner position, then all images landmark 1 also eye right corner position)

### Summary:

This idea might seem quite simple of just adding a bunch of output units to output the X,Y coordinantes of different landmarks you want to recognize.

## 4.3-3 Object detection\_slid window

use a ConvNet to perform object detection , Sliding Windows Detection Algorithm.

### 1. Car detection example

Process:

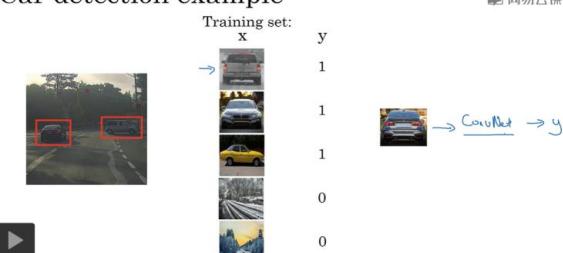
>1. create a label training set, so x and y with closely cropped examples of cars.

for the purposes in this training set, you can start off with the car closely cropped images: Meaning that x is pretty much only the object car. So, you can take a picture and crop out and just cut out anything else that's not part of a car. So you end up with the car centered in pretty much the entire image

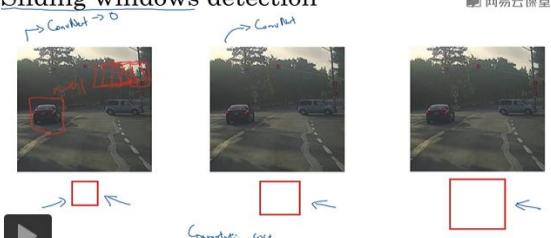
>2. Train convNet on training set, output there is car or not y=0 or 1.

>3. Use trained ConvNet classifier in sliding windows detection

### Car detection example



### Sliding windows detection



### 2. Sliding windows detection

used trained conv neural network (classifier) in sliding window detection:

idea: go through every image region of sliding window size with some stride step, and pass lots of little cropped images into convnet, have it classified 1/0 for each position and/at some stride.

#### 2.1: Process

>1: picking a certain size sliding window,

>2: putting this size red square input into conv-->conv network predict for that little red square input, contain object or not.

>3. take a second image: by moving this red square shift a little bit over ( in specified stride step) into conv network and predict contain object or not  
keep going , until: slid red square across every position of the image.

Note: slide window size could be chosen acc.to Convnet requirement.(bigger, or smaller)

>4. then repeat above process 1-3 with larger size window, with specified stride step, until run throughout the entire image until get to the end.

for the larger image region of sliding window size, resize this region into whatever input size the ConvNet is expecting, feed that to ConvNet to predict have object or not.

>5. repeat 3rd time with more large size slide window, with specified stride step

If do like above, then so long as there is a car somewhere in the image, there will be a window fed into convnet, will have output 1 for that image region. so detect there is a car there.

So this algorithm is called Sliding Windows Detection because you take these windows and slide them across the entire image and classify every square region with some stride as containing a car or not.

#### 2.2 Disadvantage: Huge computation cost:

As cropping out so many different square regions in the image and run them through Convnet independently.

(e.g: pixel 'a' in 3 cropped images: cropped1, 2, 3, then for each cropped image, pixel 'a' is computed by conv network-compute may many times repeatedly by conv same one filter.-->pixel a is computed repeatedly firstly in convnet classifier by same filter, then computed again repeatedly in different cropped image by same filter in conv net.)

>And if use a very coarse stride / big stride, then will reduce the number of windows you need to pass through the convNet, yet this coarser granularity may hurt performance.

>If use a very fine granularity or a very small stride, then a huge number of all these little regions will pass through the ConvNet, means there is a very high computational cost.

#### Summary:

>1. before the rise of Neural Networks people used to use much simpler classifiers like a simple linear classifier over hand engineer features in order to perform object detection. And in that era because each classifier was relatively cheap to compute, just a linear function, Sliding Windows Detection ran okay.it was not a bad method.

>2. but now with convnet, running a signal classification task is much more expensive, and sliding window in this way is too slow.and unless use fine granularity (little cropped image) or very small stride, will not able to localize the object that actually within the image as well.

This problem of computational cost has a pretty good solution. In particular, the Sliding Windows Object Detector can be implemented convolutionally or much more efficiently.

## 4.3-4 Convolutional implementation of sliding window

Sliding windows object detection algorithm using a convnet is too slow, could implement this algorithm convolutionally.

## 1. Turning FC layer into Conv layers

(motivation: getting result representing/corresponding to multi image position)

e.g: image  $14 \times 14 \times 3 \rightarrow \text{conv } 5 \times 5 \rightarrow \text{Max pool } (5 \times 5 \times 16) \rightarrow \text{FC } (400, 1) \rightarrow \text{FC}(400, 1) \rightarrow \text{Softmax}$ , output y (4 classes here)(4,1)

Turn FC  $\rightarrow$  Softmax into convolutional layers:

> Turn FC into convolutional layer: conv filter size:  $5 \times 5 \times 16$ , 400 filters  $\rightarrow$  result:  $1 \times 1 \times 400$  (400 channels) filter number = FC units number, filter size same as input size (height, width, channel number).

(original no weight, current have  $5 \times 5 \times 16$ ,  $\times 400$  weight parameters  $\rightarrow$  why could not just unrolled?)

> Turn 2nd FC into convolutional layer: filter  $1 \times 1 \times 400$  convolution, 400 filters. result ( $1 \times 1 \times 400$ )

(current and original weight number is  $400 \times 400$ , and still fully connection)

FC  $\rightarrow$  FC, change to cov layer using  $1 \times 1$  conv : do not change FC property.

as FC is fully connect , each input have a unique weight for each output value. while  $1 \times 1$  conv also adding unique weight to each input position for each output, different output generated by different  $1 \times 1$  conv filters-adding different, unique weights to each input positions )

> Turn softmax function into convolutional layer: filter  $1 \times 1 \times 400$ , 4 result ( $1 \times 1 \times 4$ )

(current and original weight number both is  $4 \times 400$ , and still fully connection)

Note:

>1. FC output vector (400,1) turn into actually conv. result ( $1 \times 1 \times 400$ ): 400 channels.

>2. Mathematically change FC to conv is same as FC layer:, as each of output is arbitrary linear function of  $5 \times 5 \times 16$  (filter size) activations ,from the previous layer: filter function on input position

1.2. FC  $\rightarrow$  softmax: change Fully connection (standard neural network weight matrix) into  $1 \times 1$  conv + softmax

## 2. Convolutional implementation of sliding window

### 2.1 Example: convolution implement of sliding windows

sliding Window ConNet inputs  $14 \times 14 \times 3$  image:

image  $14 \times 14 \times 3 \rightarrow \text{conv } 5 \times 5 \rightarrow \text{Max pool } (5 \times 5 \times 16) \rightarrow \text{FC } (1 \times 1 \times 400) \rightarrow \text{FC}(1 \times 1 \times 400) \rightarrow \text{Softmax}$ , output y (4 classes here)( $1 \times 4$ )

test image:  $16 \times 16 \times 3$

>1. In original sliding windows algorithm: input the blue region  $14 \times 14 \times 3$  into a convnet and run that once to generate a classification 0/1 and then slightly down (horizontally, vertically, 1 step here) to every position of the image: run this convnet four times in order to get four labels. But it turns out a lot of this computation done by these four convnets is highly duplicative.

(e.g: pixel 'a' in these 4 cropped images: cropped1, 2,3,4. then for each cropped image, pixel 'a' is computed by conv network-compute many times repeatedly by conv same one filter. $\rightarrow$ pixel a is computed repeatedly firstly in singal convnet classifier by same filter, then computed again repeatedly in different 4 cropped image by again the same filter in conv net.)

>Convolutional implementation of sliding windows: allows these 4 forward passes in the convnet to share a lot of computation.

image  $16 \times 16 \times 3 \rightarrow \text{conv } 5 \times 5 \rightarrow \text{Max pool } (5 \times 5 \times 16) \rightarrow \text{FC } (1 \times 1 \times 400) \rightarrow \text{FC}(1 \times 1 \times 400) \rightarrow \text{Softmax}$ , output y (4 classes here)( $2 \times 2 \times 4$ )

(could take it as using a sliding window same size as input image,  $\rightarrow$ but will not shrink image size to conv classifier expected size  $14 \times 14 \times 3$ , instead run conv filters on whole image  $\rightarrow$ pixel 'a' will go through only one time of conv computation-also repeated computation by filters in strides involve this pixel, but will not go through 4 times-conv computation-pixel 'a' may computed more time than separated one conv computation in cases : pixel a in cropped image-conv classifier expected size  $14 \times 14 \times 3$ , covered by filter, but filter now is outside cropped image:-then no computation. but in convolution implementation-for the same filter position, filter may still inside image, then will compute pixel a also.

-may not cause computation more, as if filter still within image, then there are may cropped image-decided by sliding window and window stride, covered pixel a and also same position filter  $\rightarrow$ compute pixel a in another cropped image, computing not less than convolution implementation.

### 2.1 Example: convolution implement of sliding windows

method: use same conv filter parameter of sliding ConvNet run on the entire test image, instead of inputting high duplicated image regions into sliding ConvNet.

Convolution implementation: instead of forcing to run four propagation on four subsets of the input image independently. Instead, it combines all four into one form of computation and shares a lot of the computation in the regions of image that are common. So all four of the 14 by 14 patches in this e.g.

(common regions of input: shared in different(at least two) cropped windows, + could cover conv classifier filter size, this region is saved from duplicated computation, go through conv classifier only once.)

### 2.2 Example\_bigger image: $28 \times 28 \times 3$

Run forward prop the same way (conv filter , pooling and FC parameters same as sliding window conv ):

image  $28 \times 28 \times 3 \rightarrow \text{conv } 5 \times 5 \rightarrow \text{Max pool } (5 \times 5 \times 16) \rightarrow \text{FC } (1 \times 1 \times 400) \rightarrow \text{FC}(1 \times 1 \times 400) \rightarrow \text{Softmax}$ , output y (4 classes here)( $8 \times 8 \times 4$ )

(note: image convolution implementation result size = image size -conv classifier expected size (classifier result size is  $1 \times 1 \times \text{channels}$ )/sliding window step +1:

e.g: image size:  $28 \times 28$ , conv classifier  $14 \times 14$ , sliding window step  $s=2 \rightarrow$ image convolution implementation result size =  $(28-14)/2 + 1 = 8$

>> sliding window step = multiplication of step of each classifier filter in conv, max, etc.; e.g above: conv stride is 2, max pool stride is 2  $\rightarrow 1 \times 2 = 2$ , is sliding window step.

### 3. Summary:

for implementing sliding windows, previously, need to crop out a region (e.g  $14 \times 14$ ) of input image, and run that through your convnet and do that for the next region over, and so on, until hopefully that one recognizes the car.

Now instead of doing it sequentially, with this convolutional implementation you can implement the entire image (e.g  $28 \times 28$ ) and convolutionally make all the predictions at the same time by one forward pass through this big convnet and hopefully have it recognize the position of the car.

This makes the whole thing much more efficient.

Yet still has one weakness: the position of the bounding boxes is not going to be too accurate.

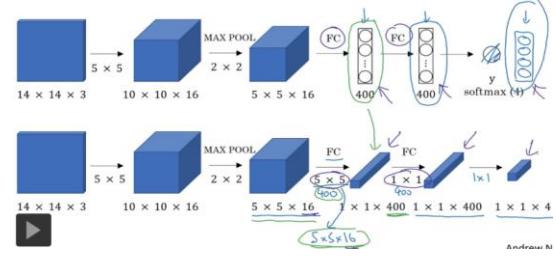
(bounding boxes position: is limited image subregion-size of classifier trained and required??

classifier is trained with cropped image, work on this size image will get one output: [pc, bx, by, ..c1, c2, c3]; $\rightarrow$  in image convolutional implementation, each result of position hwx is:  $1 \times 1 \times c$ : is one classifier result: corresponding also to the size of cropped image used for training classifier, on test image.

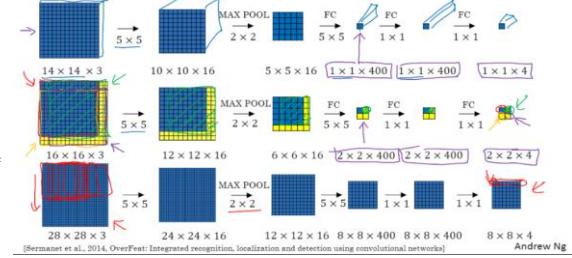
## 4.3-5 Bounding box predictions

Convolutional implementation of sliding windows is more computationally efficient, but it still has a problem of not quite outputting the most accurate bounding boxes.  
Need to improve getting more accurate bounding boxes.

## Turning FC layer into convolutional layers



## Convolution implementation of sliding windows



### 2.2 Summary:

#### 1. Principle of convolution implementation of sliding windows:

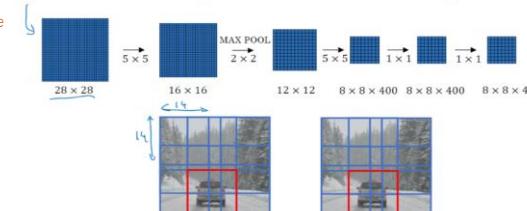
instead of forcing to run forward propagation on all subsets of the input image independently, combine all subset into one computation, and shares a lot of the computation in the regions of the image that are common.-->make all the predictions (whether there is object in subset region of image) at the same time by one computation (as the input image pass through the ConvNet).

Note:

1. the slid window (red square ( $14 \times 14 \times 3$ )) cropped test image into 4 subset (right top, right down, left top, left down) with sliding stride step is = classifier's conv filter step \* pooling step= $1 \times 2 = 2$   
(sliding window size should have relation to conv filter step?)

2. if go convnet for each subset, need go through 4 times computation for the common areas in the test image-->conv. allows to share common position computation by taking it as one image instead of cropped into 4.

## Convolution implementation of sliding windows



## 1. Output accurate bounding boxes

### 1.1 Bounding boxes accurtation problem:

With sliding windows, take discrete set of locations and run the classifier through it. In this case, none of the boxes really match up perfectly with the position of the car. And also the perfect bounding box isn't even quite square, which based on object shape.

(convolution implementation of sliding window:

1. each result position get output  $y^* = [c_1, c_2, \dots]$ , corresponding to that cropped image by sliding window, have which object.  
2. will try different sliding window size on image, and the size that cropped image have object, will take it as object size (no-have bx, by, bw, bh : limited in this cropped size).

1. from convolution implementation of sliding window result, how to match corresponding image region?  
e.g. Image 28x28 x3conv implement result is 8x8x4, if one position result have object, how to find corresponding position in original image, and region size-sliding window size (window size seems not involved in conv implementation at all)

2. why not use a trained classifier +localization: prediction output  $y^* = [pc, bx, by, bw, c_1, c_2, \dots]$  -->get box size??--> too many output, computation cost? -while YOLO may only 19x19 output vector  
)

Convolution implementing:

not try different sliding window size on image-->result hwxw: decided by image size required by trained classifier and classifier conv & max pool step, each position in hwxw, corresponding to origin input image classifier required size region, and the region position could be localized by result positon, classifier image size, and classifier step, e.g. of above, result 8x8x4/(9), first position in second row-->region in input image 28x28x3: region size: 14x14x3, region location: shift sliding wind of 14x14 down bv step2.

## 2. YOLO algorithm

YOLO: stands for you only look once.

### 2.1 Basic idea:

Place down a grid on input image, (e.g. 3x3 grid, in an actual implementation, use a finer one, like maybe a 19x19 grid, and then apply the image classification and localization algorithm to each of the image grids

(note: seems more efficient than convolutioning implementation: as this 19x19 grid sub-image has no common area to go through classifier: less duplicated computation. while in convolutioning implementation, common area may involved in conv calculation several times.)

### 2.2 Label for training

For each grid cells, specify a label Y

>1. Label Y is this eight dimensional vector:  $y = [pc; bx; by; bw; c_1; c_2; c_3]^T$

PC: 0/1 depending on whether or not there's an object in that grid cell

BX, BY, BH, BW: specify the bounding box if there is an object, associated with that grid cell.

C1, C2, C3: classes trying to recognize.

e.g. pedestrian's class, motorcycles and the background class. Then C1 C2 C3 can be the pedestrian, car and motorcycle classes.

e.g. right pic:  $y_1$  (first grid) =  $[0; ?; ?; ?]$  (no object in the first cell, and then not care other values in this vector)

e.g.: grid has object--> $y = [1; bx; by; bw; 0; 1; 0]^T$

### 2.3 YOLO algorithm example:

Use YOLO to train neural network:

input 100x100x3-->convnet-->max pool layers-->...-->eventually maps to 3x3x8 output volume.

Have training set: image - 100x100x3, and target labels Y which are 3x3x8, (labeled  $y_{bh}$ ,  $y_{bw}$  may occupy more than one grid)

-->use back propagation to train the neural network to map from any input X to this type of output volume Y.

(note: above conv. implementation of sliding window is using a trained conv classifier to predict image, classifier trained on cropped image with label y-only for classification c1, c2,... -- no, could add bounding box info.; bx, by, bh, bw, but only bounding box dimension and position limited to cropped image size that for classifier training)

YOLO algorithm advantage: neural network outputs precise bounding boxes (bx, by, bh, bw).

At test set: feed an input image X and run forward prop until you get this output Y.

then for each of the outputs, can then just read off 1 or 0: Is there an object associated with that image positions/grid. if there is an object, what object it is, and where is the bounding box for the object in that grid cell.

Note:

1. so long as do not have more than one object in each grid cell, YOLO algorithm should work okay.  
in practice, could use a much finer grid e.g. (19x19x8), this reduces the chance that there are multiple objects assigned to the same grid cell.

2. If the object spends multiple grid cells, it is assigned only to one grid cells

As the way to assign an object to grid cell is looking at the midpoint of an object and then assign that object to whichever one grid cell contains the midpoint of the object

### 3. Specify the bounding boxes : bx, by, bh, bw

for each grid cell, define left up corner coordinate (0,0), right down corner (1,1), bounding box (bx, by, bh, bw) are specified relative to the grid cell size (h, w)

bx, by: Object middle point position in grid cell: must between 0-1, grid cell whose coordinate defied as: left up corner coordinate (0,0), right down corner (1,1).

If bx/by not in range 0-1, that means in other grid, then current grid has no object pc=0, no care bx, by or other dimension output.

bx, bw: specified as fraction of grid cell size (h, w)

bx, by must between 0-1; but bh, bw may over 1 as object may occupy two or more grids.

Note: have many other parameterized methods of specifying bounding boxes, but above one conventional one works ok.

YOLO:

1. Training set:

x: image, fix all image size and set grid number;

labeled y: grid number gi, labeled y matrix:  $[g_1.pc; g_1.bx; g_1.bn; g_1.bn; g_2.pc; g_2.bn; g_2.bn; \dots; g_n.pc; g_n.bn; g_n.bn]$ ;

2. Algorithm:

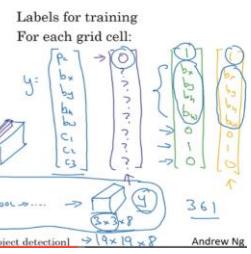
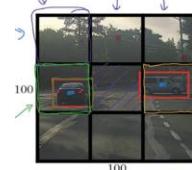
> choose architecture: image-->interception conv (1x1 conv, 1x1 conv + 3x3 conv, 1x1 conv + 5x5 conv, pooling +1x1 conv)-->1x1 conv (FC-FC turn into conv)-->output activation function-->output.

choose: hyperparameter: filter size, output activation function(as only one object, c1, c2, ...-->softmax? pc-sigmoid? bx, by, bw, bh-linear function?)

3. Cost function: for a signal example:

pc-logistic regression, box (bx, by, bh, bw)-squared error, class c1, c2...-softmax loss  $y_k * \log(h_k)$

## YOLO algorithm



### 2.2 Label for training

For each grid cells, specify a label Y

>2. for objects in image:

YOLO algorithm takes the midpoint of each objects and then assigns the object to the grid cell containing the midpoint.

Note:

1. For the grid that has part of the object yet not contains object midpoints: YOLO will pretend this grid cell has no interesting object, class label Y vector has no object

(for one grid computation: how to judge it is part of one object, and how to find midpoints if object occupy more than one grid?)

2. Total output dimension= grid number x (y dimension is 8 in this case: pc, ..c2); each grid of input image corresponds to 8 dimension output  $y_i$ .

### 2.4 YOLO algorithm features:

>1. A lot like the image classification and localization algorithm:

output the bounding box coordinates explicitly, allows the neural network to output bounding boxes any aspect ratio as well as output much more precise coordinates.-->not dictated (限制) by the stride size of sliding windows classifier.

>2. Is single Convolution implementation: do not run conv calculation on each grid separately, but run on input image only once.

Not implementing this algorithm e.g. 9 times on the 3x3 grid or if 362 times if using a 19x19 grid. Instead, this is one single convolutional implementation, where use one ConvNet with a lot of shared computation between all the computations needed for all of your 3 by 3 or all of your 19 by 19 grid cells., this is a pretty efficient algorithm

question:

1. each grid has no common region, how to share computation?  
2. image grid number: has to be fixed value for one algorithm?-algorithm trained on fixed grid number (output dimension decided by grid number), what if want to change grid number, re-train algorithm-filter number/size need to change as well?  
algorithm output dimension decide by grid number x dimension, grid number -->decide conv filter hyper parameter (filter size) - input image size. -->input image size and filter hyperparameter then is fixed by grid number. if grid number change, algorithm need to change.

>3. popular: real-time object detection: conv.calculation actually runs very fast , so this works even for real-time object detection.

note: YOLO seems a special convolutioning implementation: trained classifier, conv the while input image one time.

Trained YOLO is a special classifier: result grid number = { input (H) - classifier size (f) }

/classifier's conv/& max\_pool step +1 since each result corresponding to input image grid position, no shared position of each grid-->

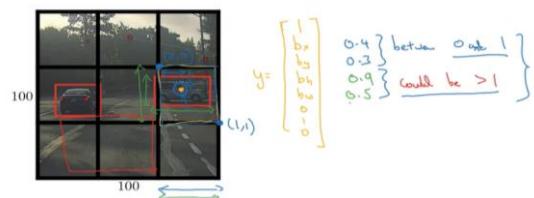
classifier step = classifier's size f

-->grid number(g) = input(H)/f

--> classifier size f = input H / grid number

different: YOLO output bounding box dimension not limited to grid/classifier image size..

## Specify the bounding boxes



## 4.3-6 Intersection over union

Intersection over union function: evaluating object detection algorithm and add another component to algorithm to make it perform better

## 1. Evaluating object detection

### Intersection over union/IoU:

Method: Compute the intersection (交集) over union (并集) of output bounding box and labeled output bounding box (计算两个边界框交集和并集之比) = size of intersection / size of union

### Threshold: If $\text{IoU} >= 0.5$ , -->object detection is correct

0.5 value: human-chosen convention, no particularly deep theoretical reason for it, the higher the IoU, the more accurate the bounding box. higher threshold, more stringent. (rarely lower than 0.5)

### IoU is one way to map localization to accuracy:

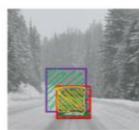
algorithm accuracy: count up the number of times an algorithm correctly detects and localizes an object where you could use a definition like this ( $\text{IoU} >= 0.5$ , correct detection and localization).

### Note:

IoU is also a way of measuring how similar two boxes are to each other.

IoU will be used in max-suppression method.

## Evaluating object localization



$$\begin{aligned} \text{Intersection over Union (IoU)} \\ = \frac{\text{size of intersection}}{\text{size of union}} \\ \text{"Correct" if } \text{IoU} \geq 0.5 \leftarrow 0.6 \leftarrow \end{aligned}$$

More generally, IoU is a measure of the overlap between two bounding boxes.

## 4.3-7 Non-max suppression

One problem of Object Detection is that algorithm may find multiple detections of the same objects. Rather than detecting an object just once, it might detect it multiple times. Non-max suppression is a way to make sure that algorithm detects each object only once.

### 1. Non-max suppression example

e.g: input image grid is 19x19, detect pedestrian, car, etc. in image.

#### Multiple detection problem:

>1. Technically this one object/car has just one midpoint, so it should be assigned just one grid cell. So technically only one of those grid cells should predict that there is specified object/a car.

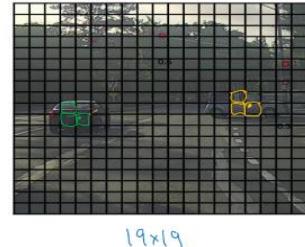
>2. In practice, one object might extend more than one grid, and these grids might take middle point of that sub-part as object middle point and make prediction detected one object. multiple grids may predict have one object-which is actually only sub-part of the same object

As we are running the image classification and localization algorithm on every grid cell, on 361 grid cells, it's possible that many of them will raise their hand and say, "my chance of thinking I have an object in it is large."

-->

Might end up with multiple detections of each object. while Non-Max suppression: will clean up these detections. So they end up with just one detection per one object, rather than multiple detections per object.

### Non-max suppression example



### Non-max suppression example

#### 2. Non-max suppression Procedure

##### for one object category:

>1. first: look at the probabilities associated with each of these detections count on the pc (is actually pc times c1, or c2 or c3)

>2. discard all box with pc <= 0.6 (think at least 0.6 chance it is an object there)

>3: take the box has the largest pc.

>4: looks at rest remaining boxes: all the ones with a high overlap - a high IoU with this largest pc box, will get suppressed.

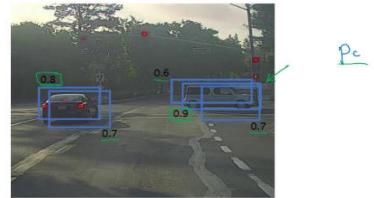
>5: go through remaining rectangles, repeat step 3, &4 above.

finally till all boxes(not discarded: pc > 0.6) deal: either be gotten rid of as suppressed rectangles, or output as the rectangles for different objects.

Non-max suppression: means output maximal probabilities classifications but suppress the close-by ones that are non-maximal.

(key point: take the largest pc as object

if object fully occupy many grids, is it possible largest pc in the grid not containing the whole object midpoint - how does this influence localization-bounding box position and size-no much influence??)



#### 2. Non-max suppression Procedure

Note: if there is more than one object category: c1, c2, c3...then independently carry out non-max suppression for each category (e.g. times for 3 categories): (pc actually is-times category) - in week exercise

If actually tried to detect more than one objects e.g. pedestrians, cars, and motorcycles, then the output vector will have more (e.g. 3) additional components (C1, C2, ...CN). And it turns out, the right thing to do is to independently carry out non-max suppression N times-N classes category number, one on each of the outputs classes.

## 3. Non-max suppression algorithm

e.g: simplified algorithm :

only do car detection: each grid output vector get rid of c1, c2, c3 : [pc; bx; by; bh; bw]T

input image: 19x19 grid --> total output 19x19 x 5

### Process:

>Each output of 19x19 grid/positions, get an output prediction [pc; bx; by; bh; bw]T

>1. Discard all the predictions of the bounding boxes with Pc less than or equal to some threshold, e.g.  $pc \leq 0.6$ . Means unless there's at least a 0.6 chance it is an object there, let's just get rid of it.

The way to think about this is for each of the 361 positions, you output a bounding box together with a probability of that bounding box being a good one.

>3. while there are any remaining bounding boxes:

>>1. Repeatedly pick the box with the highest probability, with the highest pc, and then output that as a prediction.

Commit to outputting that as a prediction for that there is an object/car there.

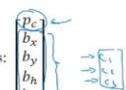
>>2. Discard any remaining box with a high overlap, a high IoU (e.g.  $> 0.5$ ), with the box that you just output in the previous step.

keep doing this while there's still any remaining boxes that you've not yet processed, until you've taken each of the boxes and either output it as a prediction, or discarded it as having too high an overlap, or too high an IoU, with one of the boxes that you have just output as your predicted position for one of the detected objects.

### Non-max suppression algorithm



Each output prediction is:



Discard all boxes with  $pc \leq 0.6$

While there are any remaining boxes:

- Pick the box with the largest pc. Output that as a prediction.
- Discard any remaining box with  $IoU \geq 0.5$  with the box output in the previous step.

Andrew Ng

## 4.3-8 Anchor boxes

One of the problems with object detection as you have seen it so far is that each of the grid cells can detect only one object. When a grid cell wants to detect multiple objects, can use the idea of anchor boxes.

## 1. Overlapping objects

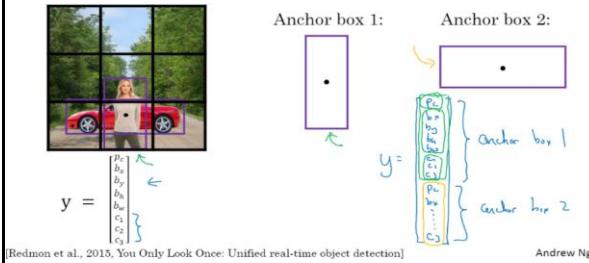
e.g.: image have two objects: person and car, and the midpoint of these two objects are in almost the same place and both of them fall into the same grid cell.

While grid cell Y outputs=  $[pc; bx; by; bh; bw; c1;c2;c3]^T$ . - 3 classes, pedestrians, cars and motorcycles,  
-->it won't be able to output two detections: only set 1 to one of 3 classes.

So have to pick one of the two detections (person or car) to output.

Anchor boxes idea: pre-define two different shapes anchor boxes

## Overlapping objects:



## 2. Anchor box algorithm:

### Summary:

>1. Previously: each object in training image is assigned to the grid cell that corresponds to that object's midpoint. The output Y is  $3 \times 3 \times 8$ , as has a  $3 \times 3$  grid.

And for each grid position, output vector which is PC, then the bounding box, and C1, C2, C3.

>2. With the anchor box,

Each object is assigned to the same grid cell as before, assigned to the grid cell that contains the object's midpoint and anchor box ( $y_1/y_2/y_3/\dots$ ) with the highest IoU with the object's shape.

e.g.: each grid has two anchor boxes, then the one has a higher IoU with object shape,, will be drawn through bounding box.

(for one grid:

1. if only have one object, in the prediction output, have box position  $bx$ ,  $by$  and size  $bh$ ,  $bw$  -->could just create a rectangular box of this size and put around object.
2. if have two object--no anchbox box: no idea how to assign  $y_1$  and  $y_2$  for which object.
3. if have two different shape object--have different shape anchbox box associated with  $y_1$  and  $y_2$ :  $y_1$  will assign to object that shape ( $bx$ , $by$ ,  
 $bh$ , $bw$ ) is higher overlap with  $y_1$  box--object box size still ( $bx$ ,  $by$ ,  $bh$ ,  $bw$ )
4. if have two same shape objects--have different shape anchbox box associated with  $y_1$  and  $y_2$ : could not tell / fix  $y_1/y_2$  for which object.
5. if have more than two shape objects-->assign more than two different anchbox box with  $y_1$ ,  $y_2$ ,  $y_3$ ... each object prediction  $y_i$  have a unique anchbox box, to assign to the similar shape object.

Anchbox box -->use shape similarity of object and  $y_i$  to locate multi object to corresponding output position  $y_i$ .

how to set anchbox box size, ratio of grid cell, how to set this ratio?

)

The object then gets assigned not just to a grid cell but to a pair. It gets assigned to grid cell comma anchor box pair. And that's how that object gets encoded in the target label.

## 3. Cases Anchor box algorithm doesn't handle well:

>1. In the same grid, objects number bigger than anchor boxes number  
e.g.: have two anchor boxes but three objects in the same grid cell

>2.In the same grid, objects have same anchor box shape

e.g.: have two objects associated with the same grid cell, but both of them have the same anchor box shape.

That's cases that this algorithm doesn't handle well. it won't happen much at all, but if it does, this algorithm doesn't have a great way of handling it, better influence some default tiebreaker for those cases.

## 4. summary:

### 4.1 Motivation of anchor boxes:

>1. As a way to deal with what happens if two objects appear in the same grid cell.

In practice, that happens quite rarely, especially if you use e.g a  $19 \times 19$  rather than a  $3 \times 3$  grid. The chance of two objects having the same midpoint out of these 361 cells, it does happen, but it doesn't happen that often.

>2. Even better motivation or even better results: that anchor boxes allows learning algorithm to specialize better(different shape objects).

In particular, if your data set has some tall, skinny objects like pedestrians, and some white objects like cars, then this allows your learning algorithm to specialize so that some of the outputs can specialize in detecting white, fat objects like cars, and some of the output units can specialize in detecting tall, skinny objects like pedestrians.

## 4.3-9 YOLO Detection Putting above together

## 2. Anchor box algorithm:

>1. Anchor box: Associate two predictions with the two anchor boxes.  
In general, you might use more anchor boxes, maybe five or even more.

>output y:  $= [y_1; y_2]^T$   
 $y_1 = [pc; bx; by; bh; bw; c1;c2;c3]^T$   
 $y_2 = [pc; bx; by; bh; bw; c1;c2;c3]^T$   
 $y_1$  associated with anchor box 1, and  $y_2$  associated with anchor box 2.

e.g:  
Pedestrian shape is more similar to anchor box 1, then can encode  $y_1$  to make prediction:  
 $pc=1$ : there is a pedestrian.  
Use ( $bx$ ,  $by$ ,  $bh$ ,  $bw$ ) to encode the bounding box size around the pedestrian, Use ( $c1$ ,  $c2$ ,  $c3$ ) to encode that the object is a pedestrian.

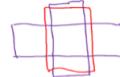
As the car shape is more similar to anchor box 2, than encode  $y_2$  to make prediction.

## Anchor box algorithm

### Previously:

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

Output y:  
 $3 \times 3 \times 8$



Andrew Ng

### With two anchor boxes:

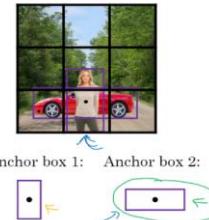
Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.

(grid cell, anchbox)

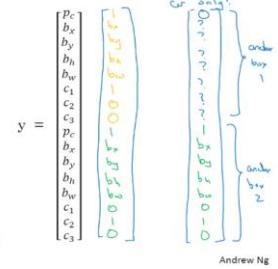
Output y:  
 $3 \times 3 \times 16$   
 $3 \times 3 \times 2 \times 8$

Andrew Ng

## Anchor box example



Anchor box 1: Anchor box 2:



Andrew Ng

## 2. Anchor box algorithm:

### E.g:

now the output Y is going to be  $3 \times 3 \times 16$ , or  $3 \times 3 \times 2 \times 8$ .

As there are now two anchor boxes and Y is eight dimensional, dimension of Y being eight was because we have three objects classes, if you have more objects then the dimension of Y would be even higher.

e.g1: image has both person and car, total output  $Y = [1; bx; by; bh; bw; 1; 0; 0; 1; bx; by; bh; bw; 0; 1; 0; JT]$

e.g1: image has only car, total output  $Y = [0; ?; ?; ?; 1; 0; 0; 1; bx; by; bh; bw; 0; 1; 0; JT]$

## 4. summary:

### 4.2 Anchor box design

>1. People used to just choose them by hand or choose maybe five or 10 anchor box shapes (then 5 or 10  $y_i$  are associated) that spans a variety of shapes that seems to cover the types of objects you seem to detect.

>2. A much more advanced version/ way to do this in one of the later YOLO research papers, is to use a K-means algorithm:

Group together two types of objects shapes you tend to get, and then to use that to select a set of anchor boxes that this most stereotypically representative of the maybe multiple, of the maybe dozens of object causes you're trying to detect.

Put all the components together to form the YOLO object detection algorithm.

### 1. Traing set:

E.g:

Objects to detect: suppose trying to train an algorithm to detect three objects: pedestrians, cars, and motorcycles. And you will need to explicitly have the full background class, so just the class labels

Anchor box: use two anchor boxes,

Grid cell: 3x3 grid

-->outputs y  $3 \times 3 \times 8 \times 2$  (-->fix algorithm architecture hyper parameter : filter size, also fix input image size)  
Eight is actually five plus the number of classes.Five is  $P_c$  and then the bounding boxes  $b_x$ ,  $b_y$ ,  $b_w$  and then C1, C2, C3. 2 anchor box.  
or either view this as  $3 \times 3 \times 2 \times 8$ .

**Traning set:** x: image, y: go through each of these nine grid cells and form the appropriate target vector y.  
e.g first grid label y=[0,?,?,?,?,?, [0,?,?,?,?,?]]  
grid has car: [0,?,?,?,?,?; 1; bx; by; bw; 0;1;0]

### 2. Train a conv net on the input:

e.g image 100x100x3, conv net would output  $3 \times 3 \times 16$

used a pre-trained classifier with cropped image? no, this is for conv implementation for sliding window.  
YOLO no pre-trained/trained classifier, just train the conv net using training set.

### 3. Making predictions

Given an image, neural network output  $3 \times 3 \times 16$  volume:  
expected output: grid does not have object:  $y=[0, \dots, 0, \dots]$   
once  $p_c=0$ , then the following 7 numbers has random value-as not care these values)  
for grid has car:  $y=[0, \dots, 1, bx, by, bh, bw, 0;1;0]$

### 4. Outputting the non-max supressed outputs

>1.If use two anchor boxes, then for each of the grid cells, get two predicted bounding boxes.

Note: some of the bounding boxes( $b_x$ ,  $b_y$ ) can go outside the height and width of the grid cell that they came from.

>2. Get rid of low probility prediction:  $p_c$ (like <0.6)

>3. for each object classes, Independently run non-max suppression to generate final predictions  
e.g three category (pedestrian, car, motorcycle)-->run non-max suppression 3 times, separately for each class.

Final output hopefully will detect all the cars and all the pedestrians in this image.

(Note: is it possible mis supressed object?: e.g. one grid cell have two objects of same class, only object shape different belong to diffrent anchox box.  
if these two objets box ( $b_x$ ,  $b_y$ ,  $b_h$ ,  $b_w$ ) have high IOU-->one object will supressed.  
these objects box compare with anchox boxes, assigne to different anchox.: anchoxes could tell these two object apart by their IOU.-Objects shape are different enough to assign to different anchox box, and to get output both objects. in this case, will these two objects shape are similar enough (IOU), to suppress the other one??-decided by anchox box shape??)

### 5. Summary:

YOLO object detection algorithm is really one of the most effective object detection algorithms, that also encompasses many of the best ideas across the entire computer vision literature that relate to object detection

## 4.3-10: Region proposal\_R-CNN (Regions with convolutional networks)

If you look at the object detection literature, there's a set of ideas called region proposals that's been very influential in computer vision as well.

Make this video optional because tend to use the region proposal algorithm a bit less often but nonetheless, it has been an influential body of work and an idea that you might come across in your own work.

### 1. Region proposal: R-CNN

Sliding windows idea: take a trained classifier and run it across all of these different windows and run the detector to see if there's a object (e.g. car, pedestrian, or maybe a motorcycle). Run sliding window algorithm convolutional: but one disadvantage is it classifies a lot of the regions where there's clearly no object.

R-CNN stands for Regions with convolutional networks or regions with CNNs: tries to pick just a few regions that makes sense to run your conv net classifier on.

So rather than running your sliding windows on every single window, instead select just a few windows to run conv net classifier on just a few window.

(sliding window: using a trained classifier (trained with cropped image) on image.  
conv classifier parameter: will not influe by image size; and image size do not need to be fixed value.(could reshape image size)

YOLO:No trained algorithm: detector(not classifier as more than one object) +localization.  
instead setting training set (X, Labeled Y-acc. to grid number, class number, , anchox box number); train a algorithm to map X to Y.)

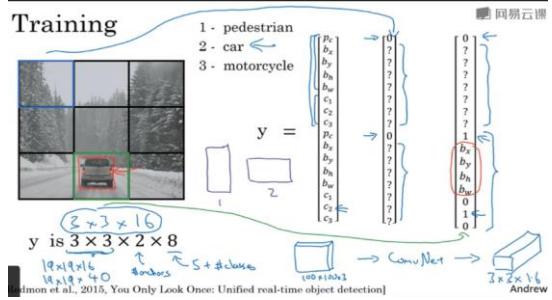
### 2. R-CNN

**Segmentation algorithm:** R-CNN use segmentation algorithm to perform the region proposals

Detail: segmentation algorithm find maybe 2000 blobs and place bounding boxes around about 2000 blobs and value classifiers on just those 2000 blobs, and this can be a much smaller number of positions on which to run conv net classifier, than to run it at every single position throughout the image.

e.g. in the Segmentation result, may finds some blobs -blu, green, etc. color over here, and might pick that bounding boxes and run classifier to see if there's some interesting there.

Note: could run conv net not just on square-shaped regions but also on tall skinny regions to try to find pedestrians or running them on your white fat regions try to find cars and running them at multiple scales as well.



### Making predictions



### Outputting the non-max suppressed outputs



- 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.

### Region proposal: R-CNN



Girshik et al., 2013, Rich feature hierarchies for accurate object detection and semantic segmentation| Andrew Ng

### Faster algorithms

→ 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 windows to classify all the proposed regions.

Faster R-CNN: Use convolutional network to propose regions.

[Girshik et al., 2013, Rich feature hierarchies for accurate object detection and semantic segmentation]  
[Girshik, 2015, Fast R-CNN]  
[Ren et al., 2016, Faster R-CNN: Towards real-time object detection with region proposal networks] Andrew Ng

### 3. Faster algorithm

Now, it turns out the R-CNN algorithm is still quite slow. So there's been a line of work to explore how to speed up this algorithm.

#### 3.1. Basic R-CNN:

>>>1. Using algorithm to find proposed regions,  
>>>2. then classify the proposal regions one at a time  
(use a trained classifier for classification + localization, like sliding window with grained classifier.  
Note: like YOLO without anchor box, and only one grid for the input image).  
>>>3. For each of the regions, output: label+bounding box  
    > Label:  
So is there a car? Is there a pedestrian? Is there a motorcycle there?  
    >also outputs a bounding box:  
so can get an accurate bounding box if there is a object in that region.

#### Basic R-CNN Note:

1. R-CNN algorithm doesn't just trust the bounding box it was given.  
It also outputs a bounding box,  $bx, by, bh, bw$ , in order to get a more accurate bounding box than whatever happened to surround the blob that the image segmentation algorithm gave it.

So it can get pretty accurate bounding boxes.

Now, one downside of the R-CNN algorithm was that it is actually quite slow.

### 4. Summary

The idea of region proposals has been quite influential in computer vision, and there are people still use these ideas.

Person opinion: not the opinion of the computer vision research committee as a whole:

Region propose is an interesting idea but that now having two steps: first proposed region and then classifier it, being able to do everything more or at the same time, similar to the YOLO -you Only Look Once algorithm , that seems to me like a more promising direction for the long term.

### 3. Faster algorithm

#### 3.2 Fast R-CNN

Fast R-CNN is basically the R-CNN algorithm but with a convolutional implementation of sliding windows to classify all the proposed regions.

The original implementation would actually classify the regions one at a time, while fast R-CNN uses a convolutional implementation of sliding windows, and this speeds up R-CNN quite a bit.  
(note: classify the regions one at a time: these regions has no common area, then no duplicated computation.-->convolutional implementation will not be faster. or ??)

Fast R-CNN algorithm problem:

The clustering step to propose the regions is still quite slow.

#### 3.3: Faster R-CNN:

uses a convolutional neural network instead of one of the more traditional segmentation algorithms to propose a blob on those regions, and that wound up running quite a bit faster than the fast R-CNN algorithm.

Yet faster R-CNN algorithm, most implementations are usually still quite a bit slower than the YOLO algorithm.

## 4-4\_ConvNet Special use

### 4.4-1 Face recognition

This week is to show a couple important special applications of conv net. We'll start the face recognition, and then go on later this week to neurosal transfer.

#### 1. Face recognition example

e.g. living people/face get access to pass door , while using picure only did not get persion to pass.

This example demoted: **face detection , as well as liveness detection** (make sure you are a live human)

liveness detection can be implemented using supervised learning as well to predict live human versus not live human.

#### 2.2: Face recognition: 1:K problem

Face recognition problem is much harder than the verification problem:

e.g if having a verification system that's 99% accurate, might not be too bad.  
but now suppose that K= 100 in a recognition system. If you apply this system to a recognition task with a 100 people in your database, you now have a hundred times of chance of making a mistake .  
then if you want an acceptable recognition error, you might actually need a verification system with maybe 99.9 or even higher accuracy to work well on the 100 data.

Note:

recognition is much harder as chances to get error is K times of face verificaiton.-->if face recognition works well , requirement on face verification accuracy should be much higher.

1. Has a database of K persons;
2. Get an input image;
3. Output ID if the image is any of the K persons (or 'not recognized')

#### Summary:

In below items we will focus on building a face verification system as a building block and then if the accuracy is high enough, then you probably use that in a recognition system as well.

### 4.4-2: one-shot learning\_face recognition

One of the challenges of face recognition is that you need to solve the **one-shot learning problem**:

What that means is that for most face recognition applications you need to be able to recognize a person given just one single image, or given just one example of that person's face. while historically, deep learning algorithms don't work well if you have only one training example.

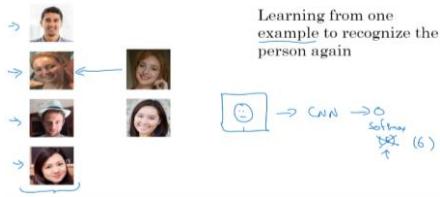
#### 1. One-shot learning

e.g: having image database of multi different people, to recognize if the person occurred now is one of them.

In the one shot learning problem, you have to learn from just one example to recognize the person again.

And you need this for most face recognition systems because you might have only one picture of each of your employees or of your team members in your employee database.

#### One-shot learning



#### Learning a "similarity" function

$d(\text{img1}, \text{img2})$  = degree of difference between images

If  $d(\text{img1}, \text{img2}) \leq \tau$       "same"  
 $> \tau$       "different"



#### 2.2: Learning a 'Similarity' function:

'Similarity' function: neural network to learn a function D: inputs two images and outputs the degree of difference between the two images.

- >1. If the two images are of the same person, D output a small number.
- >2. If the two images are of two very different people, D output a large number.

If  $d \leq t$ , --> predict same person  
if  $d > t$ , --> predict different person  
Threshold t: is a hyperparameter

This is how to address the face verification problem.

#### 2. One-shot learning approach:

2.1: Conv net: image-->conv net-->softmax, output people classes.

Input the image of the person, feed it to a ConvNet. And have it output a label, y, using a softmax unit with four outputs or maybe five outputs corresponding to each of these four persons or none of the above.

Note:

This really doesn't work well:

1. Having such a small training set, it is really not enough to train a robust neural network for this task.

2. If a new person joins your team, have to retrain the ConNet every time:

e.g.: if now you have 5 persons you need to recognize, when new one join, there should now be six outputs, have to retrain the ConvNet every time.

So instead, to make this work, learn a similarity function.

#### 3. Face recogniton

to use similarity function on face recognition problem:

with given new image:

for each image in database: use similarity function to compare new image with this database image, and output prediciton of these two images' similarity.

use similarity to compare new input image and each images in database, predict their similarity (big value or small value).

if the new image/person in database-->d will has small value, otherwise d have big value.

So long as you can learn this function d, which inputs a pair of images and tells basically, if they're the same person or different persons. Then if you have someone new join your team, you can add a new image to database, and it just works fine.

#### Summary:

Similarity function d, inputs two images, allows to address the one-shot learning problem.

### 4.4-3 Siamese network\_face recognition

The job of the function d, is to input two faces and tell you how similar or how different they are. A good way to do this is to use a **Siamese network**.

## 1. Siamese network

standard image conv net: image--->a sequence of conv layerl ----> polling ---->fully connected layers---->FC: feature vector ----->softmax, output

### 1.1: Encoding image:

Modify on standard conv net: remove softmax output layer, only take the final feature vector as an encoding of the input image  $x$ :  $f(x)$

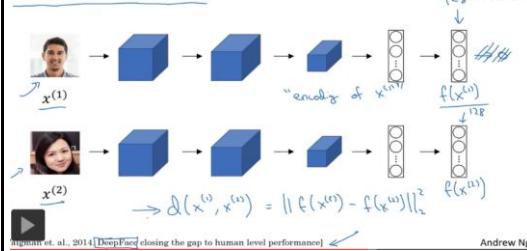
image--->a sequence of conv layerl ----> polling ---->fully connected layers---->FC: feature vector

e.g. input image  $x_1$ , then final get image feature  $f(x_1)$ .

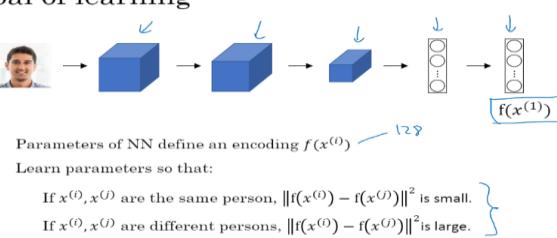
### compare two pic:

The way to build a face recognition system is then that if you want to compare two pictures: image  $x_1$ ,image  $x_2$ . What you can do is feed this second picture  $x_2$  to the same neural network with the same parameters that generated image  $x_1$  feature vector  $f(x_1)$  and get a different vector  $f(x_2)$ , which encodes this second picture.

## Siamese network



## Goal of learning



## 2. Train Siamese network

Note: two neural networks used for two images, have the same parameters.

### 2.1 Training purpose:

Train the neural network parameters so that the encoding those parameters compute results in a function  $d$  that tells whether two pictures are of the same person.

### 2.2 Training parameter and conditions to meet

Parameters of the neural network define an encoding  $f(x_i)$ .

So given any input image  $x_i$ , the neural network outputs multi dimensional encoding  $f(x_i)$ .

learn parameters to generate  $f(x_1), f(x_2)$ , meet below conditions:

If  $x_1, x_2$  are same person, distance  $\|f(x_1)-f(x_2)\|^2$  is small; condition (1)

If  $x_1, x_2$  are different persons, distance  $\|f(x_1)-f(x_2)\|^2$  is large. condition (2)

while vary the parameters in all of these layers of the neural network, will end up with different encodings  $f(x_i)$ .

--> use back propagation and vary all those parameters in order to make sure above conditions (1) and (2) are satisfied.

## 4.4-4 Triplet loss method\_for face recognition

One way to learn the parameters of the neural network so that it gives you a good encoding for your pictures of faces is to define an applied gradient descent on the **triplet loss function**.

### 1. Learning objective

#### 1.1: Triplet loss:

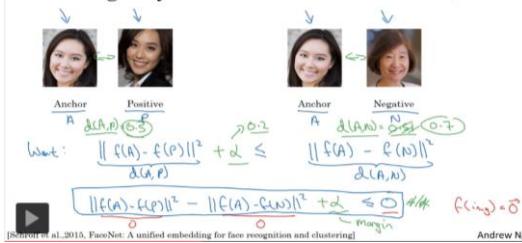
To apply the triplet loss, need to compare pairs of images.

e.g. given a pair of images, want their encodings to be similar when they are the same person. Whereas, given another pair of images, want their encodings to be quite different as they se are different persons.

**Triplet loss function:** always look at one anchor image and want distance between the anchor and the positive image (positive meaning is the same person) to be similar. Whereas want the anchor when pairs are compared to the negative example for their distances to be much further apart.

So, this is what gives rise to the term **triplet loss**, which is that you'll always be looking at three images at a time.Will be looking at an A:anchor image, a P: positive image, as well as a N:negative image.

## Learning Objective



### Triplet loss , want to meet conditions: $d(A, P) + \alpha \leq d(A, N)$

$d(A,P)$ : distance of anchor image encoding  $f(A)$  and positive image encoding  $f(P)$ .

a: **margin parameter**, is a hyperparameter: prevents a neural network from outputting the trivial solutions, by making:  
 >1. Neural network does not just output 0 for all the encoding  $f(x_i)=0$ . (then  $d=0$ );  
 >2. Neural network does not encoding every image identical to other image (then  $d=0$ )

A margin parameter pushes the anchor positive pair and the anchor negative pair further away from each other, as to satisfy  $d(A, P) + \alpha < d(A, N)$ , tend to make  $d(A,N)$  bigger, while  $d(A,P)$  smaller.

## Loss function

$$\begin{aligned} \mathcal{L}(A, P, N) &= \max \left( \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \right) \\ J &= \sum_{i=1}^m \mathcal{L}(A^{(i)}, P^{(i)}, N^{(i)}) \end{aligned}$$

Given 3 images A, P, N

Training set: 10k pictures of 1k persons

### 3. Training set:

**Training set:** have multiple pictures of the same person at least for some people

In order to define this dataset of triplets, do need some pairs of A and P. Pairs of pictures of the same person.

So for the purpose of training system, do need a dataset where have multiple pictures of the same person. If you had just one picture of each person, then you can't actually train this system.

after training, could apply the trained system to one shot learning problem where for your face recognition system, maybe you have only a single picture of someone you might be trying to recognize.

But for training set, do need to make sure having multiple images of the same person at least for some people in training set so that you can have pairs of anchor and positive images.

**Note: each triplet is independently of other triplets(A\_i, P\_i, N\_i)**

**does A need to be same one in the whole training set?**-no, each triplet is independently of other triplets (could one image A use in different triplet?-no necessary, as will choose hard training pairs, no necessary to train on multi same A, pairs.)

## 2. Loss function:

Triplet loss function is defined on triples of images. Given three images, A, P, and N, the anchor positive and negative examples. So the positive examples is of the same person as the anchor, but the negative is of a different person than the anchor.

**Loss function on a singal triplet:**  $= \max(d(A,P) - d(A,N) + \alpha, 0)$

loss =0, if  $d(A,P) - d(A,N) + \alpha < 0$ ;

loss =  $d(A,P) - d(A,N) + \alpha$ , if  $d(A,P) - d(A,N) + \alpha > 0$

By minimizing loss function, will pull  $d(A,P) - d(A,N) + \alpha <= 0$ , and does not care how negative it is once it is less than 0.

Overall cost function for neural network: sum over a training set of these individual losses on different triplets. (not average then  $1/m$ ?)

e.g: training set is 10,000 pictures with 1,000 different persons, then take 10,000 pictures and use it to generate, to select triplets (A,P,N) and then train your learning algorithm using gradient descent on cost function, which is really defined on triplets of images drawn from your training set.

#### 4. Choosing the triplets A,P,N

##### 4.1 Random choose problem: constraint $d(A,P) + \alpha \leq d(A,N)$ is easy to satisfy.

If choose A, P, and N randomly from training set subject to A and P being from the same person, and A and N being different persons, one of the problems is this constraint  $d(A,P) + \alpha \leq d(A,N)$  is very easy to satisfy.

Reason:

given two randomly chosen pictures of people, chances are A and N are much different than A and P.

if A and N are two randomly chosen different persons, then there is a very high chance that  $d(A,N)$  will be much bigger than  $d(A,P) + \alpha$ . And so, the neural network won't learn much from it.

##### 4.2: Manually choose: choose triplets that are 'hard' to train on

Expectation for all triplets is satisfying constraint  $d(A,P) + \alpha \leq d(A,N)$  be satisfied.

**Hard triplet :** is choosing A, P, N,  $d(A, P)$  is actually quite close to  $d(A, N)$ . ( $\rightarrow P, N$  image that are close)

So in this case, the learning algorithm has to try extra hard to make  $d(A,N)$  bigger and try to push  $d(A,P)$  down so that there is at least a margin of  $\alpha$  between them.

##### Effect of choosing hard triplets: increasing the computational efficiency of your learning algorithm.

If you choose your triplets randomly, then too many triplets would be really easy, and so, gradient descent won't do anything because your neural network will just get them right, pretty much all the time. And it's only by using hard triplets that the gradient descent procedure has to do some work to try to push  $d(A,N)$  further away from  $d(A,P)$ .

Parameter update:  $\theta := \theta - \alpha * \frac{\partial J}{\partial \theta}$

randomly choose triplet, easy to meat  $d(A, N) > d(A, P) + \alpha \rightarrow J = 0$ ,  $\rightarrow \frac{\partial J}{\partial \theta} = 0$ , neural network parameter  $\theta$  has no update.

hard triplet:  $d(A, N) \text{ may } < d(A, P) + \alpha \rightarrow J = d(A, P) + \alpha - d(A, N) > 0$ , function of parameter  $\rightarrow \frac{\partial J}{\partial \theta} \neq 0$ , may  $\neq 0$ ,  $\rightarrow$  neural network parameter update more.

Note: naming in deep learning in literature: 'bland' net or deep 'blank'.  
E.g: face net / deep face

#### Choosing the triplets A,P,N

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

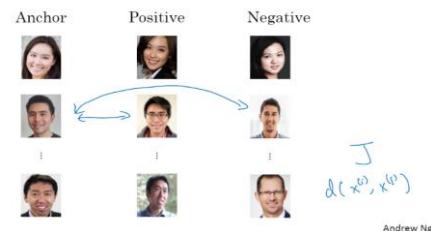
Choose triplets that're "hard" to train on.

$$\begin{aligned} d(A,P) + \alpha &\leq d(A,N) \\ d(A,P) &\approx d(A,N) \end{aligned}$$

Face Net  
Deep Face

Schroff et al., 2015, FaceNet: A unified embedding for face recognition and clustering

#### Training set using triplet loss



#### 4. Summary

This is for the triplet loss and how you can train a neural network for learning an encoding for face recognition.

To train on triplet loss, need map training set to a lot of triples: with an anchor and a positive, and a negative, where the anchor and positive are of the same person but the anchor and negative are of different persons and then use gradient descent to try to minimize the cost function  $J$ , which will have the effect of that propagating to all of the parameters of the neural network in order to learn an encoding so that  $d$  of two images will be small when these two images are of the same person, and they'll be large when these are two images of different persons.

Note:

Face recognition systems especially the large-scale commercial face recognition systems are trained on very large datasets, which are not easy to acquire. yet some of these companies have trained these large networks and posted parameters online, so it might be useful to download someone else's pre-train model, rather than do everything from scratch yourself.

#### 4.4-5 Binary classification\_face verification

Triplet Loss is one good way to learn the parameters of a conv net for face recognition. There's another way to learn these parameters, pose face recognition as a straight binary classification problem.

##### 1. Learning the similarity function

###### 1.1: Treating face recognition as Binary classification

Alternative to the triplet loss for training a system , treating face recognition just as a binary classification problem:

Train a neural network, take a pair of same Siamese Network and compute two image embeddings, maybe multi-dimensional embeddings, and then feed these embeddings to a logistic regression unit to make a prediction.

Target output will be 1 if both of these are the same persons, and 0 if both of these are of different persons.

###### 1.2: Final logistic regression:

output  $y^*$  will be a sigmoid function, applied to some set of features:  $y^* = g(f(x_1), f(x_2))$

Note:

1. Rather than just feeding in image these encodings to sigmoid function, take the differences between the encodings.

$y^* = g(\sum_{k=1}^{128} w_i |f(x_1)_k - f(x_2)_k| + b)$   
a sum over K equals 1 to 128 of the absolute value, taken element wise between the two different encodings.

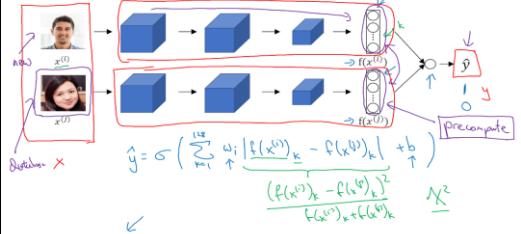
2. could take image embedding  $f(x)$  as features, and logistic regression have additional parameters  $w_i$ , and  $b$  similar to a normal logistic regression unit to each feature fed in, and then train appropriate weight on these e.g. 128 features in order to predict whether or not these two images are of the same person or of different persons. So, this will be one pretty useful way to learn to predict zero or one whether these are the same person or different persons.  $y^* = g(W^* f(x_1) - f(x_2)) + b$

Note:

There are a few other variations on how you can compute  $w_i |f(x_1)_k - f(x_2)_k|$ .  
e.g: another formula - k square formula: squared feature distance /  $|f(x_1)_k - f(x_2)_k|^2$ .

(note: triplet loss, from image feature to output d, only have one hyper parameter a-margin value.  
while binary classification loss , from image feature to output 1/0, have parameter weight  $w_i$  and b-bias could learn by back prop.)

##### Learning the similarity function



##### 1.3 Summarize:

So in this learning formulation, the input  $x$  is a pair of images, and the output  $y$  is either zero or one depending on whether you're inputting a pair of similar or dissimilar images.

Note: Same Siamese Network for the input pair of images.

computational trick:

for the database images, pre-compute their codings and store them. so instead of having to compute this embedding every single time, when the new employee walks in, could just compute this new image coding and compare it to pre-computed database encoding and then to make a prediction  $y^*$ .

As don't need to store the raw images and if have a very large database of employees, you don't need to compute these encodings every single time for every employee database.

This idea of pre- computing can save a significant computation, and both for this type of Siamese Central architecture where you treat face recognition as a binary classification problem, as well as, when you were learning encodings maybe using the Triplet Loss function as described in the last couple of videos.

##### Face verification supervised learning

$x$	$y$	
	1	"Same"
	0	"Different"
	0	
	1	

[Taigman et. al., 2014, DeepFace closing the gap to human level performance]

##### Summary:

For treating face verification supervised learning: create a training set of just pairs of image , instead of triplets images: target label is 1 When these are a pair of pictures of the same person and where the target label is 0, when these are pictures of different persons.

And use different pairs to train the Siamese network using back propagation.

Treating face verification and by extension face recognition as a binary classification problem, this works quite well as well.

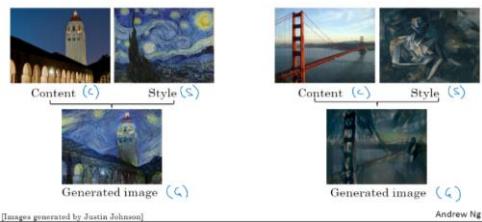
#### 4.4-6 neural style transfer

One of the most fun and exciting applications of ConvNet recently has been Neural Style Transfer.

e.g. take one image-C content image, want this picture recreated in the style of another image-S style image, Neural Style Transfer allows to generate new image-G: have C image content yet painted in S image style.

In order to implement Neural Style Transfer, you need to look at the features extracted by ConvNet at various layers, the shallow and the deeper layers of a ConvNet, get better intuition about whatever all these layers of a ConvNet is really computing.

#### Neural style transfer



Andrew Ng

#### 4.4-7 Deep convnet visualizing

Item purpose:

share visualizations that will help get intuition about what the deeper layers of a ConvNet really are doing. And this will help us think through how to implement neural style transfer as well.

##### 1. Visualizing what a deep network is learning

###### 1.1 Idea of visualizing hidden unit :

>1. Pick a hidden unit in layer l , put training set through neural network, and figure out what is the image that maximizes that particular unit's activation.

(activation result of unit i in layer l, for all training examples m :  $a[l]_i = [a[l]_i, 1, a[l]_i, 2, \dots, a[l]_i, m]$ , find the max few values  $a[l]_i ? \rightarrow$  corresponding image)

e.g. for a trained neural network like a ConvNet, pick a hidden unit in layer l, scan through training sets and find out what are the images or what are the image patches that maximize that unit's activation.

Note:

a hidden unit in layer 1, will see only a relatively small portion of the neural network. So if to visualize- plot what activated unit's activation, it makes sense to plot just a small image patches, because that's all of the image that this particular unit sees.

e.g.

If pick one hidden unit and find 9 input images that maximizes that unit's activation, nine image in the right looks like this particular hidden unit is looking for an edge or a line

###### 1.1 Idea of visualizing hidden unit :

>2. Repeat for other units

e.g.: pick a different hidden unit in layer 1 and do the same thing. Looks like this second unit, represented by another 9 image patches looks like this hidden unit is looking for a line sort of in that portion of its input region.

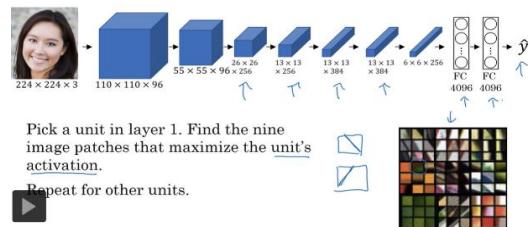
###### 1.2 Note:

1. different hidden units in the same layer, tend to activate in similar image patches: catching similar type of features.

2.Trained hidden units in shallow layers, tend to look for relatively simple features such as edge or a particular shade of color.

3. In the paper, some other more sophisticated ways of visualizing when the ConvNet is running is also used..

#### Visualizing what a deep network is learning



##### Layers visualization example:

> Layer 1 visualization: there's an edge maybe at that angle in the image;

>Layer 2 visualization: looks like it's detecting more complex shapes and patterns.

e.g. looks like it's looking for a vertical texture with lots of vertical lines. This hidden unit looks like it's highly activated when there's a rounder shape to the left part of the image. Here's one that is looking for very thin vertical lines and so on.

And so the features the second layer is detecting are getting more complicated.

>Layer 3 visualization: clearly starting to detect more complex patterns

It looks like there is a hidden unit that seems to respond highly to a rounder shape in the lower left hand portion of the image, maybe. So that ends up detecting a lot of cars, dogs and wonders is even starting to detect people. And this one looks like it is detecting certain textures like honeycomb shapes, or square shapes, this irregular texture. And some of these it's difficult to look at and manually figure out what it's detecting, but it is clearly starting to detect more complex patterns.

##### 2. Visualization deep layers: Layer 5

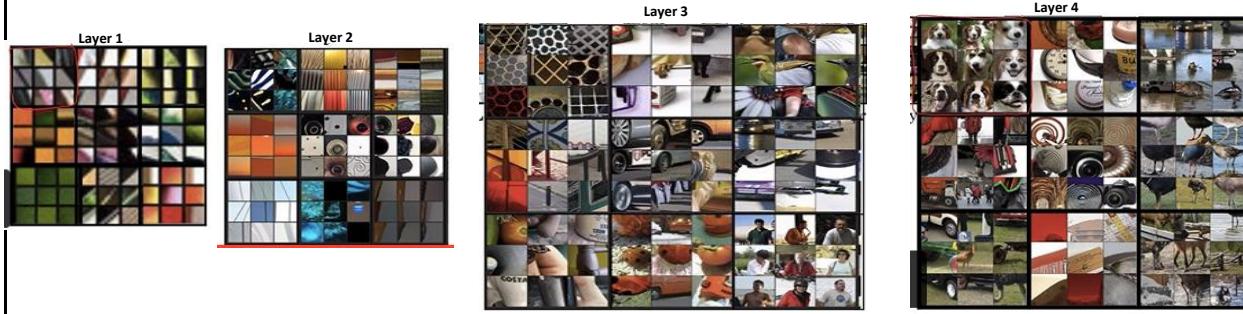
In the deeper layers, a hidden unit will see a larger region of the image.

Extreme cases, each pixel could hypothetically affect the output of these later layers of the neural network. So later units are actually seen larger image patches.

###### 2.1 Example:

Plot the image patches as the same size - see right pic. each layer have 9 units, and each units have 9 images that maximally activate it. Repeated for all layers.

If repeat this procedure, get image patches that maximally activates all the hidden units.



##### Layers visualization example:

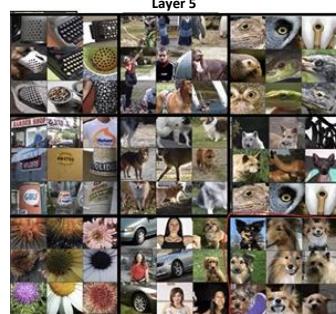
>Layer 4 visualization: features or the patterns is detected is even more complex.

It looks like this has learned almost a dog detector, but all these dogs likewise similar; there is hidden units detecting water; There are ones that look like it's actually detecting the legs of a bird and so on.

>Layer 5 visualization: is detecting even more sophisticated things.

There's also a neuron that seems to be a dog detector, but set of dogs detecting here seems to be more varied. There are ones that seem to be detecting keyboards and things with a keyboard-like texture, although maybe lots of dots against background. Also there are neurons that may be detecting text, and also ones detecting flowers.

-->Have gone a long way from detecting relatively simple things such as edges in layer 1 to textures in layer 2, up to detecting very complex objects in the deeper layers.



#### 4.4-8 cost function\_neural style transfer

Item purpose: To build a Neural Style Transfer system, define a cost function for the generated image and by minimizing this cost function, we can generate the image want.

##### 1. Neural style transfer cost function

**1.1 Neural style transfer:** given a content image C + a style image S, and the goal is to generate a new image G.

##### 1.2 Cost function J(G): purpose

Defined in order to implement neural style transfer, measures how good is a particular generated image G . And use gradient descent to minimize J of G in order to generate this image.

**1.3 Cost function part:**  $J(G) = \alpha J(C, G) + \beta J(S, G)$

Define two parts to this cost function.

>1. Content cost:  $J(C, G)$

definition: function of the content image and generated image.

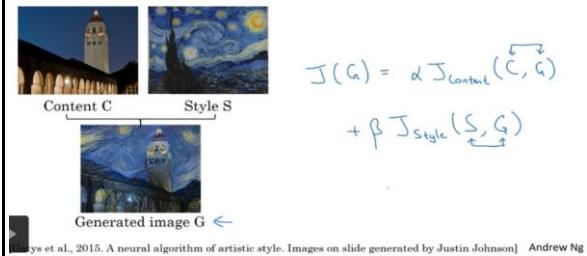
Function: measure the **similarity** of the contents of generated image to the content of the content image C.

>2. Style cost:  $J(S, G)$

definition: function of style image and generated image

Function: measure the **similarity** of the style of generated image to the style of the style image S.

##### Neural style transfer cost function



Gatys et al., 2015. A neural algorithm of artistic style. Images on slide generated by Justin Johnson | Andrew Ng

##### 1.4 Process: Algorithm have to find the cost function J (G) in order to actually generate a new image

>1. Initialize the generated image G randomly

might be  $100 \times 100 \times 3$  or  $500 \times 500 \times 3$  or whatever dimension you want it to be. (then scale dimension to Siamese conv network to generate image embedding?)

>2. Define the cost function J (G):  $J(G) = \alpha J(C, G) + \beta J(S, G)$

>3. use gradient descent to minimize cost , so can update  $G := G - d(J) / d(G)$ .

Note: in this process of running gradient descent and updating G , you're actually updating the pixel values of this image G which is a  $100 \times 100 \times 3$  maybe GRB channel image (image parameters:  $100 \times 100 \times 3 = 30,000$ ).

E.G:

Given content image and style image, initialize G randomly. As run gradient descent, minimize the cost function J(G) slowly through the pixel value so then you get slowly an image that looks more and more like your content image rendered in the style of your style image.

#### 4.4-9 content cost of cost function\_neural style transfer

##### 1. Content cost function

Overall cost function:  $J(G) = \alpha J(C, G) + \beta J(S, G)$

##### 1.1 Content cost process/Idea:

>1. Use hidden layer l to compute the content cost.

LAYER I chosen somewhere in between. It's neither too shallow nor too deep in the neural network.

>>>1. If l is a very small number, like layer 1, then it will really force algorithm to generate image to pixel values very similar to your content image.

>>>2. If use a very deep layer, then it's just asking, "Well, if there is a dog in your content image, then make sure there is a dog somewhere in the generated image."

Note: usually choose layer to be somewhere in the middle of the layers of the neural network, neither too shallow nor too deep.

>2. Use a pre-trained ConvNet () (Siamese network??), maybe a VGG network, measure similarity in content of content image C and generated image G.

>3. Let  $a[l]_C, a[l]_G$  be the activations of layer l on these two images, on the images C and G.

>4. if these two activations  $a[l]_C, a[l]_G$  are similar, then that would seem to imply that both images have similar content.

##### 1.2 Content cost definition:

Take the element-wise difference between these hidden unit activations in layer l, between when you pass in the content image compared to when you pass in the generated image, and take that squared.

$$J(C, G)_{content} = 1/2 * \| a[l]_C - a[l]_G \| ^2$$

##### Note:

1. could have a normalization constant J formular or not., so it's just 1/2, or something else.  
It doesn't really matter since this can be adjusted as well by this hyperparameter a. S

2. Both  $a[l]_C & a[l]_G$  have been unrolled into vectors.

-->content cost is really just the element-wise sum of squared differences between these two activations in layer l, between the images in C and G.

3. when later perform gradient descent on J (G) to try to find a value of G, so that the overall cost is low, this will incentivize /force the algorithm to find an image G, so that its hidden layer activations are similar to that layer activations of the content image.

#### 4.4-10 style cost of cost function\_neural style transfer

## 1. Image style:

Example: input image , use conv net to compute features that there are in different layers.

>1. Chosen some layer  $l$ , to define the measure of the style of an image.

>2. Define image style as the correlation between activations across different channels in this layer  $l$  activation. could take different channels as different hidden units in layer  $l$ , image style defined in the layer  $l$ , means each hidden units value relation with other hidden units.

### 1.1 Style intuition:

for layer  $l$  activation:  $n_h \times n_w \times n_c$  and check how correlated are the activations across different channels.

#### Example:

has five channels, look at the first two channels and see how correlated are activations in these first two channels.

#### >1. Idea:

for the same position in channel 1, and channel 2, get a pair of numbers and see when you look across all of these positions, all of these  $n_h \times n_w$  positions, how correlated are these two numbers.

This capture style: take each channel as one hidden unit in layer  $l$ .

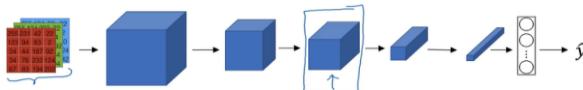
e.g. in the layer  $l$  visualizations, the red channel's max activation image is vertical texture-->red channels correspond to neurons trying to figure out if there's this little vertical texture in a particular position in the  $n_h \times n_w$ .

This yellow second channel corresponds to this neuron, which is vaguely looking for orange colored patches.

If these two channels are highly correlated, that means whenever part of first channel max activation image has this type of subtle vertical texture, that part of second channel max activation image will probably have these orange-ish tint.

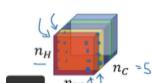
If they are uncorrelated, it means that whenever there is this vertical texture, it's probably won't have that orange-ish tint.

## Meaning of the “style” of an image



Say you are using layer  $l$ 's activation to measure “style.”

Define style as correlation between activations across channels.

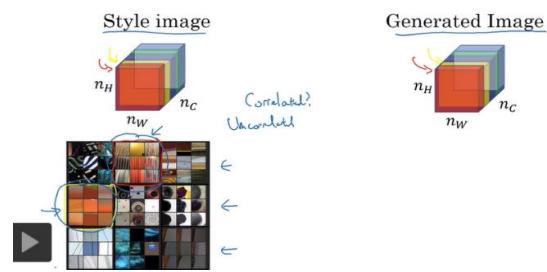


How correlated are the activations across different channels?

[Gatys et al., 2015. A neural algorithm of artistic style]

Andrew

## Intuition about style of an image



### 2. Style matrix:

layer :  $l$ ; channel:  $k$ ;  $c$  Let  $a_{i,j,k}^{[l]} = \text{activation at } (i,j,k)$ .  $G^{[l](s)}$  is  $n_c^{[l]} \times n_c^{[l]}$  matrix for layer  $l$ . activation  $a_{i,j,k}^{[l]}$ :  $a_{i,j,k}^{[l]}$

#### >1. Channel $k, k'$ correlation for style image and generated image: $G^{[l](s), k, k'}$

Measure how correlated are the activations in channel  $k$  compared to the activations in channel  $k'$ . Take a pair of activation values in same position of channel  $k, k'$  and multiply them  $a_{i,j,k}^{[l]} \cdot a_{i,j,k'}^{[l]}$  and then sum over all  $n_h \times n_w$  positions'

$$\begin{aligned} \rightarrow G^{[l](s)} &= \sum_{i=1}^{n_h} \sum_{j=1}^{n_w} a_{i,j,k}^{[l]} a_{i,j,k'}^{[l]} \\ \rightarrow G^{[l](G)} &= \sum_{i=1}^{n_h} \sum_{j=1}^{n_w} a_{i,j,k}^{[l]} a_{i,j,k'}^{[l]} \end{aligned}$$

#### >2. Style matrix for one image: $G^{[l]S}, G^{[l]G}$

take layer  $l$  activations to compute image style:  $G^{[l]S}$  is  $n_c^{[l]} \times n_c^{[l]}$  dimension.

take Channel  $k, k'$  correlation in layer  $l$  activation, go through all channels in layer  $l$  activation. As have  $n_c$  channels so have an  $n_c \times n_c$  dimensional matrix in order to measure how correlated each pair of them-channel is.

(note: style matrix is symmetric:

$[G_{11}, G_{12}, G_{13}, \dots, G_{1n_c}, \dots, G_{nn_c}]$

$[G_{21}, G_{22}, G_{23}, \dots, G_{2n_c}, \dots, G_{nn_c}]$

$\dots$

$G_{n_c 1}, G_{n_c 2}, G_{n_c 3}, \dots, G_{n_c n_c}, \dots, G_{nn_c}$

### 3. Style cost function in layer $l$ :

$$\begin{aligned} J_{\text{style}}^{[l]}(S, G) &= \frac{1}{C} \| G^{[l](s)} - G^{[l](G)} \|_F^2 \\ &= \frac{1}{(2n_h^{[l]} n_w^{[l]} n_c^{[l]})^2} \sum_{k=1}^{n_c^{[l]}} \sum_{k'=1}^{n_c^{[l]}} (G_{kk'}^{[l](s)} - G_{kk'}^{[l](G)})^2 \end{aligned}$$

**Definition:** the difference between these two matrices,  $G^{[l]S}, G^{[l]G}$ : the sum of squares of the element wise differences between these two matrices: square of  $(G^{[l]S} - G^{[l]G})$ , and sum over  $k, k'$ .  $= \text{norm}(G^{[l]S} - G^{[l]G})$

#### Note:

1. paper authors actually used normalization constant  $1 / (2 * n_h * n_w * n_c)^{1/2}$ .

But a normalization constant doesn't matter that much because this causes multiplied by some hyperparameter  $\beta$  anyway.

#### Summary:

This style cost function defined using layer  $l$ , cost function is basically the Frobenius norm between the two style matrices computed on the image  $s$  and on the image  $G$ , Frobenius on squared and divide by normalization constants, which isn't that important.

And it turns out that you get more visually pleasing results if use the style cost function from multiple different layers.

### >1. Idea:

#### Summarize:

The correlation tells which of these high level texture components tend to occur or not occur together in part of an image and the degree of correlation measuring how often these different high level features, such as vertical texture or this orange tint or other things as well, how often they occur and how often they occur together and don't occur together in different parts of an image.

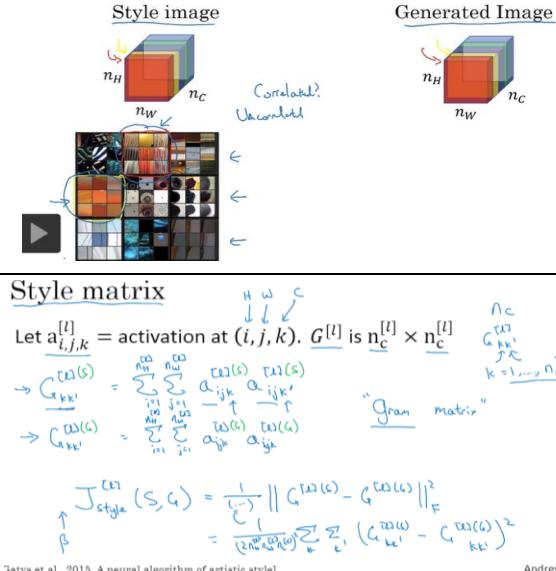
#### Generated image:

If use the degree of correlation between channels as a measure of the style, then measure the degree in generated image.

This first channel is correlated or uncorrelated with the second channel and that will tell you in the generated image how often this type of vertical texture occurs or doesn't occur with this orange-ish tint.

This gives a measure of how similar is the style of the generated image to the style of the input style image.

### Intuition about style of an image



#### Note:

1. if both of these activation tend to be large together,  $G^{[l]kk'}$  will be large. if they are uncorrelated, then  $G^{[l]kk'}$  might be small.  
(how to verify this is good representation of image style?)

2. This is actually the unnormalized cross of the areas because we're not subtracting out the mean and this is just multiplied by these elements directly.  
(no need to normalize all channels in layer  $l$ ?)

( $G^{[l]kk'}$ : kind of like conv computation: filter same size as input volume (only two channel  $k, k'$ ) here)-element wise filter and channel values and then sum, result in one value representing image feature captured by filter;  
while  $G^{[l]kk'}$  no filter, element wise channel's own values and then sum, representing image styles in this layer.).

### Style cost function

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

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

$$J(G) = \alpha J_{\text{content}}(I, G) + \beta J_{\text{style}}^{[l]}(S, G)$$

[Gatys et al., 2015. A neural algorithm of artistic style]

Andrew Ng

#### 4. Overall style cost function $J(S, G)$ :

$$J_{\text{style}}(S, G) = \sum_i \lambda^i J_{\text{style}}^{(i)}(S, G)$$

**Definition:** sum over all the different layers of the style cost function for that layer  $J[L](S, G)$ , weighted by some set of additional hyperparameters,  $\lambda[i]$ . So it allows to use different layers in a neural network: both the early ones, which measure relatively simpler low level features like edges as well as some later layers, which measure high level features and cause neural network to take both low level and high level correlations into account when computing style.

And, in the following exercise, you gain more intuition about what might be reasonable choices for this type of parameter  $\lambda[i]$  as well.

#### 5. Overall cost function: $J(G) = \alpha * J(C, G) + \beta * J(S, G)$

Could use gradient descent or a more sophisticated optimization algorithm if you want in order to try to find an image  $G$  that to minimize this cost function  $J(G)$ . and if you do that you'll be able to generate some pretty nice artwork.

### 4.4-11: 1D & 3D generalizations of models.

Have learned a lot about ConvNets, everything ranging from the architecture of the ConvNet to how to use it for image recognition, to object detection, to face recognition and neural-style transfer. Yet most of the discussion has focused on images, on sort of 2D data. It turns out that many of these ideas have learned about also apply, not just to 2D images but also to 1D data as well as to 3D data.

#### 1. convolution in 2D & 1D

##### >1. 2D convolution:

e.g: image  $14 \times 14$ , filter  $5 \times 5$ , 16 filter  $\rightarrow$  result  $10 \times 10 \times 16$   
or image  $14 \times 14 \times 3$ , -filter  $5 \times 5 \times 3$ , 16 filter  $\rightarrow$  result  $10 \times 10 \times 16$ .

Same method could apply on 1D data.

##### >2. 1D convolution:

Method for 2D convo.  $(n_h, h_w) * n_c$ , could same used in 1D  $(n_h) * n_c$

just remmove  $n_w$  dimension, same calcualtion. Just Take  $n_w=1$ .

e.g : EKG signal: heart beating vs time series. different postion (-dfferent channel) heart beating vs time

Note: filter channel still need to be same as 1D image channel.

##### >2. 1D convolution:

E.g: EKG signal- electrocardiogram. heart beating vs time, showing the voltage at each instant in time.

input is 14 (in sequence),  $\rightarrow$  convolve this with a 1 dimensional filter, rather than the  $5 \times 5$ , just have 5 ( $x1$ ) dimensional filter.

In 2D data convolution, apply the same  $5 \times 5$  feature detector across at different positions throughout the image. And end up with your  $10 \times 10$  output.

In 1D data convolution: apply 5 dimensional filter similarly in lots of different positions throughout this 1D signal.will end up 10 dimensional output ( $10 \times 1$ ).

##### Note:

If have multiple channels in 1D data, filter have same channel number.

e.g: input  $10 \times 16$  dimensional input, filter -5 dimensional with 16 channels, 32 filters, then end up  $6 \times 32$ ,

So generalization from 2D to 1D, all of 2D convolution ideas apply also to 1D data, having the same feature detector apply to a variety of positions

#### 2. 3D Convolution:

##### >1. Input: 3D block, a three dimensional input

e.g: CT scan, three dimensional model of body. so as scan through a CT scan, you can look at different slices of the human torso to see how they look .

##### >2. 3D convolution idea:

Same idea of 2D convolution, to apply a ConvNet to detect features in this three dimensional image  
e.g: input  $14 \times 14 \times 14$  x32 32 channel, filer  $5 \times 5 \times 5$ , 32 channel, 16 filter, end up with a  $6 \times 6 \times 6$  volume across 32 channels.

So 3D data can also be learned on, sort of directly using a three dimensional ConvNet. And what these filters do is really detect features across your 3D data

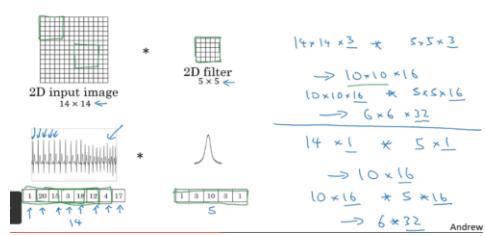
e.g2: 3D movie , different slices 3D data in time through a movie. Could use this to detect motion or people taking actions in movies.

##### >3. Note:

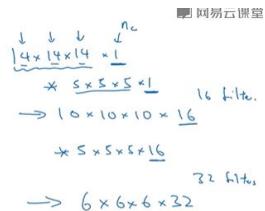
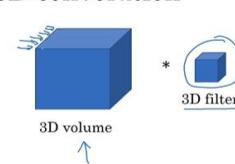
1. 3D input convolution: filter is also 3D, yet filter dimension not same as input dimension esp. in the 3rd direction:  
only filter channel is = input channel.

Apply same volume filter in different position of horizontal, vertical and volume directions.

### Convolutions in 2D and 1D



### 3D convolution



## 5.1 - Convolutional layer backward pass

### 5.1.1 - Computing dA:

This is the formula for computing  $dA$  with respect to the cost for a certain filter  $W_c$  and a given training example:

$$dA+ = \sum_{h=0}^{n_H} \sum_{w=0}^{n_W} W_c \times dZ_{hw} \quad (1)$$

Where  $W_c$  is a filter and  $dZ_{hw}$  is a scalar corresponding to the gradient of the cost with respect to the output of the conv layer Z at the  $h$ th row and  $w$ th column (corresponding to the dot product taken at the  $i$ th stride left and  $j$ th stride down). Note that at each time, we multiply the same filter  $W_c$  by a different  $dZ$  when updating  $dA$ . We do so mainly because when computing the forward propagation, each filter is dotted and summed by a different  $a\_slice$ . Therefore when computing the backprop for  $dA$ , we are just adding the gradients of all the  $a\_slices$ .

In code, inside the appropriate for-loops, this formula translates into:

```
1 da_prev_pad[vert_start:vert_end, horiz_start:horiz_end, :] += W[:, :, :, c] * dZ[i, h, w, c]
```

### 5.1.2 - Computing dW:

This is the formula for computing  $dW_c$  ( $dW_c$  is the derivative of one filter) with respect to the loss:

$$dW_c+ = \sum_{h=0}^{n_H} \sum_{w=0}^{n_W} a_{slice} \times dZ_{hw} \quad (2)$$

Where  $a_{slice}$  corresponds to the slice which was used to generate the activation  $Z_{ij}$ . Hence, this ends up giving us the gradient for  $W$  with respect to that slice. Since it is the same  $W$ , we will just add up all such gradients to get  $dW$ .

In code, inside the appropriate for-loops, this formula translates into:

```
1 dW[:, :, :, c] += a_slice * dZ[i, h, w, c]
```

### 5.1.3 - Computing db:

This is the formula for computing  $db$  with respect to the cost for a certain filter  $W_c$ :

$$db = \sum_h \sum_w dZ_{hw} \quad (3)$$

As you have previously seen in basic neural networks,  $db$  is computed by summing  $dZ$ . In this case, you are just summing over all the gradients of the conv output ( $Z$ ) with respect to the cost.

In code, inside the appropriate for-loops, this formula translates into:

```
1 db[:, :, :, c] += dZ[i, h, w, c]
```

## 5.2 Pooling layer - backward pass

Next, let's implement the backward pass for the pooling layer, starting with the MAX-POOL layer. Even though a pooling layer has no parameters for backprop to update, you still need to backpropagation the gradient through the pooling layer in order to compute gradients for layers that came before the pooling layer.

### 5.2.1 Max pooling - backward pass

Before jumping into the backpropagation of the pooling layer, you are going to build a helper function called `create_mask_from_window()` which does the following:

$$X = \begin{bmatrix} 1 & 3 \\ 4 & 2 \end{bmatrix} \rightarrow M = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} \quad (4)$$

As you can see, this function creates a "mask" matrix which keeps track of where the maximum of the matrix is. True (1) indicates the position of the maximum in  $X$ , the other entries are False (0). You'll see later that the backward pass for average pooling will be similar to this but using a different mask.

**Exercise:** Implement `create_mask_from_window()`. This function will be helpful for pooling backward.

Hints:

- `np.max()` may be helpful. It computes the maximum of an array.
- If you have a matrix  $X$  and a scalar  $x$ : `A = (X == x)` will return a matrix  $A$  of the same size as  $X$  such that:

```
1 A[i,j] = True if X[i,j] = x
2 A[i,j] = False if X[i,j] != x
```

- Here, you don't need to consider cases where there are several maxima in a matrix.

```
def create_mask_from_window(x):
    """
    Creates a mask from an input matrix x, to identify the max entry of x.

    Arguments:
    x -- Array of shape (f, f)

    Returns:
    mask -- Array of the same shape as window, contains a True at the position corresponding to the max entry of x.
    """

    ### START CODE HERE ### (≈1 line)
    mask = (x == np.max(x))
    ### END CODE HERE ###

    return mask
```

```
np.random.seed(1)
x = np.random.randn(2,3)
mask = create_mask_from_window(x)
print('x = ', x)
print("mask = ", mask)
```

```
x = [[ 1.62434536 -0.61175641 -0.52817175]
 [-1.07296862  0.86540763 -2.3015387 ]]
mask = [[ True False False]
 [False False False]]
```

Why do we keep track of the position of the max? It's because this is the input value that ultimately influenced the output, and therefore the cost. Backprop is computing gradients with respect to the cost, so anything that influences the ultimate cost should have a non-zero

gradient. So, backprop will “propagate” the gradient back to this particular input value that had influenced the cost.

### 5.2.2 - Average pooling - backward pass

In max pooling, for each input window, all the “influence” on the output came from a single input value—the max. In average pooling, every element of the input window has equal influence on the output. So to implement backprop, you will now implement a helper function that reflects this.

For example if we did average pooling in the forward pass using a 2x2 filter, then the mask you'll use for the backward pass will look like:

$$dZ = 1 \rightarrow dZ = \begin{bmatrix} 1/4 & 1/4 \\ 1/4 & 1/4 \end{bmatrix} \quad (5)$$

This implies that each position in the  $dZ$  matrix contributes equally to output because in the forward pass, we took an average.

**Exercise:** Implement the function below to equally distribute a value  $dz$  through a matrix of dimension shape. Hint

```
1 def distribute_value(dz, shape):
2     """
3         Distributes the input value in the matrix of dimension shape
4
5     Arguments:
6     dz -- input scalar
7     shape -- the shape (n_H, n_W) of the output matrix for which we want to distribute the value of dz
8
9     Returns:
10    a -- Array of size (n_H, n_W) for which we distributed the value of dz
11    """
12
13    ### START CODE HERE ###
14    # Retrieve dimensions from shape (≈1 line)
15    (n_H, n_W) = shape
16
17    # Compute the value to distribute on the matrix (≈1 line)
18    average = dz / (n_H * n_W)
19
20    # Create a matrix where every entry is the "average" value (≈1 line)
21    a = average * np.ones(shape)
22    ### END CODE HERE ###
23
24    return a
```

```
1 a = distribute_value(2, (2,2))
```

```
2 print('distributed value =', a)
```

```
1 distributed value = [[ 0.5  0.5]
2   [ 0.5  0.5]]
```











Two main types of blocks are used in a ResNet, depending mainly on whether the input/output dimensions are same or different. You are going to implement both of them.

## 2.1 - The identity block

The identity block is the standard block used in ResNets, and corresponds to the case where the input activation (say  $a^{[l]}$ ) has the same dimension as the output activation (say  $a^{[l+2]}$ ). To flesh out the different steps of what happens in a ResNet's identity block, here is an alternative diagram showing the individual steps:

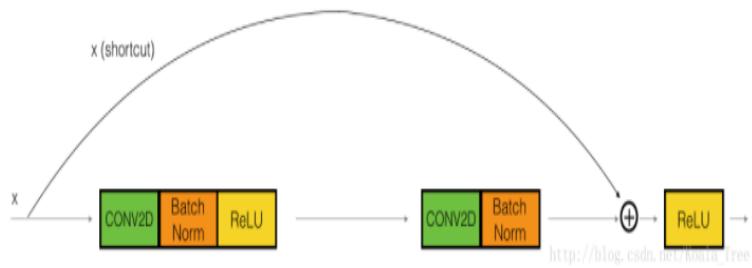
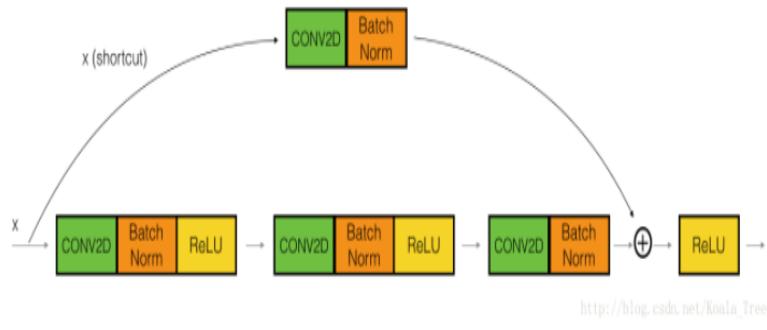


Figure 3 : Identity block. Skip connection “skips over” 2 layers.

## 2.2 - The convolutional block

You've implemented the ResNet identity block. Next, the ResNet “convolutional block” is the other type of block. You can use this type of block when the input and output dimensions don't match up. The difference with the identity block is that there is a CONV2D layer in the shortcut path:



**Figure 4 : Convolutional block**

The CONV2D layer in the shortcut path is used to resize the input  $x$  to a different dimension, so that the dimensions match up in the final addition needed to add the shortcut value back to the main path. (This plays a similar role as the matrix  $W_s$  discussed in lecture.) For example, to reduce the activation dimension's height and width by a factor of 2, you can use a  $1 \times 1$  convolution with a stride of 2. The CONV2D layer on the shortcut path does not use any non-linear activation function. Its main role is to just apply a (learned) linear function that reduces the dimension of the input, so that the dimensions match up for the later addition step.





## 2 - YOLO

YOLO ("you only look once") is a popular algorithm because it achieves high accuracy while also being able to run in real-time. This algorithm "only looks once" at the image in the sense that it requires only one forward propagation pass through the network to make predictions. After non-max suppression, it then outputs recognized objects together with the bounding boxes.

### 2.1 - Model details

First things to know:

- The **input** is a batch of images of shape  $(m, 608, 608, 3)$
- The **output** is a list of bounding boxes along with the recognized classes. Each bounding box is represented by 6 numbers  $(p_c, b_x, b_y, b_h, b_w, c)$  as explained above. If you expand  $c$  into an 80-dimensional vector, each bounding box is then represented by 85 numbers.

We will use 5 anchor boxes. So you can think of the YOLO architecture as the following: IMAGE  $(m, 608, 608, 3) \rightarrow$  DEEP CNN  $\rightarrow$  ENCODING  $(m, 19, 19, 5, 85)$ .

Lets look in greater detail at what this encoding represents.



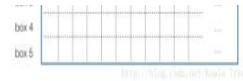


Figure 2 : Encoding architecture for YOLO

If the center/midpoint of an object falls into a grid cell, that grid cell is responsible for detecting that object.

Since we are using 5 anchor boxes, each of the 19 x19 cells thus encodes information about 5 boxes. Anchor boxes are defined only by their width and height.

For simplicity, we will flatten the last two dimensions of the shape (19, 19, 5, 85) encoding. So the output of the Deep CNN is (19, 19, 425).

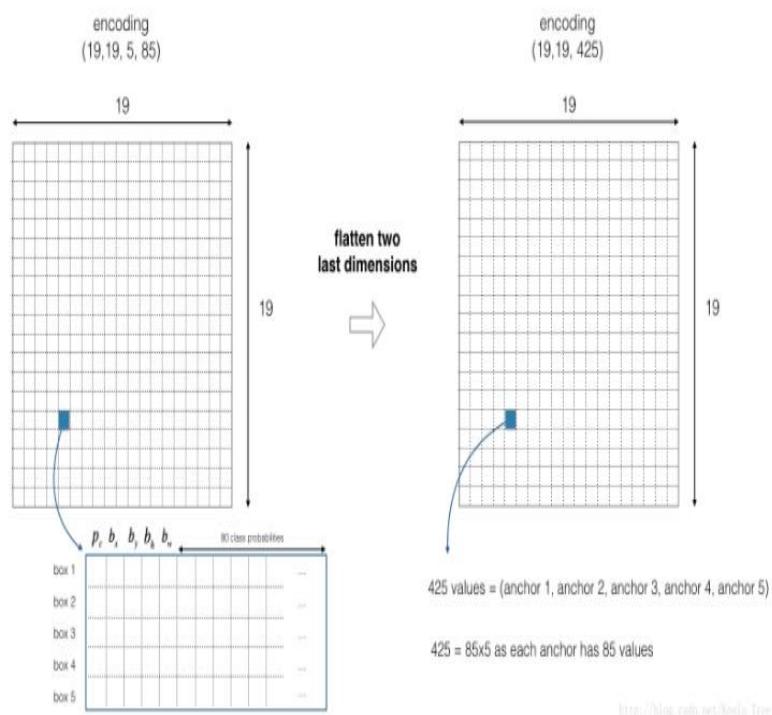


Figure 3 : Flattening the last two last dimensions

Now, for each box (of each cell) we will compute the following elementwise product and extract a probability that the box contains a certain class.





## 3 - Neural Style Transfer

We will build the NST algorithm in three steps:

- Build the content cost function  $J_{content}(C, G)$
- Build the style cost function  $J_{style}(S, G)$
- Put it together to get  $J(G) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$ .

### 3.1 - Computing the content cost

In our running example, the content image C will be the picture of the Louvre Museum in Paris. Run the code below to see a picture of the Louvre.

\* 3.1.1 - How do you ensure the generated image G matches the content of the image C?\*

As we saw in lecture, the earlier (shallower) layers of a ConvNet tend to detect lower-level features such as edges and simple textures and the later (deeper) layers tend to detect higher-level features such as more complex textures as well as object classes.

We would like the “generated” image G to have similar content as the input image C. Suppose you have chosen some layer’s activations to represent the content of an image. In practice, you’ll get the most visually pleasing results if you choose a layer in the middle of the network—neither too shallow nor too deep. (After you have finished this exercise, feel free to come back and experiment with using different layers, to see how the results vary.)

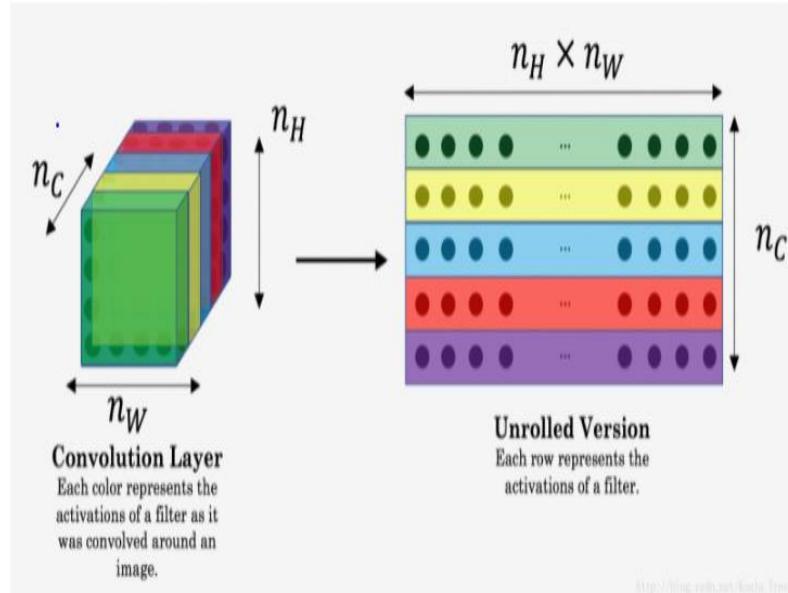
So, suppose you have picked one particular hidden layer to use. Now, set the image C as the input to the pretrained VGG network, and run forward propagation. Let  $a^{(C)}$  be the hidden layer activations in the layer you had chosen. (In lecture, we had written this as  $a^{[l]}(C)$  but here we’ll drop the superscript  $[l]$  to simplify the notation.) This will be a  $n_H \times n_W \times n_C$  tensor. Repeat this process with the image G: Set G as the input, and run forward propagation. Let

$$a^{(G)}$$

be the corresponding hidden layer activation. We will define as the content cost function as:

$$J_{content}(C, G) = \frac{1}{4 \times n_H \times n_W \times n_C} \sum_{\text{all entries}} (a^{(C)} - a^{(G)})^2 \quad (1)$$

Here,  $n_H$ ,  $n_W$  and  $n_C$  are the height, width and number of channels of the hidden layer you have chosen, and appear in a normalization term in the cost. For clarity, note that  $a^{(C)}$  and  $a^{(G)}$  are the volumes corresponding to a hidden layer's activations. In order to compute the cost  $J_{content}(C, G)$ , it might also be convenient to unroll these 3D volumes into a 2D matrix, as shown below. (Technically this unrolling step isn't needed to compute  $J_{content}$ , but it will be good practice for when you do need to carry out a similar operation later for computing the style const  $J_{style}$ .)



**Exercise:** Compute the “content cost” using TensorFlow.

**Instructions:** The 3 steps to implement this function are:

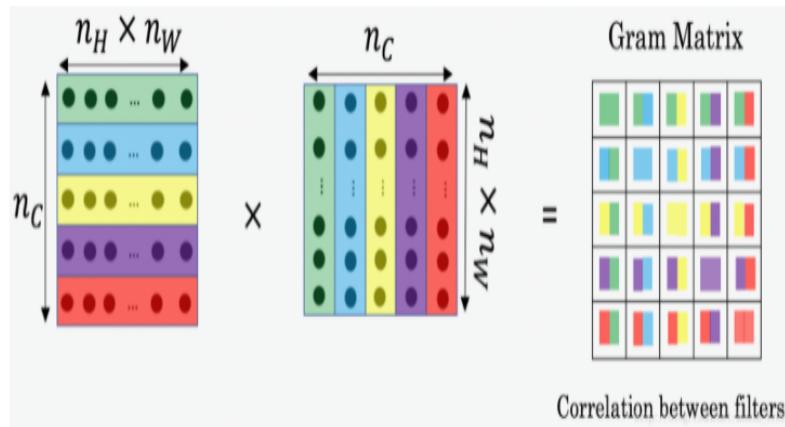
1. Retrieve dimensions from  $a_G$ :
  - To retrieve dimensions from a tensor  $X$ , use: `x.get_shape().as_list()`
2. Unroll  $a_C$  and  $a_G$  as explained in the picture above
  - If you are stuck, take a look at [Hint1](#) and [Hint2](#).
3. Compute the content cost:
  - If you are stuck, take a look at [Hint3](#), [Hint4](#) and [Hint5](#).

### 3.2.1 - Style matrix

The style matrix is also called a “Gram matrix.” In linear algebra, the Gram matrix  $G$  of a set of vectors  $(v_1, \dots, v_n)$  is the matrix of dot products, whose entries are  $G_{ij} = v_i^T v_j = np.dot(v_i, v_j)$ . In other words,  $G_{ij}$  compares how similar  $v_i$  is to  $v_j$ : If they are highly similar, you would expect them to have a large dot product, and thus for  $G_{ij}$  to be large.

Note that there is an unfortunate collision in the variable names used here. We are following common terminology used in the literature, but  $G$  is used to denote the Style matrix (or Gram matrix) as well as to denote the generated image  $G$ . We will try to make sure which  $G$  we are referring to is always clear from the context.

In NST, you can compute the Style matrix by multiplying the “unrolled” filter matrix with their transpose:



The result is a matrix of dimension  $(n_C, n_C)$  where  $n_C$  is the number of filters. The value  $G_{ij}$  measures how similar the activations of filter  $i$  are to the activations of filter  $j$ .

One important part of the gram matrix is that the diagonal elements such as  $G_{ii}$  also measures how active filter  $i$  is. For example, suppose filter  $i$  is detecting vertical textures in the image. Then  $G_{ii}$  measures how common vertical textures are in the image as a whole: If  $G_{ii}$  is large, this means that the image has a lot of vertical texture.

By capturing the prevalence of different types of features ( $G_{ii}$ ), as well as how much different features occur together ( $G_{ij}$ ), the Style matrix  $G$  measures the style of an image.

### 3.2.2 - Style cost

After generating the Style matrix (Gram matrix), your goal will be to minimize the distance between the Gram matrix of the “style” image  $S$  and that of the “generated” image  $G$ . For now, we are using only a single hidden layer  $a^{[l]}$ , and the corresponding style cost for this layer is defined as:

layer is defined as:

$$J_{style}^{[l]}(S, G) = \frac{1}{4 \times n_C^2 \times (n_H \times n_W)^2} \sum_{i=1}^{n_C} \sum_{j=1}^{n_C} (G_{ij}^{(S)} - G_{ij}^{(G)})^2 \quad (2)$$

where  $G^{(S)}$  and  $G^{(G)}$  are respectively the Gram matrices of the “style” image and the “generated” image, computed using the hidden layer activations for a particular hidden layer in the network.

**Exercise:** Compute the style cost for a single layer.

**Instructions:** The 3 steps to implement this function are:

1. Retrieve dimensions from the hidden layer activations  $a_G$ :

- To retrieve dimensions from a tensor  $X$ , use: `X.get_shape().as_list()`

2. Unroll the hidden layer activations  $a_S$  and  $a_G$  into 2D matrices, as explained in the picture above.

- You may find [Hint1](#) and [Hint2](#) useful.

3. Compute the Style matrix of the images  $S$  and  $G$ . (Use the function you had previously written.)

4. Compute the Style cost:

- You may find [Hint3](#), [Hint4](#) and [Hint5](#) useful.

### 3.2.3 Style Weights

So far you have captured the style from only one layer. We’ll get better results if we “merge” style costs from several different layers.

After completing this exercise, feel free to come back and experiment with different weights to see how it changes the generated image  $G$ . But for now, this is a pretty reasonable default:

```
1 STYLE_LAYERS = [  
2     ('conv1_1', 0.2),  
3     ('conv2_1', 0.2),  
4     ('conv3_1', 0.2),  
5     ('conv4_1', 0.2),  
6     ('conv5_1', 0.2)]
```

You can combine the style costs for different layers as follows:

$$J = \alpha \alpha + \sum \lambda[l] J[l]$$

$$J_{style}(S, G) = \sum_l \lambda^{[l]} J_{style}^{[l]}(S, G)$$

where the values for  $\lambda^{[l]}$  are given in `STYLE_LAYERS`.

We've implemented a `compute_style_cost(...)` function. It simply calls your `compute_layer_style_cost(...)` several times, and weights their results using the values in `STYLE_LAYERS`. Read over it to make sure you understand what it's doing.

### 3.3 - Defining the total cost to optimize

Finally, let's create a cost function that minimizes both the style and the content cost. The formula is:

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

## 4 - Solving the optimization problem

Finally, let's put everything together to implement Neural Style Transfer!

Here's what the program will have to do:

1. Create an Interactive Session
2. Load the content image
3. Load the style image
4. Randomly initialize the image to be generated
5. Load the VGG16 model
6. Build the TensorFlow graph:
  1. Run the content image through the VGG16 model and compute the content cost
  2. Run the style image through the VGG16 model and compute the style cost

3. Compute the total cost
  4. Define the optimizer and the learning rate
7. Initialize the TensorFlow graph and run it for a large number of iterations, updating the generated image at every step.











