## Intro to Deep Learning neural networks, CNNs, backpropagation



CSC420 David Lindell University of Toronto <u>cs.toronto.edu/~lindell/teaching/420</u> Slide credit: Babak Taati ←Ahmed Ashraf ←Sanja Fidler



# Logistics

•HW2 is out, due in 3 weeks

- Motivation
- Fully-connected Networks
- Convolutional Neural Networks
- •Training networks

•Let's take some typical tourist picture. What all do we want to recognize?



•Identification



• Scene classification: what type of scene is the picture showing?



• Classification: Is the object in the window a person, a car, etc



•Image Annotation: Which types of objects are present in the scene?



• Detection: Where are all objects of a particular class?



•Segmentation: Which pixels belong to each class of objects?



• Pose estimation: What is the pose of each object?



•Attribute recognition: Estimate attributes of the objects (color, size, etc)



•Action recognition: What is happening in the image?



• Surveillance: Why is something happening?



### Have we encountered these things before?

- Before we proceed, let's first give a shot to the techniques we already know
- Let's try detection (how?)



#### Have we encountered these things before?

- Before we proceed, let's first give a shot to the techniques we already know
- •Let's try detection (how?)
- Example techniques:
  - Template matching (remember Waldo in Lecture 3-5?)
  - Large-scale retrieval: store millions of pictures, recognize new one by finding the most similar one in database. This is a Google approach.



#### •Template matching: normalized cross-correlation with a template (filter)

Find the chair in this image

chair template





•Template matching: normalized cross-correlation with a template (filter)



template

Find the chair in this image



•Template matching: normalized cross-correlation with a template (filter)



template

Find the chair in this image



Pretty much garbage Simple template matching is not going to make it



why?

•Template matching: normalized cross-correlation with a template (filter)



Find the chair in this image



Pretty much garbage Simple template matching is not going to make it



A "popular method is that of template matching, by point to point correlation of a model pattern with the image pattern. These techniques **are inadequate for three-dimensional scene analysis for many reasons, such as occlusion, changes in viewing angle, and articulation of parts**." Nevatia & Binford, 1977.

• Upload a photo to Google image search and check if something reasonable comes out

| Googl   | e      |
|---|--------|
|   | ۹ ا    |
| Search by image<br>Search Google with an image instead of text. | ×      |
| Paste image URL ②   Upload an image                             | Search |



- Upload a photo to Google image search
- Pretty reasonable, both are Golden Gate Bridge





• Upload a photo to Google image search Let's try a typical bathtub object

| Google                              |        |  |
|-------------------------------------|--------|--|
|                                     |        |  |
| Search by image                     | ×      |  |
| Paste image URL ②   Upload an image | Search |  |

- Upload a photo to Google image search
- •A bit less reasonable, but still some striking similarity





- Make a beautiful drawing and upload to Google image search
- Can you recognize this object?

| Google  |        | M    |
|---|--------|------|
|   |        |      |
| Search by image<br>Search Google with an image instead of text. | ×      | h my |
| Paste image URL ②   Upload an image                             |        |      |
|   | Search |      |

- Make a beautiful drawing and upload to Google image search
- •Not a very reasonable result



## Why is it a Problem?

• Difficult scene conditions



[From: Grauman & Leibe]

### Why is it a Problem?

• Huge within-class variations. Recognition is mainly about modeling variation.



[Pic from: S. Lazebnik]

#### Why is it a Problem?

•Tons of classes



Biederman

•We cannot explicitly model these variations!

- •We cannot explicitly model these variations!
- •Instead our models should be relatively simple and we should learn let complexity live in the data

- •We cannot explicitly model these variations!
- •Instead our models should be relatively simple and we should learn let complexity live in the data
- •Neural networks follow this paradigm

#### Motivation

#### • Fully-connected Networks

- Convolutional Neural Networks
- •Training networks

# Image Classification

• Image classification example

#### Images



## Image Classification

• Image classification example

#### Images

3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 З 4 E Y A Y Y Y Y Y 4 4 4 9 4 4 F **7 7 7 7 7 7 7** N **9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9**  Class

"zero" "one"

"nine"

. . .

## Image Classification


• Image classification example

#### Images



#### Challenges

#### Intra-class variation

- stroke widths
- alignment
- writing styles

Challenges

• Image classification example

#### Images



• Image classification example

#### Images

ß D ູ З 4 4 4 4 4 ч ¥ S 5 5 G F フクフフ ч Π Е ¥ в **9 9 9 9** 9 9 9 

Implementation?

Can't hardcode solution!

- Data-driven approach
  - Collect training images and labels
  - Train a classifier using machine learning
  - Evaluate the classifier on unseen images

#### Implementation?













• Linear Model f(x, W) = Wx



Length of this vector is the "dimensionality" of our problem!

• Linear Model f(x, W) = Wx



In general: Wx + b



• Linear model: geometric intrepretation



Each image is a point in an Ndimensional space

- N is the number of pixels

• Linear model: geometric interpretation



$$f(x, W) = Wx$$

Computes inner product between rows of W and x!

- Each row of W is a hyperplane
- Sign of inner product tells you which side of the hyperplane
- "separates" the digits

• Linear model (visual interpretation)

Learned filters (rows of W)



• Limits of linear classifiers

Linear classifiers learn linear decision planes

What if dataset is not linearly separable?



- Linear Model f = Wx
- 2-layer MLP  $f = W_2 \max(0, W_1 x)$

- Linear Model f = W x
- 2-layer MLP  $f = W_2 \max(0, W_1 x)$
- 3-layer MLP  $f = W_3 \max(0, W_2 \max(0, W_1 x))$

- Linear Model f = W x
- 2-layer MLP  $f = W_2 \max(0, W_1 x)$
- 3-layer MLP  $f = W_3 \max(0, W_2 \max(0, W_1 x))$

Non-linearity/activation function between linear layers

- Linear Model f = W x
- 2-layer MLP  $f = W_2 \max(0, W_1 x)$
- 3-layer MLP  $f = W_3 \max(0, W_2 \max(0, W_1 x))$

Otherwise we have:

 $f = W_3 W_2 W_1 x$ 

#### **Activation Functions**

#### ...many to choose from



... ReLU is a good general-purpose choice: ReLU(x) = max(0, x)

- Linear Model f = W x
- 2-layer MLP  $f = W_2 \max(0, W_1 x)$

Back to our classification example...



- Linear Model f = W x
- 2-layer MLP  $f = W_2 \max(0, W_1 x)$

Back to our classification example...



 $x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$ 

- Linear Model f = W x
- 2-layer MLP  $f = W_2 \max(0, W_1 x)$

Back to our classification example...



• Overcomes limits of linear classifiers

- Can learn non-linear decision
  boundaries
- Complexity scales with the number of neurons/hidden layers



- More parameters is not always better!
  - Can lead to overfitting the training data
  - Performance on test data is worse



train

test



- More on classification...
  - https://cs231n.github.io/linearclassify/
  - https://amfarahmand.github.io/N N-Winter2024/

















Image: CC BY-SA Jennifer Walinga

Loose analogy!

- Neurons have activation potentials, all-or-none firing behavior
- Interconnectivity between actual neurons is dense and complicated
- Connection between neurons is complex non-linear dynamical system



#### Drawbacks of fully-connected networks



- spatial structure is destroyed
- fully-connected weights do not scale

#### Overview

#### Motivation

- Fully-connected Networks
- Convolutional Neural Networks
- •Training networks

#### Convolutional Neural Networks



- Exploit spatial structure
- Scale to large inputs with fewer parameters
- Remarkable performance for processing visual data
AlexNet & surge in popularity

2010: ImageNet Large Scale Visual Recognition Challenge

14 million labeled images

First convolutional network for image classification



AlexNet [Krizhevsky '12]

AlexNet & surge in popularity

2010: ImageNet Large Scale Visual Recognition Challenge

• 14 million labeled images

First convolutional network for image classification

0.16

2011 2012 2013

ILSVRC year [Russakovsky '15]

0.12

0.07

2014

\_ 0.28

2010

0.26

0.3

0.2

0.1

Classification error



AlexNet & surge in popularity

2010: ImageNet Large Scale Visual Recognition Challenge

• 14 million labeled images





Fully-Connected Layer





Input Image

Filter



Input Image

Filter



Input Image



Input Image



Input Image Activations



Input Image Activations



#### https://github.com/vdumoulin/conv\_arithmetic



Multiple output channels using multiple filters

Input Image Activations

#### Fully-Connected Layer



Special case of convolutional layer when filter size = input size!

#### Convolutional Neural Network



Input Image

Layer 1 Activations

Layer 2 Activations

### Input Image



### First-layer Filters





### Activations

Input Image







First-layer Filters



Image: CC BY-SA Selket



#### [Hubel & Wiesel 1959]



Simple cells in visual cortex detect edges, complex cells compose earlier responses



Dataset examples that maximize neuron outputs



#### Dataset examples that maximize neuron outputs



Dataset examples that maximize neuron outputs



Dataset examples that maximize neuron outputs

Design choices:

- filter size
- number of filters
- padding
- stride

Layer types:

- pooling
- transpose convolutions
- upsampling layers\*
- batch normalization\*
- softmax layers\*

\*no time to cover all of these layers! Check out PyTorch docs for details... e.g., https://pytorch.org/docs/stable/generated/torch.nn.BatchNorm2d.html



Number of channels

N\_out x N\_in x 5 x 5









https://github.com/vdumoulin/conv\_arithmetic

#### Convolutional Neural Network



Input Image

Layer 1 Activations

Layer 2 Activations Layer types: Pooling



e.g., max pool size=2, stride=2

•AlexNet (from UofT!): A. Krizhevsky, I. Sutskever, G. E. Hinton, ImageNet Classification with Deep Convolutional Neural Networks, NeurIPS 2012. This network won the Imagenet Challenge of 2012, and revolutionized computer vision.



[Pic adopted from: A. Krizhevsky]

- •AlexNet (from UofT!): A. Krizhevsky, I. Sutskever, G. E. Hinton, ImageNet Classification with Deep Convolutional Neural Networks, NeurIPS 2012. This network won the Imagenet Challenge of 2012, and revolutionized computer vision.
- How many parameters (weights) does this network have?



[Pic adopted from: A. Krizhevsky]

- Trained with stochastic gradient descent on two NVIDIA GPUs for about a week
- 650,000 neurons
- 60,000,000 parameters
- 630,000,000 connections
- Final feature layer: 4096-dimensional

**Convolutional layer:** convolves its input with a bank of 3D filters, then applies point-wise non-linearity

**Fully-connected layer:** applies linear filters to its input, then applies pointwise non-linearity

Figure: From http://www.image-net.org/challenges/LSVRC/2012/supervision.pdf

[Pic adopted from: A. Krizhevsky]

Image

• The trick is to not hand-fix the weights, but to train them. Train them such that when the network sees a picture of a dog, the last layer will say "dog".



[Pic adopted from: A. Krizhevsky]

•Or when the network sees a picture of a cat, the last layer will say "cat".



[Pic adopted from: A. Krizhevsky]

•Or when the network sees a picture of a boat, the last layer will say "boat"... The more pictures the network sees, the better.



Train on **lots** of examples. Millions. Tens of millions. Wait a week for training to finish.

Share your network (the weights) with others who are not fortunate enough with GPU power. [Pic adopted from: A. Krizhevsky]

# Classification

•Once trained we can do classification. Just feed in an image or a crop of the image, run through the network, and read out the class with the highest probability in the last (classification) layer.


## Overview

#### Motivation

- Fully-connected Networks
- Convolutional Neural Networks
- •Training networks

# Image Inpainting



masked input

predicted output

Image inpainting example

Training dataset:

- masked and complete image pairs
- train network to predict the complete image

#### masked images

| Λ  | Λ  | Λ | Λ  | Λ | ۸ | () | 0  | n  | $\land$ | $\cap$ | $\mathbf{c}$ | Λ  | n | $\mathbf{h}$ | 5  |
|----|----|---|----|---|---|----|----|----|---------|--------|--------------|----|---|--------------|----|
| 1  | 1  | 1 | ١  | ſ | 1 | •  | 1  | 1  | 1       | ١      | 1            | ł  | 1 | 1            | 1  |
| 2  | נ  | С | 7  | 2 | 2 | 2  | 2  | 2  | 2       | 9      | 1            | 3  | 2 | 2            | 1  |
| 2  | 2  | Z | 2  | 2 | 7 | 2  | a. | ጓ  | 2       | 3      | 2            | 7  | 7 | 2            | 2  |
| 11 | 16 | L | U  | А | Ш | 11 | J  | 11 | 11      | A.     | Ц            | Λ  | u | 17           | 11 |
| 5  | 5  | 5 | .٣ | < | C | ۲  | б  | ~  | ~       | <      | 5            | ٢  | 5 | -م           | <  |
| 1~ | ſ  | 4 | L  | 1 | ٢ | 1- | 1- | L  | /       | 4      | 1            | 1. | 1 | 1            | 1. |
| γ  | ٦  | 5 | -7 | 7 | ٦ | Н  | 7  | 2  | Π       | 7      | 7            | 7  | 7 | 7            | 7  |
| V  | đ  | ç | •  | 0 | 6 | D  | 17 | 0  |         | C      | \$           | Ś  | 0 | 0            | 4  |
| \$ | G  | a | ۵  | ۵ | 0 | 13 | a  | a  | ñ       | ۵      | 0            | A  | a | 0            | O, |

#### ground truth



Train the network to minimize the loss function

network parameters 
$$\mathcal{L}_{ heta} = rac{1}{2} \|y - \hat{y}\|_2^2$$
  
 $\theta = \{W_1, W_2\}$ 

Train the network to minimize the loss function



How do we figure out  $\theta$  ?



Gradient-based optimization

 $\theta^{(k+1)} = \theta^{(k)} - \nabla_{\theta} \mathcal{L}_{\theta}$ 



Gradient-based optimization

$$\theta^{(k+1)} = \theta^{(k)} - \nabla_{\theta} \mathcal{L}_{\theta}$$

Need to calculate the partial derivative with respect to each parameter



Generally there are 3 options

- 1. Numerical differentiation
- 2. Symbolic differentiation
- 3. "Automatic" differentiation

### Numerical Differentiation

$$\frac{\partial f(x)}{\partial x} \approx \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

Not very accurate, computationally expensive

Easy to implement! Can be used to check your analytical answers..

# Symbolic Differentiation

$$\begin{aligned} \frac{\partial \mathcal{L}_{\theta}}{\partial W_{1}} &= \frac{\partial}{\partial W_{1}} \frac{1}{2} \|y - \hat{y}\|_{2}^{2} \\ &= \frac{\partial}{\partial W_{1}} \frac{1}{2} \left( W_{2} \sigma(W_{1}x) \right)^{T} \left( W_{2} \sigma(W_{1}x) \right) \\ &= \frac{\partial}{\partial W_{1}} \frac{1}{2} \sigma(W_{1}x)^{T} W_{2}^{T} W_{2} \sigma(W_{1}x) \\ &= \dots \quad \text{chain rule, product rule...} \end{aligned}$$

Accurate, but must be manually calculated for each term Tedious!

Think about the problem as a "computational graph"

Divide and conquer using the chain rule

Enables "backpropagation" – an efficient way to take derivatives of all parameters in a computational graph

Think about the problem as a "computational graph"

Divide and conquer using the chain rule



Think about the problem as a "computational graph"

Divide and conquer using the chain rule



Think about the problem as a "computational graph"

Divide and conquer using the chain rule



Think about the problem as a "computational graph"

Divide and conquer using the chain rule



 $\frac{\partial \mathcal{L}}{\partial W_1} = \frac{\partial f}{\partial W_1} \frac{\partial g}{\partial f} \frac{\partial \hat{y}}{\partial g} \frac{\partial \mathcal{L}}{\partial \hat{y}}$ 

We can calculate analytical expressions for each of these terms and then plug in our values



| $\partial \mathcal{L}$ _    | $\partial f$              | $\partial g$            | $\partial \hat{y}$      | $\partial \mathcal{L}$ |
|-----------------------------|---------------------------|-------------------------|-------------------------|------------------------|
| $\overline{\partial w_1}$ – | $\overline{\partial w_1}$ | $\overline{\partial f}$ | $\overline{\partial g}$ | $\partial \hat{y}$     |





































What is backpropagation?



Save these intermediate values during forward computation
What is backpropagation?



Then we perform a "backward pass"



















### 1. Sample batch of images from dataset

| masked images                                     |   |
|---|---|
| indente di midige e                               | ground truth                                  |
|   | 000000000000000000                            |
|   | (           / /   /           / /             |
|   | 2222222222222222                              |
| 2 <b>2 7 7 9 7 9 7 7 7 7 7 7 7 7 7 7 7 7 7 7 </b> | 3 <b>333333</b> 33333333333                   |
| IT IL L CI A LI II CI A II A CI A N II            | 44844444444444444                             |
| <b>FFFFCCCFFF</b> FFFFF                           | 5 <b>5555555555</b> 5555555555555555555555555 |
|   | 6666666666666666666                           |
| ר ק <b>ה ה</b> ר ח <b>ח א ר ר ר ה 0 ח</b> ר ח     | <b>モフクコフ</b> フ ヤ <b>ク</b> クフ <b>フタ</b> クフフ    |
| 0 a c <b>e 0 6 0 n v e</b> c vo o n n             | 888888888888888888888888888888888888888       |
| •           | <b>9999999999999999</b>                       |

### 1. Sample batch of images from dataset



2. Run forward pass to calculate network output for each image



1. Sample batch of images from dataset

| masked images  | around truth                            |
|--|---|
| -  | giouna tratti                           |
|  | 000000000000000000                      |
|  | 111111111111111111                      |
|  | 22222222222222222                       |
| - <b>3 1 3 3 3 3 3 4 7 8 3 3 3 3 3</b> 3                   | 3333333333333333333                     |
| IT RECTALLICENTIAL AND IN                                  | 448444444444444444                      |
| K K K K K K K K K K K K K K K K K K K                      | 555555555555555555555555555555555555555 |
|  | 6666666666666666                        |
| רק <b>הר</b> ר רח <b>הח</b> א <b>רר ר רט</b> ר ה           | <b>モフクフフ</b> フ ゼ <b>クク</b> クフ <b>フ</b>  |
| <i>V &amp; C &amp; O &amp; O D D &amp; C &amp; O O A A</i> | 888888888888888888888                   |
| •                    | 999999999999999999999                   |

2. Run forward pass to calculate network output for each image



3. Run backward pass to calculate gradients with backpropagation

1. Sample batch of images from dataset

| masked images  | around truth                                     |
|--|--|
|  |  |
|  | 000000000000000000000000000000000000000          |
|  | /         / / / / /     / / / /                  |
|  | 2222222222222222                                 |
| 2 <b>3 1</b> 3 9 7 9 7 8 3 8 9 9 9                       | 3 <b>33333</b> 333333333333                      |
| IL BE CLAUTER A HA UNIT                                  | 4484444444444444                                 |
| <b><i>KEEKKEEKKEEKEE</i></b>                             | 555555555555555555555555555555555555555          |
|  | 666666666666666                                  |
| сч <b>пс</b> га <b>ччч</b> и сч                          | <b>モフクフフフ</b> セ <b>クク</b> クフ <b>フ</b>            |
| <i>va</i> c <b>e</b> e e e e e e e e e e e e e e e e e e | 8          |
| <b>BBBBBBBBBBBBB</b>                                     |  |
|  | <b>7 2 4 4 4 4 9 9 9 9 9</b> 9 9 9 9 9 9 9 9 9 9 |

2. Run forward pass to calculate network output for each image



- 3. Run backward pass to calculate gradients with backpropagation
- 4. Update parameters with stochastic gradient descent

4. Update parameters with stochastic gradient descent

$$\mathcal{L}_{\theta} = \| y - \hat{y} \|_{2}^{2}$$
$$W_{2}^{(k+1)} = W_{2}^{(k)} - \alpha \frac{\partial \mathcal{L}}{\partial W_{2}}$$
$$W_{1}^{(k+1)} = W_{1}^{(k)} - \alpha \frac{\partial \mathcal{L}}{\partial W_{1}}$$





### Takeaways

- What is a neural network? (can you write the equation for an MLP?)
- Basic building blocks/architecture of CNN
- Backpropagation, automatic differentiation, and gradient descent

### Next Time

• Embedded ethics lecture!

Supplemental Slides





Recap: vector differentiation

Scalar by Scalar

$$x, y \in \mathbb{R}$$
$$\frac{\partial y}{\partial x} \in \mathbb{R}$$



Recap: vector differentiation

Scalar by ScalarScalar by Vector $x, y \in \mathbb{R}$  $x \in \mathbb{R}^N, y \in \mathbb{R}$  $\frac{\partial y}{\partial x} \in \mathbb{R}$  $\frac{\partial y}{\partial x} \in \mathbb{R}^N$ 



**Recap: vector differentiation** 

Scalar by Scalar  $x, y \in \mathbb{R}$  $\frac{\partial y}{\partial x} \in \mathbb{R}$ 

Scalar by Vector  $x \in \mathbb{R}^N, y \in \mathbb{R}$   $x \in \mathbb{R}^N, y \in \mathbb{R}^M$  $\frac{\partial y}{\partial x} \in \mathbb{R}^N$ 

Vector by Vector  $\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M}$ 

### Recap: vector differentiation



Example 1: matrix multiply

$$\frac{\partial \hat{y}}{\partial g} = \frac{\partial}{\partial g} W_2 g$$
$$W_2 \in \mathbb{R}^{M \times N}$$
$$g \in \mathbb{R}^N$$
$$\frac{\partial \hat{y}}{\partial g} \in \mathbb{R}^{N \times M}$$

### Recap: vector differentiation



Example 1: matrix multiply

$$\frac{\partial \hat{y}}{\partial g} = \frac{\partial}{\partial g} W_2 g W_2 \in \mathbb{R}^{M \times N} g \in \mathbb{R}^N \frac{\partial \hat{y}}{\partial g} \in \mathbb{R}^{N \times M}$$

$$\frac{\partial \hat{y}}{\partial g} \in \mathbb{R}^{N \times M}$$

$$\frac{\partial \hat{y}}{\partial g} \in \mathbb{R}^{N \times M}$$

$$\frac{\partial \hat{y}}{\partial g} = \mathbb{R}^{N \times M}$$

$$\frac{\partial \hat{y}}{\partial g} = \mathbb{R}^{N \times M}$$

$$\frac{\partial \hat{y}}{\partial g} = \mathbb{R}^{N \times M}$$





Recap: vector differentiation Example 2: elementwise functions

 $h = f \odot g$  $f \in \mathbb{R}^{N}$  $g \in \mathbb{R}^{N}$  $\frac{\partial h}{\partial f} \in \mathbb{R}^{N \times N}$ 

### Recap: vector differentiation Example 2: elementwise functions

$$\begin{split} h &= f \odot g \\ f \in \mathbb{R}^{N} \\ g \in \mathbb{R}^{N} \\ \frac{\partial h}{\partial f} \in \mathbb{R}^{N \times N} \end{split} \qquad \qquad \begin{aligned} \frac{\partial h}{\partial f} &= \begin{bmatrix} \frac{\partial h_{1}}{\partial f_{1}} & \cdots & \frac{\partial h_{n}}{\partial f_{1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_{1}}{\partial f_{n}} & \cdots & \frac{\partial h_{n}}{\partial f_{n}} \end{bmatrix} \\ \frac{\partial h}{\partial f} &= \begin{bmatrix} g_{1} & & 0 \\ & \ddots & \\ 0 & & g_{n} \end{bmatrix} = \operatorname{diag}(g) \end{split}$$

### Recap: vector differentiation

Final hint: dimensions should always match up!



You should be able to calculate derivatives of each of these terms and then perform matrix multiplications without issues

Extra backpropagation example (adapted from Stanford CS231n)

$$f = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$










Extra backpropagation example



Extra backpropagation example



Extra backpropagation example

