

Deep Learning for CV

Nguyen Thi Oanh

Hanoi University of Science and Technology

oanhnt@soict.hust.edu.vn

SOICT, HUST

Content

2

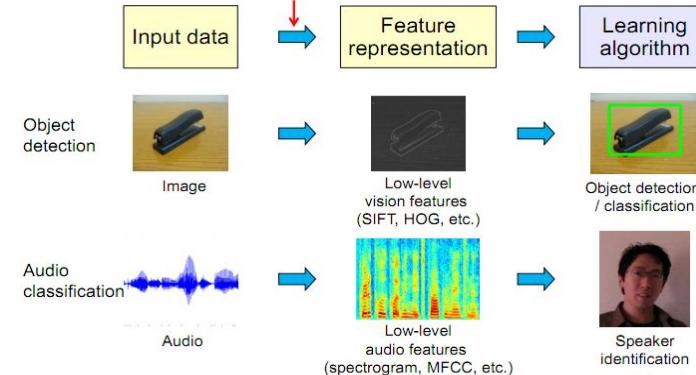
- Introduction
- ML for image classification
- CNN
 - Architectures
 - Training

Computer Vision

3

Traditional machine perception

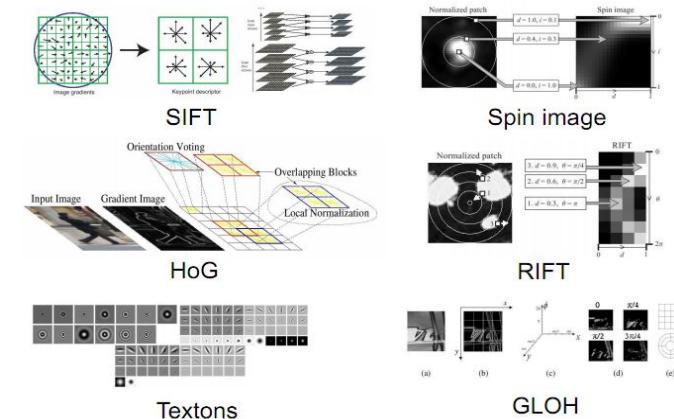
State-of-the-art:
"hand-crafting"



Computer Vision

4

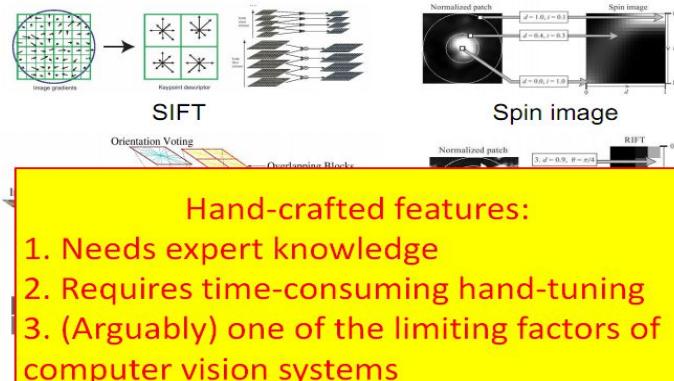
Computer vision features



Computer Vision

5

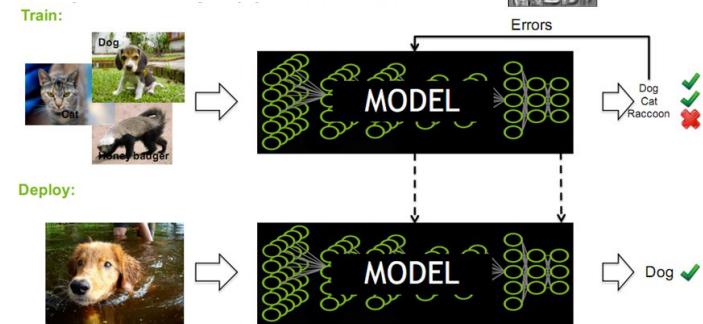
Computer vision features



Deep Learning

6

Deep learning approach



Deep Learning

7

DEEP LEARNING EVERYWHERE



INTERNET & CLOUD	MEDICINE & BIOLOGY	MEDIA & ENTERTAINMENT	SECURITY & DEFENSE	AUTONOMOUS MACHINES
Image Classification Speech Recognition Language Translation Language Processing Sentiment Analysis Recommendation	Cancer Cell Detection Diabetic Grading Drug Discovery	Video Captioning Video Search Real Time Translation	Face Detection Video Surveillance Satellite Imagery	Pedestrian Detection Lane Tracking Recognize Traffic Sign

Deep Learning

8

10 BREAKTHROUGH TECHNOLOGIES 2013

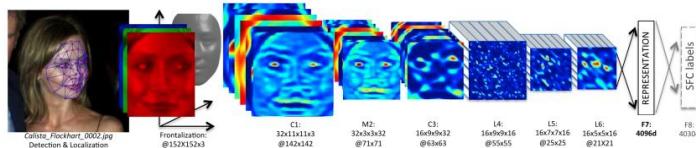
[Introduction](#) [The 10 Technologies](#) [Past Years](#)


Deep Learning vs Human

9

DeepFace: Closing the Gap to Human-Level Performance in Face Verification

Yaniv Taigman Ming Yang Marc'Aurelio Ranzato Lior Wolf
Facebook AI Research Menlo Park, CA, USA Tel Aviv University
yaniv, mingyang, ranzato@fb.com wolf@cs.tau.ac.il



DeepFace 2014

Closing the Gap to Human Level Performance in Face Verification

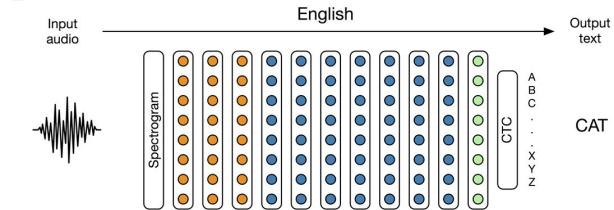
Accuracy

DeepFace: 97.35%

Human: 97.5%

Deep Learning vs Human

10



Deep Speech 2015

Baidu has developed a voice system that can recognize English and Mandarin speech **better than people**, in some cases.

"For short phrases, out of context, we seem to be **surpassing human levels of recognition**" - Andrew Ng

Deep Learning vs Human

11

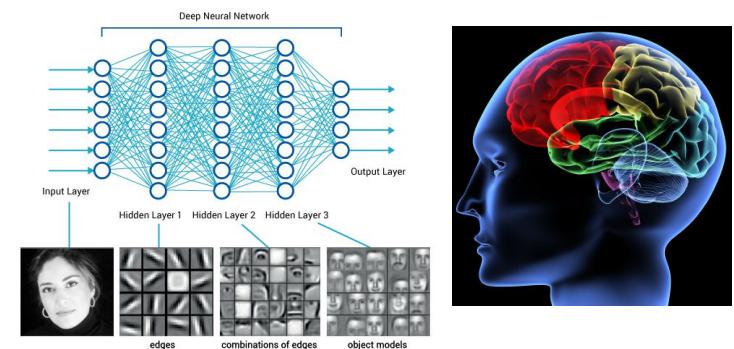


AlphaGo vs Human: 9-1

What is Deep Learning?

12

Deep learning is a branch of machine learning based on a set of algorithms that attempt to model high-level abstractions in data by using a deep graph with multiple processing layers, composed of multiple non-linear transformations.



Content

13

- Introduction
- ML for image classification
- CNN
 - Architectures
 - Training

Convolution

14

$$\begin{array}{|c|c|c|} \hline 12 & 3 & 19 \\ \hline 25 & 10 & 1 \\ \hline 9 & 7 & 17 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 133 & 75 \\ \hline 100 & 101 \\ \hline \end{array}$$

$$f[n, m] * h[n, m] = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} f[k, l] h[n - k, m - l]$$

Why they are useful

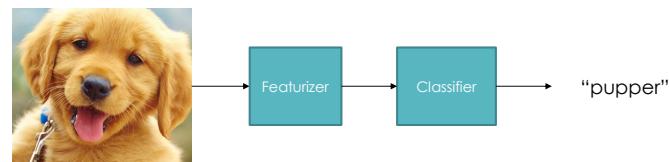
15

Allow us to find **interesting insights/features** from images!

$$\begin{array}{c} \text{Image of a mechanical part} \\ * \end{array} \begin{array}{|c|c|c|} \hline 0 & -\frac{1}{2} & 0 \\ \hline 0 & 0 & 0 \\ \hline 0 & \frac{1}{2} & 0 \\ \hline \end{array} = \begin{array}{c} \text{Image of the same part with highlights} \end{array}$$

Recall Image Classification...

16



Allow us to use features to put **images in categories**!

Wait a Minute...

17

Convolution = Image -> Features

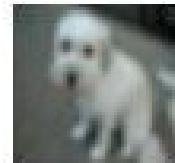
Classification Algorithm = Features -> Category

→ Let's put 'em together!



Feature Extractor

19



32x32
"Airplane
Filter"

=

●

In Specific...

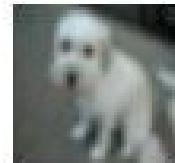
18

Let's build a **convolution-based** classification algorithm for the CIFAR-10 dataset (10 classes, 32x32 images):



Feature Extractor

20



32x32
"Airplane
Filter"

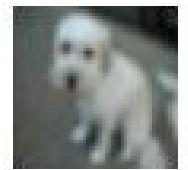
=

●

"probability" of
the image
being an
airplane

Feature Extractor

21



32x32
"Airplane
Filter"

"probability" of
the image
being an
airplane

$$= \text{circle} \longrightarrow$$



Feature Extractor

22



32x32
"Airplane
Filter"

*This is not really a probability
but a score, because it can
be less than 0 and greater
than 1*

"probability" of
the image being
an airplane

$$= \text{circle} \longrightarrow$$



Feature Extractor

23



32x32
"Automobile
Filter"

"probability" of
the image being
an automobile

$$= \text{circle} \longrightarrow$$



Feature Extractor

24



32x32 "Bird
Filter"

"probability" of
the image
being a bird

$$= \text{circle} \longrightarrow$$



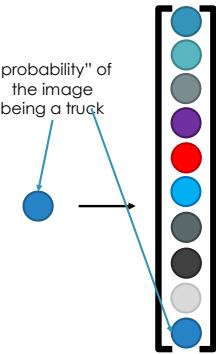
Feature Extractor

25



$*$ **32x32 "Truck Filter"**

"probability" of the image being a truck



Classifier

26

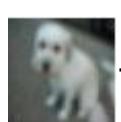
$$c_{pred} = \arg \max(\text{vector})$$

We predict the class that has the highest probability!



The Whole Shebang

27



Image

32x32x10 Conv Block

Feature Extractor

\hat{y}

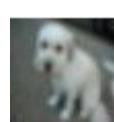
Prediction

 \hat{y}

27

The Whole Shebang

28



Image

32x32x10 Conv Block

Feature Extractor

\hat{y}

Prediction

 \hat{y}

28

Classifier

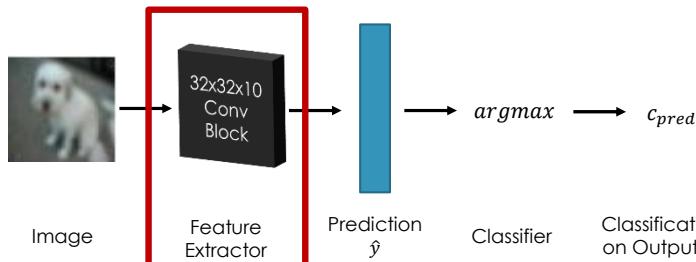
Classifier

Classification Output

The Whole Shebang

29

How ???



Reframing convolution

30

$$\begin{bmatrix} 12 & 21 \\ 18 & 31 \end{bmatrix} * \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 12 \\ 21 \\ 18 \\ 31 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$$

Reframed Feature Extractor

31

$$\begin{bmatrix} \text{Image} \end{bmatrix} * \begin{bmatrix} \text{32x32} \\ \text{"Airplane Filter"} \end{bmatrix}$$

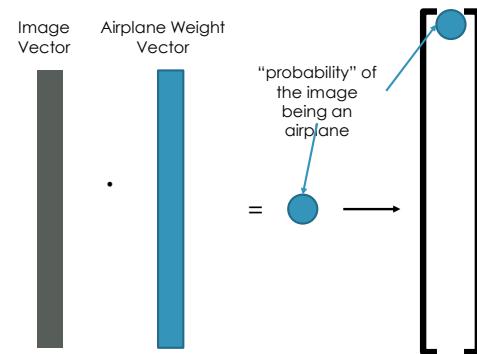
Reframed Feature Extractor

32

$$\begin{bmatrix} \text{Image Vector} \end{bmatrix} * \begin{bmatrix} \text{32x32} \\ \text{"Airplane Filter"} \end{bmatrix} = \begin{bmatrix} \text{Airplane Weight Vector} \end{bmatrix} \cdot \begin{bmatrix} \text{Image Vector} \end{bmatrix}$$

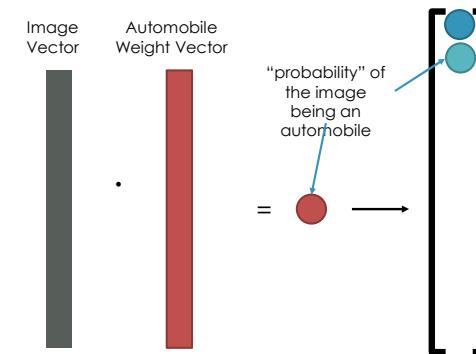
New Feature Extractor

33



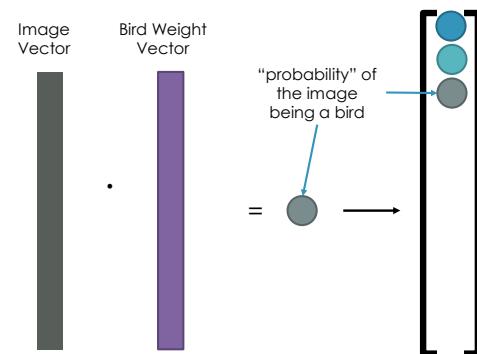
New Feature Extractor

34



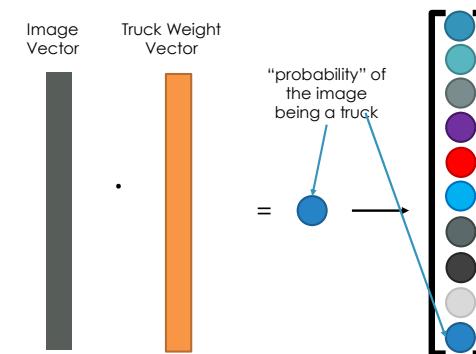
New Feature Extractor

35



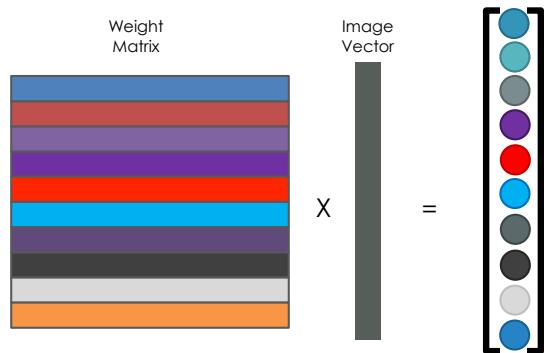
New Feature Extractor

36



New Feature Extractor

37



New Feature Extractor

38

$$Wx = \hat{y}$$

W : the (10×1024) matrix of weight vectors

x : the (1024×1) image vector

\hat{y} : the (10×1) vector of class "probabilities"

New Feature Extractor

39

This simple computation
is called a *fully-*
connected layer!

$$Wx = \hat{y}$$

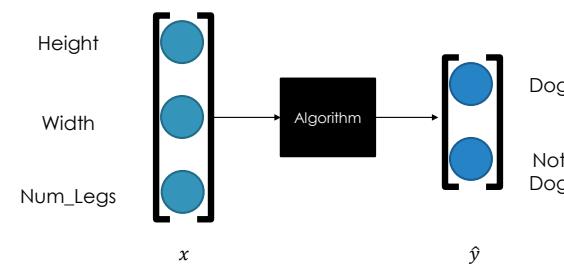
W : the (10×1024) matrix of weight vectors

x : the (1024×1) image vector

\hat{y} : the (10×1) vector of class "probabilities"

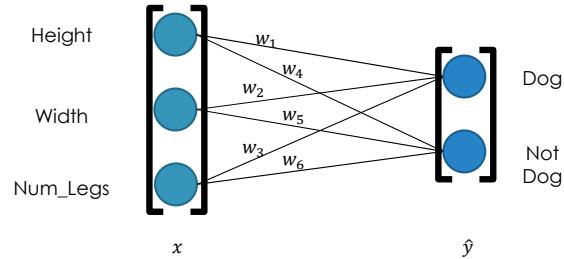
Aside: Fully-Connected Neural Networks

40



Aside: Fully-Connected Neural Networks

41



Aside: Fully-Connected Neural Networks

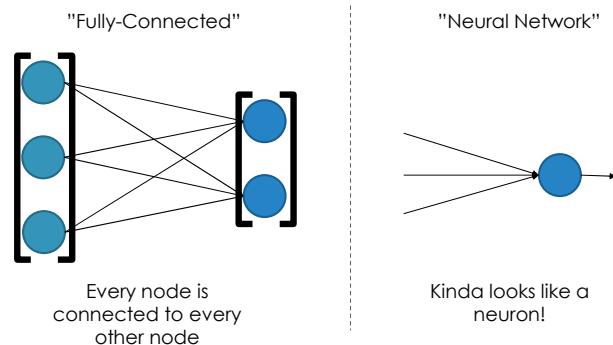
42

$$\begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \end{bmatrix} \cdot \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} \hat{y} \end{bmatrix}$$

$Wx = \hat{y}$

Aside: Fully-Connected Neural Networks

43



New Feature Extractor

44

$$Wx = \hat{y}$$

W : the (10×1024) matrix of weight vectors

x : the (1024×1) image vector

\hat{y} : the (10×1) vector of class "probabilities"

New Feature Extractor

45

$$Wx = \hat{y}$$

W : the (10x1024) matrix of weight vectors

x : the (1024x1) image vector

\hat{y} : the (10x1) vector of class "probabilities"?

Class Probability Vector

46

- Must have values between 0 and 1

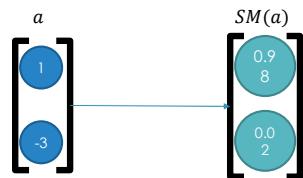
- Must sum to 1

- There's no guarantee either requirement is satisfied!

$$\hat{y} = Wx$$

Softmax Function

47



$$\text{Softmax: } a(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Class Probability Vector

48

- Must have values between 0 and 1

- Must have values between 0 and 1

- Must sum to 1

- Must sum to 1

$$\hat{y} = Wx$$

$$\hat{y} = SM(Wx)$$

System so far...

49

- Feature extractor:

$$\hat{y} = SM(Wx)$$

- Classifier:

$$c_{pred} = \arg \max(\hat{y})$$

System so far...

50

- Feature extractor:

$$\hat{y} = SM(\boxed{W}x)$$

- Classifier:

$$c_{pred} = \arg \max(\hat{y})$$

System so far...

51

- Feature extractor:

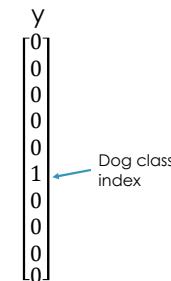


- Classifier:

Using the label

52

Let's compare our prediction with the real answer!
For each image, we have the label y which tells us the true class:



Key Insight:

53

We want:

$$\arg \max(\hat{y}) = \arg \max(y)$$

Key Insight:

54

We want:

$$\arg \max(\hat{y}) = \arg \max(y)$$

Which we can accomplish by:

$$W^* = \arg \min_W \left(- \sum_{x,y} \log(p_c) \right)$$

Where p_c is the probability of the true class in \hat{y}

Cross-Entropy Loss

55

Our loss function represents how bad we are currently doing:

$$L = -\log(p_c)$$

Examples:

$$p_c = 0 \rightarrow L = -\log(0) = \infty$$

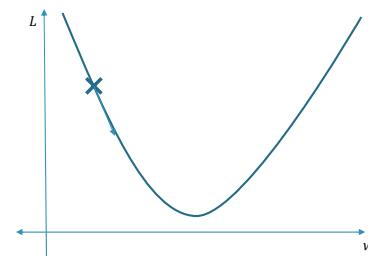
$$p_c = 0.1 \rightarrow L = -\log(0.1) = 2.3$$

$$p_c = 0.9 \rightarrow L = -\log(0.9) = 0.1$$

$$p_c = 1 \rightarrow L = -\log(1) = 0$$

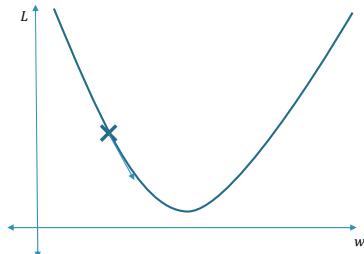
Minimizing Loss

56



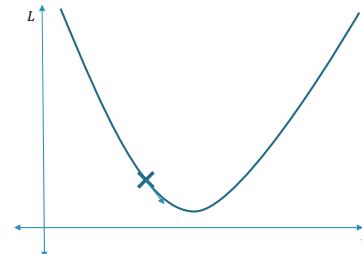
Minimizing Loss

57



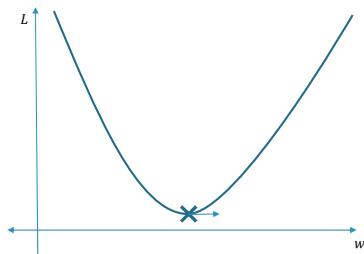
Minimizing Loss

58



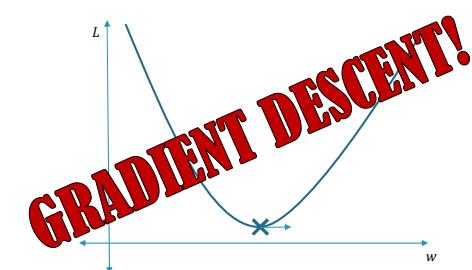
Minimizing Loss

59



Minimizing Loss

60



Gradient Descent Pseudocode

61

```
for i in {0, ..., num_epochs}:
    for x, y in data:
         $\hat{y} = SM(Wx)$ 
         $L = CE(\hat{y}, y)$ 
         $\frac{dL}{dW} = ???$ 
         $W := W - \alpha \frac{dL}{dW}$ 
```

Getting the Gradient

62

$$z = Wx$$

$$L = SCE(z, y)$$

$$\frac{dL}{dW} = \frac{dL}{dz} \frac{dz}{dW}$$

Getting the Gradient

63

$$z = Wx$$

$$L = SCE(z, y)$$

$$\frac{dL}{dW} = \frac{dL}{dz}(x)$$

Getting the Gradient

64

$$z = Wx$$

$$L = SCE(z, y)$$

$$\frac{dL}{dW} = (SM(z) - y)(x)$$

Getting the Gradient

65

$$z = Wx$$

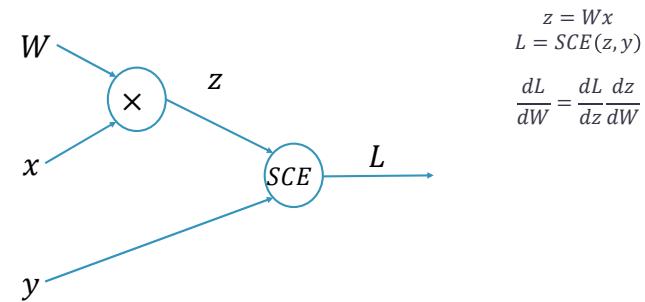
$$L = SCE(z, y)$$

$$\frac{dL}{dW} = (SM(z) - y)(x^T)$$

BACKPROPAGATION!

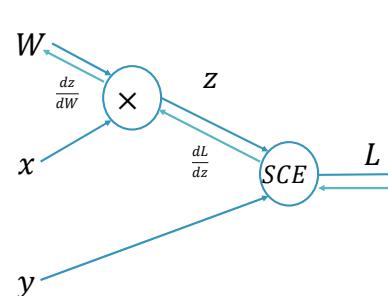
What is Backprop?

66



What is Backprop?

67



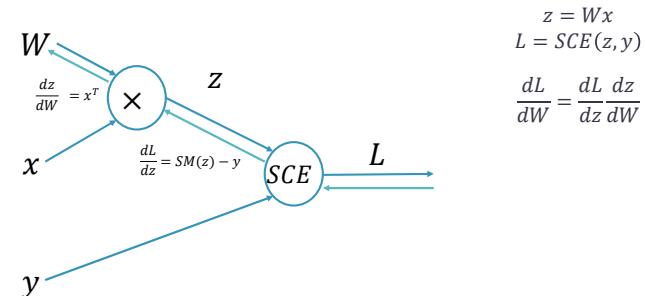
$$z = Wx$$

$$L = SCE(z, y)$$

$$\frac{dL}{dW} = \frac{dL}{dz} \frac{dz}{dW}$$

What is Backprop?

68



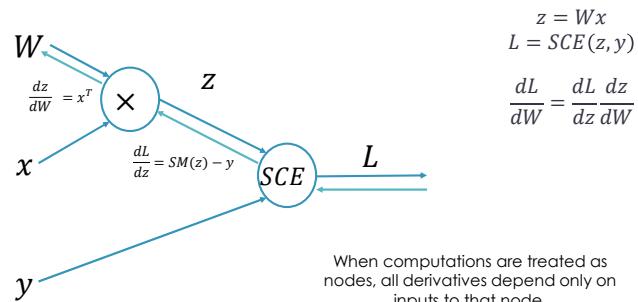
$$z = Wx$$

$$L = SCE(z, y)$$

$$\frac{dL}{dW} = \frac{dL}{dz} \frac{dz}{dW}$$

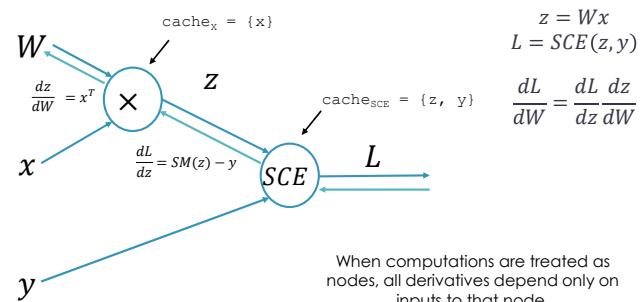
What is Backprop?

69



What is Backprop?

70



Gradient Descent Pseudocode (Updated)

71

```
for i in {0, ..., num_epochs}:
    for x, y in data:
         $\hat{y} = SM(Wx)$ 
         $L = CE(\hat{y}, y)$ 
         $\frac{dL}{dW} = \text{backprop}(L)$ 
         $W := W - \alpha \frac{dL}{dW}$ 
```

Gradient Descent Pseudocode (Updated)

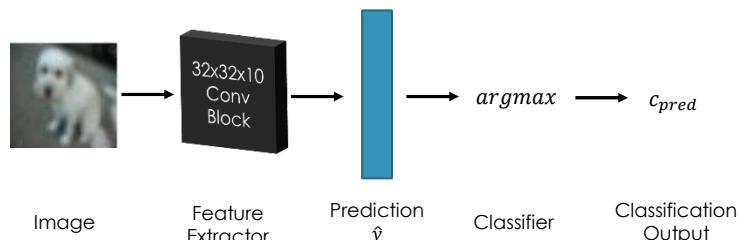
72

```
for i in {0, ..., num_epochs}:
    for x, y in data:
         $\hat{y} = SM(Wx)$ 
         $L = CE(\hat{y}, y)$ 
         $\frac{dL}{dW} = \text{backprop}(L)$ 
         $W := W - \alpha \frac{dL}{dW}$ 
```

MACHINE LEARNING!

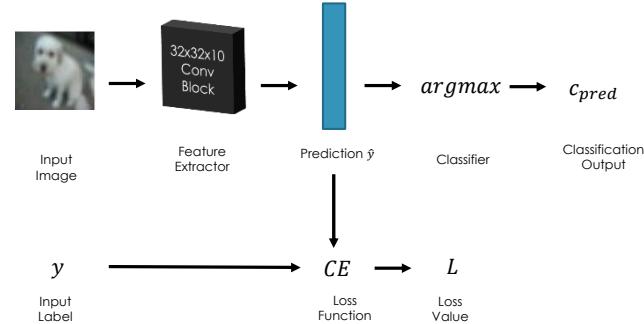
Simple Classification System

73



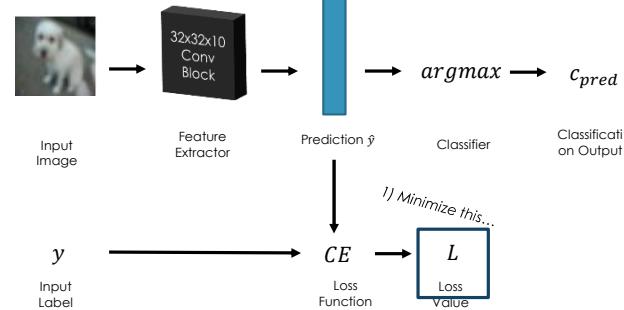
Simple Classification System (modified)

74



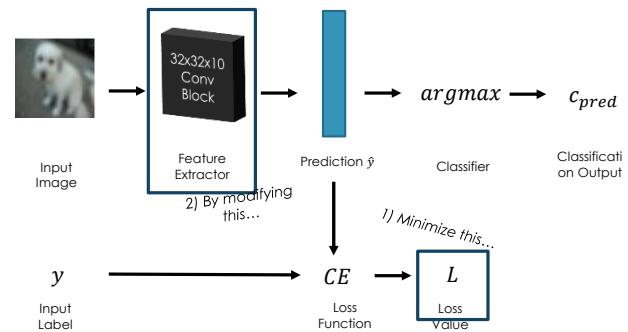
Simple Classification System (modified)

75



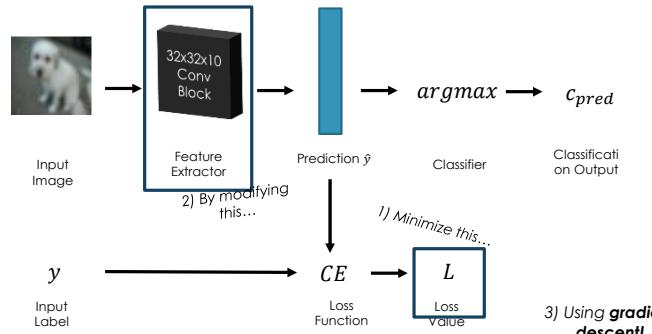
Simple Classification System (modified)

76



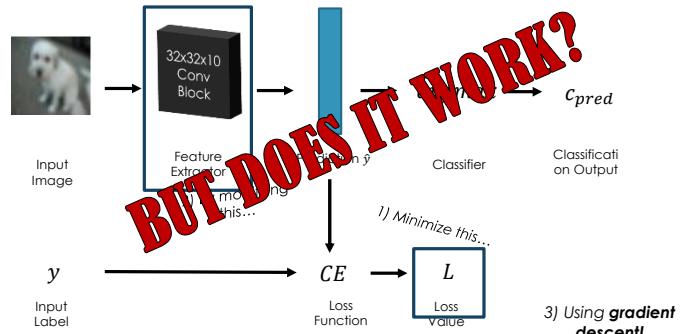
Simple Classification System (modified)

77



Simple Classification System (modified)

78



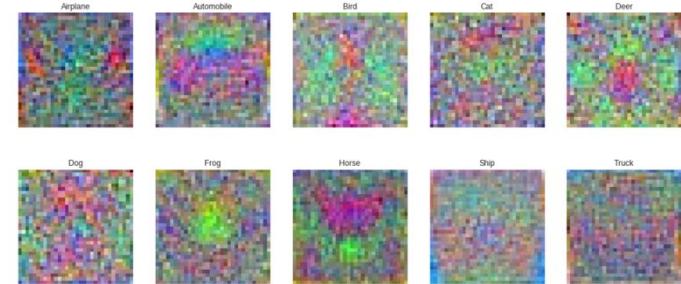
Simple System's Performance

79

- ~40% accuracy on CIFAR-10 test
 - Best class: Truck (~60%)
 - Worst class: Horse (~16%)
- Check out the model at: <https://tinyurl.com/cifar10>
- What about the filters? What do they look like?

Visualizing the Filters

80



Content

81

- Introduction
- Computer Vision & ML
- CNN
 - Architectures
 - Training

Convolutions

82

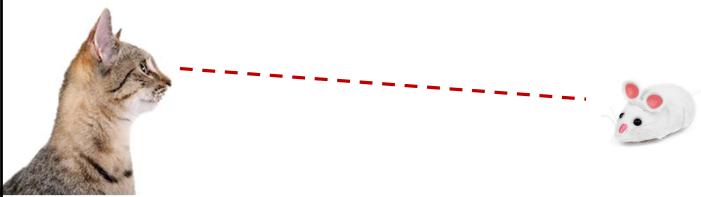
Building a stronger convolution-based feature extractor ??

Convolutions = Insights

More Convolutions = More Insights?

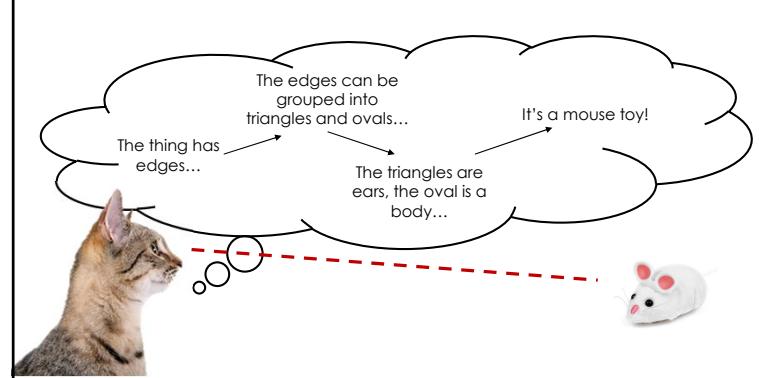
Recall Hubel and Weisel...

83



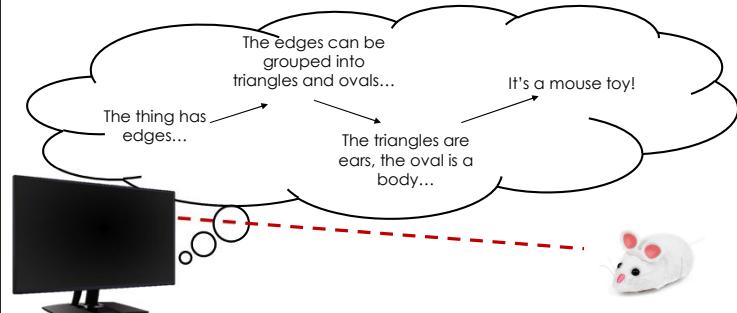
Recall Hubel and Weisel...

84



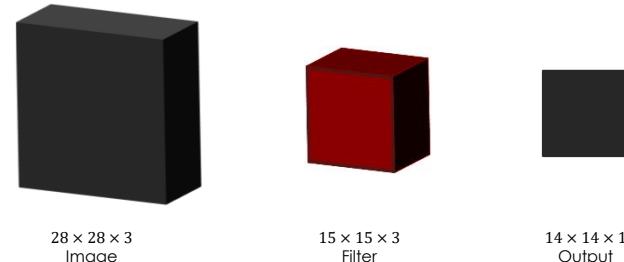
Recall Hubel and Weisel...

85



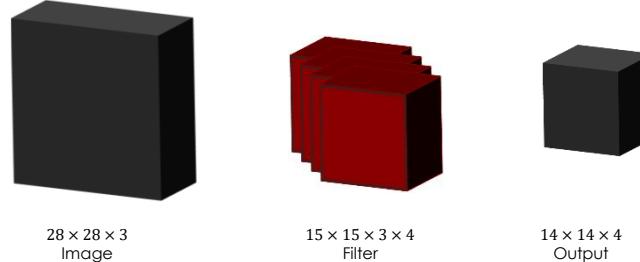
Convolutions Across Channels

86



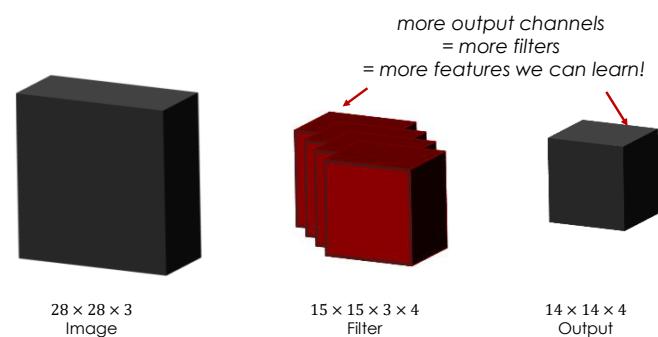
Convolutions Across Channels

87



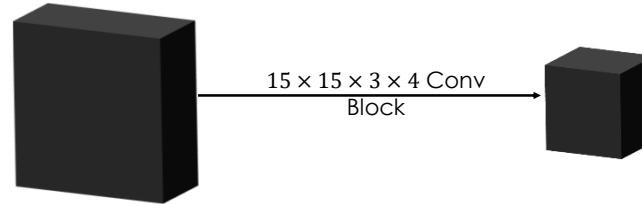
Convolutions Across Channels

88



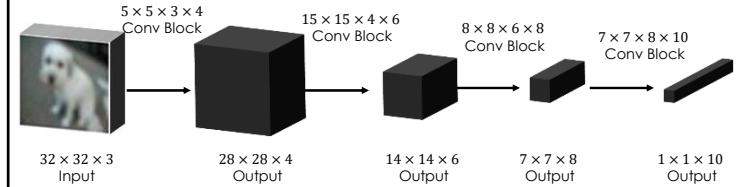
Convolutions Across Channels

89



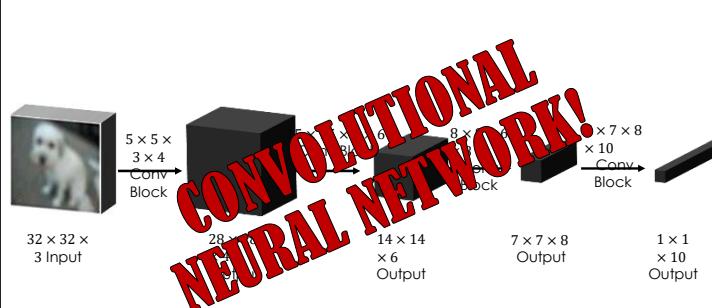
Stacking Convolutions

90



Stacking Convolutions

91



Convolutional Neural Networks (ConvNets)

92

- Neural networks which involve the stacking of multiple convolutional layers to produce output
- Often times end in fully-connected layers as the “classifier”

Why Do They Work So Well?



Why Do They Work So Well?



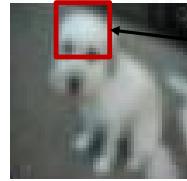
Why Do They Work So Well?



Why Do They Work So Well?



Why Do They Work So Well?

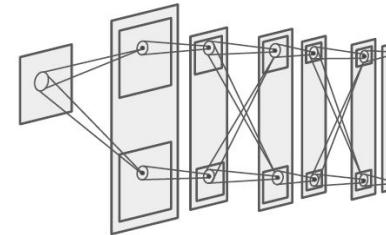


This is the neural network's "receptive field"—it's able to see!

History of ConvNets

A bit of history:

Neocognitron
[Fukushima 1980]



"sandwich" architecture (SCSCSC...)
simple cells: modifiable parameters
complex cells: perform pooling

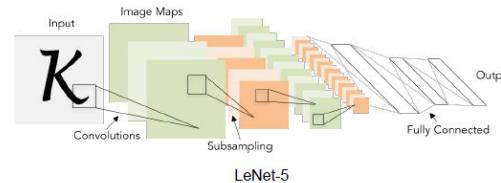
History of ConvNets

99

A bit of history:

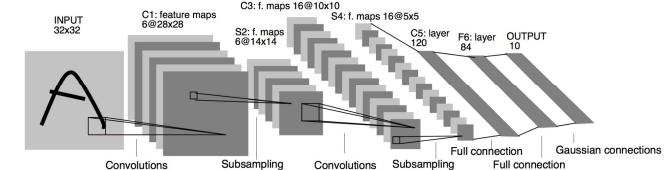
Gradient-based learning applied to document recognition

[LeCun, Bottou, Bengio, Haffner 1998]



History of ConvNets

100

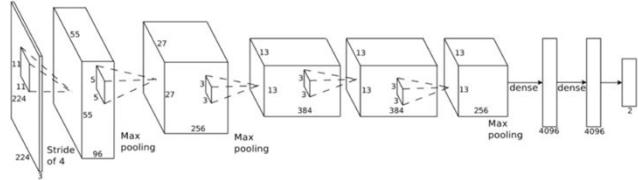


LeNet – 1998

MNIST is 0-9
0.7% test accuracy on MNIST

History of ConvNets

101

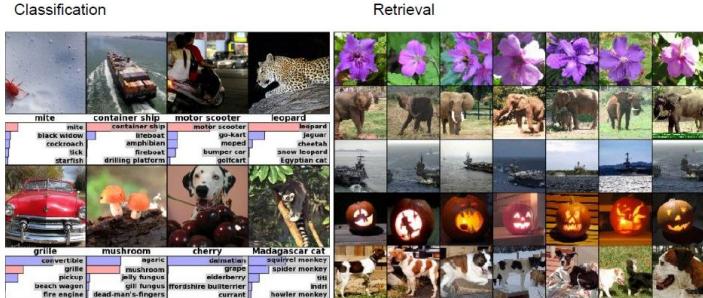


AlexNet – 2012

History of ConvNets

104

Fast-forward to today: ConvNets are everywhere



Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

History of ConvNets

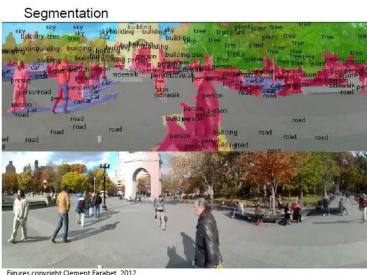
105

Fast-forward to today: ConvNets are everywhere



Figures copyright Shaoqing Ren, Kaiming He, Ross Girshick, Jian Sun, 2015. Reproduced with permission.

[Faster R-CNN: Ren, He, Girshick, Sun 2015]



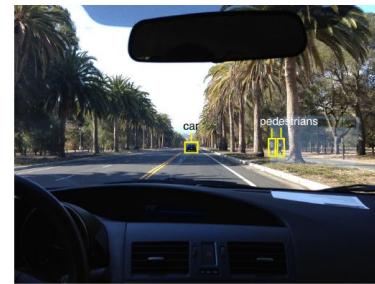
Figures copyright Clement Farabet, 2012. Reproduced with permission.

[Farabet et al., 2012]

History of ConvNets

106

Fast-forward to today: ConvNets are everywhere



self-driving cars



NVIDIA Tesla line
(these are the GPUs on rye01.stanford.edu)

This image by GBPublic_PR is licensed under CC-BY 2.0

Note that for embedded systems a typical setup would involve NVIDIA Tegras, with integrated GPU and ARM-based CPU cores.

History of ConvNets

107

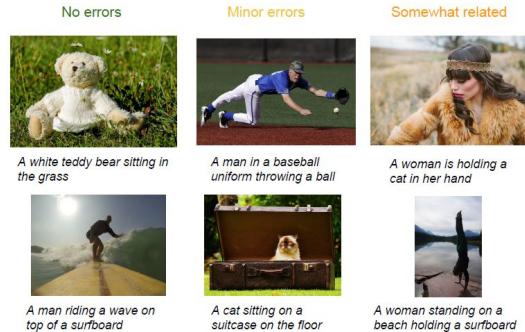


Image Captioning

[Vinyals et al., 2015]
[Karpathy and Fei-Fei, 2015]

All images are CC0 Public domain:
<https://creativecommons.org/publicdomain/zero/1.0/>
<https://creativecommons.org/publicdomain/zero/1.0/>
<https://creativecommons.org/licenses/by-sa/4.0/>
<https://creativecommons.org/licenses/by-sa/4.0/>
<https://creativecommons.org/licenses/by-sa/4.0/>

Captions generated by Justin Johnson using NeuralTalk2

History of ConvNets

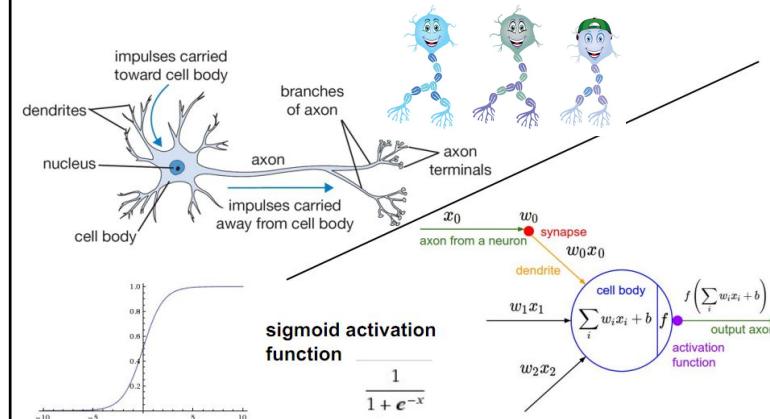
108



NN vs CNN

109

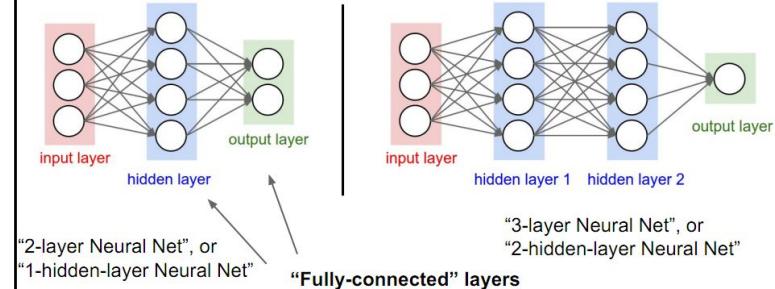
Neuron



Neural Networks

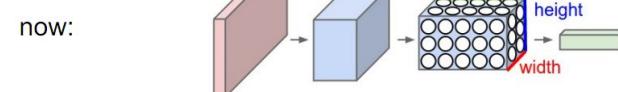
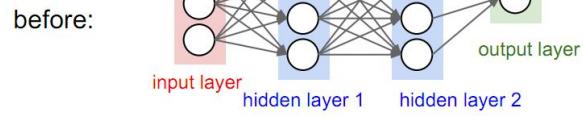
111

Neural Networks: Architectures



ConvNets

112



ConvNets

113

Convolutional Neural Networks
are just Neural Networks BUT:
Local connectivity

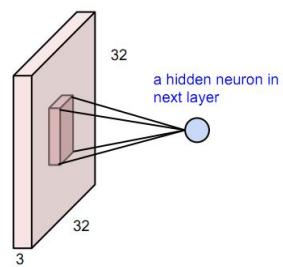
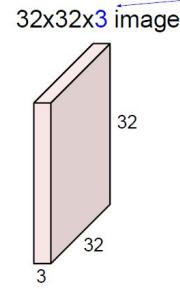


image: 32x32x3 volume
before: full connectivity: 32x32x3 weights
now: one neuron will connect to, e.g. 5x5x3 chunk and only have 5x5x3 weights.
note that connectivity is:
- local in space (5x5 inside 32x32)
- but full in depth (all 3 depth channels)

ConvNets

114

Convolution Layer



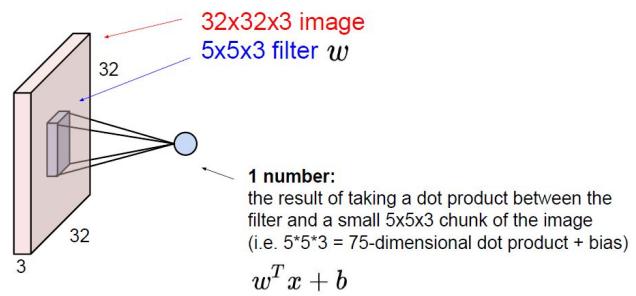
Filters always extend the full
depth of the input volume

Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

ConvNets

115

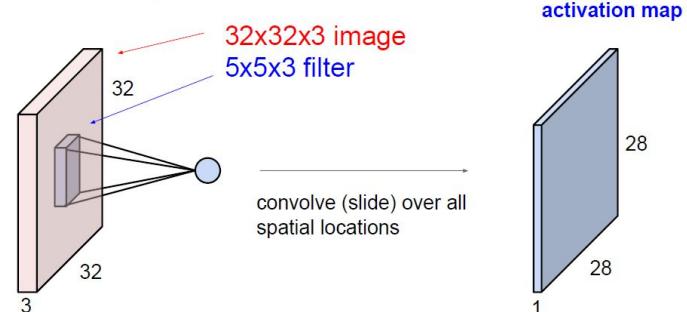
Convolution Layer



ConvNets

116

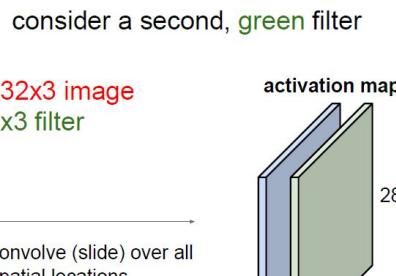
Convolution Layer



ConvNets

117

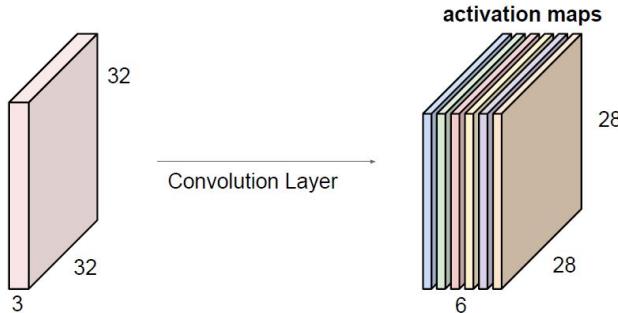
Convolution Layer



ConvNets

118

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

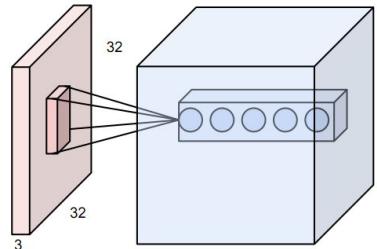


We stack these up to get a "new image" of size 28x28x6!

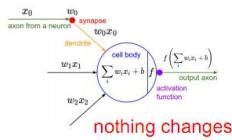
ConvNets

119

Convolutional Neural Networks
are just Neural Networks BUT:
Local connectivity



These form a single
[1 x 1 x depth]
"depth column" in the
output volume

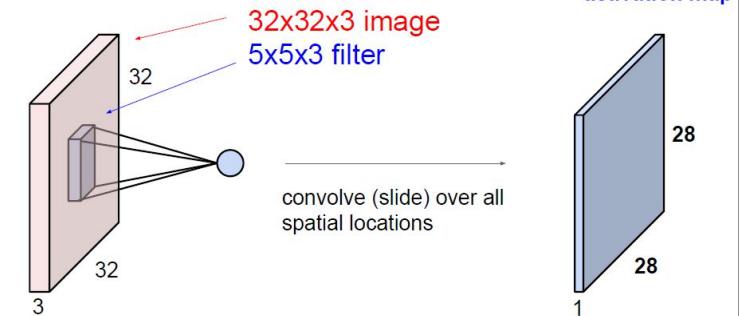


nothing changes

ConvNets

120

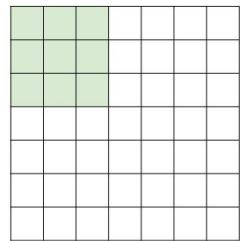
A closer look at spatial dimensions:



ConvNets

12
1

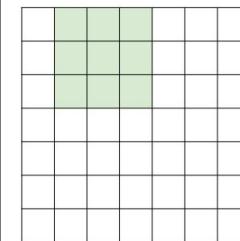
Replicate this column of hidden neurons across
space, with some **stride**.



ConvNets

12
2

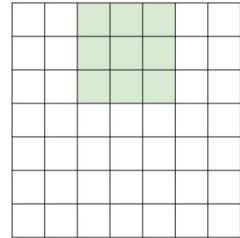
Replicate this column of hidden neurons across
space, with some **stride**.



ConvNets

12
3

Replicate this column of hidden neurons across space, with some **stride**.

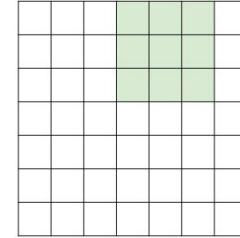


7x7 input
assume 3x3 connectivity, stride 1

ConvNets

12
4

Replicate this column of hidden neurons across space, with some **stride**.

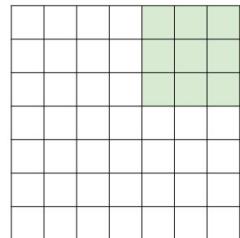


7x7 input
assume 3x3 connectivity, stride 1

ConvNets

12
5

Replicate this column of hidden neurons across space, with some **stride**.

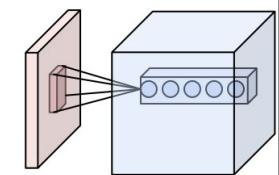


7x7 input
assume 3x3 connectivity, stride 1
=> **5x5 output**

ConvNets

126

Examples time:



Input volume: **32x32x3**
Receptive fields: **5x5, stride 1**
Number of neurons: **5**

Output volume: $(32 - 5) / 1 + 1 = 28$, so: **28x28x5**
How many weights for each of the 28x28x5 neurons?

ConvNets

127

In practice: Common to zero pad the border

(in each channel)

e.g. input 7x7
neuron with receptive field 3x3, stride 1
pad with 1 pixel border => what is the output?

7x7 => preserved size!

in general, common to see stride 1, size F, and zero-padding with $(F-1)/2$.
(Will preserve input size spatially)

ConvNets

12
8

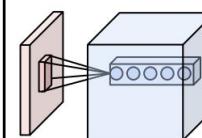
There's one more problem...

Assume input **[32 x 32 x3]**

30 neurons with receptive fields 5x5, applied at stride 1/pad1:

=> Output volume: [32 x 32 x 30] (3)

Each neuron has $5 \times 5 \times 3$ (=75) weights
=> Number of weights in such layer: $30720 \times 75 \approx 3$ million :)



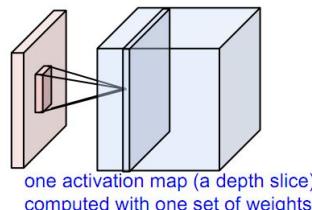
— Example trained weights
IDEA: let's not learn the same
thing across all spatial locations

ConvNets

12
9

These layers are called **Convolutional Layers**

1. Connect neurons only to local receptive fields
 2. Use the same neuron weight parameters for neurons in each “depth slice” (i.e. across spatial positions)



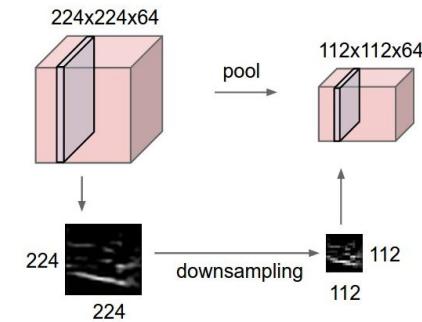
one activation map (a depth slice)
computed with one set of weights

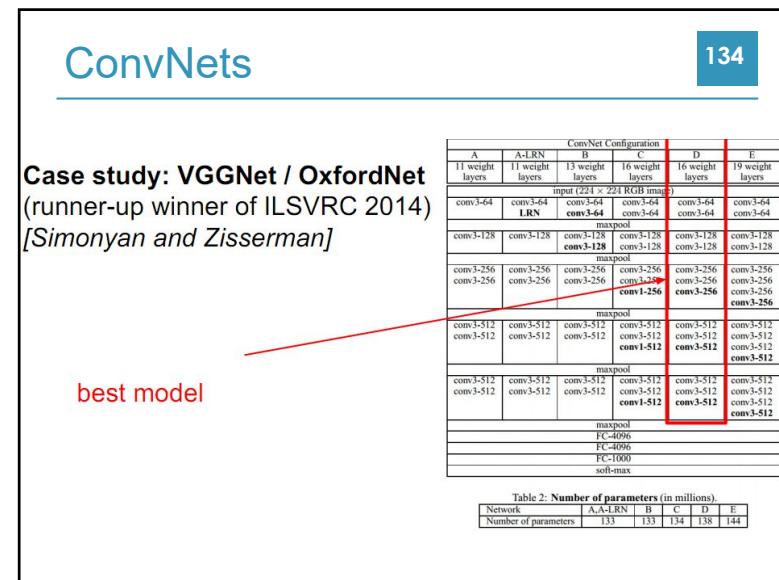
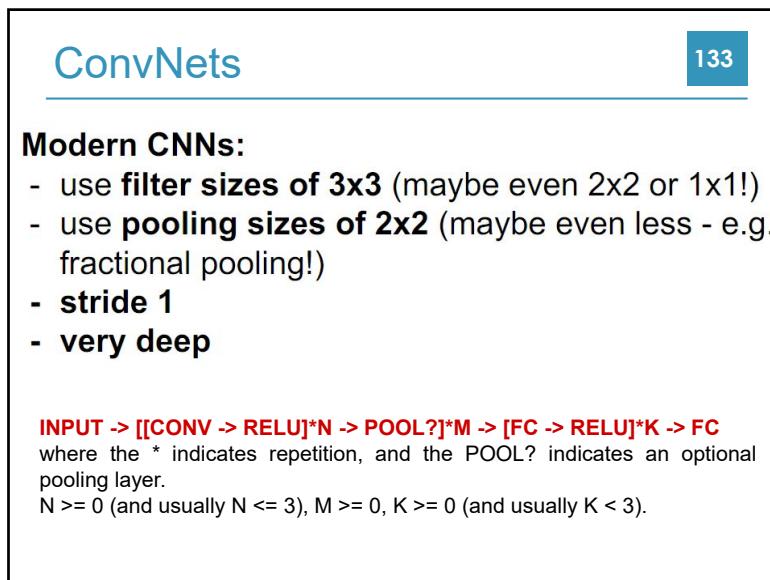
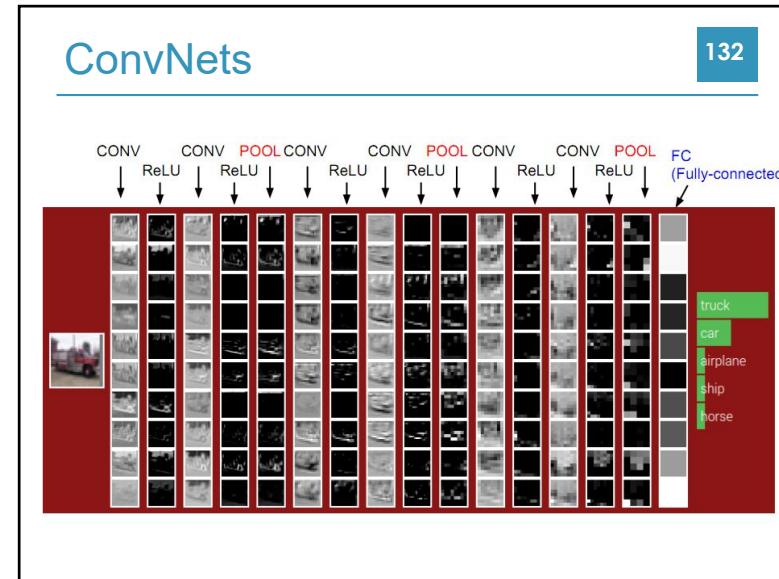
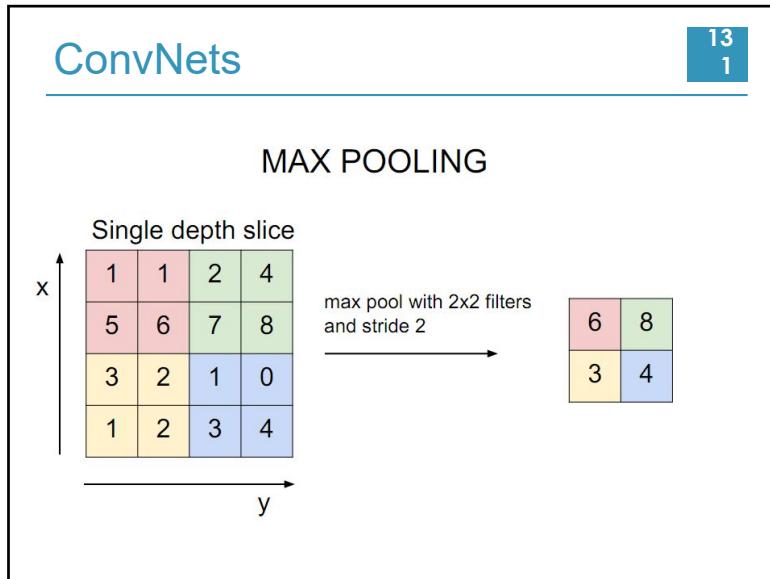
ConvNets

13
0

In ConvNet architectures, **Conv** layers are often followed by **Pool** layers

- convenience layer: makes the representations smaller and more manageable without losing too much information. Computes MAX operation (most common)





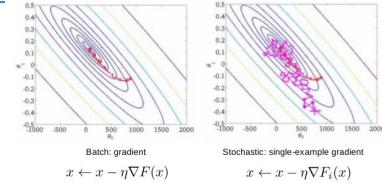
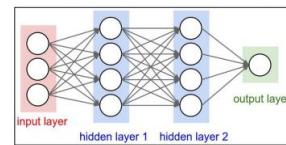
Training ConvNets

135

Mini-batch SGD

Loop:

1. **Sample** a batch of data
2. **Forward** prop it through the graph, get loss
3. **Backprop** to calculate the gradients
4. **Update** the parameters using the gradient



Training ConvNets

136

Activation Functions

Sigmoid

$$\sigma(x) = 1/(1 + e^{-x})$$

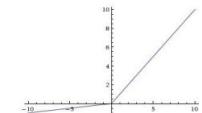
tanh

$$\tanh(x)$$

ReLU

$$\max(0, x)$$

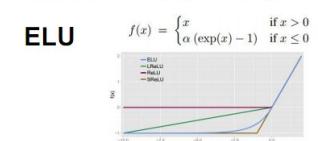
Leaky ReLU
 $\max(0.1x, x)$



Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

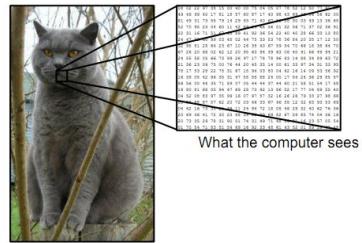


Data Augmentation

137

Data Augmentation

- i.e. simulating “fake” data
- explicitly encoding image transformations that shouldn’t change object identity.

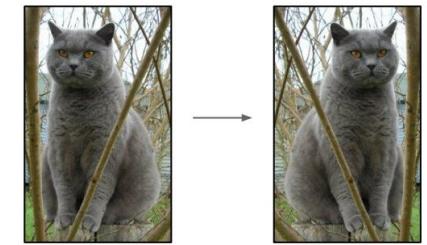


Data Augmentation

138

Data Augmentation

1. Flip horizontally



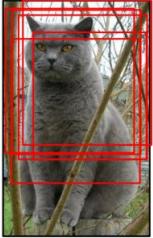
Data Augmentation

139

Data Augmentation

2. Random crops/scales

Sample these during training
(also helps a lot during test time)



e.g. common to see even up to 150 crops used

Data Augmentation

140

Data Augmentation

3.

Random mix/combinations of :

- translation
- rotation
- stretching
- shearing,
- lens distortions, ... (go crazy)

Data Augmentation

141

Data Augmentation

4. Color jittering

(maybe even contrast jittering, etc.)

- Simple: Change contrast small amounts, jitter the color distributions, etc.
- Vignette,... (go crazy)



Transfer learning

142

“You need a lot of data if you want to train/use CNNs”

Transfer learning

143

Transfer Learning

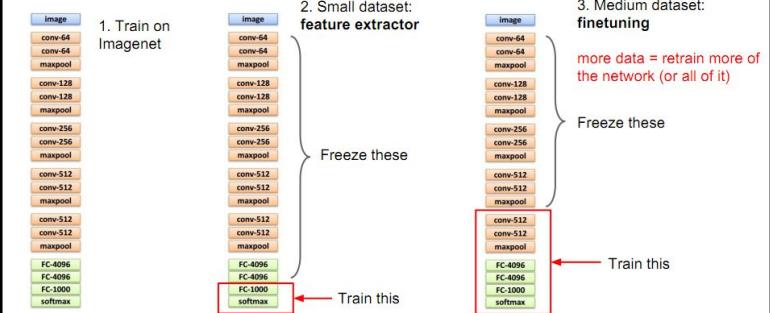
“You need a lot of data if you want to train/learn CNNs”

BUSTED

Transfer learning

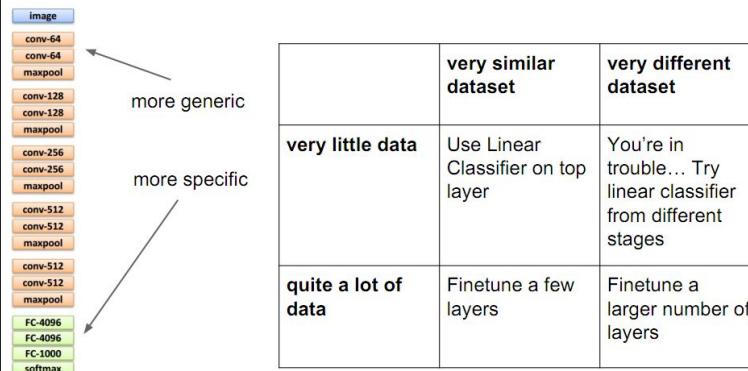
144

Transfer Learning with CNNs



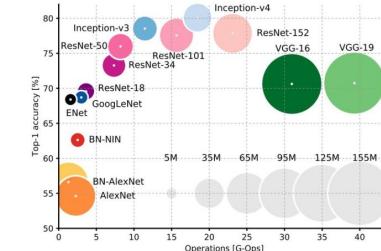
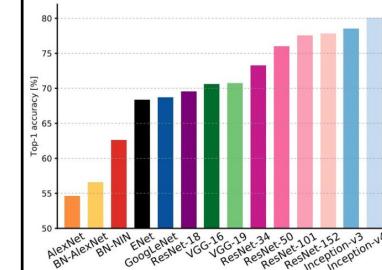
Transfer learning

145



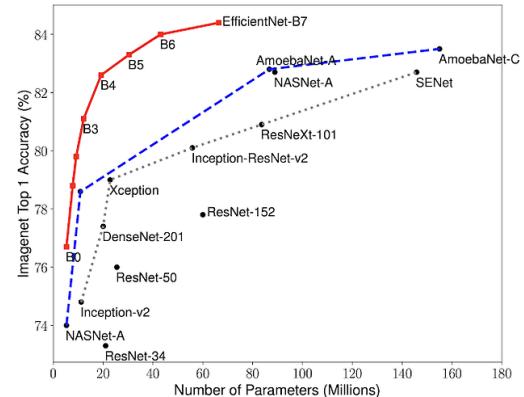
CNN Performances

146



CNN Performances

147



[EfficientNet: Improving Accuracy and Efficiency through AutoML and Model Scaling](#), 5/2019

What is CNN dev?

14
8

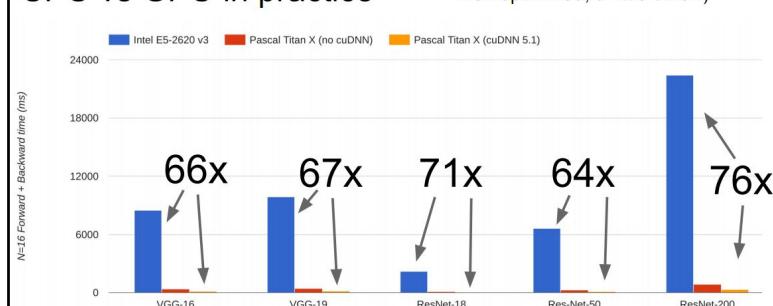
- Define the objective
 - What is the input/output?
 - What is the loss/objective function?
- Create the architecture
 - How many conv layers?
 - What size are the convolutions?
 - How many fully-connected layers?
- Define hyperparameters
 - What is the learning rate?
- Train and evaluate
 - How did we do?
 - How can we do better?

CPU vs GPU

14
9

CPU vs GPU in practice

(CPU performance not well-optimized, a little unfair)



CPU vs GPU

150

	Cores	Clock Speed	Memory	Price	Speed
CPU (Intel Core i7-7700k)	4 (8 threads with hyperthreading)	4.2 GHz	System RAM	\$339	~540 GFLOPs FP32
GPU (NVIDIA GTX 1080 Ti)	3584	1.6 GHz	11 GB GDDR5 X	\$699	~11.4 TFLOPs FP32
TPU NVIDIA TITAN V	5120 CUDA, 640 Tensor	1.5 GHz	12GB HBM2	\$2999	~14 TFLOPs FP32 ~112 TFLOP FP16
TPU Google Cloud TPU	?	?	64 GB HBM	\$6.50 per hour	~180 TFLOP

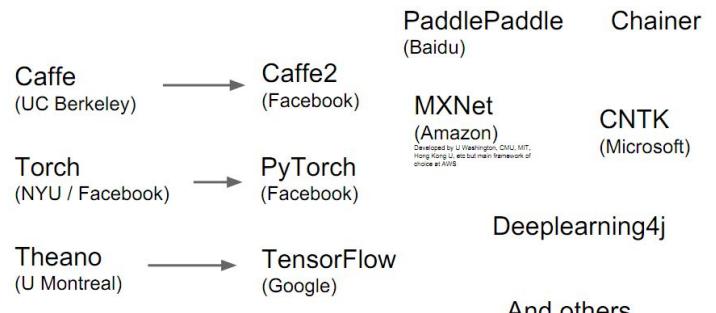
CPU: Fewer cores, but each core is much faster and much more capable; great at sequential tasks

GPU: More cores, but each core is much slower and "dumber"; great for parallel tasks

TPU: Specialized hardware for deep learning

A zoo of frameworks

15
1



References

152

- Fei-Fei Li & Justin Johnson & Serena Yeung. *CS231n – 2019 - lecture 4, lecture 5*
- Shubhang Desai, Ranjay Krishna, and Juan Carlos Niebles. *Lecture: Computer Vision and Machine Learning*. CS131 Stanford Vision and Learning Lab