

NVIDIA Deep Learning



Fundamentals



Raymond Ptucha,
Rochester Institute of Technology, NVIDIA Deep
Learning Institute



ICFHR 2018

*The 16th International Conference on
Frontiers in Handwriting Recognition*

August 5 - 8, 2018 • **Niagara Falls, USA**

Fair Use Agreement

This agreement covers the use of all slides in this document, please read carefully.

- You may freely use these slides, if:
 - You send me an email telling me the conference/venue/company name in advance, and which slides you wish to use.
 - You receive a positive confirmation email back from me.
 - My name (Ptucha) appears on each slide you use.

(c) Raymond Ptucha, rwpeec@rit.edu

Agenda

- Part I- Intuition and Theory
 - 9:00-9:45am: Introduction
 - 9:45-10:30am: Convolutional Neural Networks
- 10:30-10:45pm: Break
- Part II- Hands on
 - 10:45am-Noon: Hands-on exercises

Navigating to Qwiklabs

- Have you registered for NVIDIA account yet?
 1. Navigate to: <https://nvlabs.qwiklab.com>
 2. Login or create a new account
- Select event:
 - ICFHC CV Ambassador Workshop
- Then select class:
 - Image Classification with DIGITS
- You can return to class for up to 60 days

Machine Learning



- Machine learning is giving computers the ability to analyze, generalize, think/reason/behave like humans.
- Machine learning is transforming medical research, financial markets, international security, and generally making humans more efficient and improving quality of life.
- Inspired by the mammalian brain, deep learning is machine learning on steroids- bigger, faster, better!

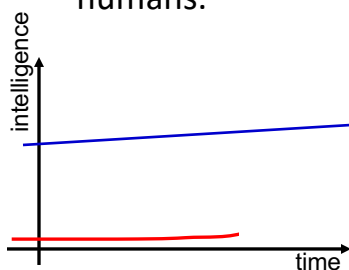


Ptucha '18

6

The point of Singularity

- The point of singularity is when computers become smarter than humans.



— Evolution of biology
— Advancement of technology

Ptucha '18

7

Unleashing of Intelligence



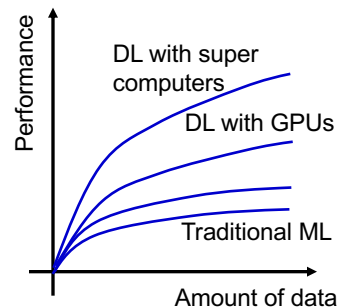
- Machines will slowly match, then quickly surpass human capabilities.
- Today it is exciting/scary/fun to drive next to an autonomous car.
- Tomorrow it may be considered irresponsible for a human to relinquish control from a car that has faster reaction times, doesn't drink/text/get distracted/tired, and is communicating with surrounding vehicles and objects.

Ptucha '18

9

Why is AI (Deep Learning) Just Now Becoming Practical in Many Day-to-Day Situations?

- Availability of data;
- Sustained advances in hardware capabilities (including GPUs running machine learning workloads);
- Omnipresent connectivity;
- Lower cost and power consumption;
- Sustained advances in algorithmic learning techniques.



Hot trend:
High performance
architecture experts
teaming up with deep
learning experts

Ptucha '18

10

2017: The Year of AI:

The Wall Street Journal, Forbes, and Fortune



NEC Face Recognition

Turn: -0.7932
 Engine: 0.99999
 Fitness: 9.78920



Alex Simmons
 @AlexIGNUK

SONY Playstation Virtual Reality

Generation: 1

Evolutionary Reinforcement Learning

Ptucha '18

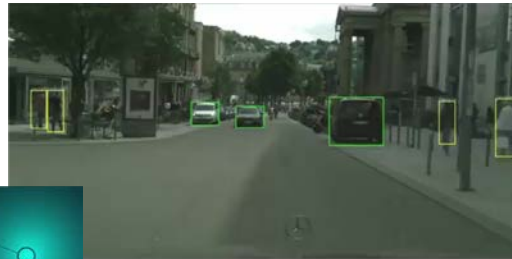
13

2017: The Year of AI:

The Wall Street Journal, Forbes, and Fortune



DeepBach



NVIDIA Autonomous Car
 Detection & Segmentation



YOLO v2 Object Detection

Ptucha '18

14

Some Things to Look for in 2018

-
-
-
-

CelebA-HQ
1024 x 1024

Progressive growing

CelebA-HQ
1024 x 1024

Latent space interpolations

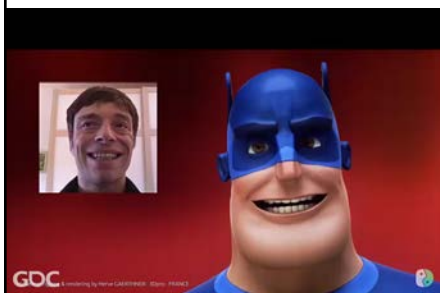
http://research.nvidia.com/sites/default/files/pubs/2017-10_Progressive-Growing-of/karras2017gan-paper.pdf

Ptucha '18

19

Some Things to Look for in 2018

Faceshift GDC



Apple iPhone X, Animoji Yourself



Ptucha '18

20

Some Things to Look for in 2018

NVIDIA DRIVE
Autonomous Vehicle Platform
October 10, 2017



NVIDIA Drive

Ptucha '18

21

The Human Brain



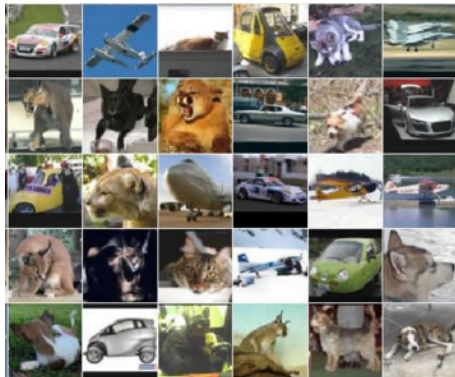
- We've learned more about the brain in the last 5 years than we have learned in the last 5000 years!
- It controls every aspect of our lives, but we still don't understand exactly how it works.

Ptucha '18

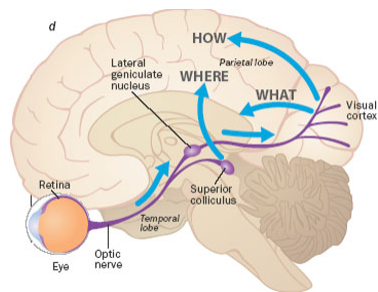
22

The Brain on Pattern Recognition

- Airplane, Cat, Car, Dog



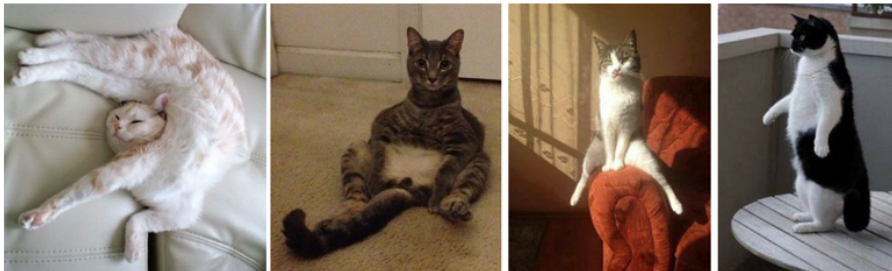
STL-10 dataset



<http://thebraingeek.blogspot.com/2012/08/blindsight.html>

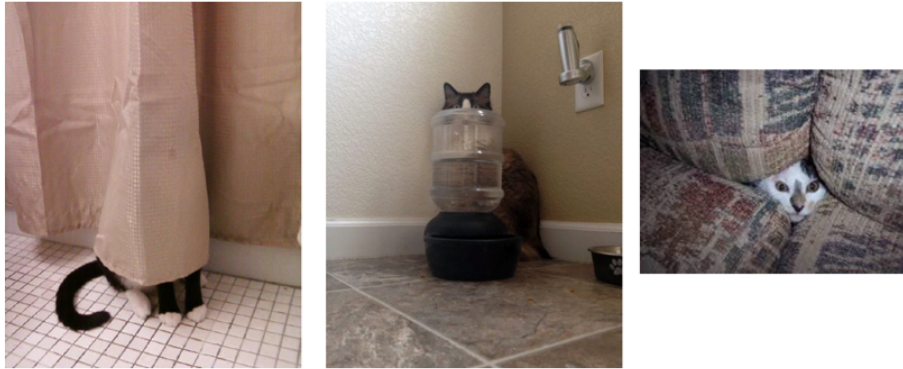
The Brain on Pattern Recognition

Despite Changes in Deformation:



The Brain on Pattern Recognition

Despite Changes in Occlusion:

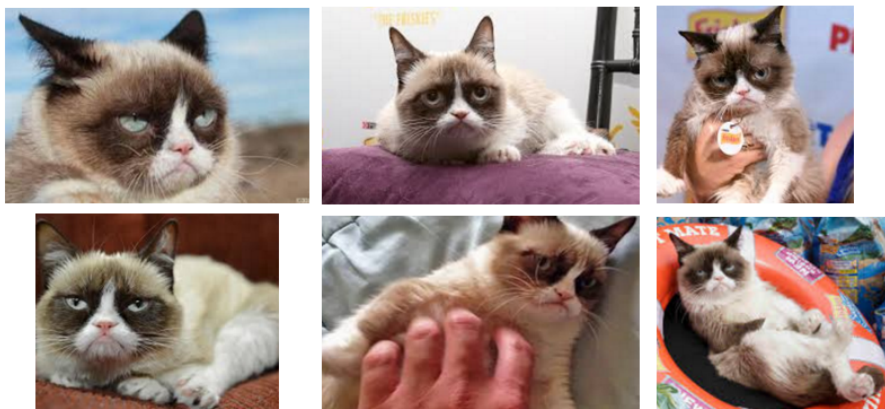


Ptucha '18

27

The Brain on Pattern Recognition

Despite Changes in Size, Pose, Angle:



Tardar Sauce "Grumpy Cat"

Ptucha '18

28

The Brain on Pattern Recognition

Despite Changes in Background Clutter:



Ptucha '18

29

The Brain on Pattern Recognition

Despite Changes in Class Variation...

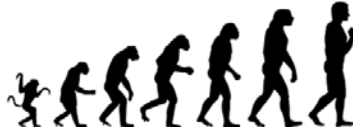


Ptucha '18

30

Teaching Computers to See

- It took evolution 540M years to develop the marvel of the eye-brain.



- Lets say a child collects a new image every 200msec.
- By age 3, this child has processed over 100M images.



$(5 \text{ images/sec})(60 \text{ sec/min})(60 \text{ min/hr})(12 \text{ hr/day})(365 \text{ days/yr})(3 \text{ yrs}) = 236 \text{ M}$

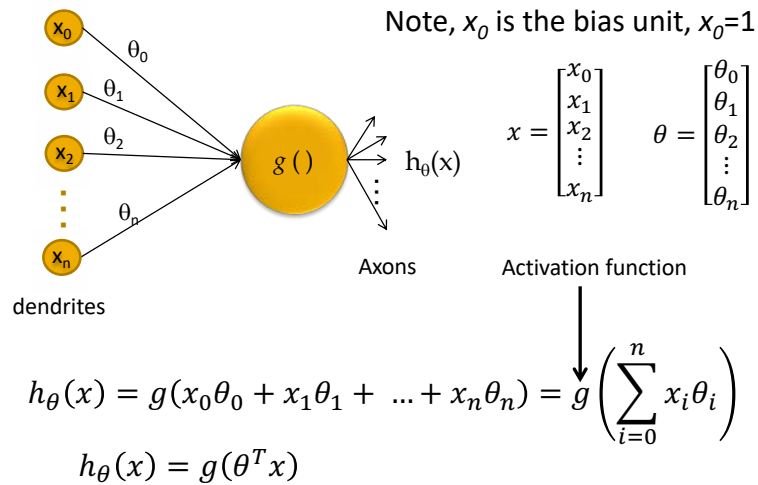
- Today's computers can do this in a few days...

Neural Nets on Pattern Recognition

- Instead of trying to code simple intuitions/rules on what makes an airplane, car, cat, and dog...
- We feed neural networks a large number of training samples, and it will automatically learn the rules!
- We will learn the magic behind this today!



Artificial Neuron

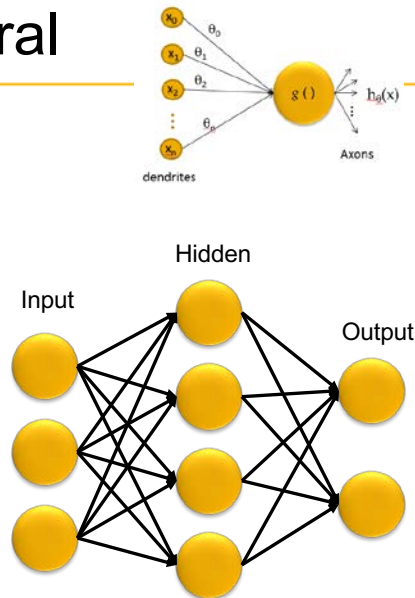


Ptucha '18

34

Artificial Neural Networks

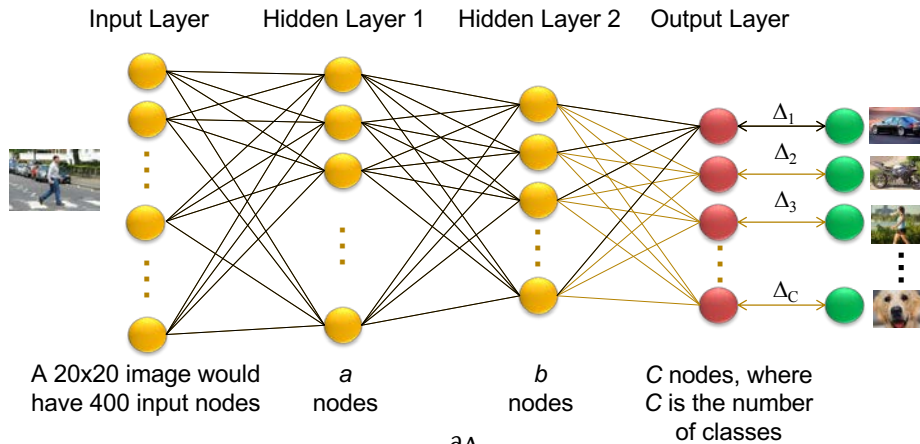
- Artificial Neural Network (ANN) – A network of interconnected nodes that “mimic” the properties of a biological network of neurons



Ptucha '18

35

4-Layer ANN Fully Connected Topology

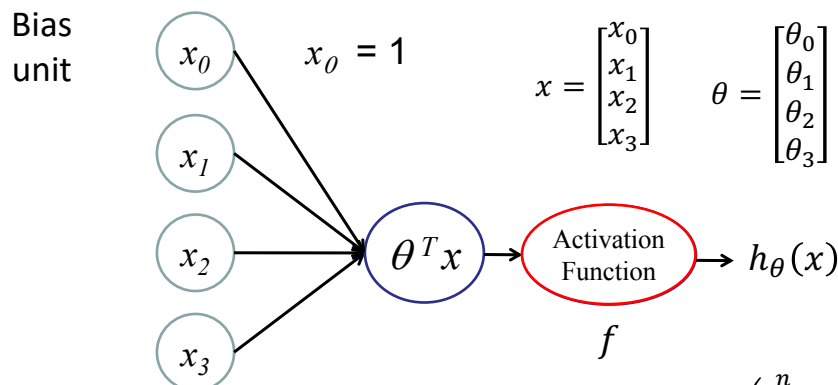


Backpropagation (~1985) uses $\frac{\partial \Delta}{\partial w}$ for learning
 Learning happens in the weights- each line is a weight.

Ptucha '18

36

Neuron Model



$$h_\theta(x) = g(x_0\theta_0 + x_1\theta_1 + \dots + x_n\theta_n) = g\left(\sum_{i=0}^n x_i\theta_i\right)$$

$$h_\theta(x) = g(\theta^T x)$$

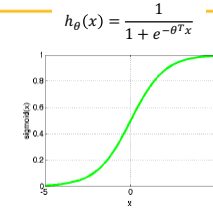
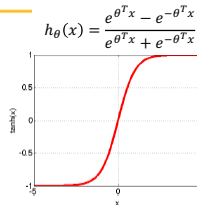
Ptucha '18

37

Activation Function Comparison

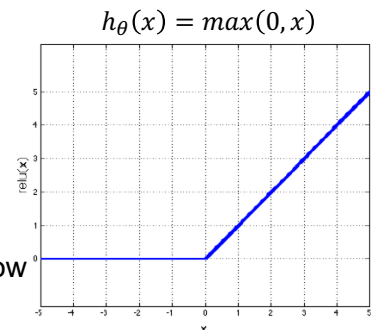
- **Tanh**
- **Sigmoid**

Gradient of both saturates at zero. Sigmoid also non-zero centered, so in practice tanh performs better.



- **Rectified Linear Units (ReLU)**

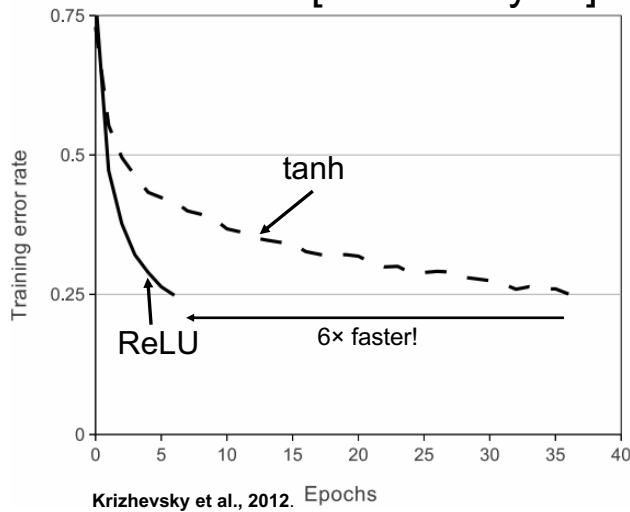
- Better for high dynamic range
- Faster learning
- Overall better result
- Neurons can “die” if allowed to grow unconstrained



Ptucha '18

39

Tanh vs. ReLU on CIFAR-10 dataset [Krizhevsky'12]



ReLU reaches 25% error 6× faster!
 Note: Learning rates optimized for each, no regularization, four layer CNN.

Krizhevsky et al., 2012. Epochs

Ptucha '18

40

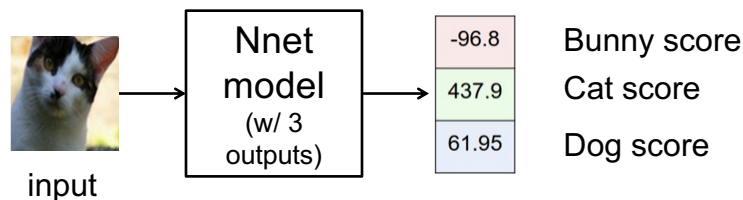
Where Do Weights Come From?

- The weights in a neural network need to be learned such that the errors are minimized.
- Just like logistic regression, we can write a cost function.
- Similar to gradient descent, we can write an iterative procedure to update weights, with each iteration decreasing our cost.
- These iterative methods may be less efficient than a direct analytical solution, but are easier to generalize.

Ptucha '18

41

Multiclass Loss Functions



- The input image scores highest against cat, but is also somewhat similar to dog.
- How do we assign a loss function?

Ptucha '18

47

Activation Function of Output Layer

- Sigmoid returns 0 or 1 for each output node.
- What if you wanted a confidence interval?
- Use a linear activation function for regression: $a^{(l)}=z^{(l)}$
- Softmax often used for classification:

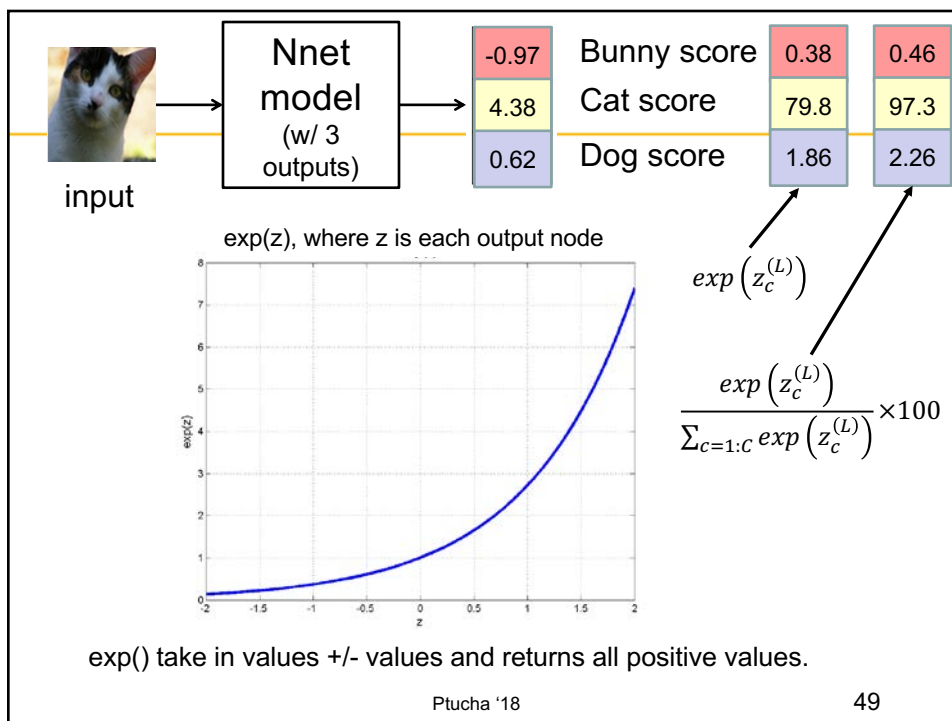
$$a_c^{(L)} = h_{\theta}(x)_c = g(z_c^{(L)}) = \frac{\exp(z_c^{(L)})}{\sum_{c=1:C} \exp(z_c^{(L)})}$$

← $\exp()$ of each output node
← Sum of all output nodes

- **Note: Only the output layer activation function changes- all hidden layer nodes activation functions would be the sigmoid/tanh/ReLU function.**

Ptucha '18

48



Most Common Loss Functions

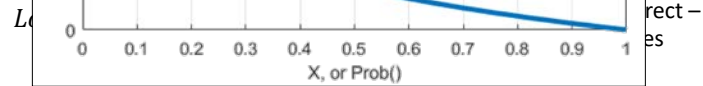
- The cost function we previously used was a direct copy from logistic regression and works great for binary classification.

- For multi-class, there are two popular data loss methods:

1. Cross-entropy loss, which uses softmax:

$$Loss^{(i)} = -\log\left(\frac{\exp(out_{y_i}^{(i)})}{\sum_{c=1:C} \exp(out_c^{(i)})}\right)$$

2. Multiclass SVM Loss (Weston Watkins formulation):



Most Common Loss Functions

- The cost function we previously used was a direct copy from logistic regression and works great for binary classification.

- For multi-class, there are two popular data loss methods:

1. Cross-entropy loss, which uses softmax:

$$Loss^{(i)} = -\log\left(\frac{\exp(out_{y_i}^{(i)})}{\sum_{c=1:C} \exp(out_c^{(i)})}\right)$$

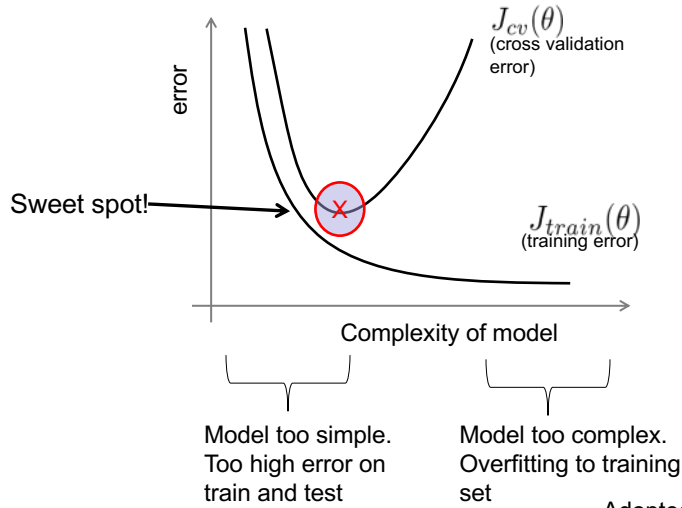
Loss for sample $i = \frac{\exp(\text{output of GT node})}{\text{Sum of exp(output) of all nodes}}$

2. Multiclass SVM Loss (Weston Watkins formulation):

$$Loss^{(i)} = \sum_{j \neq y_i} \max(0, out_j - out_{y_i} + \Delta)$$

Sum of incorrect - correct classes

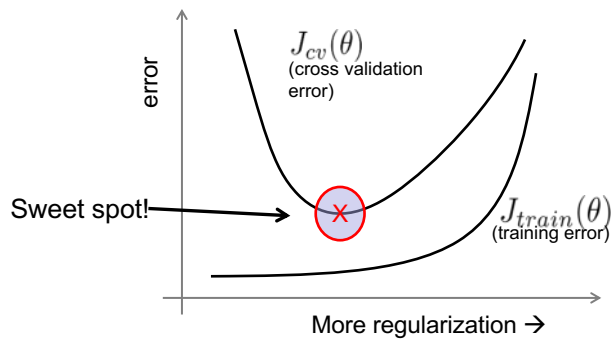
Bias (underfit) vs. Variance (overfit) errors



Ptucha '18

Adopted from:
Andrew Ng, ML class 64

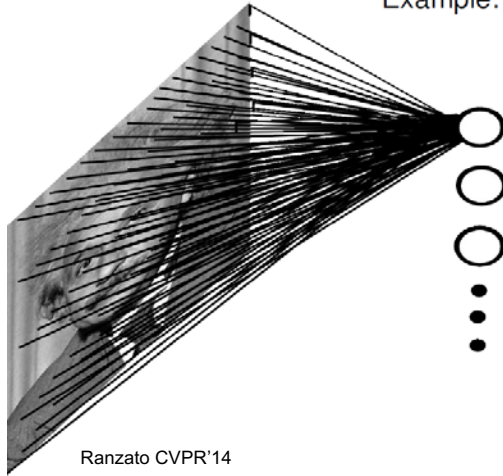
Regularization Tuning



Ptucha '18

Adopted from:
Andrew Ng, ML class 65

Fully Connected Layers?



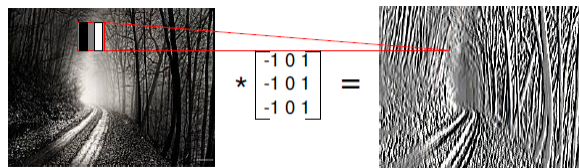
- Example:
- 200×200 pixel image.
 - 40K input fully connected to 40K hidden (or output) layer.
 - 1.6 billion weights!
 - Generally don't have enough training samples to learn that many weights.

Ranzato CVPR'14

Ptucha '18

66

Convolution Filter



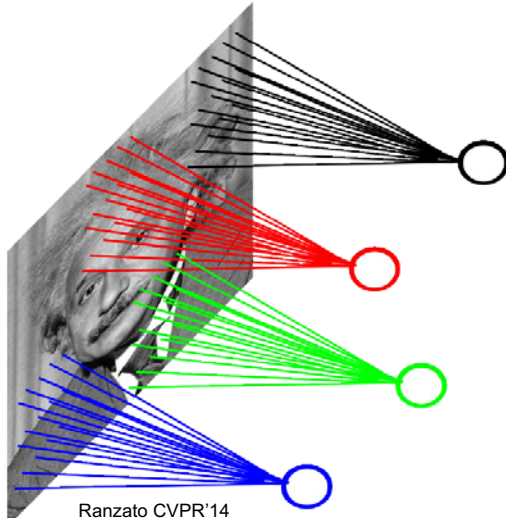
Ranzato CVPR'14

- Convolution filters apply a transform to an image.
- The above filter detects vertical edges.

Ptucha '18

67

Locally Connected Layer

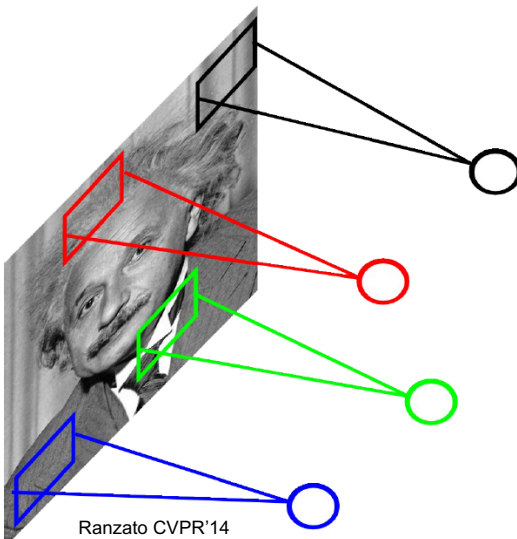


- 200×200 pixel image.
- 40K input.
- Four 10×10 filters, each fully connected
- $40K \times 10 \times 10 \times 4 = 16M$ weights....getting better!

Ptucha '18

68

Locally Connected Layer



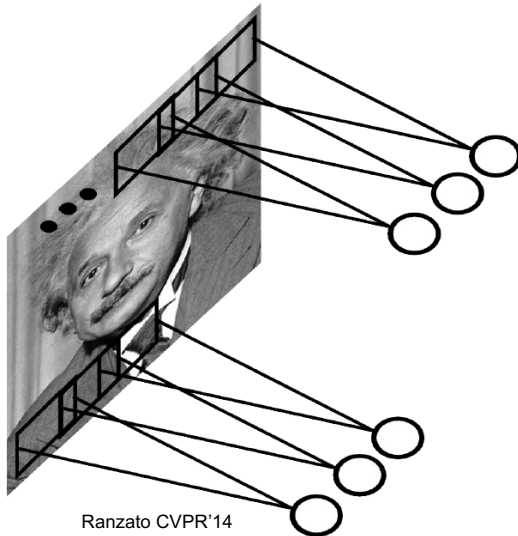
- 200×200 pixel image.
- 40K input.
- Four 10×10 filters, each fully connected
- $40K \times 10 \times 10 \times 4 = 16M$ weights....getting better!

- Can we formulate so each filter has similar statistics across all locations?

Ptucha '18

69

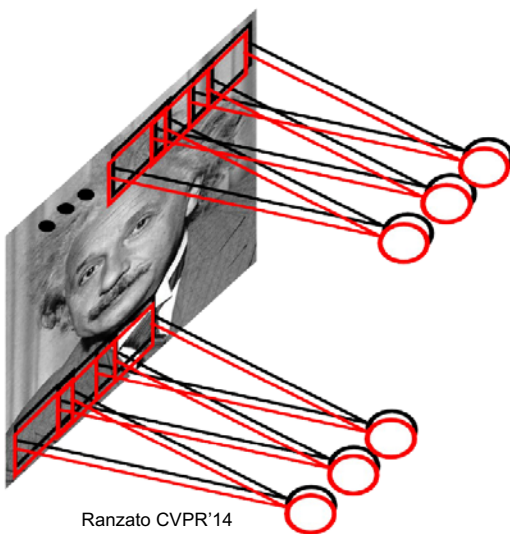
Convolution Layer



- 200×200 pixel image.
- 40K input.
- Four 10×10 filters, each fully connected
- $40K \times 10 \times 10 \times 4 = 16M$ weights....getting better!
- Require each filter has same statistics across all locations.
- Learn filters.

70

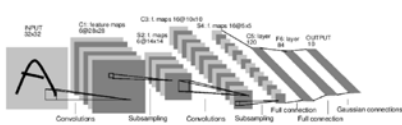
Convolution Layer



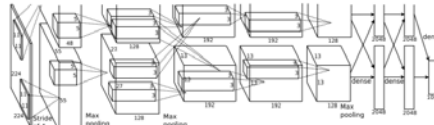
- 200×200 pixel image.
- 40K input.
- Four 10×10 filters, each fully connected
- $40K \times 10 \times 10 \times 4 = 16M$ weights....getting better!
- Require each filter has same statistics across all locations.
- Learn filters.
- To learn four filters we have $4 \times 10 \times 10 = 400$ parameters- great!

71

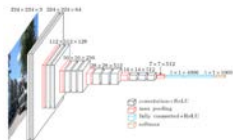
Many Flavors of Convolution Neural Networks (CNNs)...



LeNet-5, LeCun 1989



AlexNet, Krizhevsky 2012



VGGNet, Simonyan 2014



GoogLeNet (Inception), Szegedy 2014



ResNet, He 2015



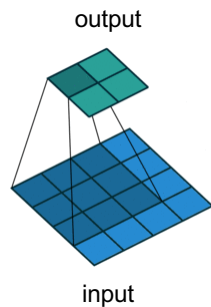
DenseNet, Huang 2017

Ptucha '18

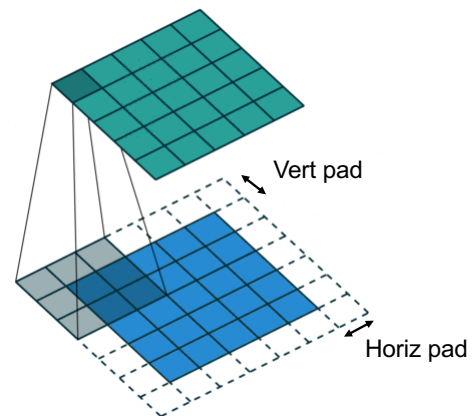
72

Image Convolution

By padding $(\text{filterWidth}-1)/2$, output image size matches input image size



3x3 filter sliding over input image

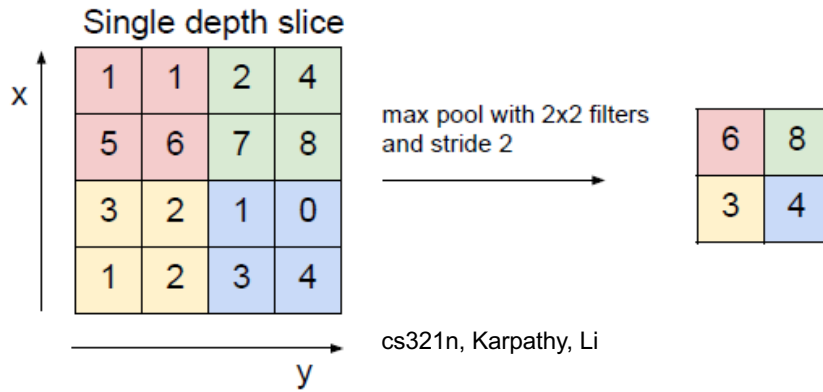


https://github.com/vdumoulin/conv_arithmetic

Ptucha '18

73

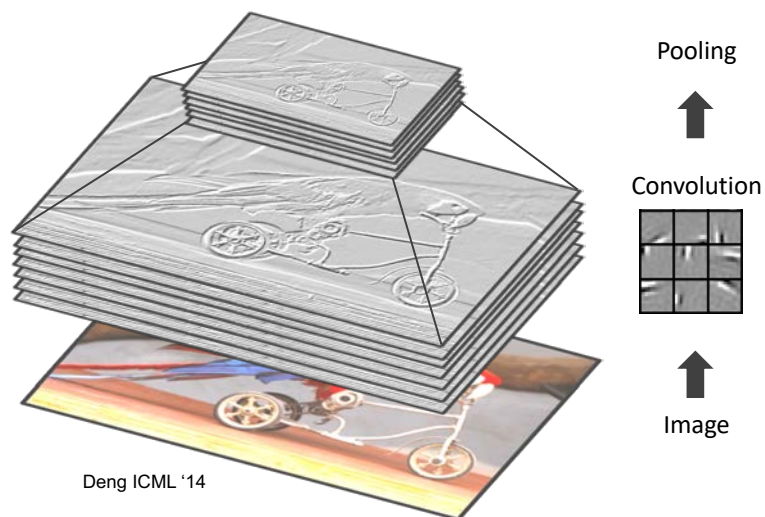
Max Pooling- Reducing the Size of an Image



Ptucha '18

74

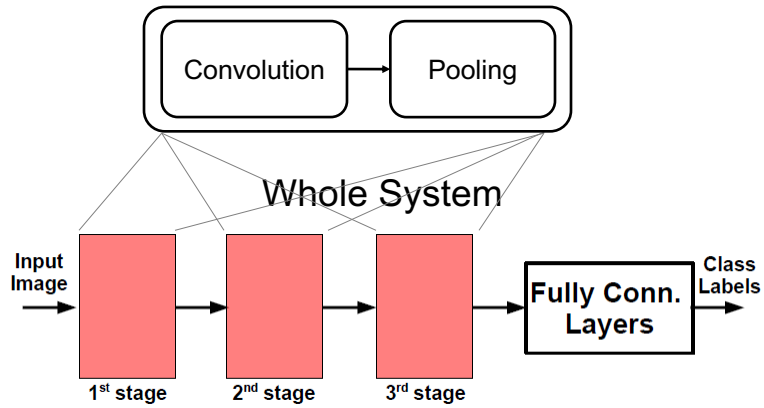
Convolution Neural Network (CNN) Building Block



Ptucha '18

75

Putting it All Together



Ptucha '18

76

Learning Filters

32 Learned Filters, each 5×5

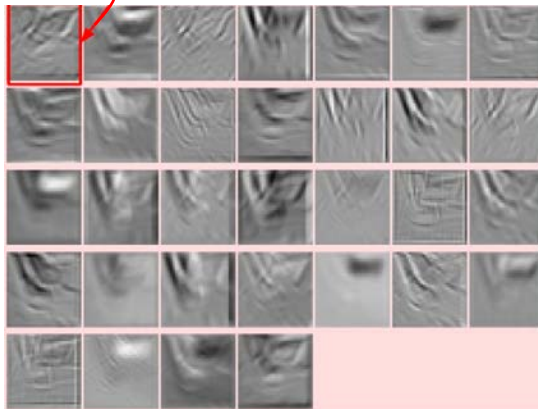


32 Filtered images, each is 28×28

Input image
28×28



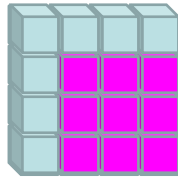
Use zero padding



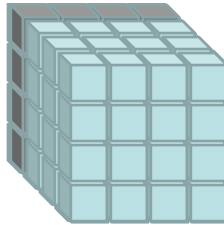
Ptucha '18

80

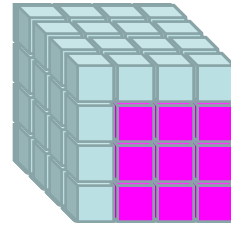
Filters



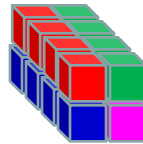
3x3 filter



3x3 filter



3x3x4 filter



Ptucha '18

81

Learning Filters

32 Learned Filters, each 5x5x3

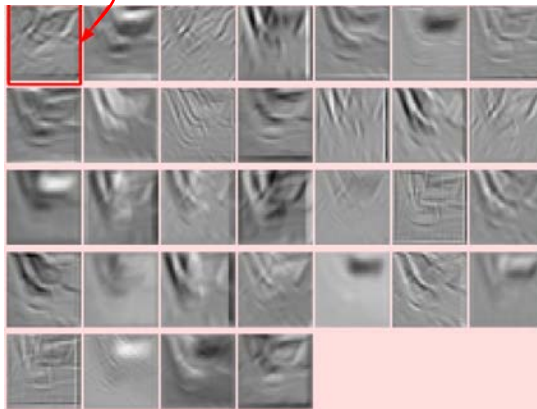


32 Filtered images, each is 28x28x1

Input image
28x28x3



Use zero padding



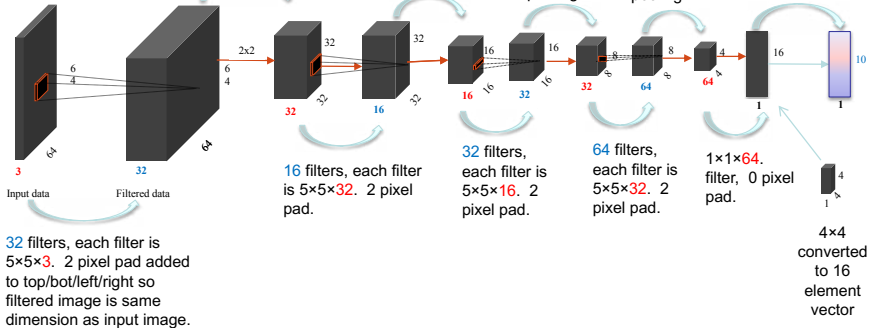
Ptucha '18

82

CNN Architecture

(Not so) Toy Example

Input RGB image:
64x64x3 pixels



Ptucha '18

83

Case Study

Case study: VGGNet / OxfordNet
(runner-up winner of ILSVRC 2014)
[Simonyan and Zisserman]

best model

ConvNet Configurations					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 x 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256	conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256
maxpool					
conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512
maxpool					
conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

cs321n, Karpathy, Li

Table 2. Number of parameters (in millions).

Network	A	A-LRN	B	C	D	E
Number of parameters	135	135	134	138	144	

Ptucha '18

89

Case Study

INPUT: [224x224x3] memory: 224*224*3=150K params: 0 (not counting biases)

CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728

CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*64)*64 = 36,864

POOL2: [112x112x64] memory: 112*112*64=800K params: 0

CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*64)*128 = 73,728

CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*128)*128 = 147,456

POOL2: [56x56x128] memory: 56*56*128=400K params: 0

CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*128)*256 = 294,912

CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824

CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824

POOL2: [28x28x256] memory: 28*28*256=200K params: 0

CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*256)*512 = 1,179,648

CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296

CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296

POOL2: [14x14x512] memory: 14*14*512=100K params: 0

CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296

CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296

CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296

POOL2: [7x7x512] memory: 7*7*512=25K params: 0

FC: [1x1x4096] memory: 4096 params: 7*7*512*4096 = 102,760,448

FC: [1x1x4096] memory: 4096 params: 4096*4096 = 16,777,216

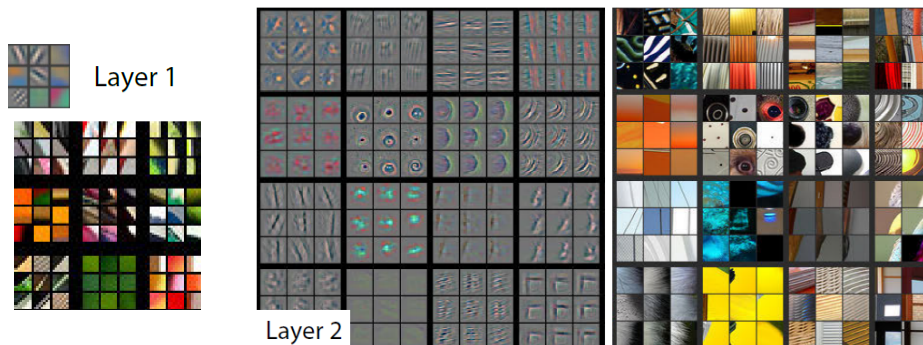
FC: [1x1x1000] memory: 1000 params: 4096*1000 = 4,096,000

ConvNet Configuration		
B	C	D
13 weight layers	16 weight layers	16 weight layers
<p>Note: Most memory in early layers</p>		
maxpool		
conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256
	conv1-256	conv3-256
maxpool		
conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512
	conv1-512	conv3-512
<p>Note: Most parameters in FC layers</p>		
FC-4096		
FC-1000		
soft-max		

Ptucha '18

90

CNN Visualization

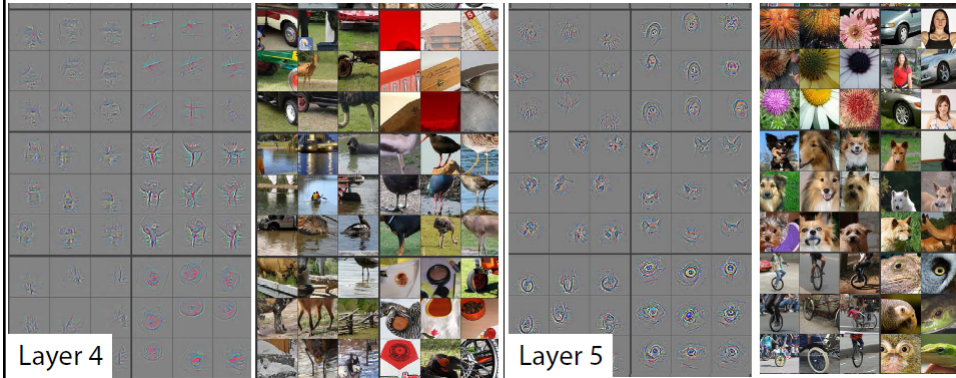


Zeiler, Fergus, 2014

Ptucha '18

91

CNN Visualization



Zeiler, Fergus, 2014

Ptucha '18

92

CNN as Vector Representation

Typical CNN Architecture



Input Image



2D Plot of fc8 Feature Vector

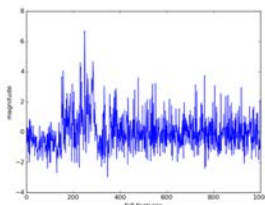
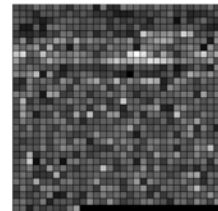


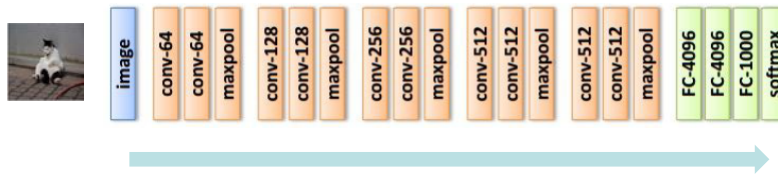
Image of fc8 Feature Vector



Ptucha '18

95

CNN as Vector Representation

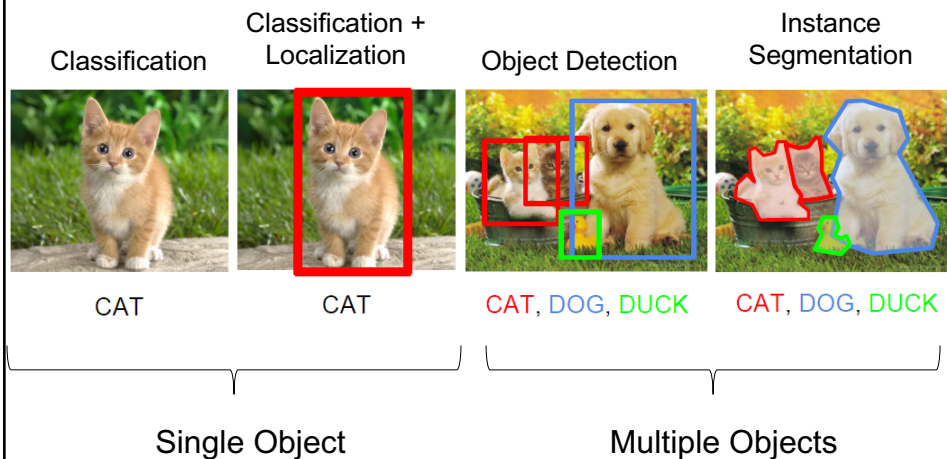


- As it turns out, these fully connected layers are excellent descriptors of the input image!
- For example, you can pass images through a pre-trained CNN, then take the output from a FC layer (image2vec) as input to a SVM classifier.
- Images in this vector space generally have the property that similar images are close in this latent representation.

Ptucha '18

96

Vision Tasks



Ptucha '18

108

Classification vs. Classification + Localization

Classification

Input: Image
Output: Class label
Evaluation metric: Accuracy



→ CAT

Classification + Localization

Input: Image
Output: Class label, Box coordinates
Evaluation metric:
 Intersection over Union (IoU)



→ (CAT,x,y,w,h)

Ptucha '18

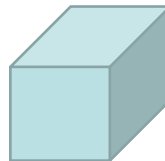
109

Classification with Localization

(0,0)



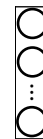
→



→

...

→



Lets allow a few classes:

1. Car
2. Truck
3. Pedestrian
4. Motorcycle

- For now, lets assume one object per image.
- Each object has $\{x, y, w, h\}$
- For this image, object location $\{x, y, w, h\} = \{0.3, 0.6, 0.4, 0.3\}$

Image from: deeplearning.ai, C4W3L01

Ptucha '18

110

Classification with Localization

Four classes:

1. Car
2. Truck
3. Pedestrian
4. Motorcycle

Localization $\{x, y, w, h\}$



Define y label: $y =$

$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \\ C_1 \\ C_2 \\ C_3 \\ C_4 \end{bmatrix}$$

Probability of an object
 Bounding box location
 0/1 for each class

Image from: deeplearning.ai, C4W3L01

Ptucha '18

111

Classification with Localization

Four classes:

1. Car
2. Truck
3. Pedestrian
4. Motorcycle

Localization $\{x, y, w, h\}$



Cost function (squared error):

$$Loss = \sum_{i=1}^9 (\hat{y}_i - y_i)^2 \quad \text{If } y_i=1$$

$$Loss = (\hat{y}_1 - y_1)^2 \quad \text{If } y_1=0$$

$$y = \begin{bmatrix} 1 \\ 0.3 \\ 0.6 \\ 0.4 \\ 0.3 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \\ C_1 \\ C_2 \\ C_3 \\ C_4 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \quad ? = \text{don't care}$$

Image from: deeplearning.ai, C4W3L01

Ptucha '18

112

Classification with Localization

Four classes:

1. Car
2. Truck
3. Pedestrian
4. Motorcycle

Localization $\{x, y, w, h\}$

Alternate cost function:

- y_1 can be logistic loss
- $y_2 \rightarrow y_5$ can be squared error
- $y_6 \rightarrow y_9$ can be softmax cross entropy



$$y = \begin{bmatrix} 1 \\ 0.3 \\ 0.6 \\ 0.4 \\ 0.3 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \\ C_1 \\ C_2 \\ C_3 \\ C_4 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

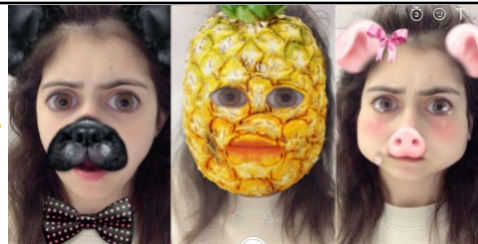
? = don't care

Image from: deeplearning.ai, C4W3L01

Ptucha '18

113

Snapchat Facewarp?



- Traditional approach:

Viola Jones
Face Detection

Search for actual point locations
using Mahalanobis distance



Repeat ~3-5x

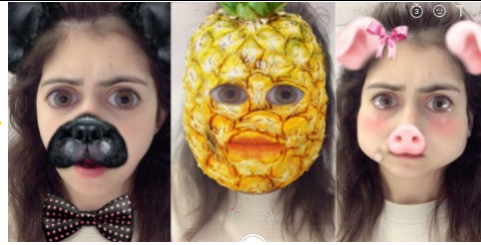
Average eye and 82
facial feature points

Restrict based on
PCA statistics

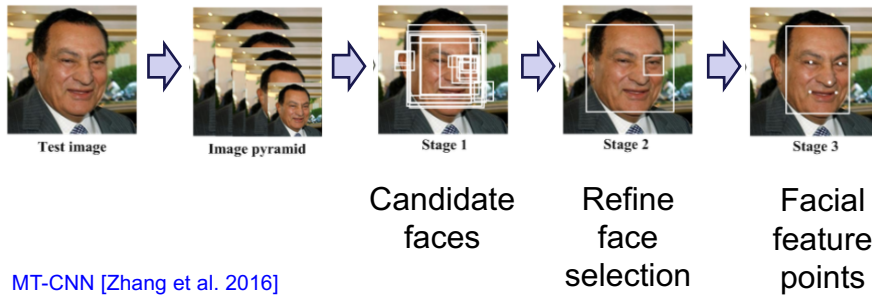
Ptucha '18

114

Snapchat Facewarp?



- Deep Learning approach:

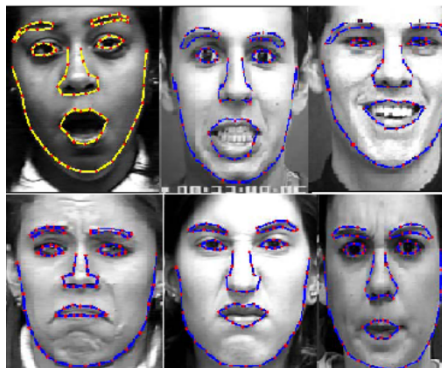


Ptucha '18

115

Localization

- Facial feature



Each face has 68 points, so CNN would output:

Face? {
 pt1X
 pt1Y
 pt2X
 pt2Y
 .
 .
 .
 pt68X
 pt68Y

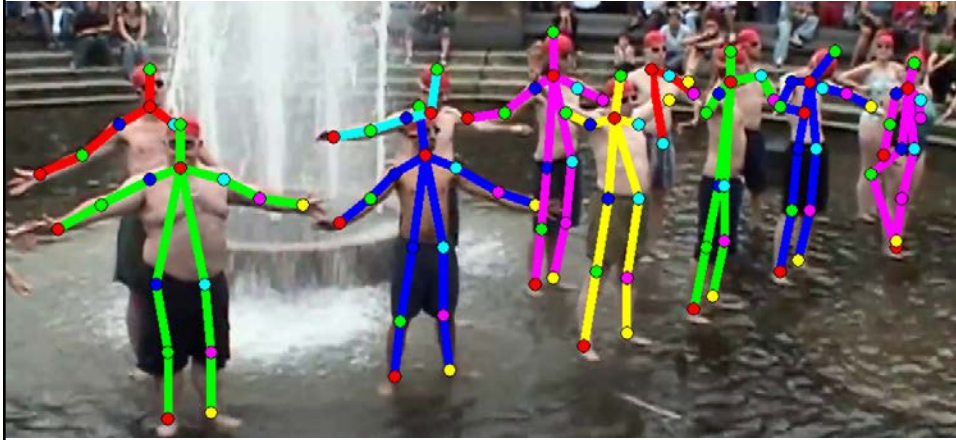
137 outputs

Of course, need GT for thousands of faces to train model.

ptucha Ptucha '18

116

Can do same with Body Pose...



Pishchulin et al. CVPR'16

Ptucha '18

117

Object Detection

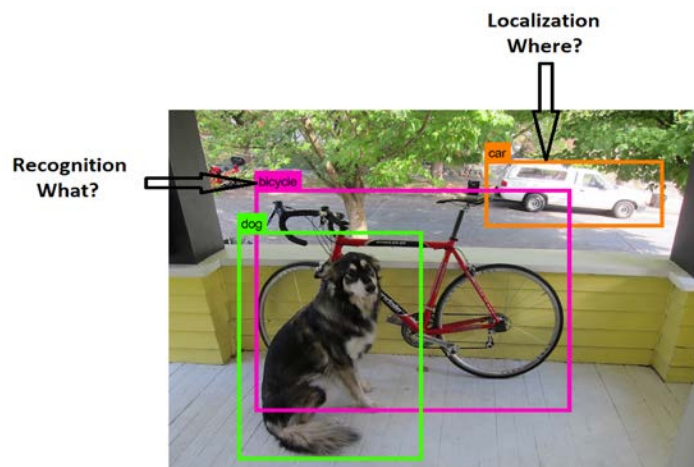


Image Credit: Redmon, Joseph, et al [4]

Ptucha '18

118

More than one object per image?

Training set:



y

1

1

1

0

0

Car detection example

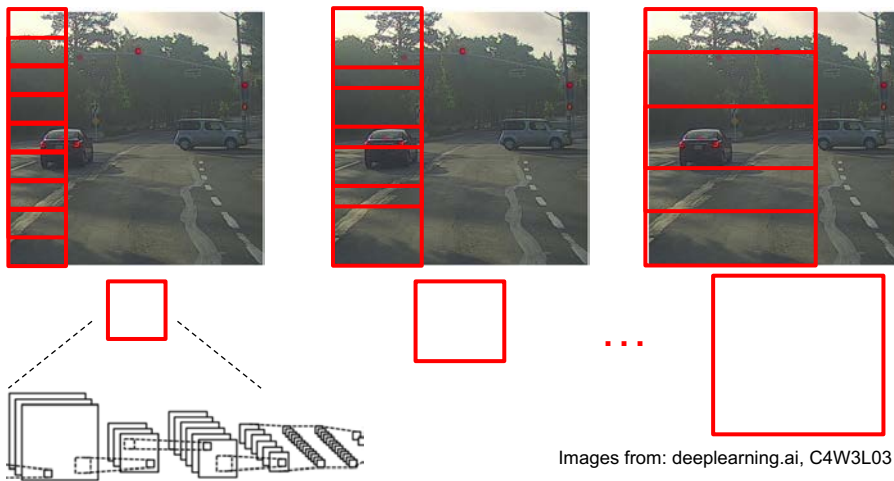


Images from: deeplearning.ai, C4W3L03

Ptucha '18

121

Sliding Window Detection



Images from: deeplearning.ai, C4W3L03

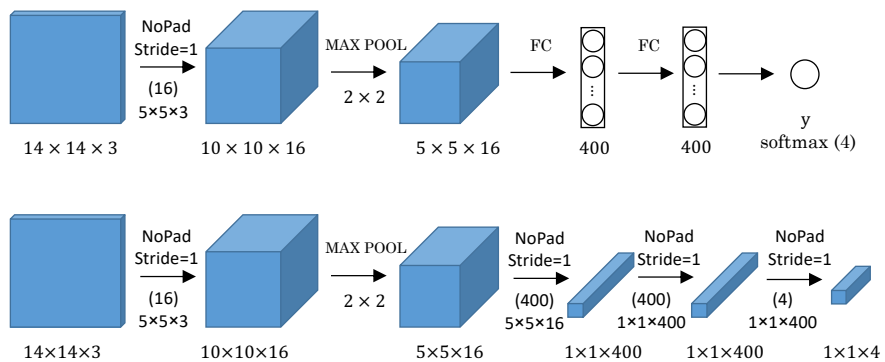
ptucha

Ptucha '18

122

122

Computing FC layers with Convolution

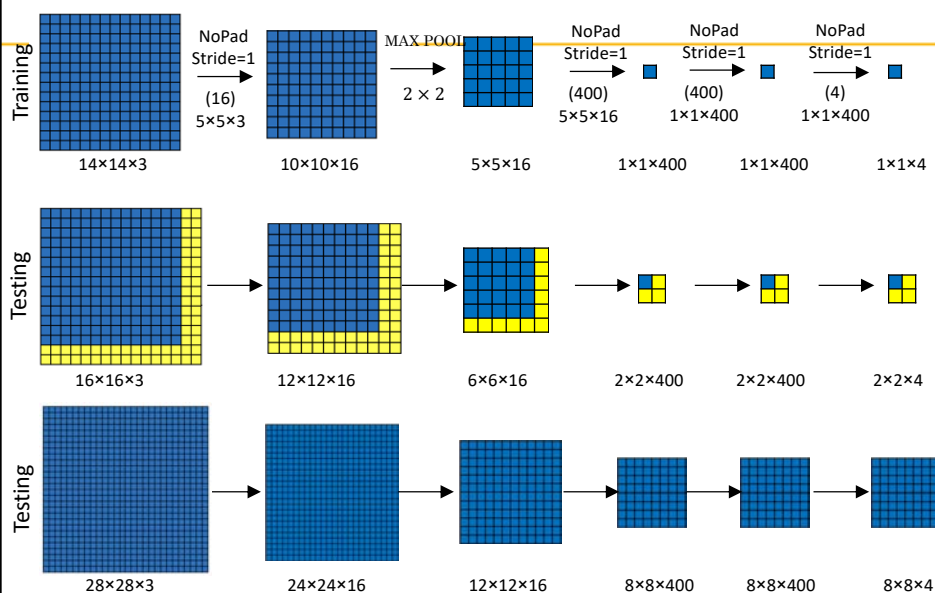


Images from: deeplearning.ai, C4W3L04

Ptucha '18

123

Replacing Sliding Windows w/Fully Convolutional CNNs



Images from: deeplearning.ai, C4W3L04

Ptucha '18

124

Replacing Sliding Windows w/Fully Convolutional CNNs

Sliding window approach:
Sequentially evaluate one window at a time

Fully convolutional approach:
Evaluate 64 regions at once

$28 \times 28 \times 3$ $24 \times 24 \times 16$ $12 \times 12 \times 16$ $8 \times 8 \times 400$ $8 \times 8 \times 400$ $8 \times 8 \times 4$
Images from: deeplearning.ai, C4W3L04 Ptucha '18 125

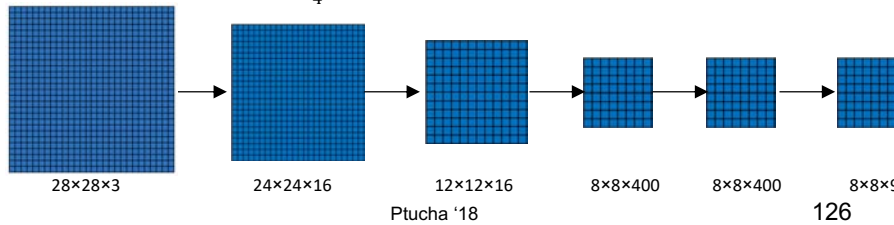
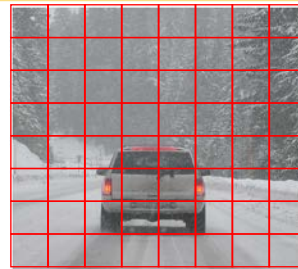
Replacing Sliding Windows w/Fully Convolutional CNNs

- Can think of this as evaluating 8×8 grid, where each of the 64 cells is independently checked for an object:

Each cell has a y label:

$$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \\ C_1 \\ C_2 \\ C_3 \\ C_4 \end{bmatrix}$$

Prob. of an object
 Object location
 0/1 for each class
 (Four classes in this example)



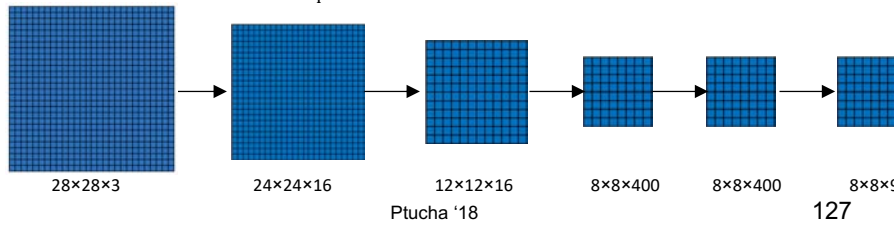
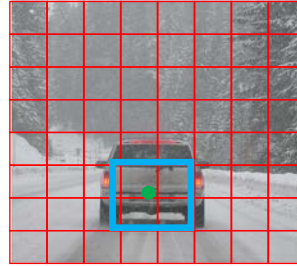
Replacing Sliding Windows w/Fully Convolutional CNNs

- Overlay GT of object
- Cell where centroid lie is responsible.

Each cell has a y label:

$$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \\ C_1 \\ C_2 \\ C_3 \\ C_4 \end{bmatrix}$$

Prob. of an object
Object location
0/1 for each class
(Four classes in this example)



Replacing Sliding Windows w/Fully Convolutional CNNs

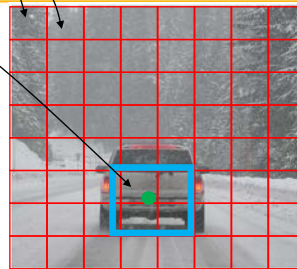
$$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \\ C_{pers} \\ C_{car} \\ C_{truck} \\ C_{motorc} \end{bmatrix}$$

$$y_1 = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

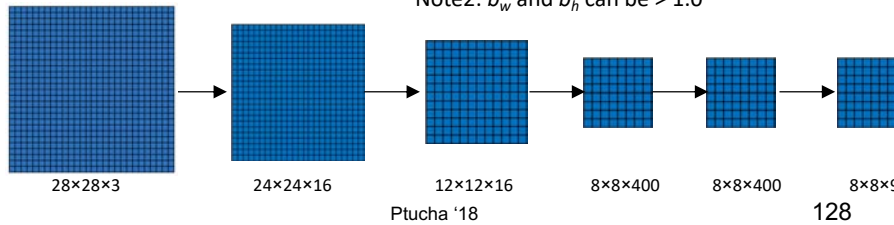
$$y_2 = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

$$\dots$$

$$y_{44} = \begin{bmatrix} 1 \\ 0.8 \\ 0.9 \\ 2.0 \\ 1.8 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$



Note1: cell upper left (0,0); cell lower right (1,1)
Note2: b_w and b_h can be > 1.0





deeplearning.ai | COURSERA



Andrew Ng, 2017

<https://www.deeplearning.ai/>

Deep Learning Specialization, Five courses:

1. Neural Networks and Deep Learning
2. Improving Deep Neural Networks
3. Structured Machine Learning Projects
4. Convolutional Neural Networks
5. Sequence Models

Ptucha '18

133



CS231n: Convolutional Neural Networks for Visual Recognition



Fei-Fei Li



Justin Johnson



Serena Young

Li, Johnson, Yeung 2017

<http://cs231n.stanford.edu/>

Ptucha '18

134



DEEP
LEARNING
INSTITUTE

www.nvidia.com/dli



<https://www.rit.edu/mil>



Raymond W. Ptucha

Assistant Professor, Computer Engineering
Director, Machine Intelligence Laboratory
Rochester Institute of Technology

Email: rwpeec@rit.edu
Phone: +1 (585) 797-5561
Office: GLE (09) 3441

Ptucha '18

135