Learning Based Vision I

Computer Vision Fall 2018 Columbia University

Homeworks

- Homework 1 grades back
 - Median 37/40, Std 3.6
- Homework 2 due a few minutes ago
- Homework 3 out today
 - Due October 22nd

Project

- Project Proposals due October 31
- Pick one of our suggested projects, or pitch your own
- Must use something in this course
- Groups of 2 strongly recommended
- We'll give you Google Cloud credits once you turn in your project proposal
- Details here: http://w4731.cs.columbia.edu/project



We need translation invariance

Lots of useful linear filters...

High order Gaussian derivatives

 g_{x^6}



Gabor

And many more...



We need translation and scale invariance

Lots of image pyramids...



Gaussian Pyr

Laplacian Pyr

And many more: QMF, steerable, ...



We need ...

What is the best representation?

- All the previous representation are manually constructed.
- Could they be learnt from data?

A brief history of Neural Networks



Where did this all start?

Hubel & Wiesel (1962)

Insights about early image processing in the brain.

Simple cells detect local features

Complex cells pool local features in a retinotopic neighborhood



The perceptron



Slide credit: Deva Ramanan



Percepton

If input features are binary, can model logical "and"s and "or"s



What about "xors"?



Minsky and Papert, Perceptrons, 1972





FOR BUYING OPTIONS, START HERE	
Select Shipping Destination	÷

Paperback | \$35.00 Short | £24.95 | ISBN: 9780262631112 | 308 pp. | 6 x 8.9 in | December 1987

Perceptrons, expanded edition

An Introduction to Computational Geometry

By Marvin Minsky and Seymour A. Papert

Overview

Perceptrons - the first systematic study of parallelism in computation - has remained a classical work on threshold automata networks for nearly two decades. It marked a historical turn in artificial intelligence, and it is required reading for anyone who wants to understand the connectionist counterrevolution that is going on today.

Artificial-intelligence research, which for a time concentrated on the programming of ton Neumann computers, is swinging back to the idea that intelligence might emerge from the activity of networks of neuronlike entities. Minsky and Papert's book was the first example of a mathematical analysis carried far enough to show the exact limitations of a class of computing machines that could seriously be considered as models of the brain. Now the new developments in mathematical tools, the recent interest of physicists in the theory of disordered matter, the new insights into and psychological models of how the brain works, and the evolution of fast computers that can simulate networks of automata have given *Perceptrons* new importance.

Witnessing the swing of the intellectual pendulum, Minsky and Papert have added a new chapter in which they discuss the current state of parallel computers, review developments since the appearance of the 1972 edition, and identify new research directions related to connectionism. They note a central theoretical challenge facing connectionism: the challenge to reach a deeper understanding of how "objects" or "agents" with individuality can emerge in a network. Progress in this area would link connectionism with what the authors have called "society theories of mind."



Parallel Distributed Processing (PDP), 1986



XOR problem



PDP authors pointed to the backpropagation algorithm as a breakthrough, allowing multi-layer neural networks to be trained. Among the functions that a multi-layer network can represent but a single-layer network cannot: the XOR function.



LeCun conv nets, 1998

PROC. OF THE IEEE, NOVEMBER 1998



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.

Demos: http://yann.lecun.com/exdb/lenet/index.html

Slide credit: Antonio Torralba

 $\overline{7}$



Fig. 13. Examples of unusual, distorted, and noisy characters correctly recognized by LeNet-5. The grey-level of the output label represents the penalty (lighter for higher penalties).



http://pub.clement.farabet.net/ecvw09.pdf

Neural networks to recognize handwritten digits? yes

Neural networks for tougher problems? not really

NIPS 2000

- NIPS, Neural Information Processing Systems, is the premier conference on machine learning. Evolved from an interdisciplinary conference to a machine learning conference.
- For the NIPS 2000 conference:
 - <u>title words predictive of paper acceptance</u>:
 "Belief Propagation" and "Gaussian".
 - <u>title words predictive of paper rejection</u>:
 "Neural" and "Network".



Krizhevsky, Sutskever, and Hinton, NIPS 2012



ImageNet Classification 2012

- Krizhevsky et al. -- 16.4% error (top-5)
- Next best (non-convnet) 26.2% error



Krizhevsky, Sutskever, and Hinton, NIPS 2012



Test Nearby images, according to NN features



Krizhevsky, Sutskever, and Hinton, NIPS 2012

Slide credit: Antonio Torralba

25









http://www.deeplearningbook.org/

By Ian Goodfellow, Yoshua Bengio and Aaron Courville

November 2016

"Classical" Recognition

Speech recognition (early 90's – 2011)



Object recognition (2006 – 2012)



End-to-end

Deep learning can be summarized as learning both the representation and the classifier out of it

• Fixed engineered features (or kernels) + trainable classifier



End-to-end

In deep learning we have multiple stages of non linear feature transformation



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Slide credit: LeCun
Artificial Neural Networks



Artificial Neural Networks



 $x_i \in \mathbb{R}^{H \times 1}$ $W_i \in \mathbb{R}^{D \times H}$ $b_i \in \mathbb{R}^{D \times 1}$

$$x_{i+1} = f\left(W_i x_i + b_i\right)$$

Bio/Artificial Neurons



Slide credit: Andrej Karpathy

Bio/Artificial Networks

The ventral (recognition) pathway in the visual cortex has multiple stages



Artificial Neural Networks



 $x_i \in \mathbb{R}^{H \times 1}$ $W_i \in \mathbb{R}^{D \times H}$ $b_i \in \mathbb{R}^{D \times 1}$

$$x_{i+1} = f\left(W_i x_i + b_i\right)$$

How do we find W and b?







Kitchens







Forests



Theatres





Bedrooms







Loss Functions



The objective min $\sum_{\theta} \sum_{i} \mathscr{L}(f(x_i; \theta), y_i)$ of learning:

Common Loss Functions

Squared error:

$$\mathscr{L}(x, y) = \|x - y\|_2^2$$

Hinge loss:

$$\mathcal{L}(x, y) = \max(0, 1 - x \cdot y)$$

Cross entropy:

$$\mathscr{L}(x, y) = -\sum_{i} y_{i} \log x_{i}$$























θ Sometimes millions of dimensions

Flavors of Gradient Descent

Gradient descent:

$$\theta_{t+1} = \theta_t + \alpha \frac{\delta \mathscr{L}}{\delta \theta_t}$$

Gradient descent with momentum:

$$z_{t+1} = \beta z_t + \frac{\delta \mathscr{L}}{\delta \theta_t}$$
$$\theta_{t+1} = \theta_t + \alpha z_{t+1}$$

Stochastic gradient descent:

$$\theta_{t+1} = \theta_t + \alpha \mathbb{E}_x \left[\frac{\delta \mathscr{L}}{\delta \theta_t} \right]$$

Are we done?



Wikipedia

Regularization



The objective of learning:

$$\min_{\theta} \sum_{i} \mathscr{L}\left(f(x_{i}), y_{i}\right) + \lambda R(\theta)$$

Common $R(\theta) = \|\theta\|_2^2$ regularization:

Artificial Neural Networks



 $x_i \in \mathbb{R}^{H \times 1}$ $W_i \in \mathbb{R}^{D \times H}$ $b_i \in \mathbb{R}^{D \times 1}$

$$x_{i+1} = f\left(W_i x_i + b_i\right)$$

$$\min_{\theta} \sum_{j} \mathscr{L}\left(x_{last}^{j}, y^{j}\right)$$

Where we are headed



- Let $\mathbf{x}_{\mathbf{L}}$ be the output of the Lth layer.
- We then apply a loss given a target label **y**
- Ultimate goal: compute **dz/dw**

Composable units



$$g(\mathbf{x}_0, \mathbf{w}_1, \ldots) = f_L(f_{L-1}(f_{L-2}(\ldots, \mathbf{w}_{L-2}), \mathbf{w}_{L-1}), \mathbf{w}_L)$$

= $f_L(\cdot, \mathbf{w}_L) \circ f_{L-1}(\cdot, \mathbf{w}_{L-1}) \ldots$

[Shorthand notation for *composition* that we'll use in next slides]

Slide credit: Deva Ramanan

Review: Chain Rule

 $\frac{\delta}{\delta x} \left[f \circ g(x) \right] = f' \left(g(x) \right) g'(x)$





 $\frac{dz}{d\mathbf{w}_l} = \frac{d}{d\mathbf{w}_l} \left[\ell_{\mathbf{y}} \circ f_L(\cdot; \mathbf{w}_L) \circ \dots \circ f_2(\cdot; \mathbf{w}_2) \circ f_1(\mathbf{x}_0; \mathbf{w}_1) \right]$





$$\frac{dz}{d\mathbf{w}_l} = \frac{d}{d\mathbf{w}_l} \left[\ell_{\mathbf{y}} \circ f_L(\cdot; \mathbf{w}_L) \circ \dots \circ f_2(\cdot; \mathbf{w}_2) \circ f_1(\mathbf{x}_0; \mathbf{w}_1) \right]$$

$$\frac{dz}{d\mathbf{w}_l} = \frac{dz}{d(\operatorname{vec} \mathbf{x}_L)^{\top}} \frac{d\operatorname{vec} \mathbf{x}_L}{d(\operatorname{vec} \mathbf{x}_{L-1})^{\top}} \cdots \frac{d\operatorname{vec} \mathbf{x}_{l+1}}{d(\operatorname{vec} \mathbf{x}_l)^{\top}} \frac{d\operatorname{vec} \mathbf{x}_l}{d\mathbf{w}_l^{\top}}$$

We can add *any* functional module f so long as its interface provides: (1) derivative of output wrt to input

(2) derivative of output wrt parameter

Slide credit: Deva Ramanan



 $\frac{\partial y}{\partial w}$: needed to compute gradiant updates for this block

 $\frac{\partial y}{\partial x}$: needed to compute gradiant updates next block down stream

Computational savings: cache intermediate gradients





Slide credit: Deva Ramanan

Artificial Neural Networks



 $x_i \in \mathbb{R}^{H \times 1}$ $W_i \in \mathbb{R}^{D \times H}$ $b_i \in \mathbb{R}^{D \times 1}$

$$x_{i+1} = f(W_i x_i + b_i)$$

$$\min_{\theta} \sum_{j} \mathscr{L}\left(x_{last}^{j}, y^{j}\right)$$

Non-linearities: sigmoid



- Interpretation as firing rate of neuron
- Bounded between [0,1]
- Saturation for large +ve,-ve inputs
- Gradients go to zero
- Outputs centered at 0.5 (poor conditioning)
- Not used in practice

Non-linearities: tanh

$$f(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$



- Bounded between [-1,+1]
- Saturation for large +ve,-ve inputs
- Gradients go to zero
- Outputs centered at 0
- Preferable to sigmoid

tanh(x) = 2 sigmoid(2x) - 1

Non-linearities: rectified linear (ReLU)

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



- Unbounded output (on positive side)
- Efficient to implement:

$$f'(a) = \frac{df}{da} = \begin{cases} 0 & a < 0\\ 1 & a \ge 0 \end{cases}$$

- Also seems to help convergence (see 6x speedup vs tanh in [Krizhevsky et al.])
- Drawback: if strongly in negative region, unit is dead forever (no gradient).
 - Default choice: widely used in current models. Slide credit: Antonio Torralba
Non-linearities: Leaky ReLU

$$f(a) = \begin{cases} \max(0,a) & a > 0\\ a\min(0,a) & a < 0 \end{cases}$$



- where α is small (e.g. 0.02)
- Efficient to implement:

$$f'(a) = \frac{df}{da} = \begin{cases} -a & a < 0\\ 1 & a > 0 \end{cases}$$

- Also known as probabilistic ReLU (PReLU)
- Has non-zero gradients everywhere (unlike ReLU)
- α can also be learned (see Kaiming He et al. 2015).

Multilayer Perceptron (MLP)



 $x_i \in \mathbb{R}^{H \times 1}$ $W_i \in \mathbb{R}^{D \times H}$ $b_i \in \mathbb{R}^{D \times 1}$

$$x_{i+1} = f\left(W_i x_i + b_i\right)$$

$$\min_{\theta} \sum_{j} \mathscr{L}\left(x_{last}^{j}, y^{j}\right)$$

Convolutional Networks

PROC. OF THE IEEE, NOVEMBER 1998

1

Gradient-Based Learning Applied to Document Recognition

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner



Stack together convolution and pooling (avg + subsample) operations.

Convolutional Layer



$$x_i \in \mathbb{R}^{W \times H \times D}$$

 $x_{i+1} \in \mathbb{R}^{W \times H \times 1}$

Convolutional Layer



$$x_i \in \mathbb{R}^{W \times H \times D}$$

 $x_{i+1} \in \mathbb{R}^{W \times H \times 1}$

Convolutional Layer



 $x_{i+1} \in \mathbb{R}^{W \times H \times K}$

A better visualization of AlexNet



Max Pooling

Fast way to resize an "image"

What's the derivative?





y

Х

max pool with 2x2 filters and stride 2

6	8
3	4



All filter dimensions 3x3 except fc6 (which uses 7x7)

Inception / resnet



Convolution Pooling Softmax Other





Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.





Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

Get to know your units



	-		•				•	1.000	1
					in:	100			3
•				*	nin in	ou the			
			un ann			0			æ
	-		•			N.			ni san
100	-				a.	100.000			•
-			•	••	•			•	
100	0	• •	••••	N. S.	••				
	AV AV		•	•••	-		· .		• •
• •	· ·		10 M						

96 Units in conv1

Gabor wavelets



Slide credit: Antonio Torralba

Fourier transform of a Gabor wavelet



$$\psi_{c}(x,y) = e^{-\frac{x^{2} + y^{2}}{2\sigma^{2}}} \cos(2\pi u_{0}x)$$







Slide credit: Antonio Torralba



Fig. 5. Top row: illustrations of empirical 2-D receptive field profiles measured by J. P. Jones and L. A. Palmer (personal communication) in simple cells of the cat visual cortex. Middle row: best-fitting 2-D Gabor elementary function for each neuron, described by (10). Bottom row: residual error of the fit, indistinguishable from random error in the Chisquared sense for 97 percent of the cells studied.





Top Activated Images Interpretation: lamp



Top Activated Images Interpreta





conv5 unit 79 car (object) IoU=0.13



conv5 unit 107 road (object)





Slide credit: Bolei Zhou

conv5 unit 144 mountain (object) IoU=0.13



conv5 unit 200 mountain (object)





Slide credit: Bolei Zhou

House



conv5_3 unit 243



conv5 unit 36

GoogLeNet inception_4e unit 789

loU=0.137

res5c unit 1410

IoU=0.142

IoU=0.053

IoU=0.070

ResNet



Airplane

conv5 unit 13

IoU=0.101



conv5_3 unit 151

IoU=0.150



inception_4e unit 92

loU=0.164



res5c unit 1243

IoU=0.172



Emergence of Interpretable Units during Training





Optimization

Appears to be a significant hurdle for training

Deep Residual Learning for Image Recognition



Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer "plain" networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

Intuition: bias deep models to behave like shallow models during learning



Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.



Batch normalization

[Ioffe et al]

Intuition: build second-order behaviour into SGD by normalizing variables (zero-mean, identity covariance) before nonlinearity



Many (if not most) contemporary networks make use of this

Drop-out regularization



Intuition: we should really train a family of models with different architectures and average their predictions (c.f. model averaging from machine learning)

Practical implementation: learn a single "superset" architecture that randomly removes nodes (by randomly zero'ing out activations) during gradient updates