



10-301/10-601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Convolutional Neural Networks (CNNs)

Matt Gormley, Henry Chai, Hoda Heidari

Lecture 17

Mar. 20, 2024

Reminders

- **Homework 6: Learning Theory & Generative Models**
 - **Out: Mon, Mar 18**
 - **Due: Sun, Mar 24 at 11:59pm**

THE BIG PICTURE

ML Big Picture

Learning Paradigms:

What data is available and when? What form of prediction?

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

Theoretical Foundations:

What principles guide learning?

- probabilistic
- information theoretic
- evolutionary search
- ML as optimization

Problem Formulation:

What is the structure of our output prediction?

boolean	Binary Classification
categorical	Multiclass Classification
ordinal	Ordinal Classification
real	Regression
ordering	Ranking
multiple discrete	Structured Prediction
multiple continuous	(e.g. dynamical systems)
both discrete & cont.	(e.g. mixed graphical models)

Facets of Building ML Systems:

How to build systems that are robust, efficient, adaptive, effective?

1. Data prep
2. Model selection
3. Training (optimization / search)
4. Hyperparameter tuning on validation data
5. (Blind) Assessment on test data

Big Ideas in ML:

Which are the ideas driving development of the field?

- inductive bias
- generalization / overfitting
- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

Application Areas

Key challenges?

NLP, Speech, Computer Vision, Robotics, Medicine, Search

Classification and Regression: The Big Picture

Recipe for Machine Learning

1. Given data $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$
2. (a) Choose a decision function $h_{\theta}(\mathbf{x}) = \dots$
(parameterized by θ)
(b) Choose an objective function $J_{\mathcal{D}}(\theta) = \dots$
(relies on data)
3. Learn by choosing parameters that optimize the objective $J_{\mathcal{D}}(\theta)$

$$\hat{\theta} \approx \underset{\theta}{\operatorname{argmin}} J_{\mathcal{D}}(\theta)$$

4. Predict on new test example \mathbf{x}_{new} using $h_{\theta}(\cdot)$

$$\hat{y} = h_{\theta}(\mathbf{x}_{\text{new}})$$

Optimization Method

- Gradient Descent: $\theta \rightarrow \theta - \gamma \nabla_{\theta} J(\theta)$
- SGD: $\theta \rightarrow \theta - \gamma \nabla_{\theta} J^{(i)}(\theta)$
for $i \sim \text{Uniform}(1, \dots, N)$
where $J(\theta) = \frac{1}{N} \sum_{i=1}^N J^{(i)}(\theta)$
- mini-batch SGD
- closed form
 1. compute partial derivatives
 2. set equal to zero and solve

Decision Functions

- Perceptron: $h_{\theta}(\mathbf{x}) = \operatorname{sign}(\theta^T \mathbf{x})$
- Linear Regression: $h_{\theta}(\mathbf{x}) = \theta^T \mathbf{x}$
- Discriminative Models: $h_{\theta}(\mathbf{x}) = \underset{y}{\operatorname{argmax}} p_{\theta}(y | \mathbf{x})$
 - Logistic Regression: $p_{\theta}(y = 1 | \mathbf{x}) = \sigma(\theta^T \mathbf{x})$
 - Neural Net (classification):
 $p_{\theta}(y = 1 | \mathbf{x}) = \sigma((\mathbf{W}^{(2)})^T \sigma((\mathbf{W}^{(1)})^T \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$
- Generative Models: $h_{\theta}(\mathbf{x}) = \underset{y}{\operatorname{argmax}} p_{\theta}(\mathbf{x}, y)$
 - Naive Bayes: $p_{\theta}(\mathbf{x}, y) = p_{\theta}(y) \prod_{m=1}^M p_{\theta}(x_m | y)$

Objective Function

- MLE: $J(\theta) = - \sum_{i=1}^N \log p(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$
- MCLE: $J(\theta) = - \sum_{i=1}^N \log p(\mathbf{y}^{(i)} | \mathbf{x}^{(i)})$
- L2 Regularized: $J'(\theta) = J(\theta) + \lambda \|\theta\|_2^2$
(same as Gaussian prior $p(\theta)$ over parameters)
- L1 Regularized: $J'(\theta) = J(\theta) + \lambda \|\theta\|_1$
(same as Laplace prior $p(\theta)$ over parameters)

Backpropagation and Deep Learning

Convolutional neural networks (CNNs) and **recurrent neural networks (RNNs)** are simply fancy computation graphs (aka. hypotheses or decision functions).

Our recipe also applies to these models and (again) relies on the **backpropagation algorithm** to compute the necessary gradients.

BACKGROUND: COMPUTER VISION

Example: Image Classification

- ImageNet LSVRC-2011 contest:
 - **Dataset:** 1.2 million labeled images, 1000 classes
 - **Task:** Given a new image, label it with the correct class
 - **Multiclass** classification problem
- Examples from <http://image-net.org/>

Bird

Warm-blooded egg-laying vertebrates characterized by feathers and forelimbs modified as wings

2126 pictures

92.85% Popularity Percentile

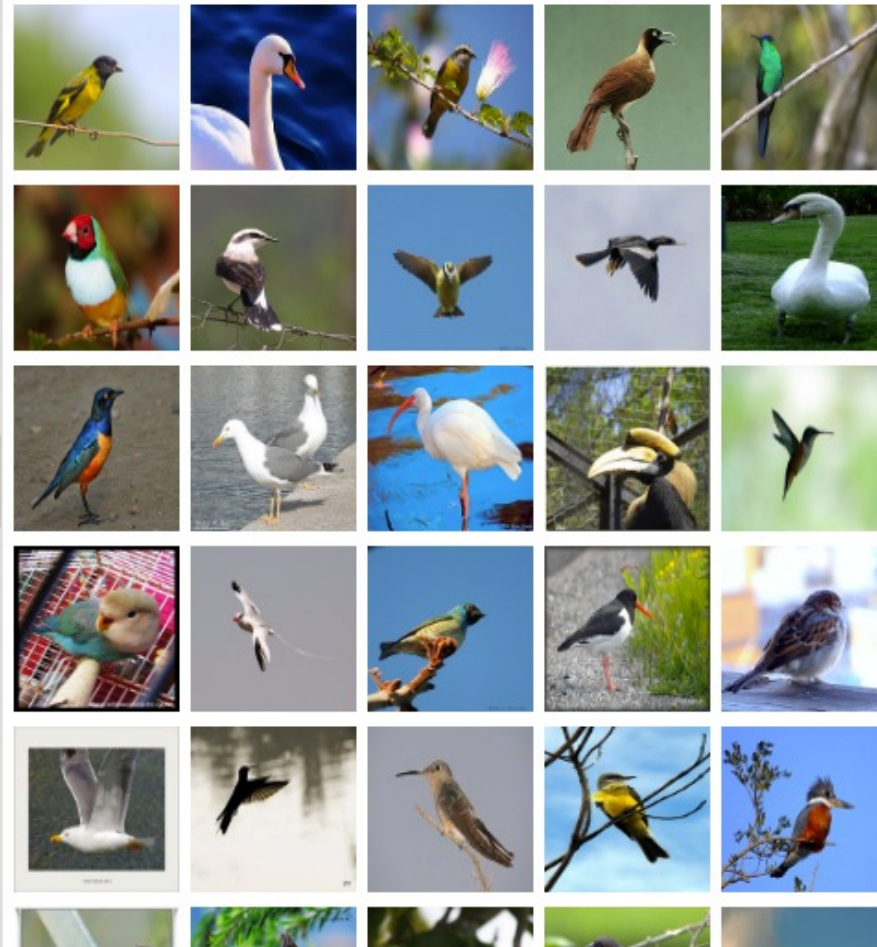


- marine animal, marine creature, sea animal, sea creature (1)
- scavenger (1)
- biped (0)
- predator, predatory animal (1)
- larva (49)
- acrodont (0)
- feeder (0)
- stunt (0)
- chordate (3087)**
 - tunicate, urochordate, urochord (6)
 - cephalochordate (1)
 - vertebrate, craniate (3077)
 - mammal, mammalian (1169)
 - bird (871)**
 - dickeybird, dickey-bird, dickybird, dicky-bird (0)
 - cock (1)
 - hen (0)
 - nester (0)
 - night bird (1)
 - bird of passage (0)
 - protoavis (0)
 - archaeopteryx, archeopteryx, Archaeopteryx lithographi
 - Sinornis (0)
 - Ibero-mesornis (0)
 - archaeornis (0)
 - ratite, ratite bird, flightless bird (10)
 - carinate, carinate bird, flying bird (0)
 - passerine, passeriform bird (279)
 - nonpasserine bird (0)
 - bird of prey, raptor, raptorial bird (80)
 - gallinaceous bird, gallinacean (114)

Treemap Visualization

Images of the Synset

Downloads



German iris, *Iris kochii*

Iris of northern Italy having deep blue-purple flowers; similar to but smaller than *Iris germanica*

469 pictures

49.6% Popularity Percentile



- ... halophyte (0)
- ... succulent (39)
- ... cultivar (0)
- ... cultivated plant (0)
- ... weed (54)
- ... evergreen, evergreen plant (0)
- ... deciduous plant (0)
- ... vine (272)
- ... creeper (0)
- ... woody plant, ligneous plant (1868)
- ... geophyte (0)
- ... desert plant, xerophyte, xerophytic plant, xerophile, xerophilic mesophyte, mesophytic plant (0)
- ... aquatic plant, water plant, hydrophyte, hydrophytic plant (11)
- ... tuberous plant (0)
- ... bulbous plant (179)
 - ... iridaceous plant (27)
 - ... iris, flag, fleur-de-lis, sword lily (19)
 - ... bearded iris (4)
 - ... Florentine iris, orris, *Iris germanica florentina*, *Iris German iris, Iris germanica* (0)
 - ... **German iris, *Iris kochii*** (0)
 - ... Dalmatian iris, *Iris pallida* (0)
 - ... beardless iris (4)
 - ... bulbous iris (0)
 - ... dwarf iris, *Iris cristata* (0)
 - ... stinking iris, gladdon, gladdon iris, stinking gladwyn, Persian iris, *Iris persica* (0)
 - ... yellow iris, yellow flag, yellow water flag, *Iris pseudo*
 - ... dwarf iris, vernal iris, *Iris verna* (0)
 - ... blue flag, *Iris versicolor* (0)

Treemap Visualization

Images of the Synset

Downloads



Court, courtyard

An area wholly or partly surrounded by walls or buildings; "the house was built around an inner court"

165 pictures

92.61% Popularity Percentile



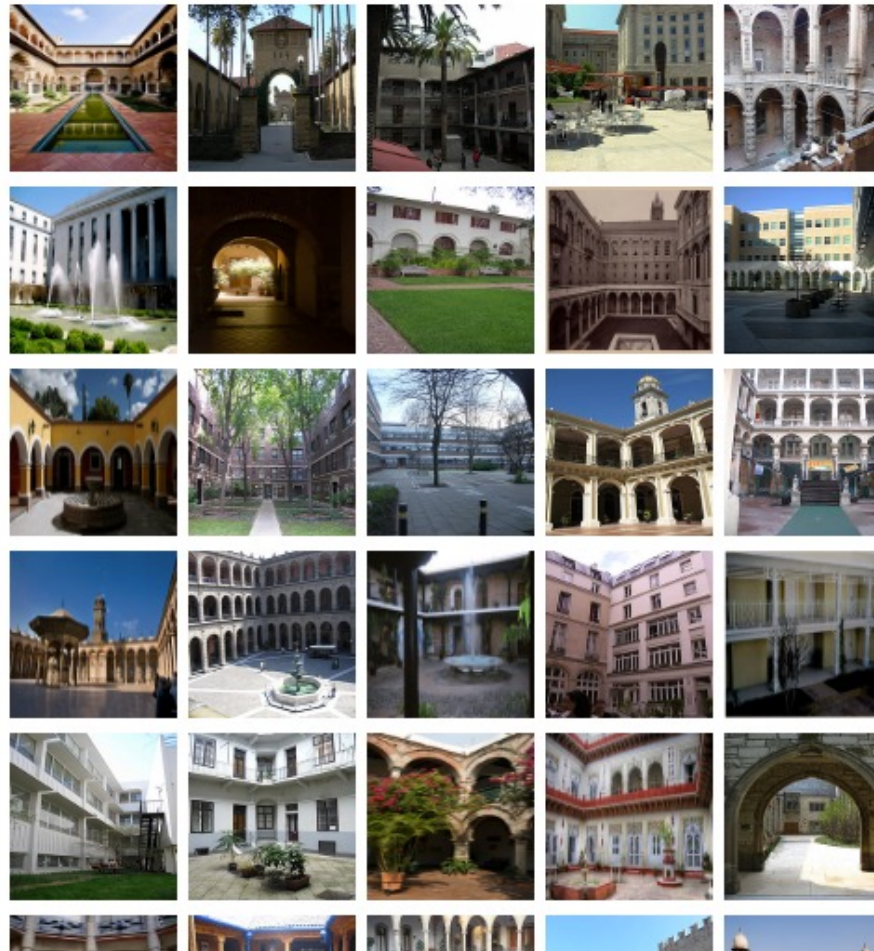
Numbers in brackets: (the number of synsets in the subtree).

- ImageNet 2011 Fall Release (32326)
 - plant, flora, plant life (4486)
 - geological formation, formation (175)
 - natural object (1112)
 - sport, athletics (176)
 - artifact, artefact (10504)
 - instrumentality, instrumentation (5494)
 - structure, construction (1405)
 - airdock, hangar, repair shed (0)
 - altar (1)
 - arcade, colonnade (1)
 - arch (31)
 - area (344)
 - aisle (0)
 - auditorium (1)
 - baggage claim (0)
 - box (1)
 - breakfast area, breakfast nook (0)
 - bullpen (0)
 - chancel, sanctuary, bema (0)
 - choir (0)
 - corner, nook (2)
 - court, courtyard (6)
 - atrium (0)
 - bailey (0)
 - cloister (0)
 - food court (0)
 - forecourt (0)
 - narvis (0)

Treemap Visualization

Images of the Synset

Downloads

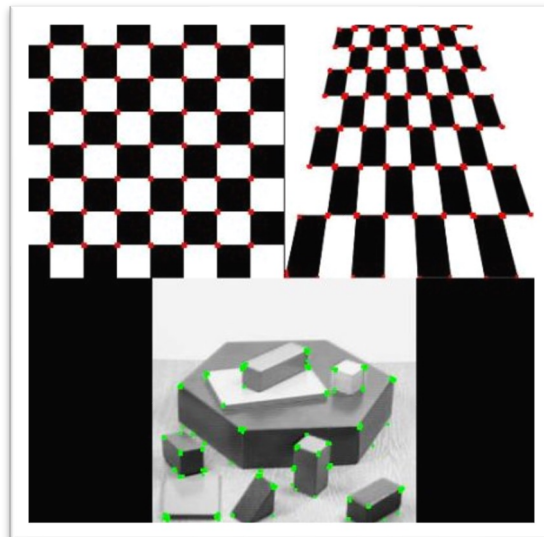


Feature Engineering for CV

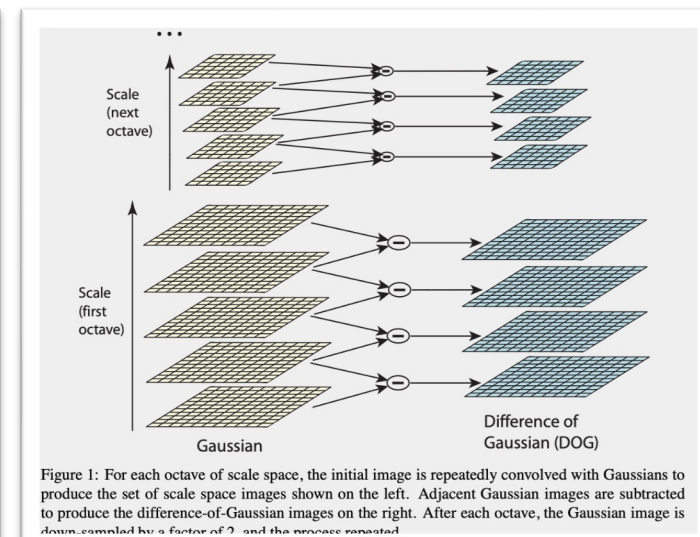
Edge detection (Canny)



Corner Detection (Harris)



Scale Invariant Feature Transform (SIFT)



Example: Image Classification

CNN for Image Classification

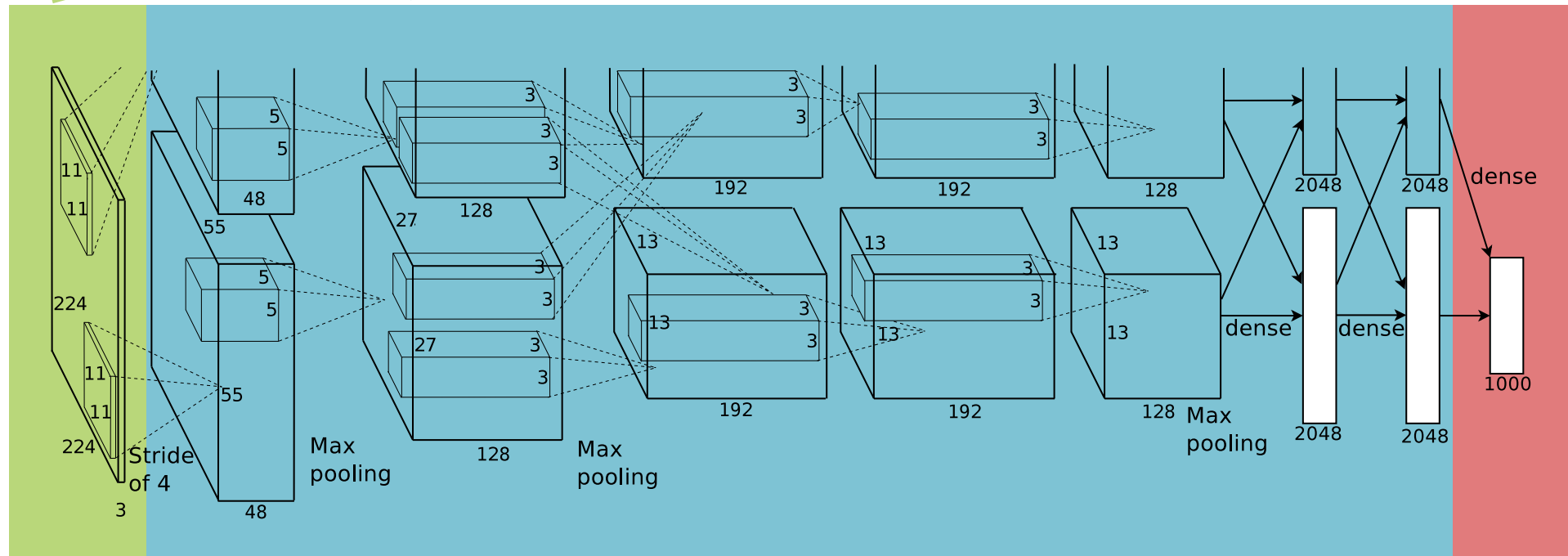
(Krizhevsky, Sutskever & Hinton, 2012)

15.3% error on ImageNet LSVRC-2012 contest

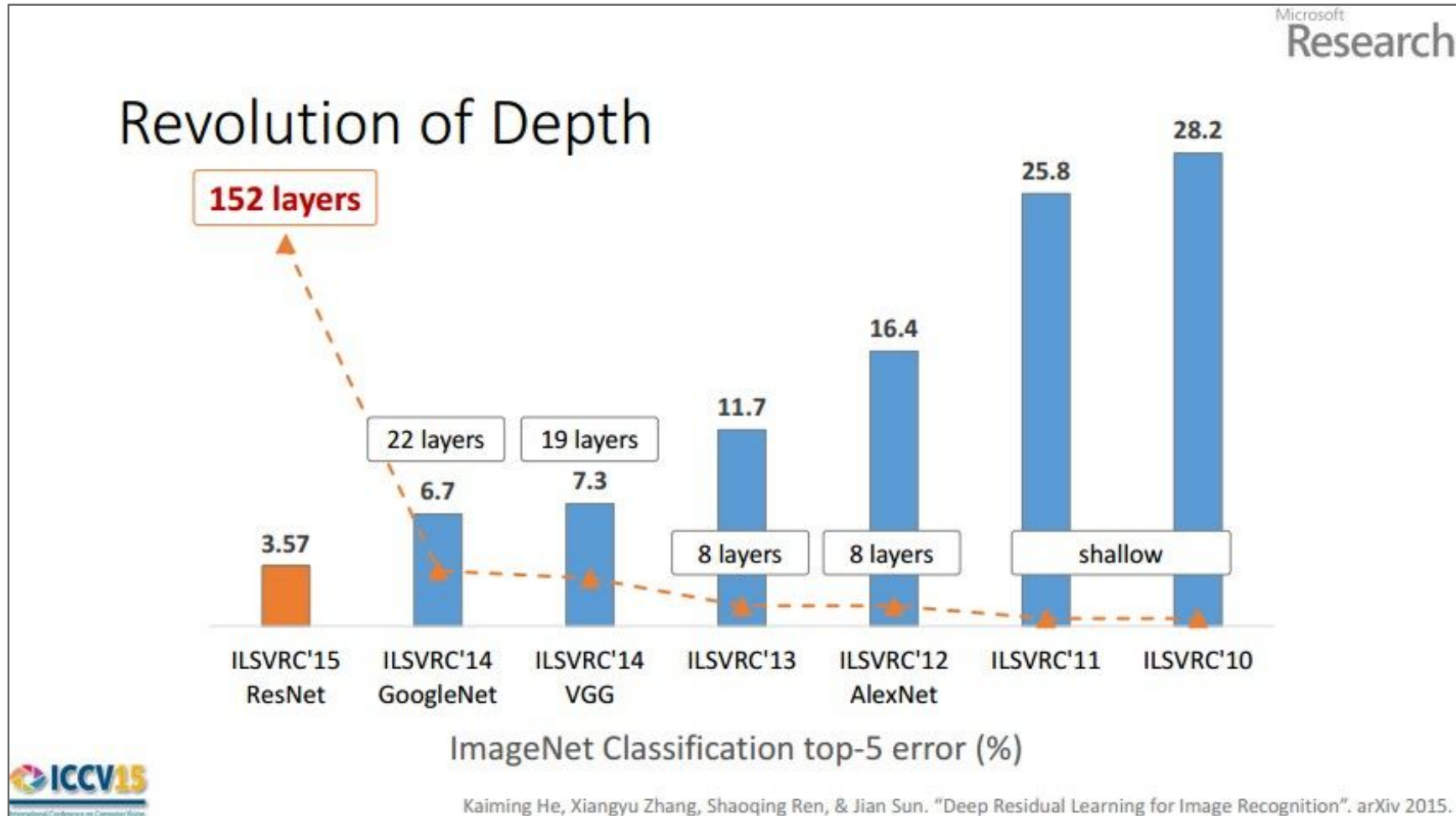
Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax



CNNs for Image Recognition



Feed-forward Neural Networks for Computer Vision

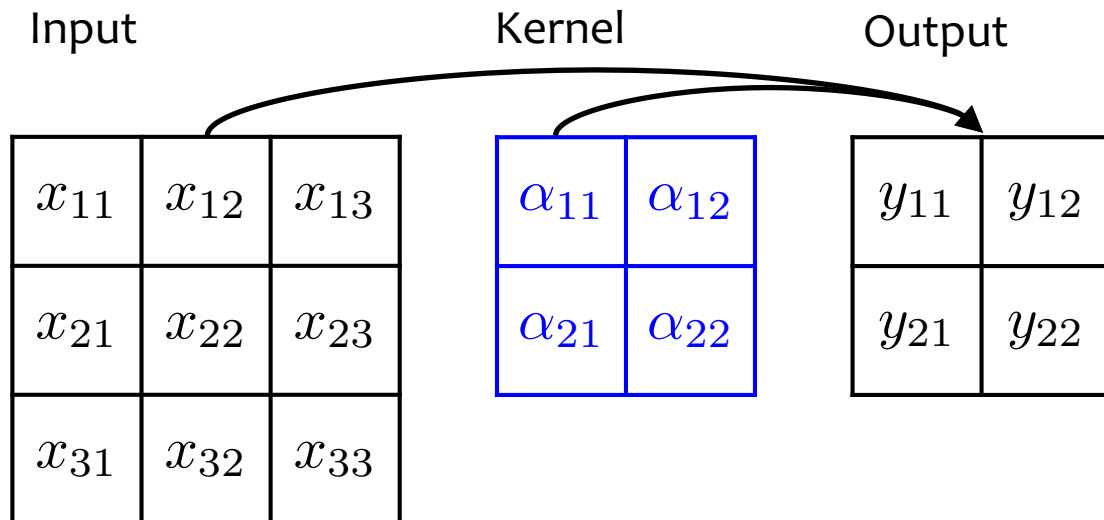
Feed-forward Neural Networks for Computer Vision

CONVOLUTION

2D Convolution

- Basic idea:
 - Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
 - Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation
- Key point:
 - Different convolutions extract different types of low-level “features” from an image
 - All that we need to vary to generate these different features is the weights of F

Example: 1 input channel, 1 output channel



$$y_{11} = \alpha_{11}x_{11} + \alpha_{12}x_{12} + \alpha_{21}x_{21} + \alpha_{22}x_{22} + \alpha_0$$

$$y_{12} = \alpha_{11}x_{12} + \alpha_{12}x_{13} + \alpha_{21}x_{22} + \alpha_{22}x_{23} + \alpha_0$$

$$y_{21} = \alpha_{11}x_{21} + \alpha_{12}x_{22} + \alpha_{21}x_{31} + \alpha_{22}x_{32} + \alpha_0$$

$$y_{22} = \alpha_{11}x_{22} + \alpha_{12}x_{23} + \alpha_{21}x_{32} + \alpha_{22}x_{33} + \alpha_0$$

2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	1	1
0	1	0

Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	1	1
0	1	0

Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

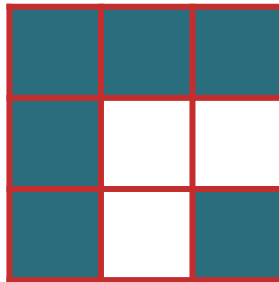
2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution



Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

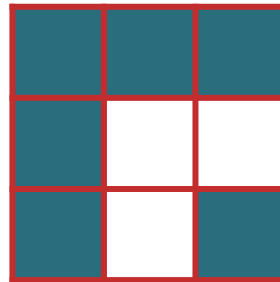
2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution



Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

			0	0	0	0
	1	1	1	1	1	0
	1		0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

Convolved Image

3				

2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	0	0	0	0
0	0	1	1	1	1	0
0	0	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	0	0
0	0	0

Convolved Image

3	2	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0				0	0
0	1		1	1	1	0
0	1		0		0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

Convolved Image

3	2	2		

2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	1	1	1	0
0	1	1	1	1	1	0
0	1	0	1	1	1	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

1	1	1
1	0	0
1	0	1

Convolved Image

3	2	2	3	
	1			1
			1	1
		1	1	1
	1	1	1	1

2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	0			
0	1	1	1		1	0
0	1	0	0		0	
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

Convolved Image

3	2	2	3	1

2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	0	0	0	0
0	0	0	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	0	0
0	0	0

Convolved Image

3	2	2	3	1
2	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	0	0	0	0
0	0	0	0	1	1	0
0	0	0	0	1	0	0
0	0	0	0	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	0	0
0	0	0

Convolved Image

3	2	2	3	1
2	0			

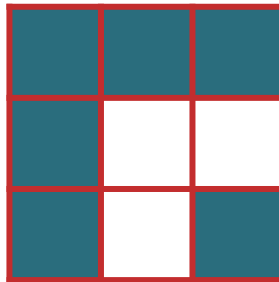
2D Convolution

- Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
- Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution



Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

Padding

Suppose you want to preserve the size of the original input image in your convolved image.

You can accomplish this by padding your input image with zeros.

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Identity Convolution

0	0	0
0	1	0
0	0	0

Convolved Image

1	1	1	1	1
1	0	0	1	0
1	0	1	0	0
1	1	0	0	0
1	0	0	0	0

Padding

Suppose you want to preserve the size of the original input image in your convolved image.

You can accomplish this by padding your input image with zeros.

Input Image

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	1	1	1	1	1	0	0
0	0	1	0	0	1	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

Identity Convolution

0	0	0
0	1	0
0	0	0

Convolved Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Kernels for Image Processing

A **convolution matrix** (aka. kernel) is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	1	1	1	1	1	0	0
0	0	1	0	0	1	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

Identity
Convolution

0	0	0
0	1	0
0	0	0

Convolved Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Kernels for Image Processing

A **convolution matrix** (aka. kernel) is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	1	1	1	1	1	0	0
0	0	1	0	0	1	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

Blurring
Convolution

.1	.1	.1
.1	.2	.1
.1	.1	.1

Convolved Image

.1	.2	.3	.3	.3	.2	.1
.2	.4	.5	.5	.5	.4	.1
.3	.4	.2	.3	.6	.3	.1
.3	.5	.4	.4	.2	.1	0
.3	.5	.6	.2	.1	0	0
.2	.4	.3	.1	0	0	0
.1	.1	.1	0	0	0	0

Kernels for Image Processing

A **convolution matrix** (aka. kernel) is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	1	1	1	1	1	0	0
0	0	1	0	0	1	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

Vertical
Edge
Detector

-1	0	1
-1	0	1
-1	0	1

Convolved Image

-1	-1	0	0	0	1	1
-2	-1	1	-1	0	2	1
-3	-1	1	-1	1	2	1
-3	-1	2	0	1	1	0
-3	-1	2	1	1	0	0
-2	-1	2	1	0	0	0
-1	0	1	0	0	0	0

Kernels for Image Processing

A **convolution matrix** (aka. kernel) is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	1	1	1	1	1	0	0
0	0	1	0	0	1	0	0	0
0	0	1	0	1	0	0	0	0
0	0	1	1	0	0	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

Horizontal
Edge
Detector

-1	-1	-1
0	0	0
1	1	1

Convolved Image

-1	-2	-3	-3	-3	-2	-1
-1	-1	-1	-1	-1	-1	0
0	1	1	2	2	2	1
0	-1	-1	0	1	1	0
0	0	1	1	1	0	0
1	2	2	1	0	0	0
1	1	1	0	0	0	0

Convolution Examples

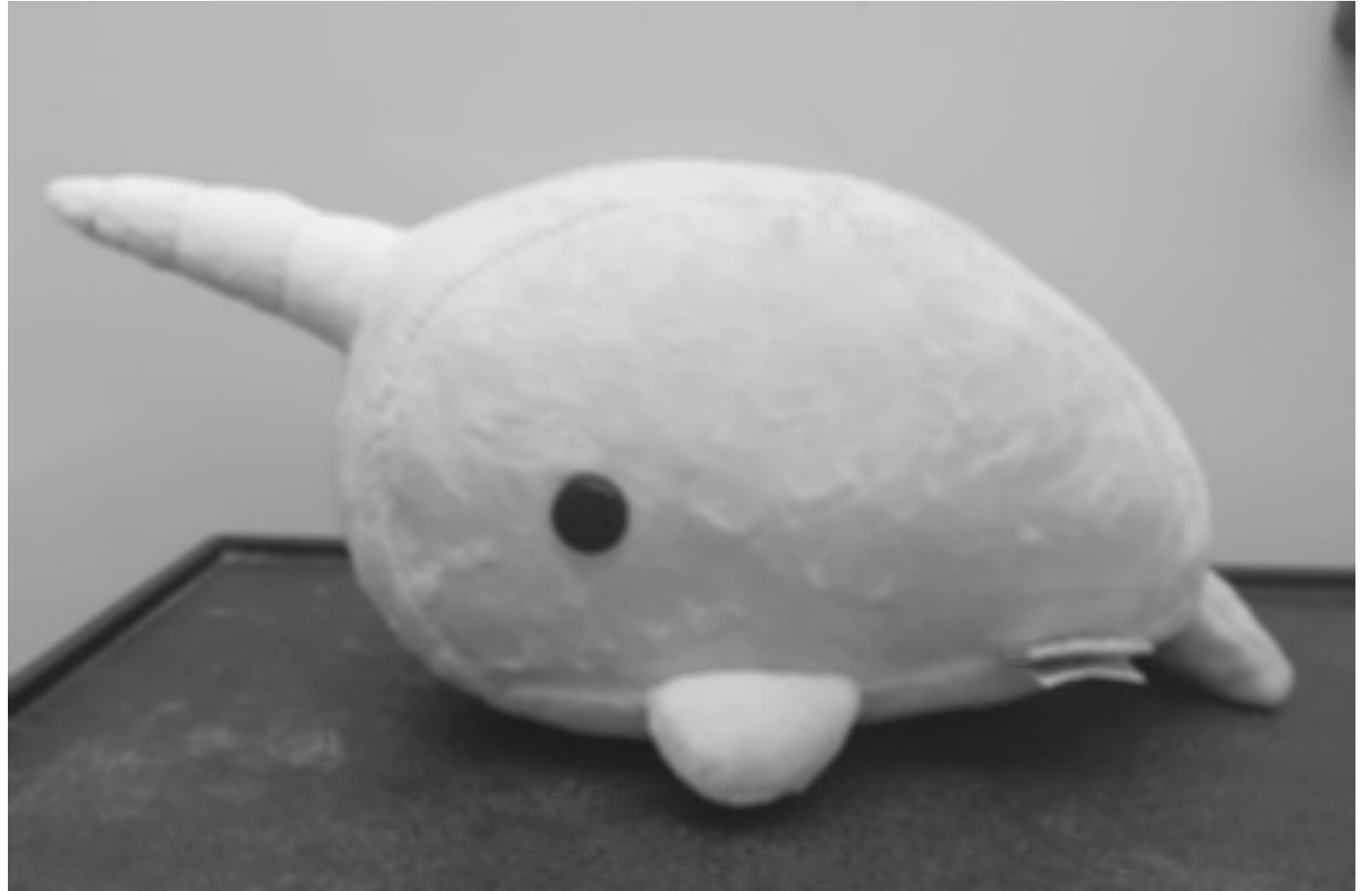
Original
Image



Convolution Examples

Smoothing
Convolution

$1/9$	$1/9$	$1/9$
$1/9$	$1/9$	$1/9$
$1/9$	$1/9$	$1/9$



Convolution Examples

Gaussian
Blur

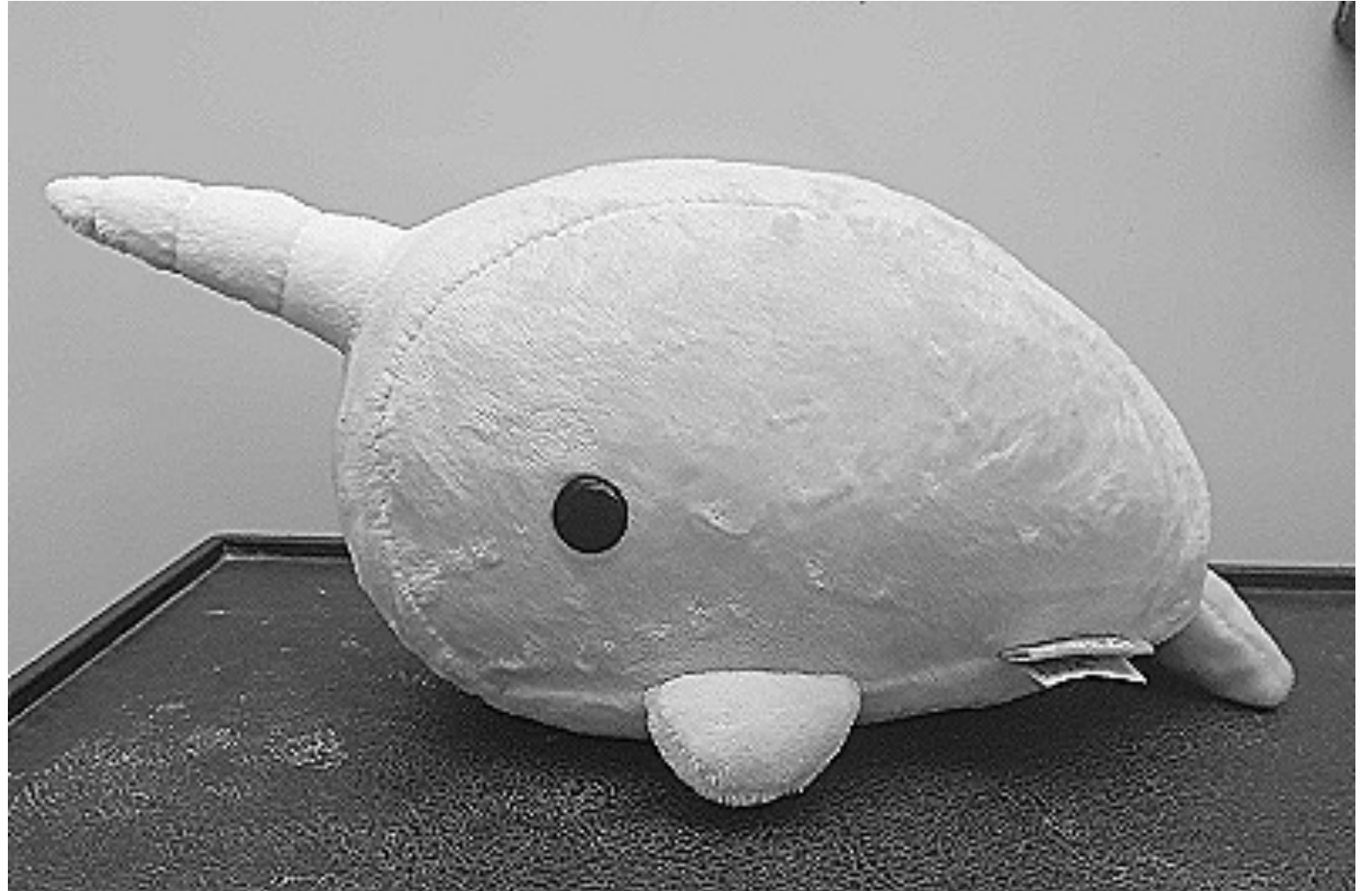
.01	.04	.06	.04	.01
.04	.19	.25	.19	.04
.06	.25	.37	.25	.06
.04	.19	.25	.19	.04
.01	.04	.06	.04	.01



Convolution Examples

Sharpening
Kernel

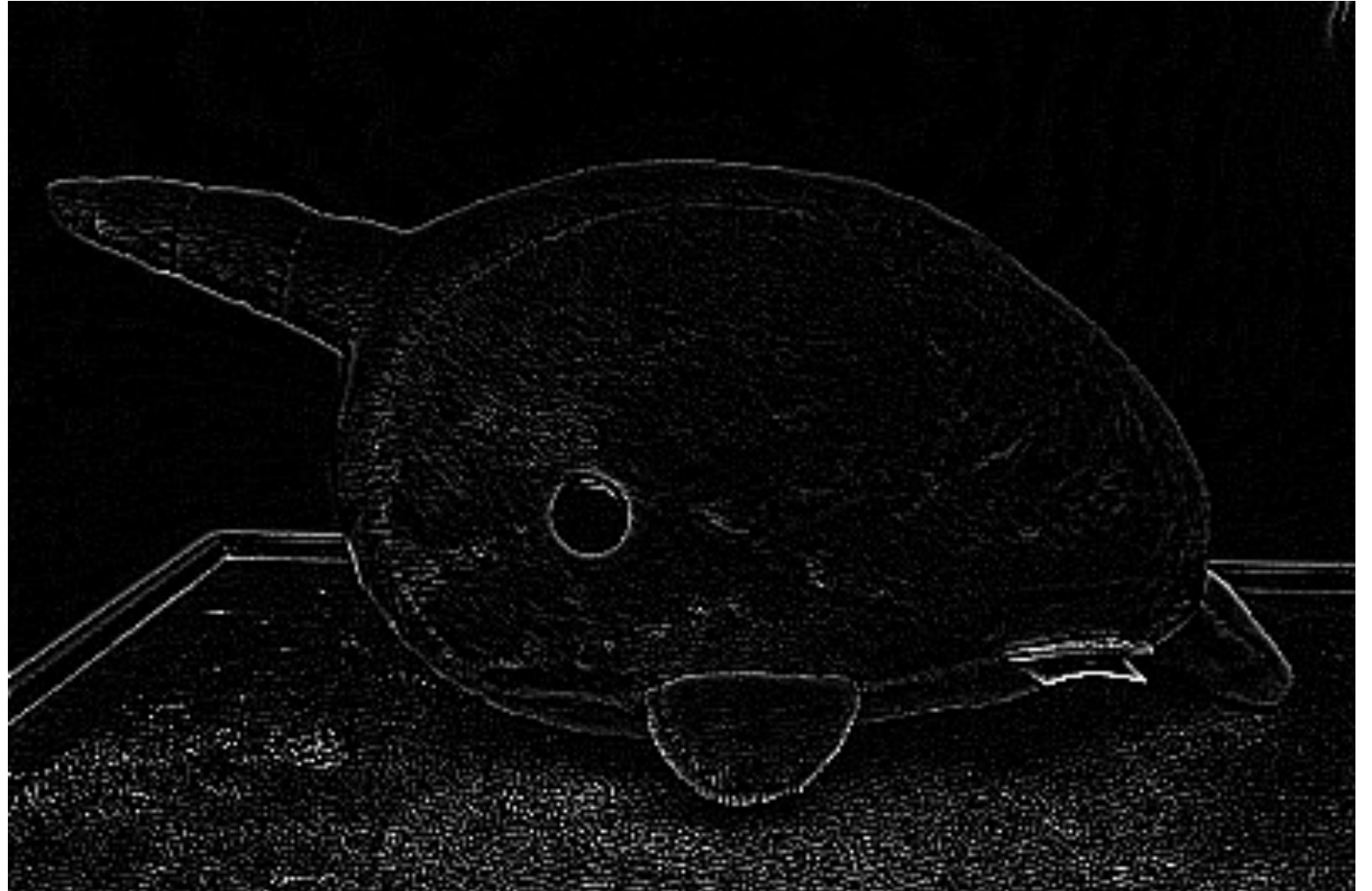
0	-1	0
-1	5	-1
0	-1	0



Convolution Examples

Edge
Detector

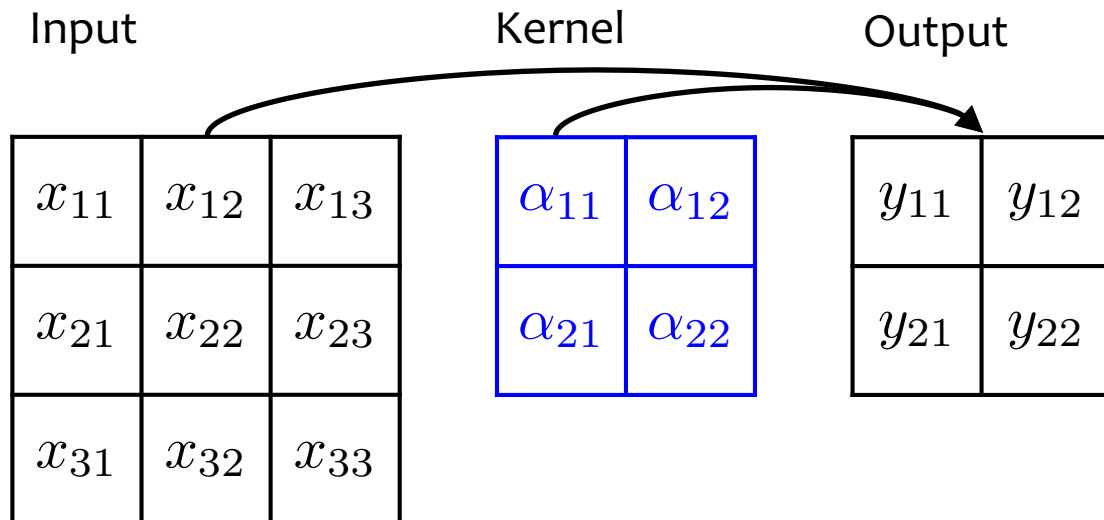
-1	-1	-1
-1	8	-1
-1	-1	-1



2D Convolution

- Basic idea:
 - Pick a 2x2 matrix F of weights (called a kernel or convolution matrix)
 - Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation
- Key point:
 - Different convolutions extract different types of low-level “features” from an image
 - All that we need to vary to generate these different features is the weights of F

Example: 1 input channel, 1 output channel



$$y_{11} = \alpha_{11}x_{11} + \alpha_{12}x_{12} + \alpha_{21}x_{21} + \alpha_{22}x_{22} + \alpha_0$$

$$y_{12} = \alpha_{11}x_{12} + \alpha_{12}x_{13} + \alpha_{21}x_{22} + \alpha_{22}x_{23} + \alpha_0$$

$$y_{21} = \alpha_{11}x_{21} + \alpha_{12}x_{22} + \alpha_{21}x_{31} + \alpha_{22}x_{32} + \alpha_0$$

$$y_{22} = \alpha_{11}x_{22} + \alpha_{12}x_{23} + \alpha_{21}x_{32} + \alpha_{22}x_{33} + \alpha_0$$

DOWNSAMPLING

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3		

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3		

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	0

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	0
1		

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	0
1	0	0

Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	0
1	0	0

Downsampling by Averaging

- Downsampling by averaging is a special case of convolution where the weights are fixed to a uniform distribution
- The example below uses a stride of 2

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

$1/4$	$1/4$
$1/4$	$1/4$

Convolved Image

$3/4$	$3/4$	$1/4$
$3/4$	$1/4$	0
$1/4$	0	0

Max-Pooling

- Max-pooling with a stride > 1 is another form of downsampling
- Instead of averaging, we take the max value within the same range as the equivalently-sized convolution
- The example below uses a stride of 2

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Max-pooling

$x_{i,j}$	$x_{i,j+1}$
$x_{i+1,j}$	$x_{i+1,j+1}$

$$y_{ij} = \max(x_{ij}, x_{i,j+1}, x_{i+1,j}, x_{i+1,j+1})$$

Max-Pooled Image

1	1	1
1	1	0
1	0	0

CONVOLUTIONAL NEURAL NETS

A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

1. • Convolutional Neural Networks (CNNs) provide another form of **decision function**
• Let's see what they look like...

2. CHOOSE EACH OF THESE.

- Decision function

$$\hat{y} = f_{\theta}(x_i)$$

- Loss function

$$\ell(\hat{y}, y_i) \in \mathbb{R}$$

4. Train with SGD:

(Take small steps opposite the gradient)

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(x_i), y_i)$$

Convolutional Layer

CNN key idea:
Treat convolution matrix as parameters and learn them!

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0



Learned Convolution

θ_{11}	θ_{12}	θ_{13}
θ_{21}	θ_{22}	θ_{23}
θ_{31}	θ_{32}	θ_{33}

Convolved Image

.4	.5	.5	.5	.4
.4	.2	.3	.6	.3
.5	.4	.4	.2	.1
.5	.6	.2	.1	0
.4	.3	.1	0	0

Convolutional Neural Network (CNN)

- Typical layers include:
 - Convolutional layer
 - Max-pooling layer
 - Fully-connected (Linear) layer
 - ReLU layer (or some other nonlinear activation function)
 - Softmax
- These can be arranged into arbitrarily deep topologies

Architecture #1: LeNet-5

PROC. OF THE IEEE, NOVEMBER 1998

7

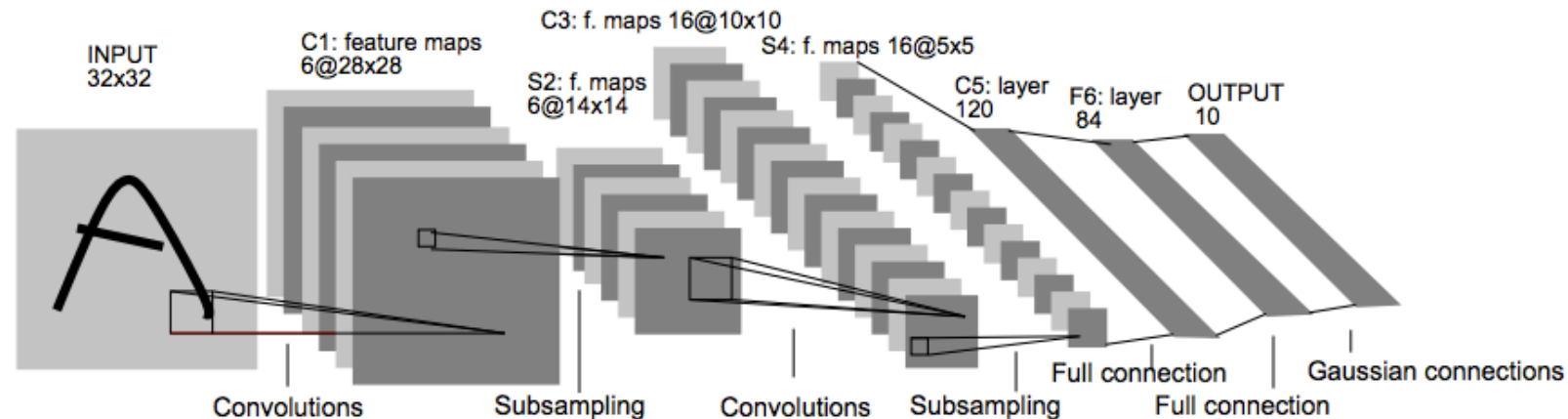


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

TRAINING CNNs

A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

(take small steps
opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of the following:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

- Q: Now that we have the CNN as a decision function, how do we compute the gradient?
- A: Backpropagation of course!

opposite the gradient)


$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

SGD for CNNs

Example: Simple CNN Architecture

Given \mathbf{x} , \mathbf{y}^* and parameters $\boldsymbol{\theta} = [\boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{W}]$

$$J = \ell(\mathbf{y}, \mathbf{y}^*)$$

$$\mathbf{y} = \text{softmax}(\mathbf{z}^{(5)})$$

$$\mathbf{z}^{(5)} = \text{linear}(\mathbf{z}^{(4)}, \mathbf{W})$$

$$\mathbf{z}^{(4)} = \text{relu}(\mathbf{z}^{(3)})$$

$$\mathbf{z}^{(3)} = \text{conv}(\mathbf{z}^{(2)}, \boldsymbol{\beta})$$

$$\mathbf{z}^{(2)} = \text{max-pool}(\mathbf{z}^{(1)})$$

$$\mathbf{z}^{(1)} = \text{conv}(\mathbf{x}, \boldsymbol{\alpha})$$

Algorithm 1 Stochastic Gradient Descent (SGD)

- 1: Initialize $\boldsymbol{\theta}$
 - 2: **while** not converged **do**
 - 3: Sample $i \in \{1, \dots, N\}$
 - 4: Forward: $\mathbf{y} = h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})$,
 - 5: $J(\boldsymbol{\theta}) = \ell(\mathbf{y}, \mathbf{y}^{(i)})$
 - 6: Backward: Compute $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
 - 7: Update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
-

LAYERS OF A CNN

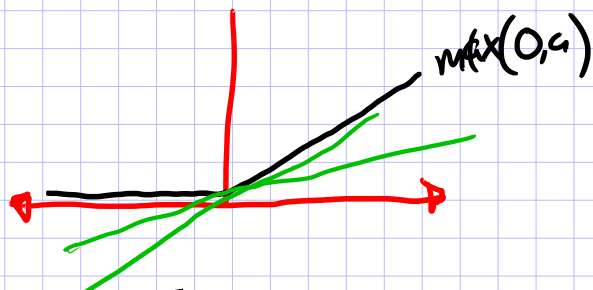
ReLU Layer

ReLU Layer Input: $\vec{x} \in \mathbb{R}^k$ Output: $\vec{y} \in \mathbb{R}^k$

Forward:

$$\vec{y} = \sigma(\vec{x}) \leftarrow \text{element-wise}$$

$$\sigma(a) = \max(0, a)$$



Backward:

$$\frac{dJ}{dx_i} = \frac{dJ}{dy_i} \frac{dy_i}{dx_i}$$

where $\frac{dy_i}{dx_i} = \begin{cases} 1 & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$

subderivative

Softmax Layer

Input: $\mathbf{x} \in \mathbb{R}^K$, Output: $\mathbf{y} \in \mathbb{R}^K$

Forward: for each i ,

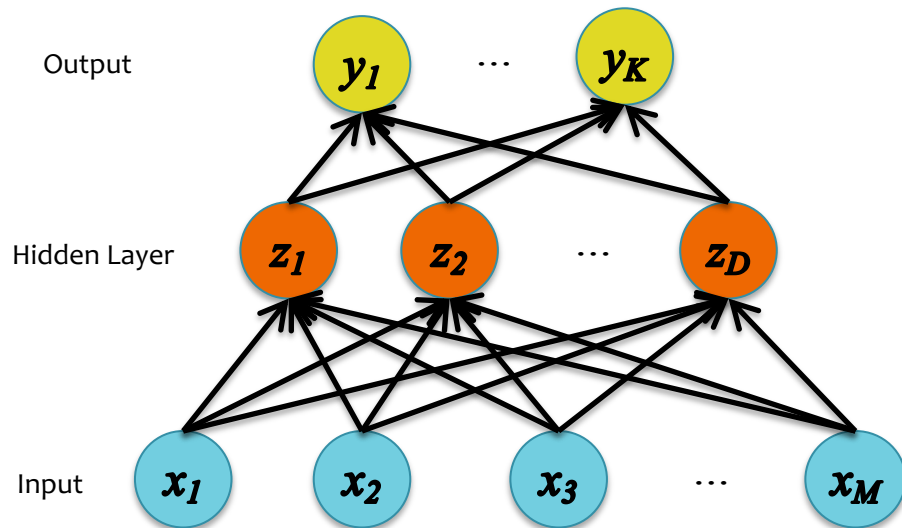
$$y_i = \frac{\exp(x_i)}{\sum_{k=1}^K \exp(x_k)}$$

Backward: for each j ,

$$\frac{\partial J}{\partial x_j} = \sum_{i=1}^K \frac{\partial J}{\partial y_i} \frac{\partial y_i}{\partial x_j}$$

where

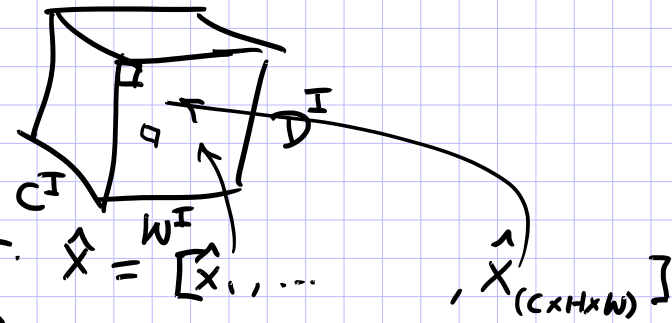
$$\frac{\partial y_i}{\partial x_j} = \begin{cases} y_i(1 - y_i) & \text{if } i = j \\ -y_i y_j & \text{otherwise} \end{cases}$$



Fully-Connected Layer

Fully Connected Layer (w/ tensor input)

- Suppose input is a 3D Tensor: $X =$



- stretch out into a long vector. $\hat{x} = [x_1, \dots]$

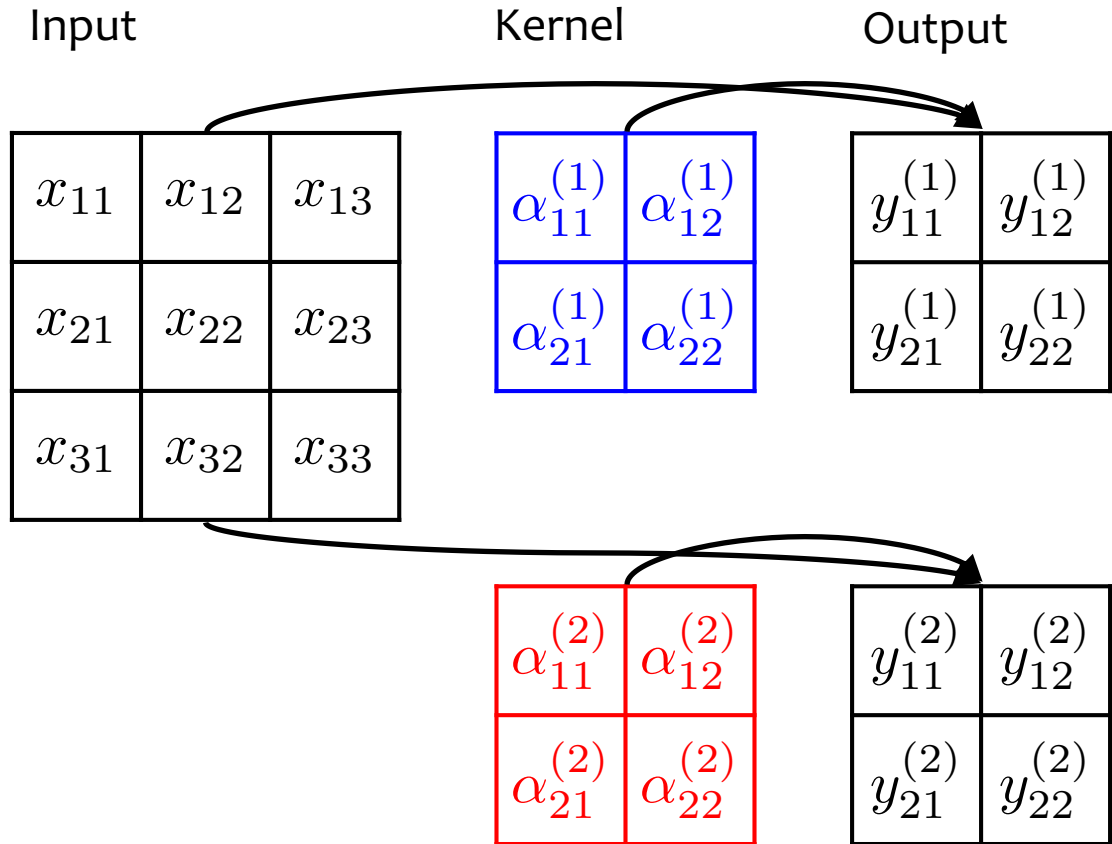
- then standard linear layer:

$$y = \alpha^T \hat{x} + \alpha_0 \quad \text{where } \alpha \in \mathbb{R}^{A \times B}$$

$|\hat{x}| = A, |y| = B$

2D Convolution

Example: 1 input channel, 2 output channels



$$y_{11}^{(1)} = \alpha_{11}^{(1)} x_{11} + \alpha_{12}^{(1)} x_{12} + \alpha_{21}^{(1)} x_{21} + \alpha_{22}^{(1)} x_{22} + \alpha_0^{(1)}$$

$$y_{12}^{(1)} = \alpha_{11}^{(1)} x_{12} + \alpha_{12}^{(1)} x_{13} + \alpha_{21}^{(1)} x_{22} + \alpha_{22}^{(1)} x_{23} + \alpha_0^{(1)}$$

$$y_{21}^{(1)} = \alpha_{11}^{(1)} x_{21} + \alpha_{12}^{(1)} x_{22} + \alpha_{21}^{(1)} x_{31} + \alpha_{22}^{(1)} x_{32} + \alpha_0^{(1)}$$

$$y_{22}^{(1)} = \alpha_{11}^{(1)} x_{22} + \alpha_{12}^{(1)} x_{23} + \alpha_{21}^{(1)} x_{32} + \alpha_{22}^{(1)} x_{33} + \alpha_0^{(1)}$$

$$y_{11}^{(2)} = \alpha_{11}^{(2)} x_{11} + \alpha_{12}^{(2)} x_{12} + \alpha_{21}^{(2)} x_{21} + \alpha_{22}^{(2)} x_{22} + \alpha_0^{(2)}$$

$$y_{12}^{(2)} = \alpha_{11}^{(2)} x_{12} + \alpha_{12}^{(2)} x_{13} + \alpha_{21}^{(2)} x_{22} + \alpha_{22}^{(2)} x_{23} + \alpha_0^{(2)}$$

$$y_{21}^{(2)} = \alpha_{11}^{(2)} x_{21} + \alpha_{12}^{(2)} x_{22} + \alpha_{21}^{(2)} x_{31} + \alpha_{22}^{(2)} x_{32} + \alpha_0^{(2)}$$

$$y_{22}^{(2)} = \alpha_{11}^{(2)} x_{22} + \alpha_{12}^{(2)} x_{23} + \alpha_{21}^{(2)} x_{32} + \alpha_{22}^{(2)} x_{33} + \alpha_0^{(2)}$$

Convolution of a Color Image

- Color images consist of 3 floats per pixel for RGB (red, green blue) color values
- Convolution must also be 3-dimensional

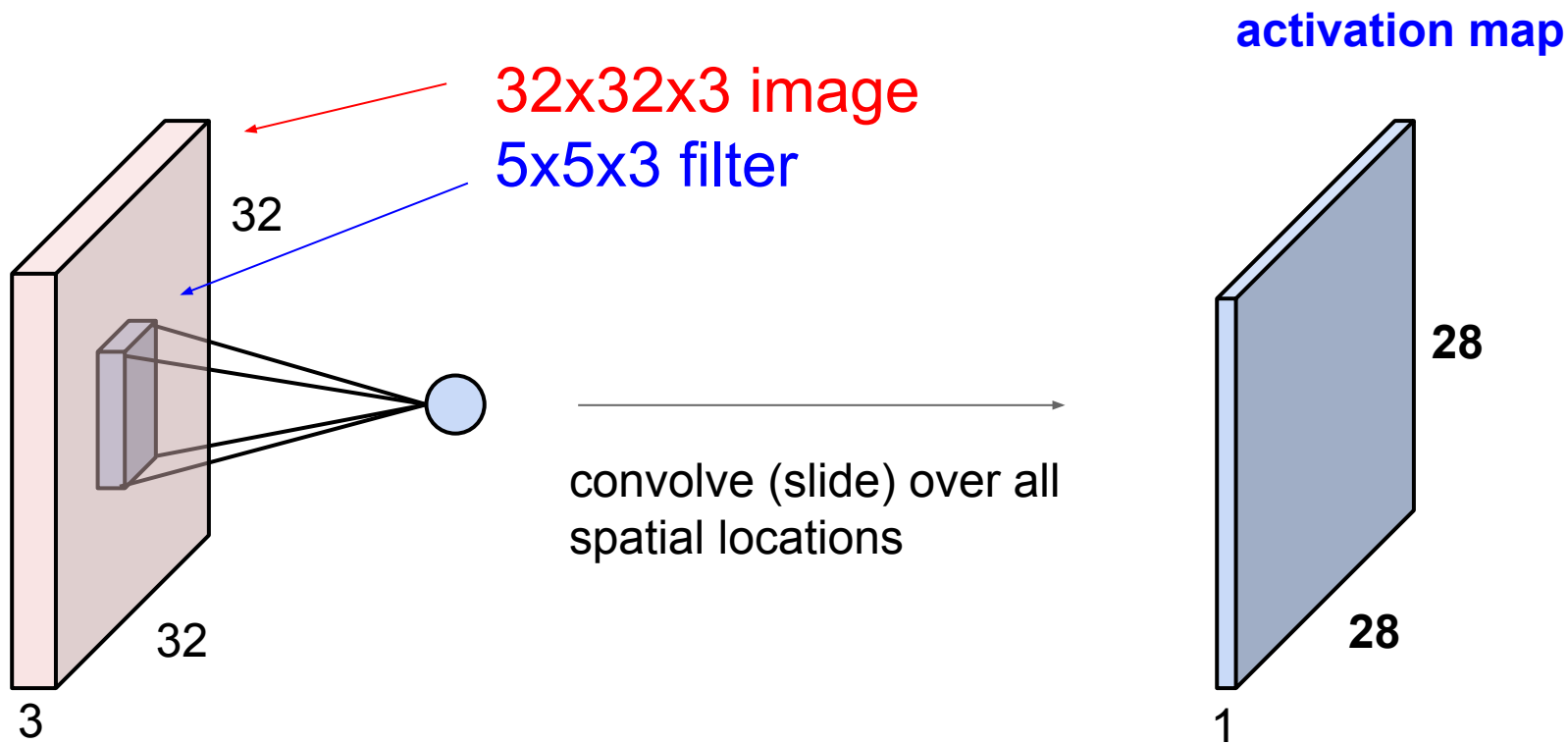


Figure from Fei-Fei Li & Andrej Karpathy & Justin Johnson (CS231N)

Animation of 3D Convolution

<http://cs231n.github.io/convolutional-networks/>

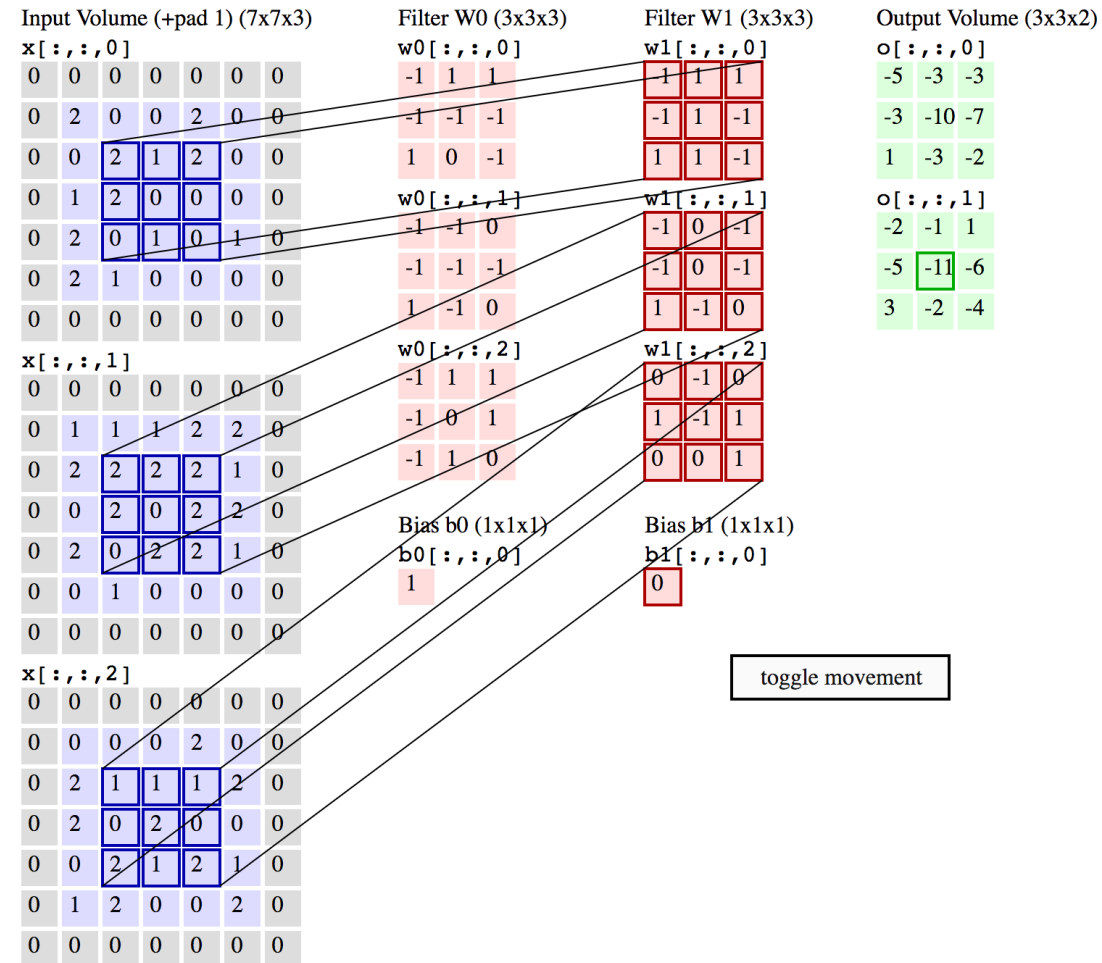
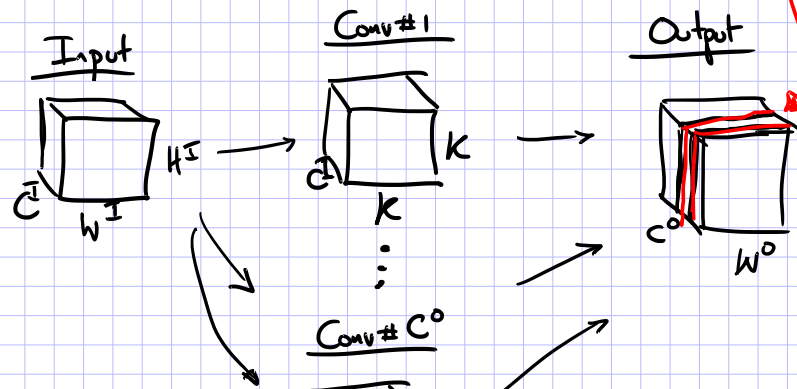


Figure from Fei-Fei Li & Andrej Karpathy & Justin Johnson (CS231N)

Convolutional Layer

Ex: C^I input channels, C^O output channels

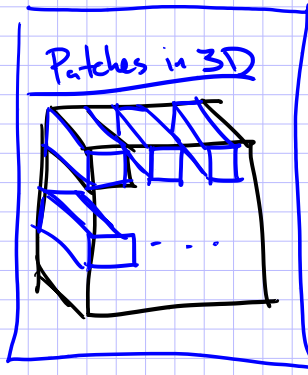


i -th slice is output from j -th convolution matrix

$$H^O = \lfloor (H^I + 2p - K) / s + 1 \rfloor$$

$$W^O = \lfloor (W^I + 2p - K) / s + 1 \rfloor$$

where p = # pixels of padding on input
 K = size of conv. matrix
 s = stride length



Forward:

$$y_{ij}^{(k)} = \alpha_0^{(k)} + \sum_{c=1}^{C^I} \sum_{q=1}^K \sum_{r=1}^K \alpha_{qr}^{(k)} x_{mn}^{(c)}$$

where $m = s(i-1) + q$
 $n = s(j-1) + r$

Backward:

$$\frac{dJ}{d\alpha_0^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{d\alpha_0^{(k)}}$$

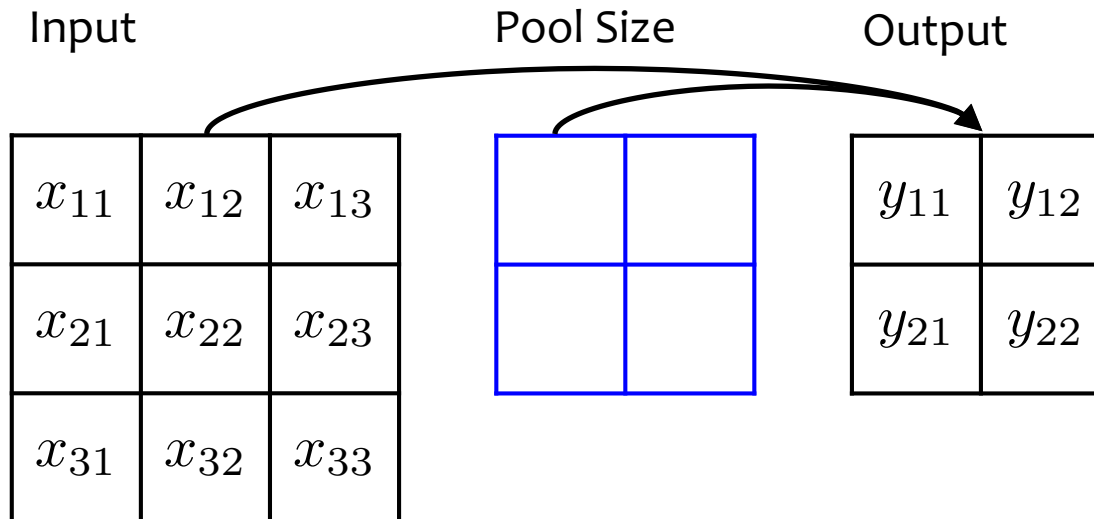
$$\frac{dJ}{d\alpha_{qr}^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{d\alpha_{qr}^{(k)}}$$

$$\frac{dJ}{dx_{mn}^{(c)}} = \sum_i \sum_j \sum_k \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{dx_{mn}^{(c)}}$$

just some calculus

Max-Pooling Layer

Example: 1 input channel, 1 output channel, stride of 1



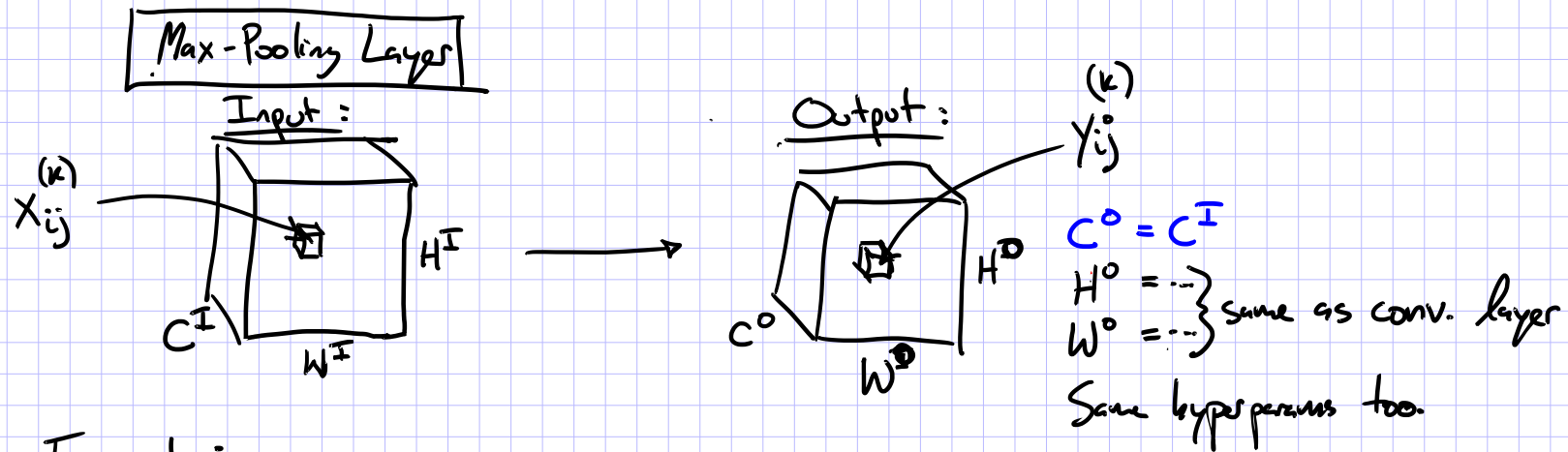
$$y_{11} = \max(x_{11}, x_{12}, x_{21}, x_{22})$$

$$y_{12} = \max(x_{12}, x_{13}, x_{22}, x_{23})$$

$$y_{21} = \max(x_{21}, x_{22}, x_{31}, x_{32})$$

$$y_{22} = \max(x_{22}, x_{23}, x_{32}, x_{33})$$

Max-Pooling Layer



Forward:

$$Y_{ij}^{(k)} = \max_{\substack{q \in \{1, \dots, k\} \\ r \in \{1, \dots, k\}}} X_{mn}^{(k)} \text{ where } \begin{cases} m = s(i-1) + q \\ n = s(j-1) + r \end{cases}$$

Backward:

$$\frac{dJ}{dx_{mn}^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{dx_{mn}^{(k)}}$$

Subderivatives

- + $\max()$ is not differentiable, but subdifferentiable.
- + There are a set of derivatives and we can just choose one for SGD.

$$y = \max(a, b)$$

$$\Rightarrow \frac{dJ}{da} = \frac{dJ}{dy} \frac{dy}{da} \text{ where } \frac{dy}{da} = \begin{cases} 1 & \text{if } a > b \\ 0 & \text{otherwise} \end{cases}$$

Convolutional Neural Network (CNN)

- Typical layers include:
 - Convolutional layer
 - Max-pooling layer
 - Fully-connected (Linear) layer
 - ReLU layer (or some other nonlinear activation function)
 - Softmax
- These can be arranged into arbitrarily deep topologies

Architecture #1: LeNet-5

PROC. OF THE IEEE, NOVEMBER 1998

7

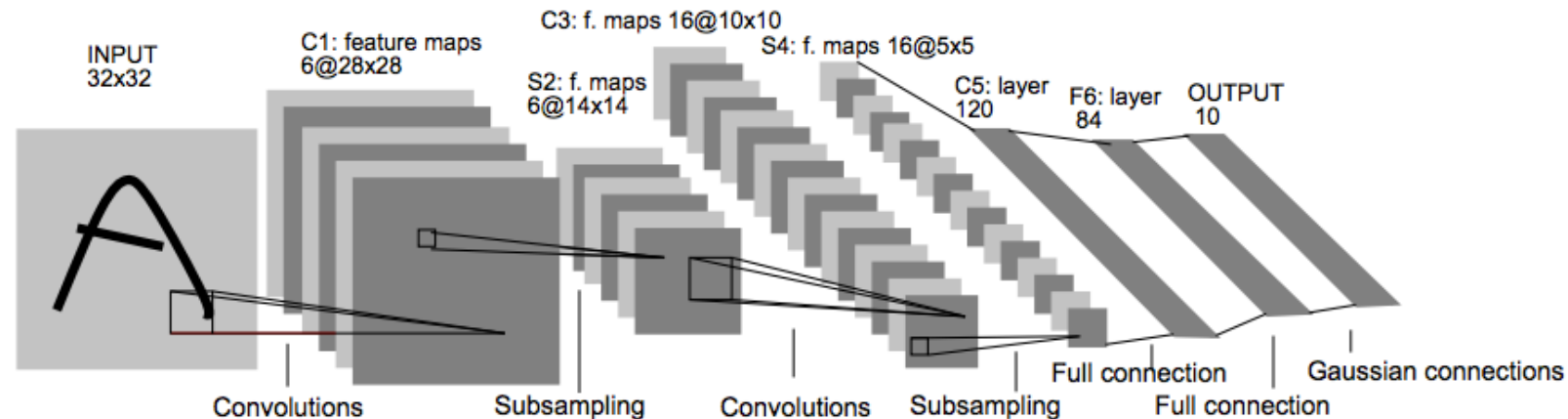


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Architecture #2: AlexNet

CNN for Image Classification

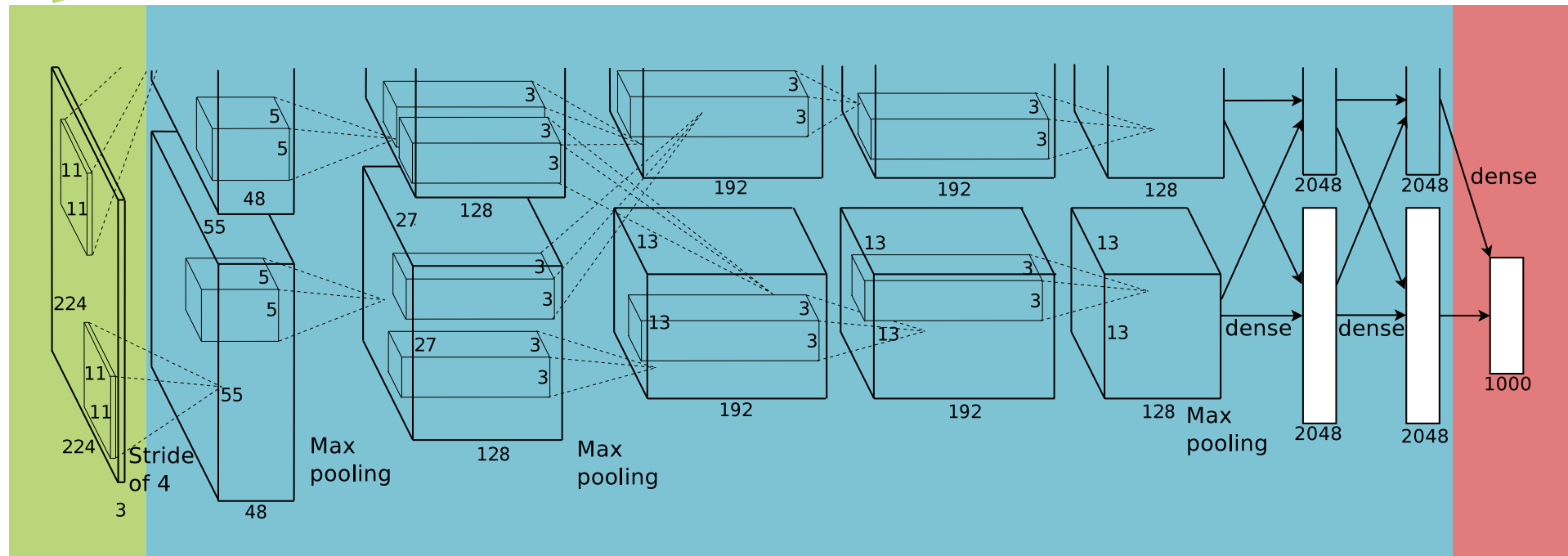
(Krizhevsky, Sutskever & Hinton, 2012)

15.3% error on ImageNet LSVRC-2012 contest

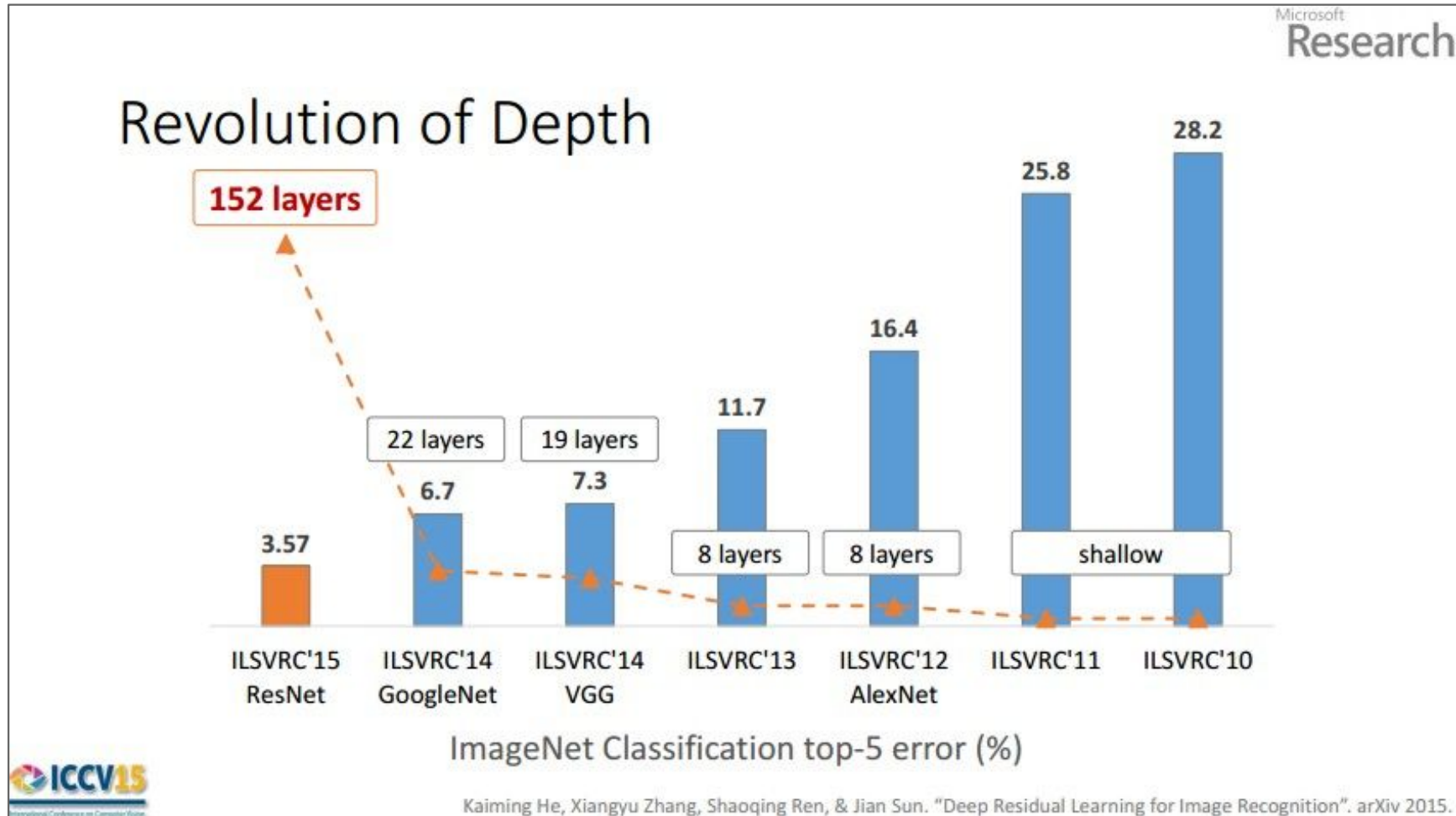
Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax

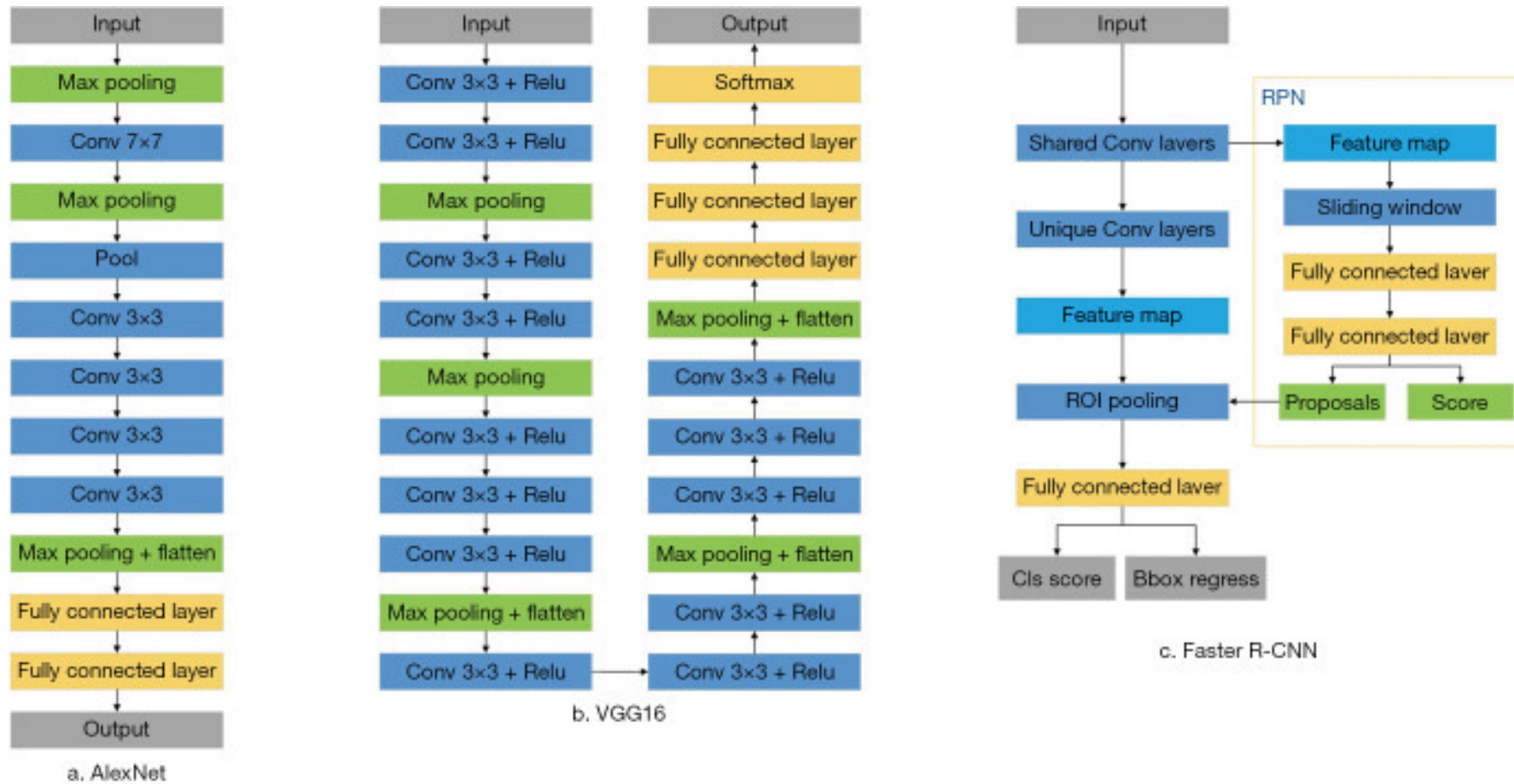


CNNs for Image Recognition



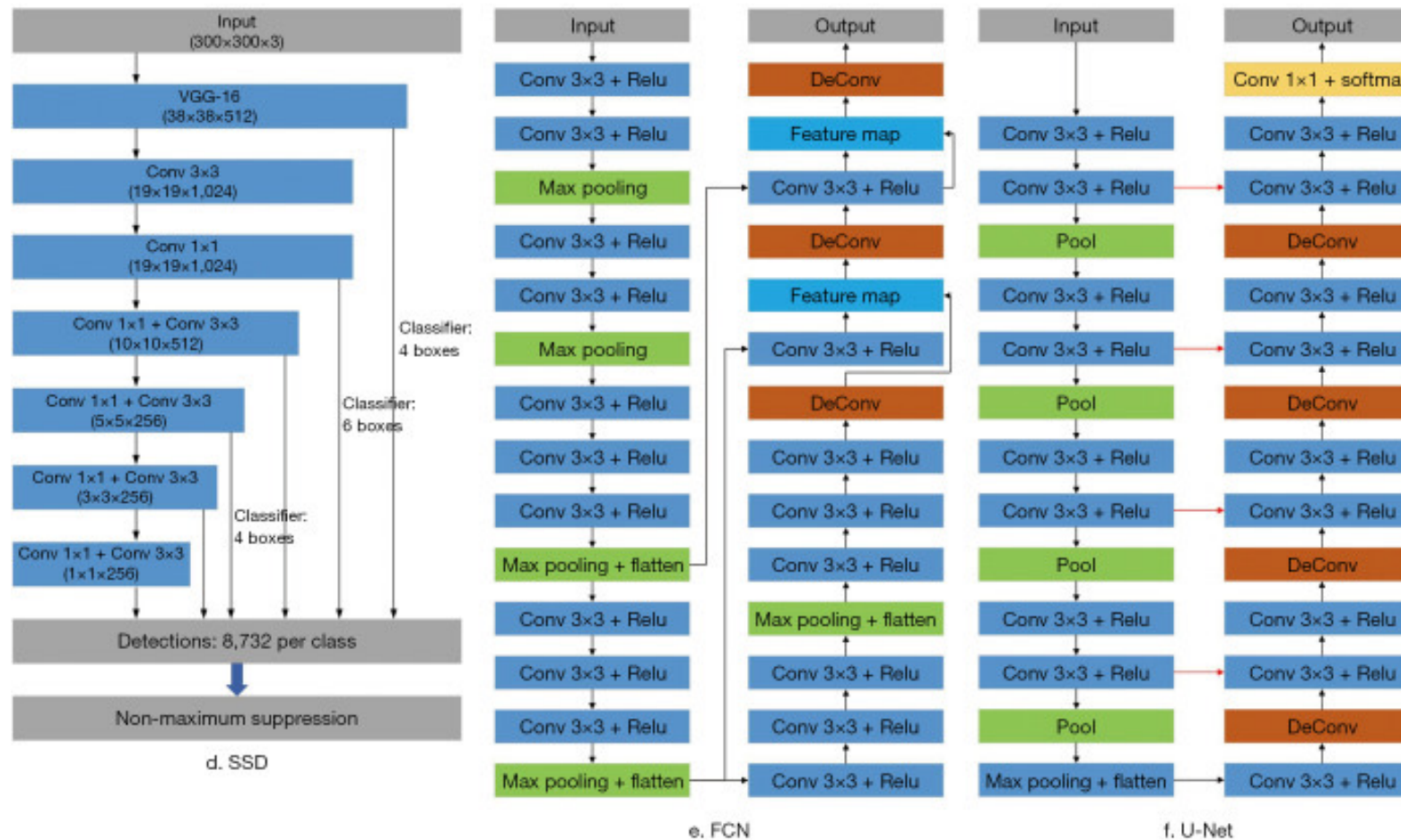
Convolutional Neural Network (CNN)

Typical Architectures



Convolutional Neural Network (CNN)

Typical Architectures



Convolutional Neural Network (CNN)

Typical Architectures

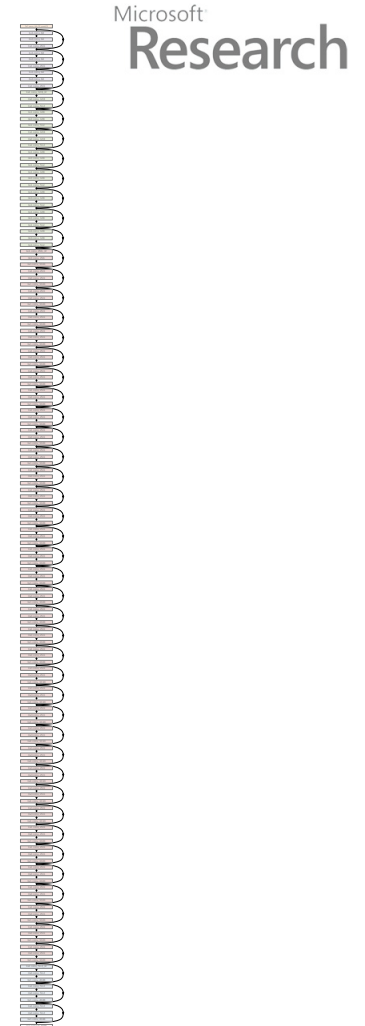
AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



ResNet, 152 layers
(ILSVRC 2015)



In-Class Poll

Question:

Why do many layers used in computer vision *not have* location specific parameters?

Answer:

Convolutional Layer

For a convolutional layer, how do we pick the kernel size (aka. the size of the convolution)?

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

2x2
Convolution

θ_{11}	θ_{12}
θ_{21}	θ_{22}

3x3
Convolution

θ_{11}	θ_{12}	θ_{13}
θ_{21}	θ_{22}	θ_{23}
θ_{31}	θ_{32}	θ_{33}

4x4
Convolution

θ_{11}	θ_{12}	θ_{13}	θ_{14}
θ_{21}	θ_{22}	θ_{23}	θ_{24}
θ_{31}	θ_{32}	θ_{33}	θ_{34}
θ_{41}	θ_{42}	θ_{43}	θ_{44}

- A small kernel can only see a very small part of the image, but is fast to compute
- A large kernel can see more of the image, but at the expense of speed

CNN VISUALIZATIONS

Visualization of CNN

https://adamharley.com/nn_vis/cnn/2d.html

60 FPS (1-61) Draw your number here 0123456789

Layer visibility

Input layer	Hide
Convolution layer 1	Hide
Downsampling layer 1	Hide
Convolution layer 2	Hide
Downsampling layer 2	Hide
Fully-connected layer 1	Hide
Fully-connected layer 2	Hide
Output layer	Hide

Downsampled drawing: 4
First guess: 7
Second guess: 8

The visualization shows a large digit '4' drawn on a canvas. Below the canvas is a grid of feature maps for each layer. The layers are: Input layer (the original drawing), Convolution layer 1, Downsampling layer 1, Convolution layer 2, Downsampling layer 2, Fully-connected layer 1, Fully-connected layer 2, and Output layer. The output layer shows the digit '4' with a probability of 0.123456789.

MNIST Digit Recognition with CNNs (in your browser)

<https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

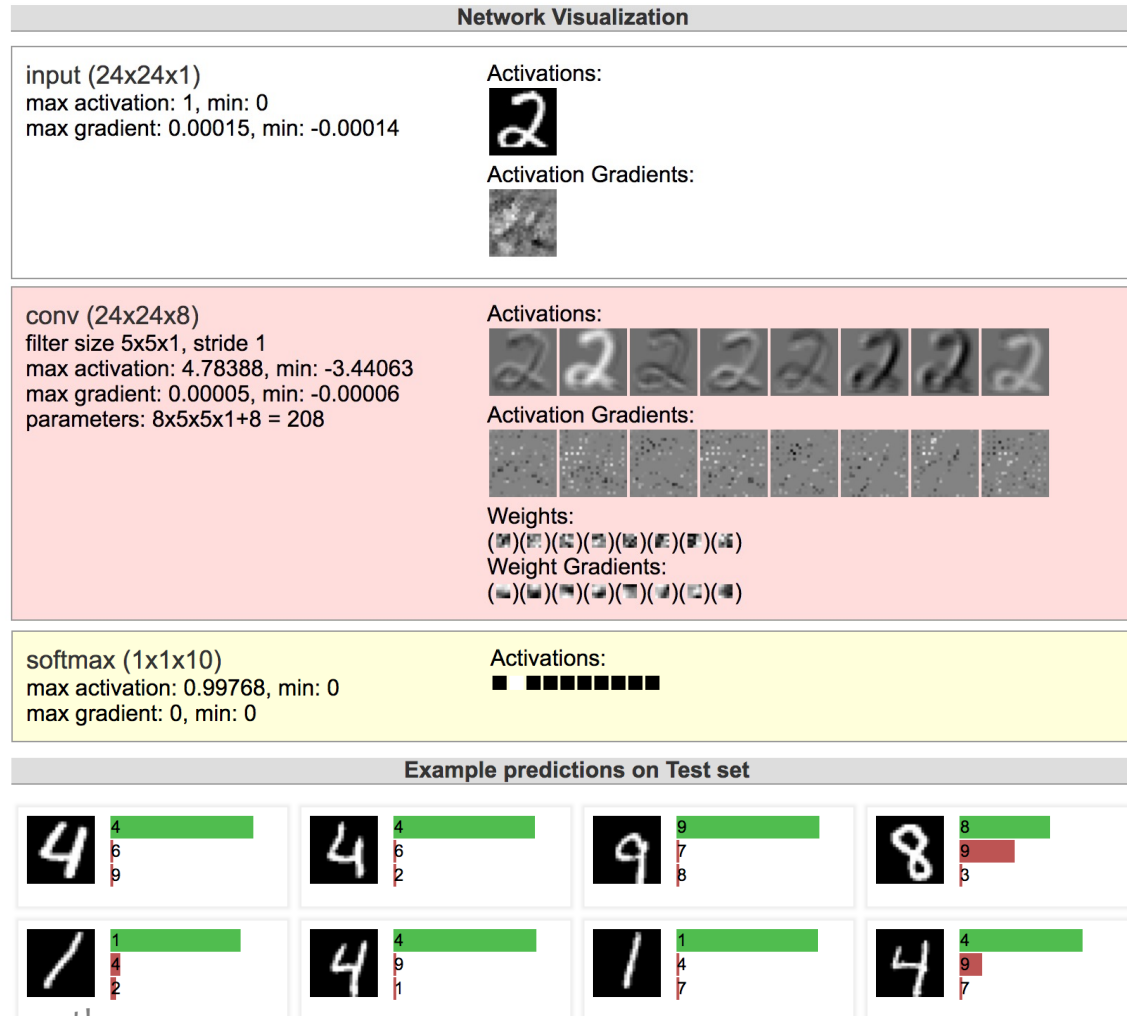


Figure from Andrej Karpathy

CNN Summary

CNNs

- Are used for all aspects of **computer vision**, and have won numerous pattern recognition competitions
- Able learn **interpretable features** at different levels of abstraction
- Typically, consist of **convolution** layers, **pooling** layers, **nonlinearities**, and **fully connected** layers

Deep Learning Objectives

You should be able to...

- Implement the common layers found in Convolutional Neural Networks (CNNs) such as linear layers, convolution layers, max-pooling layers, and rectified linear units (ReLU)
- Explain how the shared parameters of a convolutional layer could learn to detect spatial patterns in an image
- Describe the backpropagation algorithm for a CNN
- Identify the parameter sharing used in a basic recurrent neural network, e.g. an Elman network
- Apply a recurrent neural network to model sequence data
- Differentiate between an RNN and an RNN-LM