



Lecture 1: Introduction

Fei-Fei Li & Justin Johnson & Serena Yeung

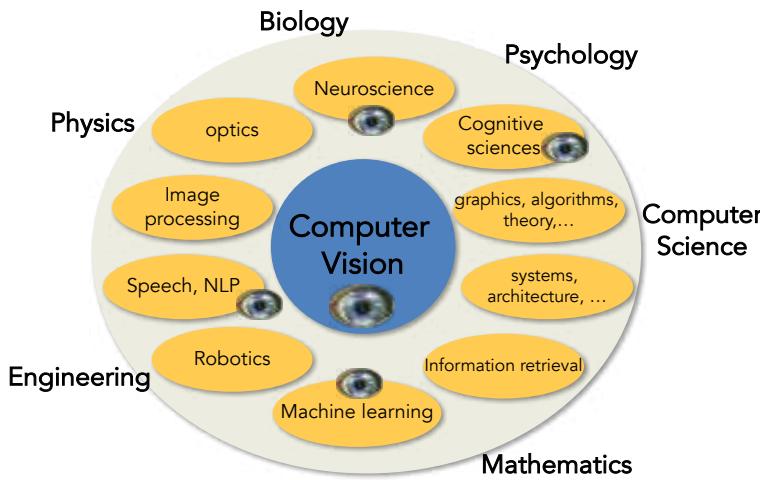
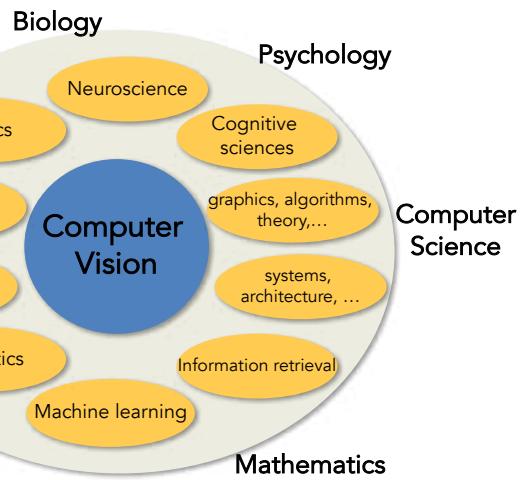
Lecture 1 - 1

4/3/2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 2

4/3/2018



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 3

4/3/2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 4

4/3/2018

Related Courses @ Stanford

- CS131: Computer Vision: Foundations and Applications
 - Fall 2017, Juan Carlos Niebles and Ranjay Krishna
 - Undergraduate introductory class
- CS231a: Computer Vision, from 3D Reconstruction to Recognition
 - Winter 2018, Professor Silvio Savarese
 - Core computer vision class for seniors, masters, and PhDs
 - Image processing, cameras, 3D reconstruction, segmentation, object recognition, scene understanding; not just deep learning
- CS 224n: Natural Language Processing with Deep Learning
 - Winter 2018, Richard Socher
- CS 230: Deep Learning
 - Spring 2018, Prof. Andrew Ng and Kian Katanforoosh
- **CS231n: Convolutional Neural Networks for Visual Recognition**
 - This course, Prof. Fei-Fei Li & Justin Johnson & Serena Yeung
 - Focusing on applications of deep learning to computer vision

Today's agenda

- A brief history of computer vision
- CS231n overview

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 5

4/3/2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 6

4/3/2018

Evolution's Big Bang



543 million years, B.C. *vision in animals*

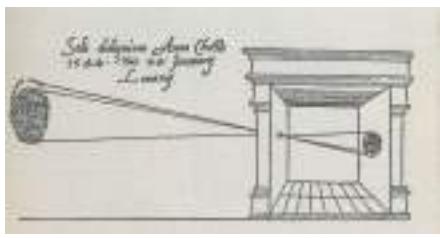
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 7

4/3/2018

Camera Obscura

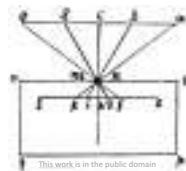
Gemma Frisius, 1545



Encyclopedie, 18th Century



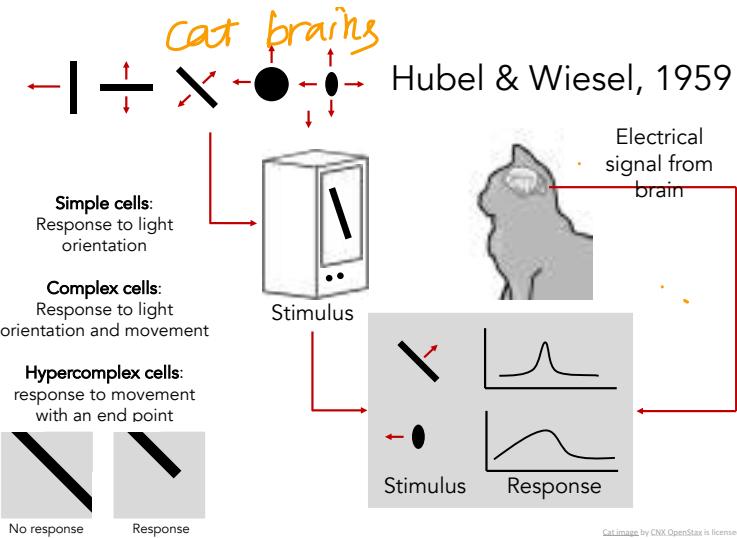
Leonardo da Vinci, 16th Century AD



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 8

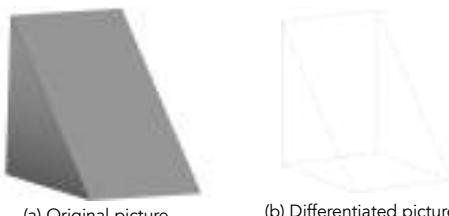
4/3/2018



Block world

Larry Roberts, 1963

early fbs



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 9

4/3/2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 10

4/3/2018

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
PROJECT MAC

Artificial Intelligence Group
Vision Issue, No. 100.

July 1, 1968

THE SUMMER VISION PROJECT

Seymour Papert

The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real landmark in the development of "pattern recognition".

VISION



David Marr

Summer of 1970
in support of
Project MAC

David Marr, 1970s

Fei-Fei Li & Justin Johnson & Serena Yeung

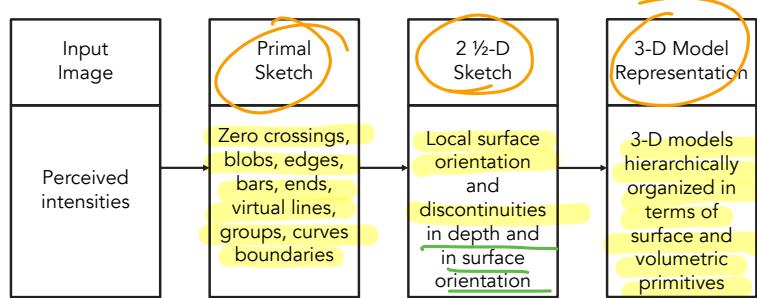
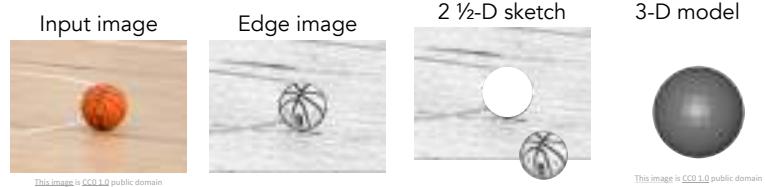
Lecture 1 - 11

4/3/2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 12

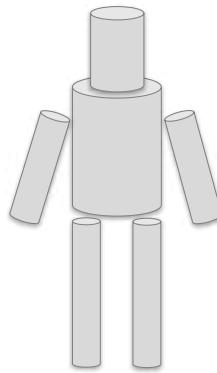
4/3/2018



Stages of Visual Representation, David Marr, 1970s

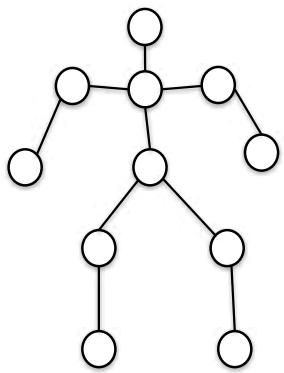
• Generalized Cylinder

Brooks & Binford, 1979



• Pictorial Structure

Fischler and Elschlager, 1973



group pixels



David Lowe, 1987

Normalized Cut (Shi & Malik, 1997)

Image is CC BY 3.0

Image is public domain

Image is CC-BY 2.0; changes made



Face Detection, Viola & Jones, 2001

face detection

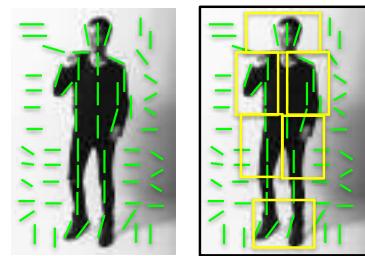
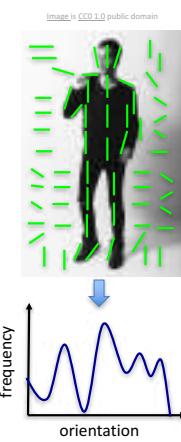
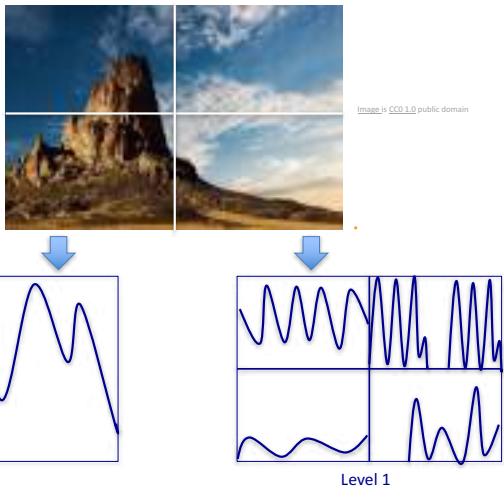


feature-based object recognition



"SIFT" & Object Recognition, David Lowe, 1999

identify visual features



Deformable Part Model
Felzenswalb, McAllester, Ramanan, 2009

Histogram of Gradients (HoG)
Dalal & Triggs, 2005

PASCAL Visual Object Challenge (20 object categories)

[Everingham et al. 2006-2012]



detect images

IMAGENET

www.image-net.org

22K categories and 14M images

recognize objects

- Animals
- Plants
- Structures
- Person
- Bird
- Tree
- Artifact
- Scenes
- Fish
- Flower
- Tools
- Indoor
- Mammal
- Food
- Appliances
- Geological Formations
- Invertebrate
- Materials
- Structures
- Sport Activities

IMAGENET Large Scale Visual Recognition Challenge

The Image Classification Challenge:

1,000 object classes

1,431,167 images



Output:
Scale
T-shirt
Steel drum
Drumstick
Mud turtle



Output:
Scale
T-shirt
Giant panda
Drumstick
Mud turtle



complex → easy to overfit

IMAGENET Large Scale Visual Recognition Challenge

The Image Classification Challenge:

1,000 object classes

1,431,167 images

error rate

drops significantly



CNN

Today's agenda

- A brief history of computer vision
- CS231n overview

CS231n focuses on one of the most fundamental problems of visual recognition –
image classification



Image by US Army is licensed under CC BY 2.0



Image is CC0 1.0 public domain



Image by Karpelhow is licensed under CC BY-SA 3.0



Image by Christina C. is licensed under CC BY-SA 4.0

There are many visual recognition problems that are related to image classification, such as
object detection, image captioning

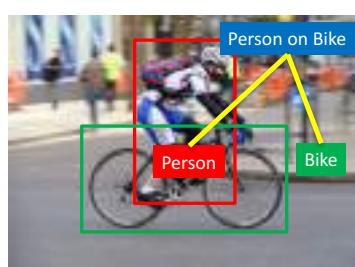


This image is licensed under CC BY-NC-SA 2.0 changes made

- Object detection
- Action classification
- Image captioning
- ...

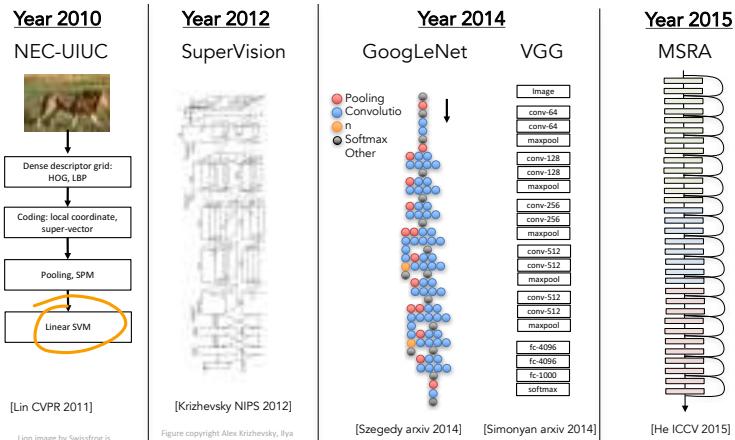


Person
Hammer



Convolutional Neural Networks (CNN) have become an important tool for object recognition

IMAGENET Large Scale Visual Recognition Challenge



Convolutional Neural Networks (CNN)
were not invented overnight

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 31

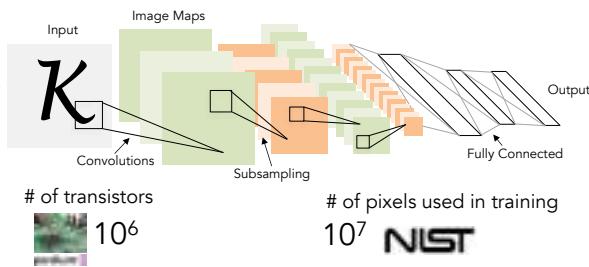
4/3/2018

Fei-Fei Li & Justin Johnson & Serena Yeung

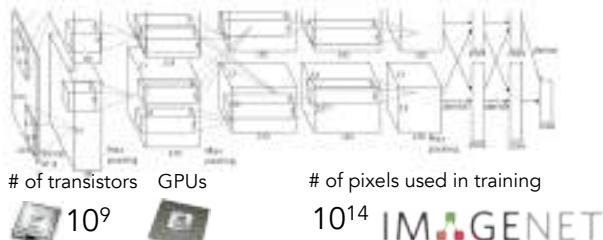
Lecture 1 - 32

4/3/2018

1998
LeCun et al.



2012
Krizhevsky et al.



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 33

4/3/2018

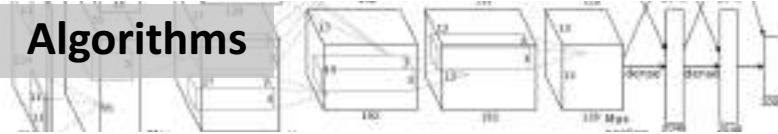
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 34

4/3/2018

Ingredients for Deep Learning

Algorithms



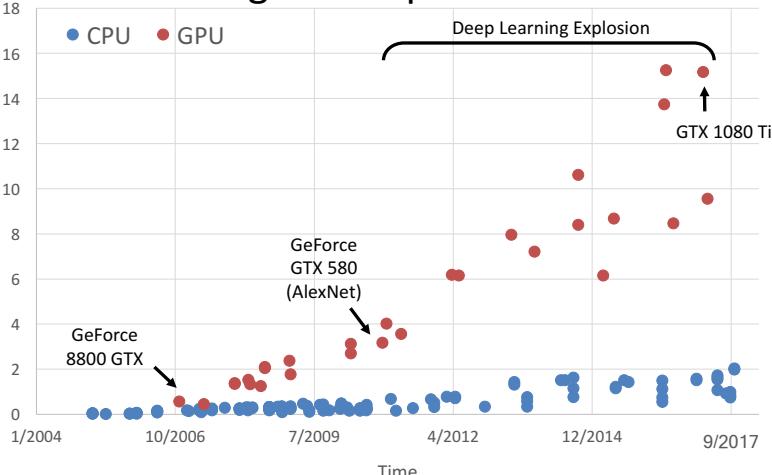
Data



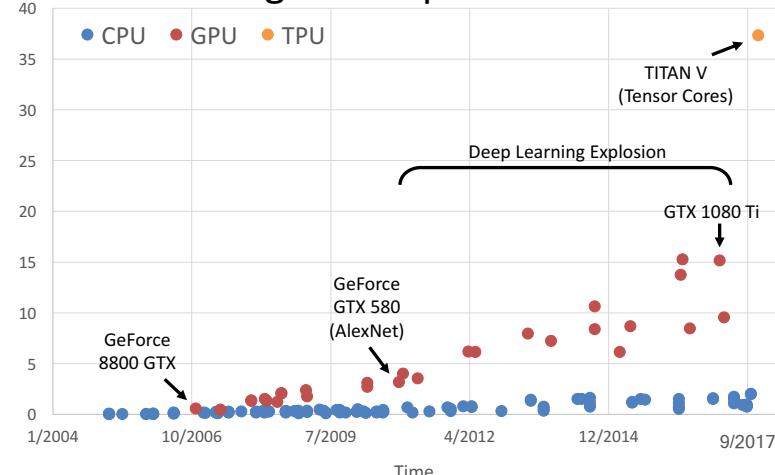
Computation



GigaFLOPs per Dollar



GigaFLOPs per Dollar



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 35

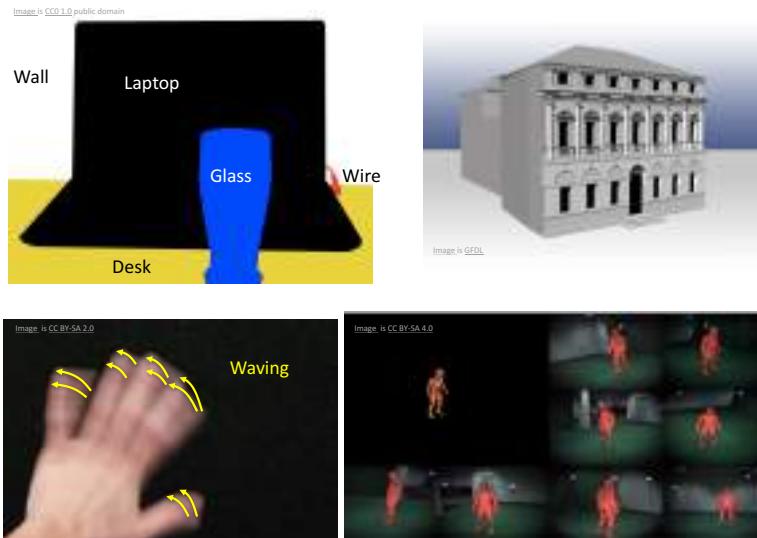
4/3/2018

Fei-Fei Li & Justin Johnson & Serena Yeung

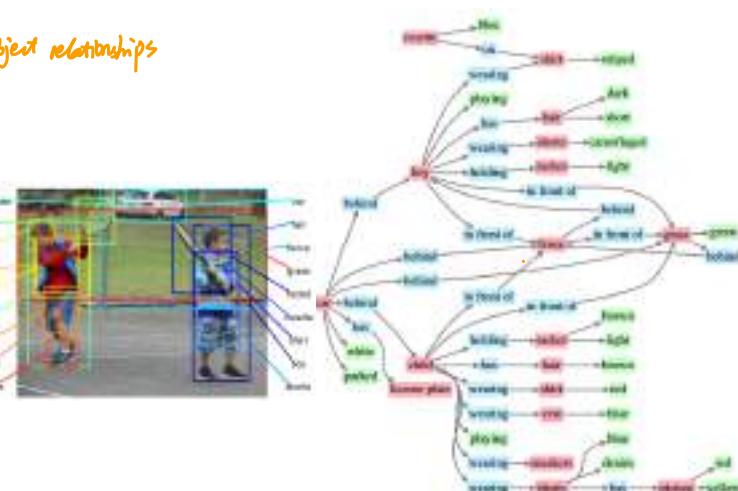
Lecture 1 - 36

4/3/2018

The quest for visual intelligence
goes far beyond object recognition...



object relationships



$PT = 500ms$

Some kind of game or fight. Two groups of two men? The man on the left is throwing something. Outdoors seemed like because i have an impression of grass and maybe lines on the grass? That would be why I think perhaps a game, rough game though, more like rugby than football because they pairs weren't in pads and helmets, though I did get the impression of similar clothing, maybe some trees? in the background. (Subject: SM)



This image is copyright-free United States government work

Example credit: Andrej Karpathy



Outside border images, clockwise, starting from top left:
Image by Ivan Cultural_Globe is licensed under CC BY 2.0; changes made
Image by US Government is in the public domain
Image by US Government is in the public domain
Image by Glogster is licensed under CC BY-SA 3.0; changes made
Image by US Government is in the public domain
Image by US Government is in the public domain

Inside four images, clockwise, starting from top left:
Image is CC 1.0 public domain
Image by Tucan is licensed under CC BY 2.0; changes made
Image by Intuitive Surgical, Inc. is licensed under CC BY-SA 3.0; changes made
Image by Oyunardi Zorigtbaatar is licensed under CC BY-SA 4.0

Who we are

Instructors



Teaching Assistants



Fei-Fei Li & Justin Johnson & Serena Yeung

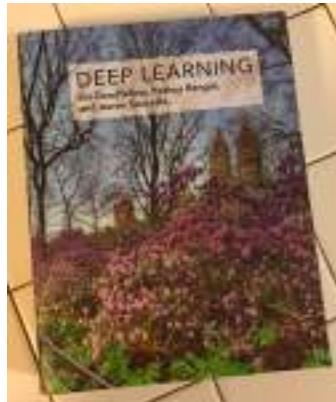
Lecture 1 - 43

4/3/2018

How to Contact Us

- Course Website: <http://cs231n.stanford.edu/>
 - Syllabus, lecture slides, links to assignment downloads, etc
- Piazza: <http://piazza.com/stanford/spring2018/cs231n>
 - Use this for most communication with course staff
 - Ask questions about homework, grading, logistics, etc
 - Use private questions if you want to post code
- Gradescope
 - For turning in homework and receiving grades
- Canvas
 - For watching lecture videos

Optional Textbook



- Deep Learning by Goodfellow, Bengio, and Courville
- Free online

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 45

4/3/2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 44

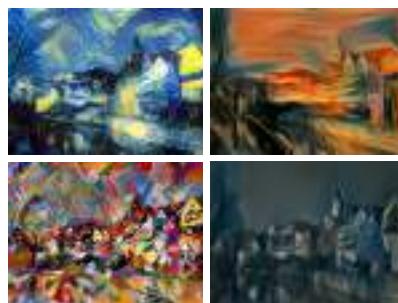
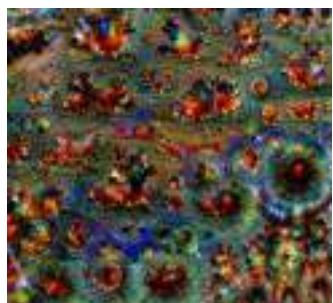
4/3/2018

Our philosophy

- Thorough and Detailed.
 - Understand how to write from scratch, debug and train convolutional neural networks.
- Practical.
 - Focus on practical techniques for training these networks at scale, and on GPUs (e.g. will touch on distributed optimization, differences between CPU vs. GPU, etc.) Also look at state of the art software tools such as TensorFlow, and PyTorch
- State of the art.
 - Most materials are new from research world in the past 1-3 years. Very exciting stuff!

Our philosophy (cont'd)

- Fun.
 - Some fun topics such as Image Captioning (using RNN)
 - Also DeepDream, NeuralStyle, etc.



Pre-requisite

- Proficiency in Python, some high-level familiarity with C/C++
 - All class assignments will be in Python (and use numpy), but some of the deep learning libraries we may look at later in the class are written in C++.
 - A Python tutorial available on course website
- College Calculus, Linear Algebra
- Equivalent knowledge of CS229 (Machine Learning)
 - We will be formulating cost functions, taking derivatives and performing optimization with gradient descent.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 47

4/3/2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 1 - 48

4/3/2018

Grading Policy

- 3 Problem Sets: **15% x 3 = 45%**
- Midterm Exam: **20%**
- Course Project: **35%**
 - Project Proposal: 1%
 - Milestone: 2%
 - Poster: 2%
 - Project Report: 30%
- Late policy
 - 4 free late days – use up to 2 late days per assignment
 - Afterwards, 25% off per day late
 - No late days for project report

Collaboration Policy

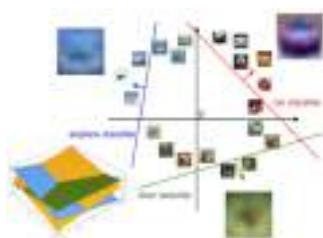
- We follow the [Stanford Honor Code](#) and the [CS Department Honor Code](#) – read them!
- **Rule 1:** Don't look at solutions or code that are not your own; everything you submit should be your own work
- **Rule 2:** Don't share your solution code with others; however discussing ideas or general strategies is fine and encouraged
- **Rule 3:** Indicate in your submissions anyone you worked with
- Turning in something late / incomplete is better than violating the honor code

Next Time: Image Classification

K-Nearest Neighbor



Linear Classifier



References

- Hubel, David H., and Torsten N. Wiesel. "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex." *The Journal of physiology* 160.1 (1962): 106. [\[PDF\]](#)
- Roberts, Lawrence Gilman. "Machine Perception of Three-dimensional Solids." Diss. Massachusetts Institute of Technology, 1963. [\[PDF\]](#)
- Marr, David. "Vision." The MIT Press, 1982. [\[PDF\]](#)
- Brooks, Rodney A., and Creiner, Russell and Binford, Thomas O. "The ACRONYM model-based vision system." In Proceedings of the 6th International Joint Conference on Artificial Intelligence (1979): 105-113. [\[PDF\]](#)
- Fischler, Martin A., and Robert A. Elschlager. "The representation and matching of pictorial structures." *IEEE Transactions on Computers* 22.1 (1973): 67-92. [\[PDF\]](#)
- Lowe, David G., "Three-dimensional object recognition from single two-dimensional images," *Artificial Intelligence*, 31, 3 (1987), pp. 355-395. [\[PDF\]](#)
- Shi, Jianbo, and Jitendra Malik. "Normalized cuts and image segmentation." *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 22.8 (2000): 888-905. [\[PDF\]](#)
- Viola, Paul, and Michael Jones. "Rapid object detection using a boosted cascade of simple features." *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*. Vol. 1. IEEE, 2001. [\[PDF\]](#)
- Lowe, David G. "Distinctive image features from scale-invariant keypoints." *International Journal of Computer Vision* 60.2 (2004): 91-110. [\[PDF\]](#)
- Lazebnik, Svetlana, Cordelia Schmid, and Jean Ponce. "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories." *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*. Vol. 2. IEEE, 2006. [\[PDF\]](#)

- Dalal, Navneet, and Bill Triggs. "Histograms of oriented gradients for human detection." *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. Vol. 1. IEEE, 2005. [\[PDF\]](#)
- Felzenszwalb, Pedro, David McAllester, and Deva Ramanan. "A discriminatively trained, multiscale, deformable part model." *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*. IEEE, 2008 [\[PDF\]](#)
- Everingham, Mark, et al. "The pascal visual object classes (VOC) challenge." *International Journal of Computer Vision* 88.2 (2010): 303-338. [\[PDF\]](#)
- Deng, Jia, et al. "Imagenet: A large-scale hierarchical image database." *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, 2009. [\[PDF\]](#)
- Russakovsky, Olga, et al. "Imagenet Large Scale Visual Recognition Challenge." *arXiv:1409.0575*. [\[PDF\]](#)
- Lin, Yuanqing, et al. "Large-scale image classification: fast feature extraction and SVM training." *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011. [\[PDF\]](#)
- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems*. 2012. [\[PDF\]](#)
- Szegedy, Christian, et al. "Going deeper with convolutions." *arXiv preprint arXiv:1409.4842* (2014). [\[PDF\]](#)
- Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition." *arXiv preprint arXiv:1409.1556* (2014). [\[PDF\]](#)
- He, Kaiming, et al. "Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition." *arXiv preprint arXiv:1406.4729* (2014). [\[PDF\]](#)
- LeCun, Yann, et al. "Gradient-based learning applied to document recognition." *Proceedings of the IEEE* 86.11 (1998): 2278-2324. [\[PDF\]](#)
- Fei-Fei, Li, et al. "What do we perceive in a glance of a real-world scene?." *Journal of vision* 7.1 (2007): 10. [\[PDF\]](#)

Lecture 2: Image Classification pipeline

Administrative: Piazza

For questions about midterm, poster session, projects, etc, use Piazza!

SCPD students: Use your @stanford.edu address to register for Piazza; contact scpd-customerservice@stanford.edu for help.

Administrative: Assignment 1

Out yesterday, due 4/18 11:59pm

- K-Nearest Neighbor
- Linear classifiers: SVM, Softmax
- Two-layer neural network
- Image features

Administrative: Friday Discussion Sections

Fridays 12:30pm - 1:20pm in Skilling Auditorium

Hands-on tutorials, with more practical detail than main lecture

Check course website for schedule:

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

This Friday: Python / numpy / Google Cloud setup

Administrative: Python + Numpy

CS231n Convolutional Neural Networks for Visual Recognition

Python Numpy Tutorial

This tutorial was contributed by Justin Johnson.

We will use the Python programming language for all assignments in this course. Python is a general-purpose programming language on its own but with the help of a few popular libraries (in this case, numpy) it becomes a powerful environment for scientific computing.

We expect that many of you will have some experience with Python and numpy. For the rest of you, this section will serve as a quick crash course both on the Python programming language and on the use of Python for scientific computing.

<http://cs231n.github.io/python-numpy-tutorial/>

Administrative: Google Cloud

We will be using Google Cloud in this class

We will be distributing coupons to all enrolled students

See our tutorial here for walking through Google Cloud setup:
<http://cs231n.github.io/gce-tutorial/>

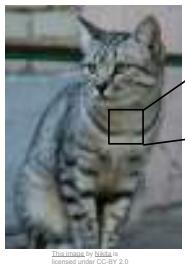
Image Classification: A core task in Computer Vision



(assume given set of discrete labels)
{dog, cat, truck, plane, ...}

→ cat

The Problem: Semantic Gap



An image is just a big grid of numbers between [0, 255]:

e.g. 800 x 600 x 3
(3 channels RGB)

pixels

Challenges: Viewpoint variation



All pixels change when the camera moves!

Challenges: Illumination



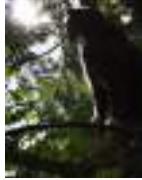
This image is CC0 1.0 public domain



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain

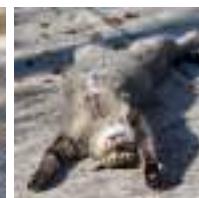


This image is CC0 1.0 public domain

Challenges: Deformation



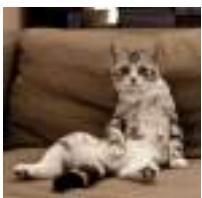
This image by Umberto Salvagnin is licensed under [CC-BY-2.0](#)



This image by Umberto Salvagnin is licensed under [CC-BY-2.0](#)



This image by [vanhoek](#) is licensed under [CC0 1.0](#)



This image by [Zon](#) is licensed under [CC0 1.0](#)

Challenges: Occlusion



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain



This image by [jonsso](#) is licensed under [CC-BY-2.0](#)

Challenges: Background Clutter



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain

Challenges: Intraclass variation



This image is CC0 1.0 public domain

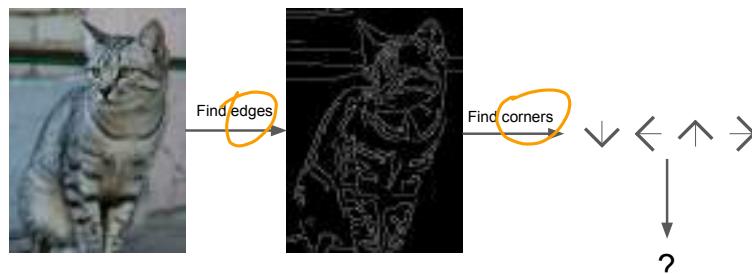
An image classifier

```
def classify_image(image):
    # Some magic here!
    return class_label
```

Unlike e.g. sorting a list of numbers,

no obvious way to hard-code the algorithm for
recognizing a cat, or other classes.

Attempts have been made



John Canny, "A Computational Approach to Edge Detection", IEEE TPAMI 1986

Machine Learning: Data-Driven Approach

1. Collect a **dataset** of images and labels
2. Use Machine Learning to **train a classifier**
3. Evaluate the classifier on new images

```
def train(images, labels):
    # Machine learning
    return model

def predict(model, test_images):
    # Use model to predict labels
    return test_labels
```

Example training set



First classifier: Nearest Neighbor

```
def train(images, labels):
    # Machine learning
    return model

def predict(model, test_images):
    # Use model to predict labels
    return test_labels
```

- Memorize all data and labels
- Predict the label of the most similar training image

Example Dataset: CIFAR10

10 classes
50,000 training images
10,000 testing images



Alex Krizhevsky, "Learning Multiple Layers of Features from Tiny Images", Technical Report, 2009.

Example Dataset: CIFAR10

10 classes
50,000 training images
10,000 testing images



Alex Krizhevsky, "Learning Multiple Layers of Features from Tiny Images", Technical Report, 2009.

Test images and nearest neighbors



Distance Metric to compare images

L1 distance:

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

Mahalon

test image				training image				pixel-wise absolute-value differences			
58	32	10	18	10	20	24	17	46	12	14	1
90	23	128	133	8	10	89	100	82	13	39	33
24	26	178	200	12	16	178	170	12	10	0	30
2	0	255	220	4	32	233	112	2	32	22	108

= add

```

import numpy as np
class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where rows are examples. Y is L-dimensional vector of size N """
        self.X = X
        self.y = y

    def predict(self, X):
        """ X is B x D where rows are examples. To predict class, for each row in X, compare to all training rows and return index of the closest example """
        if len(X.shape) != 2:
            raise ValueError('X must be a [N x D] matrix')
        else:
            X = X.reshape(-1, X.shape[1])
        distances = np.sqrt(np.sum((X - self.X)**2, axis=1))
        min_index = np.argmin(distances)
        return self.y[min_index]

```

Nearest Neighbor classifier

```

import numpy as np
class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where rows are examples. Y is L-dimensional vector of size N """
        self.X = X
        self.y = y

    def predict(self, X):
        """ X is B x D where rows are examples. To predict class, for each row in X, compare to all training rows and return index of the closest example """
        if len(X.shape) != 2:
            raise ValueError('X must be a [N x D] matrix')
        else:
            X = X.reshape(-1, X.shape[1])
        distances = np.sqrt(np.sum((X - self.X)**2, axis=1))
        min_index = np.argmin(distances)
        return self.y[min_index]

```

Nearest Neighbor classifier

Memorize training data

```

import numpy as np
class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where rows are examples. Y is L-dimensional vector of size N """
        self.X = X
        self.y = y

    def predict(self, X):
        """ X is B x D where rows are examples. To predict class, for each row in X, compare to all training rows and return index of the closest example """
        if len(X.shape) != 2:
            raise ValueError('X must be a [N x D] matrix')
        else:
            X = X.reshape(-1, X.shape[1])
        distances = np.sqrt(np.sum((X - self.X)**2, axis=1))
        min_index = np.argmin(distances)
        return self.y[min_index]

```

Nearest Neighbor classifier

```

import numpy as np
class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where rows are examples. Y is L-dimensional vector of size N """
        self.X = X
        self.y = y

    def predict(self, X):
        """ X is B x D where rows are examples. To predict class, for each row in X, compare to all training rows and return index of the closest example """
        if len(X.shape) != 2:
            raise ValueError('X must be a [N x D] matrix')
        else:
            X = X.reshape(-1, X.shape[1])
        distances = np.sqrt(np.sum((X - self.X)**2, axis=1))
        min_index = np.argmin(distances)
        return self.y[min_index]

```

Nearest Neighbor classifier

Q: With N examples,
how fast are training
and prediction?

For each test image:
Find closest train image
Predict label of nearest image

```

In [1]: In [1]: In [1]: In [1]:
class NearestNeighbor:
    def __init__(self, X, y):
        self.X = X
        self.y = y

    def classify(self, X, k=1):
        """ X is N x D where rows are examples; y is 1-dimensional of size N;
            X[i] is a D-dimensional point. Compute distances between test points and
            nearest neighbor training points. Return a list of labels for each test point """
        dists = self.compute_distances_no_loops(X)
        return self.predict_labels(dists, k=k)

    def compute_distances_no_loops(self, X):
        """ X is N x D where rows are examples; y is 1-dimensional of size N;
            X[i] is a D-dimensional point. Compute distances between test points and
            nearest neighbor training points. Return a list of labels for each test point """
        dists = np.sqrt(np.sum(self.X**2, axis=1)[:, np.newaxis] + np.sum(X**2, axis=1) - 2 * np.dot(self.X, X.T))
        return dists

```

Nearest Neighbor classifier

Q: With N examples, how fast are training and prediction?

A: Train O(1), predict O(N)

```

In [1]: In [1]: In [1]: In [1]:
class NearestNeighbor:
    def __init__(self, X, y):
        self.X = X
        self.y = y

    def classify(self, X, k=1):
        """ X is N x D where rows are examples; y is 1-dimensional of size N;
            X[i] is a D-dimensional point. Compute distances between test points and
            nearest neighbor training points. Return a list of labels for each test point """
        dists = self.compute_distances_no_loops(X)
        return self.predict_labels(dists, k=k)

    def compute_distances_no_loops(self, X):
        """ X is N x D where rows are examples; y is 1-dimensional of size N;
            X[i] is a D-dimensional point. Compute distances between test points and
            nearest neighbor training points. Return a list of labels for each test point """
        dists = np.sqrt(np.sum(self.X**2, axis=1)[:, np.newaxis] + np.sum(X**2, axis=1) - 2 * np.dot(self.X, X.T))
        return dists

```

Nearest Neighbor classifier

Q: With N examples, how fast are training and prediction?

A: Train O(1), predict O(N)

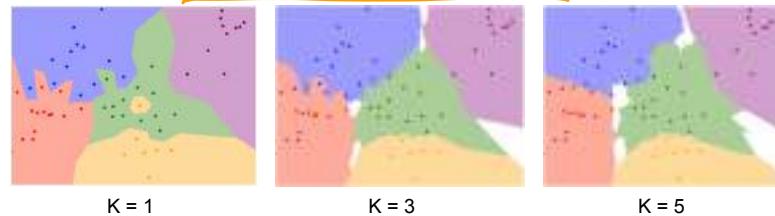
This is bad: we want classifiers that are fast at prediction; slow for training is ok

What does this look like?



K-Nearest Neighbors

Instead of copying label from nearest neighbor, take **majority vote** from K closest points



K = 1

K = 3

K = 5

What does this look like?



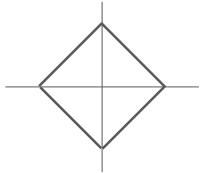
What does this look like?



K-Nearest Neighbors: Distance Metric

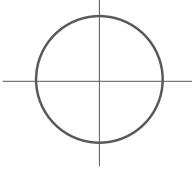
L1 (Manhattan) distance

$$d_1(f_1, f_2) = \sum_p |f_1^p - f_2^p|$$



L2 (Euclidean) distance

$$d_2(f_1, f_2) = \sqrt{\sum_p (f_1^p - f_2^p)^2}$$



K-Nearest Neighbors: Distance Metric

L1 (Manhattan) distance

$$d_1(f_1, f_2) = \sum_p |f_1^p - f_2^p|$$



K = 1

L2 (Euclidean) distance

$$d_2(f_1, f_2) = \sqrt{\sum_p (f_1^p - f_2^p)^2}$$



K = 1

K-Nearest Neighbors: Demo Time



<http://vision.stanford.edu/teaching/cs231n-demos/knn/>

Hyperparameters

What is the best value of **k** to use?
What is the best **distance** to use?

These are **hyperparameters**: choices about
the algorithm that we set rather than learn

Hyperparameters

What is the best value of **k** to use?
What is the best **distance** to use?

These are **hyperparameters**: choices about
the algorithm that we set rather than learn

Very problem-dependent.
Must try them all out and see what works best.

Setting Hyperparameters

Idea #1: Choose hyperparameters
that work best on the data

Your Dataset

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: K = 1 always works perfectly on training data

Your Dataset

*↑ focus
unseen data*

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: K = 1 always works perfectly on training data

Your Dataset

Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on test data

train

test

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: K = 1 always works perfectly on training data

Your Dataset

Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on test data

BAD: No idea how algorithm will perform on new data

train

test

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: K = 1 always works perfectly on training data

Your Dataset

Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on test data

BAD: No idea how algorithm will perform on new data

train

test

Idea #3: Split data into **train**, **val**, and **test**; choose hyperparameters on val and evaluate on test

Better!

train

validation

test

Setting Hyperparameters

Your Dataset

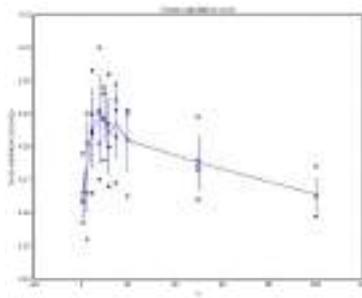
Idea #4: Cross-Validation: Split data into **folds**, try each fold as validation and average the results

fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test

Useful for small datasets, but not used too frequently in deep learning

Setting Hyperparameters

Example of 5-fold cross-validation for the value of k.



Each point: single outcome.

The line goes through the mean, bars indicated standard deviation

(Seems that k ≈ 7 works best for this data)

k-Nearest Neighbor on images never used.

- Very slow at test time
- Distance metrics on pixels are not informative

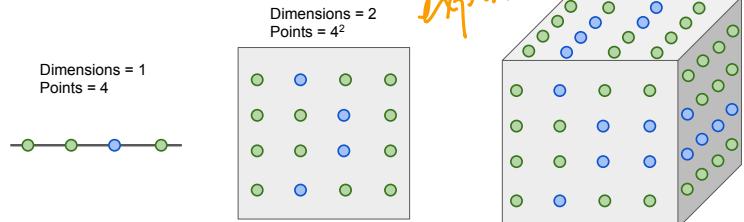


(all 3 images have same L2 distance to the one on the left)

Original image is CC0 public domain

k-Nearest Neighbor on images never used.

- Curse of dimensionality



K-Nearest Neighbors: Summary

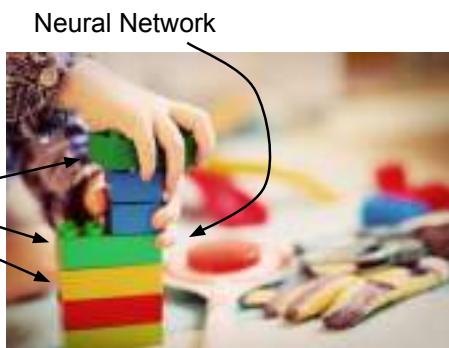
In **Image classification** we start with a **training set** of images and labels, and must predict labels on the **test set**

The **K-Nearest Neighbors** classifier predicts labels based on nearest training examples

Distance metric and K are **hyperparameters**

Choose hyperparameters using the **validation set**; only run on the test set once at the very end!

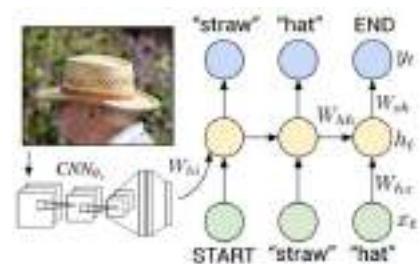
Linear Classification



Two young girls are playing with lego toy.
Boy is doing backflip on wakeboard



Man in black shirt is playing guitar.
Construction worker in orange safety vest is working on road.



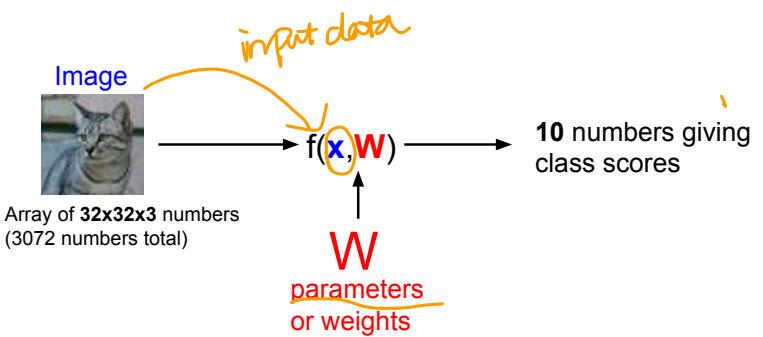
Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015
Figures copyright IEEE, 2015. Reproduced for educational purposes.

Recall CIFAR10

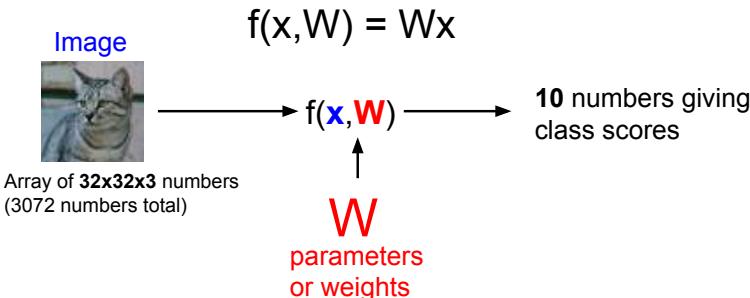


50,000 training images
each image is 32x32x3
10,000 test images.

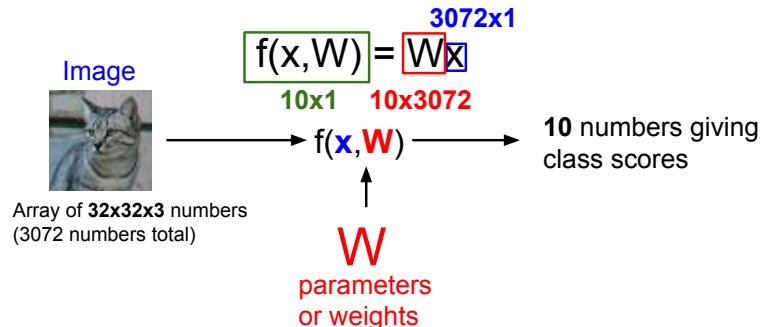
Parametric Approach



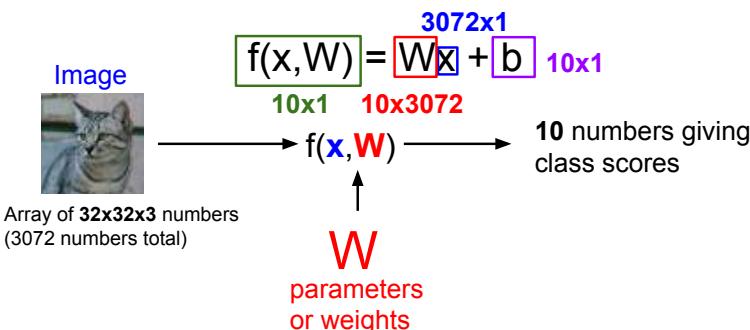
Parametric Approach: Linear Classifier



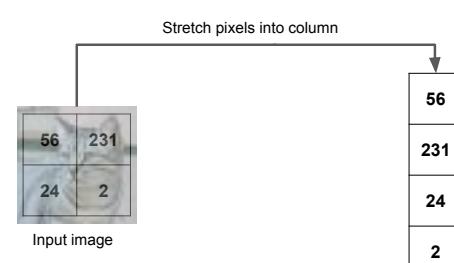
Parametric Approach: Linear Classifier



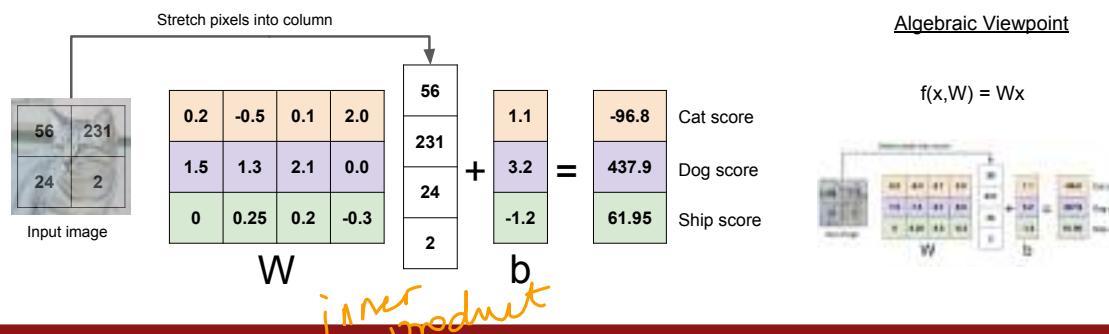
Parametric Approach: Linear Classifier



Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Example with an image with 4 pixels, and 3 classes (cat/dog/ship) Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



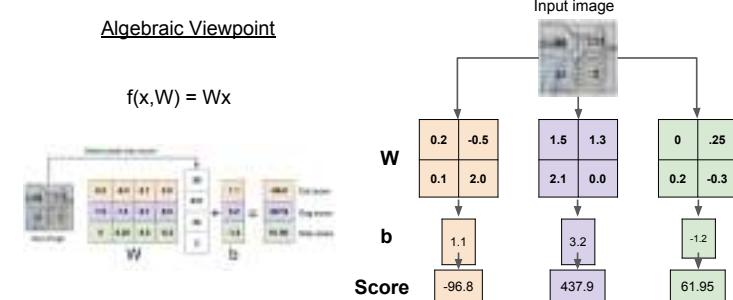
Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 2 - 56

April 5, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 57

April 5, 2018

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Fei-Fei Li & Justin Johnson & Serena Yeung

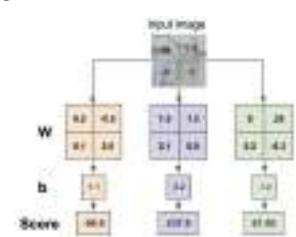
Lecture 2 - 58

April 5, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 59

April 5, 2018

Interpreting a Linear Classifier



Interpreting a Linear Classifier: Visual Viewpoint



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 60

April 5, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 61

April 5, 2018

Interpreting a Linear Classifier: Geometric Viewpoint

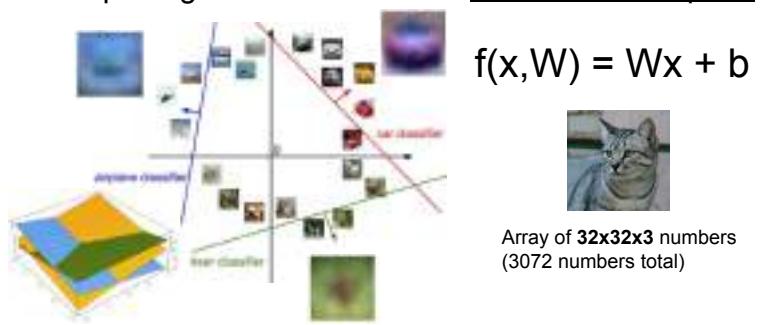
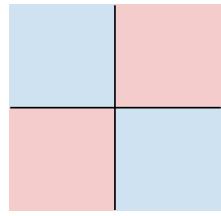


Image by Nata is licensed under CC BY 2.0

Hard cases for a linear classifier

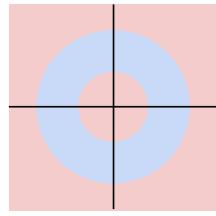
Class 1:
First and third quadrants

Class 2:
Second and fourth quadrants



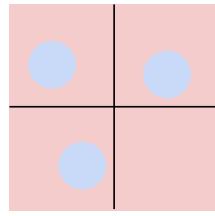
Class 1:
 $1 \leq L_2 \text{ norm} \leq 2$

Class 2:
Everything else



Class 1:
Three modes

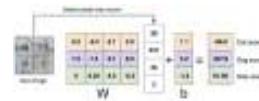
Class 2:
Everything else



Linear Classifier: Three Viewpoints

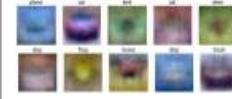
Algebraic Viewpoint

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



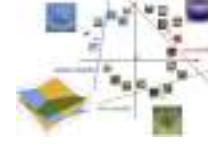
Visual Viewpoint

One template per class



Geometric Viewpoint

Hyperplanes cutting up space



So far: Defined a (linear) score function $f(x, W) = Wx + b$

Example class scores for 3 images for some W :



How can we tell whether this W is good or bad?

	-3.45	-0.51	3.42
cat	-8.87	6.04	4.64
car	0.09	5.31	2.65
frog	2.9	-4.22	5.1
bus	4.48	-4.19	2.64
dog	8.02	3.58	5.55
horse	3.78	4.49	-4.34
train	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
bike	-0.72	-2.93	6.14

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

Coming up:

- Loss function
- Optimization
- ConvNets!

(quantifying what it means to have a “good” W)

(start with random W and find a W that minimizes the loss)

(tweak the functional form of f)

Administrative: Live Questions

We'll use Zoom to take questions from remote students live-streaming the lecture

Check Piazza for instructions and meeting ID:

<https://piazza.com/class/jdmurnqexkt47x?cid=108>

Lecture 3: Loss Functions and Optimization

Administrative: Office Hours

Office hours started this week, schedule is on the course website:

http://cs231n.stanford.edu/office_hours.html

Areas of expertise for all TAs are posted on Piazza:

<https://piazza.com/class/jdmurnqexkt47x?cid=155>

Administrative: Assignment 1

Assignment 1 is released:

<http://cs231n.github.io/assignments2018/assignment1/>

Due **Wednesday April 18, 11:59pm**

Administrative: Google Cloud

You should have received an email yesterday about claiming a coupon for Google Cloud; make a private post on Piazza if you didn't get it

There was a problem with @cs.stanford.edu emails; this is resolved

If you have problems with coupons: **Post on Piazza**

DO NOT email me, DO NOT email Prof. Phil Levis

Administrative: SCPD Tutors

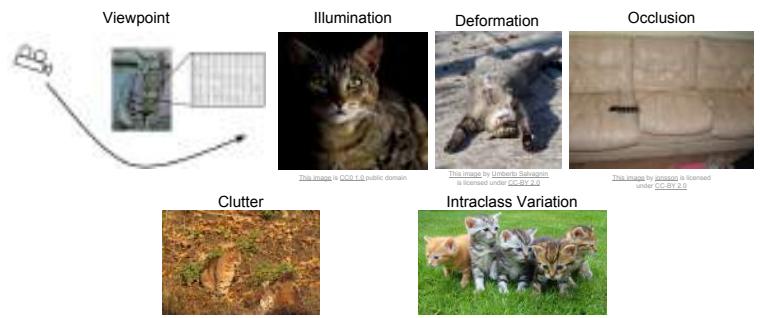
This year the SCPD office has hired tutors specifically for SCPD students taking CS231N; you should have received an email about this yesterday (4/9/2018)

Administrative: Poster Session

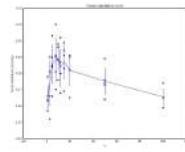
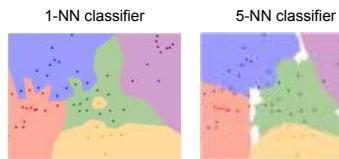
Poster session will be Tuesday June 12 (our final exam slot)

Attendance is mandatory for non-SCPD students; if you don't have a legitimate reason for skipping it then you forfeit the points for the poster presentation

Recall from last time: Challenges of recognition

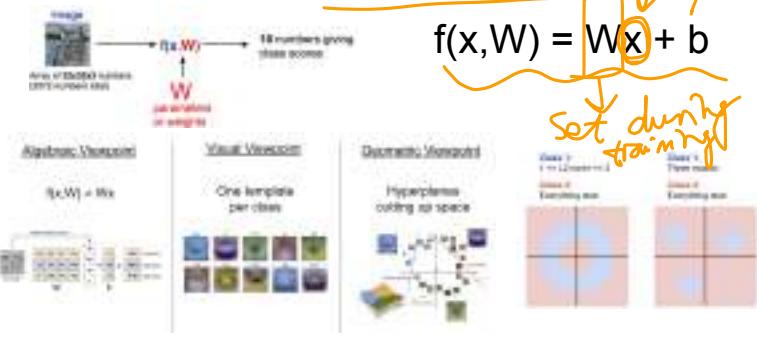


Recall from last time: data-driven approach, kNN



train	test
train	validation

Recall from last time: Linear Classifier



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

Input

Set during training

Recall from last time: Linear Classifier

TODO:



	-3.45	-0.51	3.42
automobile	-8.87	8.54	8.64
bird	0.29	3.32	2.63
cat	2.9	-4.22	5.1
deer	3.88	-4.19	2.84
dog	6.02	3.58	3.55
frog	3.78	6.49	-6.98
horse	1.04	-4.17	-1.2
ship	-8.36	-2.09	-4.79
truck	-8.12	-3.93	8.18

- Define a **loss function** that quantifies our unhappiness with the scores across the training data.
- Come up with a way of efficiently finding the parameters that minimize the loss function. (**optimization**)

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



A **loss function** tells how good our current classifier is

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where x_i is image and y_i is (integer) label

Loss over the dataset is a sum of loss over examples:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



Multiclass SVM loss:

Given an example (x_i, y_i) where x_i is the image and where y_i is the (integer) label

and using the shorthand for the scores vector: $s_j = f(x_i, W)_j$

score of correct category

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

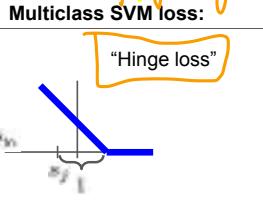
score of first

sum over incorrect integer

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



Solving many "Hinge loss"



$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,

and using the shorthand for the
scores vector: $\mathbf{s} = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,

and using the shorthand for the
scores vector: $\mathbf{s} = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$= \max(0, 5.1 - 3.2 + 1)$$

$$+ \max(0, -1.7 - 3.2 + 1)$$

$$= \max(0, 2.9) + \max(0, -3.9)$$

$$= 2.9 + 0$$

$$= 2.9$$

cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9		

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,

and using the shorthand for the
scores vector: $\mathbf{s} = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$= \max(0, 1.3 - 4.9 + 1)$$

$$+ \max(0, 2.0 - 4.9 + 1)$$

$$= \max(0, -2.6) + \max(0, -1.9)$$

$$= 0 + 0$$

$$= 0$$

cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9		0

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,

and using the shorthand for the
scores vector: $\mathbf{s} = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$= \max(0, 2.2 - (-3.1) + 1)$$

$$+ \max(0, 2.5 - (-3.1) + 1)$$

$$= \max(0, 6.3) + \max(0, 6.6)$$

$$= 6.3 + 6.6$$

$$= 12.9$$

cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,

and using the shorthand for the
scores vector: $\mathbf{s} = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Loss over full dataset is average:

$$L = \frac{1}{N} \sum_{i=1}^N L_i$$

$$L = (2.9 + 0 + 12.9)/3$$

$$= 5.27$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,
and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q: What happens to
loss if car scores
change a bit?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,
and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q2: what is the
min/max possible
loss?



Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,
and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q3: At initialization W
is small so all $s \approx 0$.
What is the loss?
Classes - 1

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,
and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q4: What if the sum
was over all classes?
(including $j = y_i$)

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,
and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q5: What if we used
mean instead of
sum?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Multiclass SVM loss:

Given an example (x_i, y_i)
where x_i is the image and
where y_i is the (integer) label,
and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q6: What if we used
mean instead of
sum?

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Multiclass SVM Loss: Example code

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

```
def L_i_vectorized(x, y, W):
    scores = W.dot(x)
    margins = np.maximum(0, scores - scores[y] + 1)
    margins[y] = 0
    loss_i = np.sum(margins)
    return loss_i
```

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

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

E.g. Suppose that we found a W such that $L = 0$. Is this W unique?

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

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

E.g. Suppose that we found a W such that $L = 0$. Is this W unique?

No! $2W$ is also has $L = 0!$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9		
	0		

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Before:

$$\begin{aligned} &= \max(0, 1.3 - 4.9 + 1) \\ &\quad + \max(0, 2.0 - 4.9 + 1) \\ &= \max(0, -2.6) + \max(0, -1.9) \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

With W twice as large:

$$\begin{aligned} &= \max(0, 2.6 - 9.8 + 1) \\ &\quad + \max(0, 4.0 - 9.8 + 1) \\ &= \max(0, -6.2) + \max(0, -4.8) \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

Regularization

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

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

E.g. Suppose that we found a W such that $L = 0$. Is this W unique?

No! $2W$ is also has $L = 0!$

How do we choose between W and $2W$?

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i; W), y_i)}_{\text{Data loss}}$$

Data loss: Model predictions should match training data

Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions should match training data}} + \lambda R(W)$$

Regularization: Prevent the model from doing *too well* on training data

Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions should match training data}} + \lambda R(W)$$

Regularization: Prevent the model from doing *too well* on training data

λ = regularization strength (hyperparameter)

Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions should match training data}} + \lambda R(W)$$

Regularization: Prevent the model from doing *too well* on training data

λ = regularization strength (hyperparameter)

Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions should match training data}} + \lambda R(W)$$

Regularization: Prevent the model from doing *too well* on training data

λ = regularization strength (hyperparameter)

Simple examples

L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

Simple examples

L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

More complex:

Dropout

Batch normalization

Stochastic depth, fractional pooling, etc

Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss: Model predictions should match training data}} + \lambda R(W)$$

Regularization: Prevent the model from doing *too well* on training data

λ = regularization strength (hyperparameter)

Regularization: Expressing Preferences

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

L2 Regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

Why regularize?

- Express preferences over weights
- Make the model *simple* so it works on test data
- Improve optimization by adding curvature

$$w_1^T x = w_2^T x = 1$$

Regularization: Expressing Preferences

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

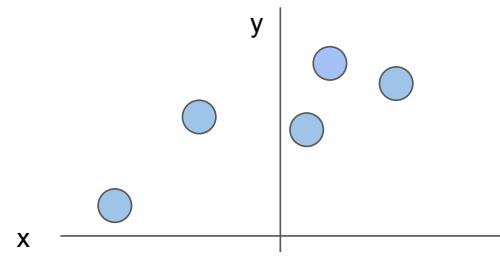
L2 Regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

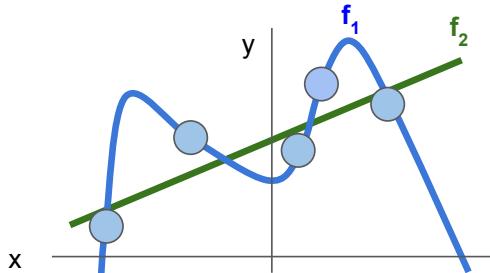
L2 regularization likes to "spread out" the weights

$$w_1^T x = w_2^T x = 1$$

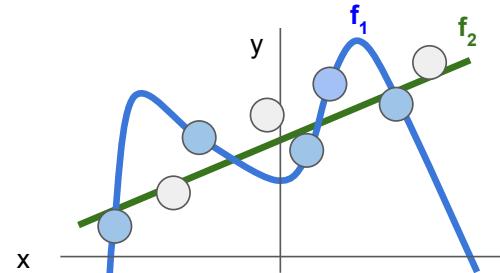
Regularization: Prefer Simpler Models



Regularization: Prefer Simpler Models



Regularization: Prefer Simpler Models



Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as probabilities

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as probabilities

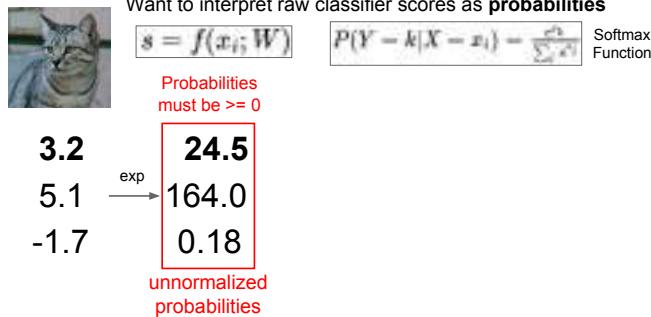
$$s = f(x_i; W) \quad P(Y=k|X=x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

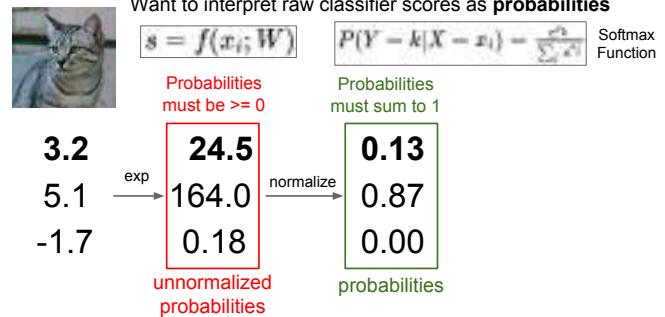
cat	3.2
car	5.1
frog	-1.7

cat	3.2
car	5.1
frog	-1.7

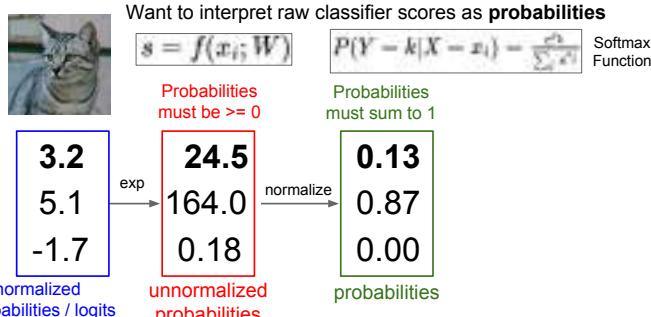
Softmax Classifier (Multinomial Logistic Regression)



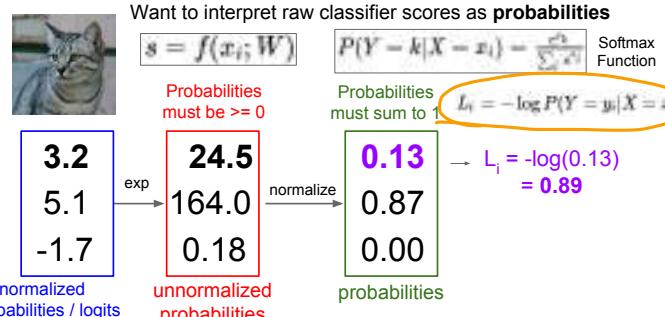
Softmax Classifier (Multinomial Logistic Regression)



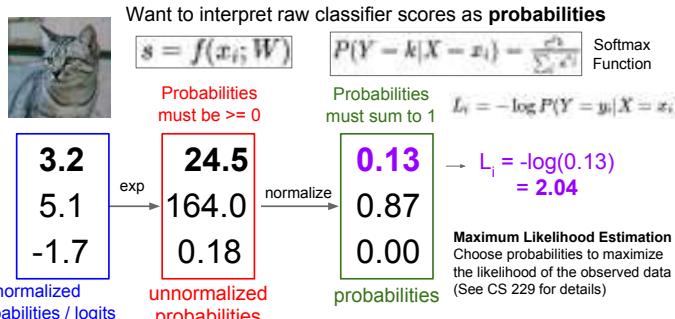
Softmax Classifier (Multinomial Logistic Regression)



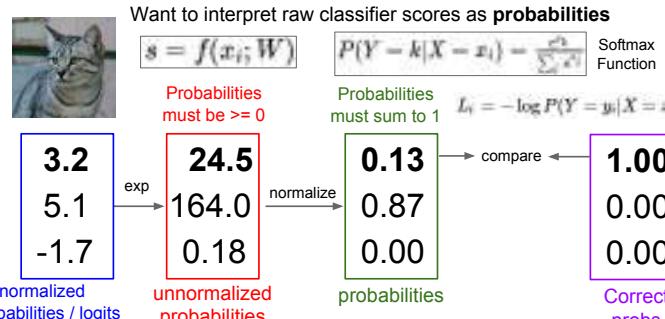
Softmax Classifier (Multinomial Logistic Regression)



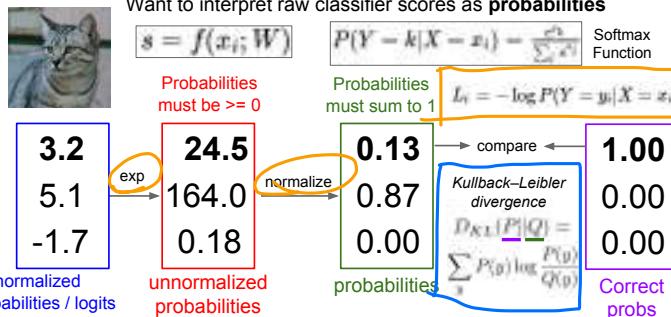
Softmax Classifier (Multinomial Logistic Regression)



Softmax Classifier (Multinomial Logistic Regression)



Softmax Classifier (Multinomial Logistic Regression)

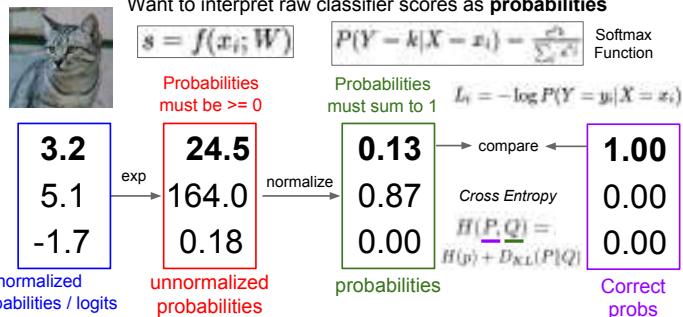


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 51

April 10, 2018

Softmax Classifier (Multinomial Logistic Regression)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 52

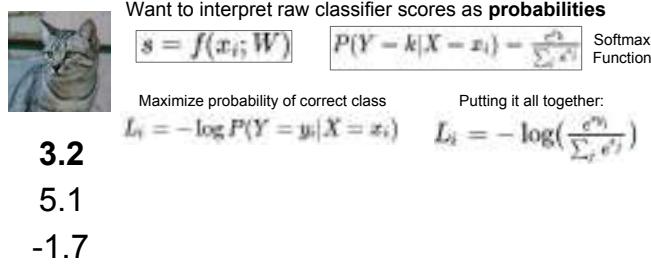
April 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 51

April 10, 2018

Softmax Classifier (Multinomial Logistic Regression)

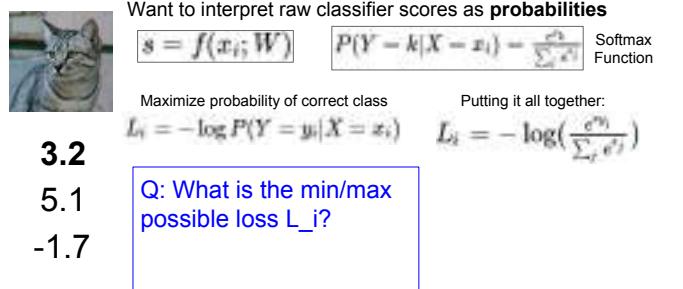


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 53

April 10, 2018

Softmax Classifier (Multinomial Logistic Regression)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 54

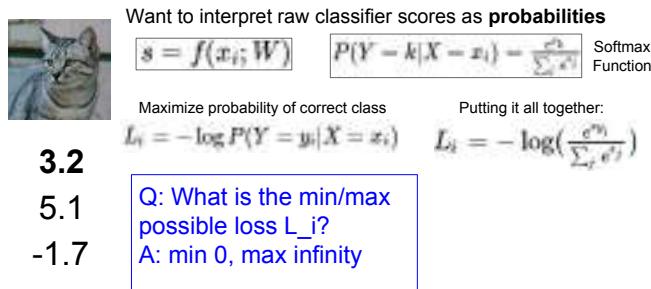
April 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 53

April 10, 2018

Softmax Classifier (Multinomial Logistic Regression)

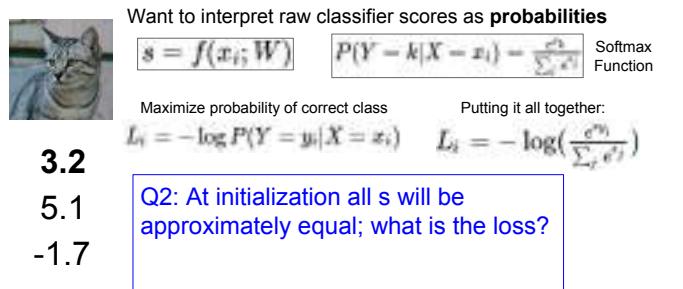


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 55

April 10, 2018

Softmax Classifier (Multinomial Logistic Regression)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 56

April 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 55

April 10, 2018

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y=k|X=x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

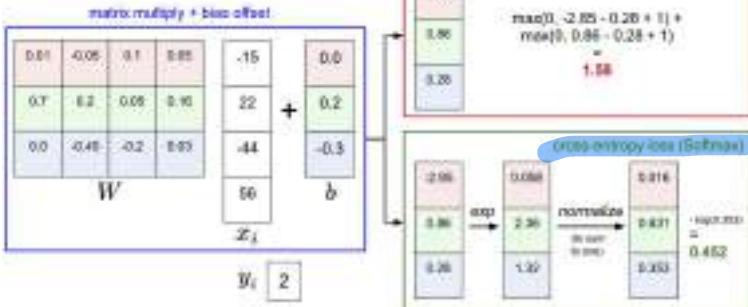
Maximize probability of correct class

$$L_i = -\log P(Y=y_i|X=x_i) \quad L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

cat	3.2
car	5.1
frog	-1.7

Q2: At initialization all s will be approximately equal; what is the loss?
A: $\log(C)$, eg $\log(10) \approx 2.3$

Softmax vs. SVM



Softmax vs. SVM

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Softmax vs. SVM

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

assume scores:

$$[10, -2, 3]$$

$$[10, 9, 9]$$

$$[10, -100, -100]$$

and $y_i = 0$

Q: Suppose I take a datapoint and I jiggle a bit (changing its score slightly). What happens to the loss in both cases?

Recap

- We have some dataset of (x, y)
- We have a **score function**: $s = f(x; W) = Wx$ e.g.
- We have a **loss function**:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

SVM

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + R(W)$$

Full loss

Recap

- We have some dataset of (x, y)
- We have a **score function**: $s = f(x; W) = Wx$ e.g.
- We have a **loss function**:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

SVM

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + R(W)$$

Full loss

How do we find the best W ?

optimisation

Optimization



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 63

April 10, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 64

April 10, 2018



Strategy #1: A first very bad idea solution: **Random search**

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 65

April 10, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 66

April 10, 2018

Lets see how well this works on the test set...

Strategy #2: Follow the slope

```

# Assume X_test is (3072 x 30000), Y_true (30000 x 1)
scores = Xtest.dot(Xte_cols) # n x 10000, the class scores for all test samples
# Find the index with max score in each column (the predictor class)
Xte_predict = np.argmax(scores, axis=0)
# and calculate accuracy (fraction of predictions that are correct)
sp.mean(Xte_predict == Yte)
# 0.770458563886

```

15.5% accuracy! not bad!
(SOTA is ~95%)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 67

April 10, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 68

April 10, 2018

Strategy #2: Follow the slope

In 1-dimension, the derivative of a function:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

current W:

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,...]
loss 1.25347

gradient dW:

[?,
?,
?,
?,
?,
?,
?,
?,
?,...]

In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension

The slope in any direction is the dot product of the direction with the gradient
The direction of steepest descent is the **negative gradient**

loss 1.25347

current W:	W + h (first dim):	gradient dW:	current W:	W + h (first dim):	gradient dW:
[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[0.34 + 0.0001 , -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25322	[?, ?, ?, ?, ?, ?, ?, ?, ?,...]	[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[0.34 + 0.0001 , -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25322	[-2.5, ?, ?, (1.25322 - 1.25347)/0.0001 = -2.5 [?, ?, ?,...]

current W:	W + h (second dim):	gradient dW:	current W:	W + h (second dim):	gradient dW:
[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[0.34, -1.11 + 0.0001 , 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25353	[-2.5, ?, ?, ?, ?, ?, ?, ?, ?,...]	[0.34, -1.11 + 0.0001 , 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[0.34, -1.11 + 0.0001 , 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25353	[-2.5, 0.6, ?, ?, (1.25353 - 1.25347)/0.0001 = 0.6 [?, ?, ?,...]

current W:	W + h (third dim):	gradient dW:	current W:	W + h (third dim):	gradient dW:
[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[0.34, -1.11, 0.78 + 0.0001 , 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[-2.5, 0.6, ?, ?, ?, ?, ?, ?, ?,...]	[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[0.34, -1.11, 0.78 + 0.0001 , 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[-2.5, 0.6, 0 , ?, ?, $\frac{(1.25347 - 1.25347)/0.0001}{0.0001} = 0$, ?,...]

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 3 - 75 April 10, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 3 - 76 April 10, 2018

current W:	W + h (third dim):	gradient dW:	
[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[0.34, -1.11, 0.78 + 0.0001 , 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[-2.5, 0.6, 0 , ?, ?, ?, ?, ?,...]	This is silly. The loss is just a function of W: $L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2$ $L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$ $s = f(x; W) = Wx$ want $\nabla_W L$

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 3 - 77 April 10, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 3 - 78 April 10, 2018

This is silly. The loss is just a function of W:

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$s = f(x; W) = Wx$$

want $\nabla_W L$



Use calculus to compute an
analytic gradient

current W:	gradient dW:
[0.34, -1.11, 0.78, 0.12, 0.55, 2.81, -3.1, -1.5, 0.33,...] loss 1.25347	[-2.5, 0.6, 0, 0.2, 0.7, -0.5, 1.1, 1.3, -2.1,...]

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 3 - 79 April 10, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 3 - 80 April 10, 2018

Gradient Descent

In summary:

- Numerical gradient: approximate, slow, easy to write
- Analytic gradient: exact, fast, error-prone

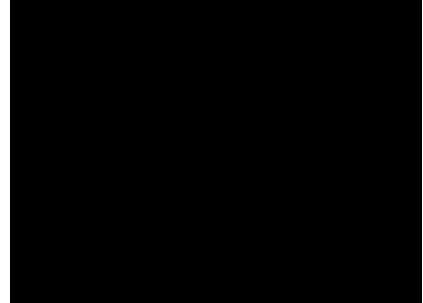
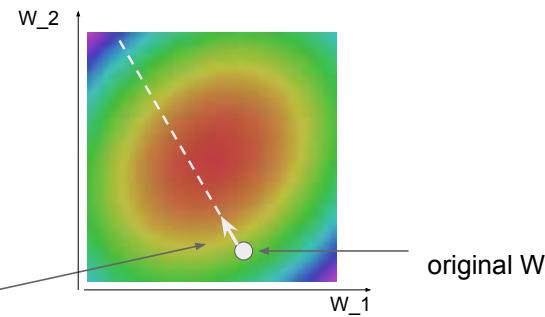
=>

In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a gradient check.

```
# Vanilla Gradient Descent
while True:
    weights_grad = evaluate_gradient(loss_fn, data, weights)
    weights += -step_size * weights_grad # perform parameter update
```

update weights hyper-parameter
 (learning rate)

negative gradient direction



Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

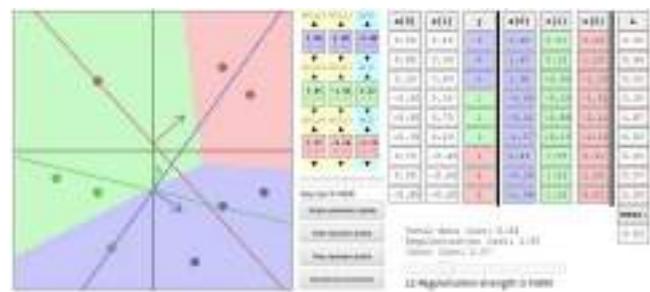
Full sum expensive when N is large!

Approximate sum using a minibatch of examples (sample)
32 / 64 / 128 common

Vanilla Stochastic Gradient Descent:

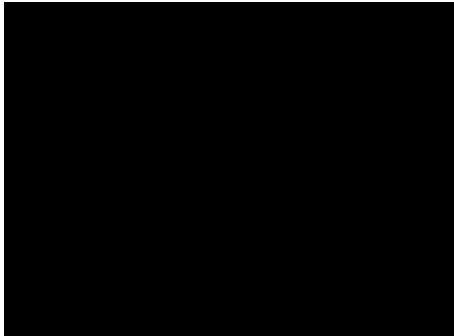
```
while True:
    data_batch = sample_training_data(data, 32) # sample 32 examples
    weights_grad = evaluate_gradient(loss_fn, data_batch, weights)
    weights += -step_size * weights_grad # perform parameter update
```

Interactive Web Demo time....



<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>

Interactive Web Demo time....



Aside: Image Features



$$f(x) = Wx$$

Class scores



Aside: Image Features

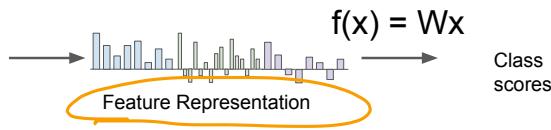
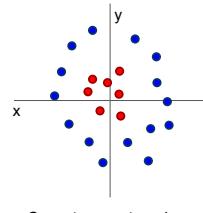
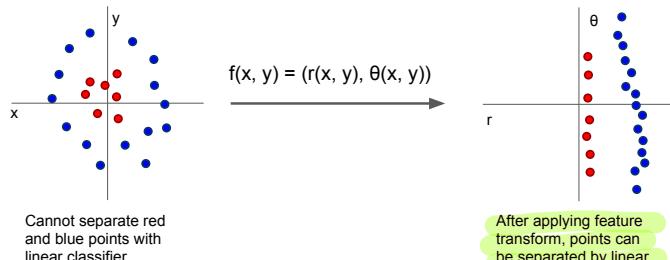


Image Features: Motivation



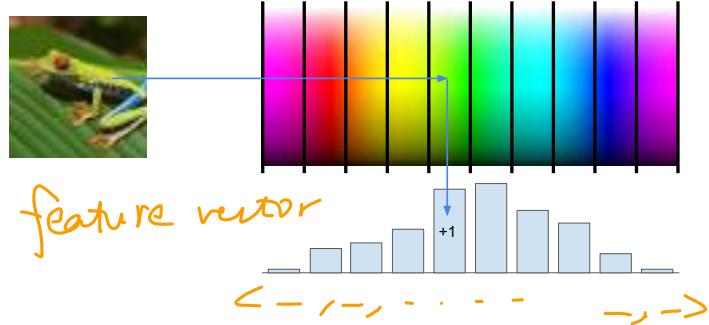
Cannot separate red
and blue points with
linear classifier

Image Features: Motivation

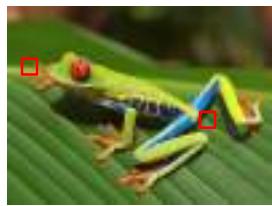


Cannot separate red
and blue points with
linear classifier

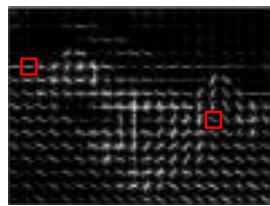
Example: Color Histogram



Example: Histogram of Oriented Gradients (HoG)



Divide image into 8x8 pixel regions
Within each region quantize edge direction into 9 bins



Example: 320x240 image gets divided into 40x30 bins; in each bin there are 9 numbers so feature vector has $30 \times 40 \times 9 = 10,800$ numbers

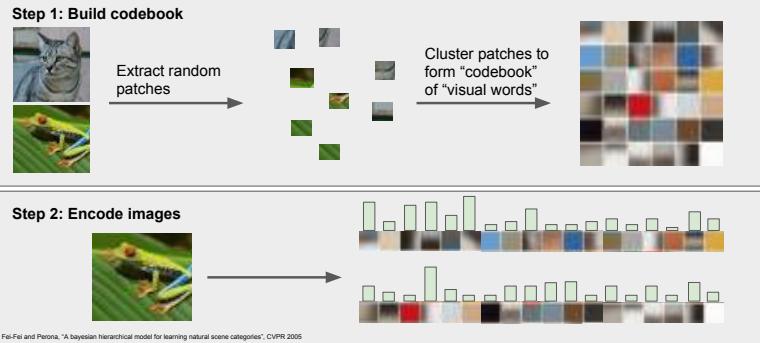
Lowe, "Object recognition from local scale-invariant features", ICCV 1999
Dalal and Triggs, "Histograms of oriented gradients for human detection", CVPR 2005

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 93

April 10, 2018

Example: Bag of Words

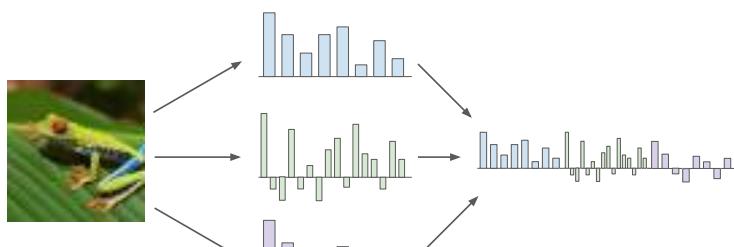


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 94

April 10, 2018

Aside: Image Features

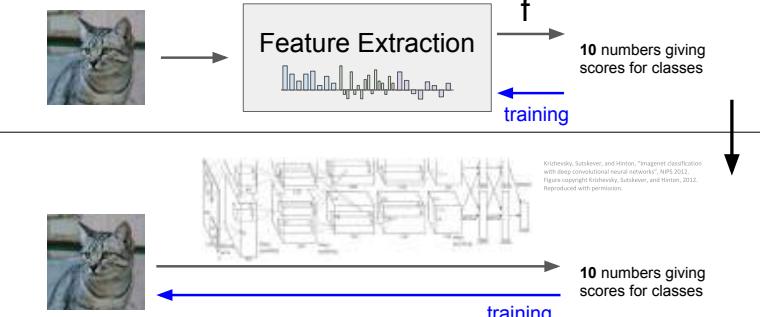


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 95

April 10, 2018

Image features vs ConvNets



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 96

April 10, 2018

Next time:

Introduction to neural networks

Backpropagation

Lecture 4: Backpropagation and Neural Networks

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 3 - 97

April 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 1

April 12, 2018

Administrative

Assignment 1 due Wednesday April 18, 11:59pm

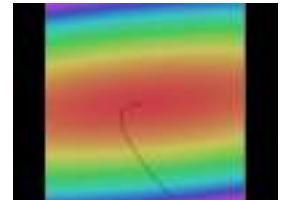
Administrative

All office hours this week will use queuestatus

Where we are...

$$\begin{aligned} s &= f(x; W) = Wx && \text{scores function} \\ L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) && \text{SVM loss} \\ L &= \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2 && \text{data loss + regularization} \\ \text{want } \nabla_W L & \text{ minimise "L"} \end{aligned}$$

Optimization



```
# Vanilla Gradient Descent
while True:
    weights -= grad * evaluate_gradient(loss_fn, data, weights)
    weights += step_size * weights.grad # prevent overflow update
```

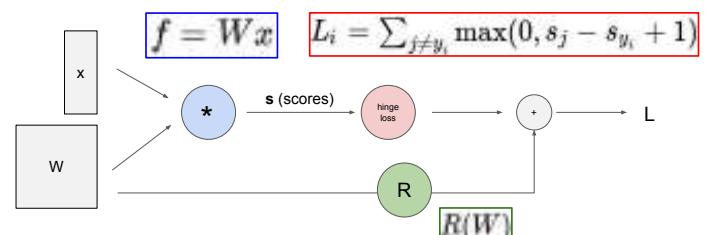
Gradient descent

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

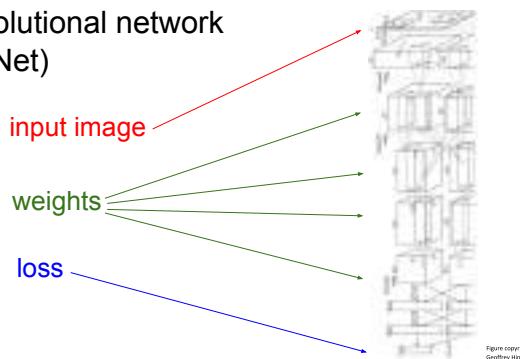
Numerical gradient: slow (:), approximate (:), easy to write :()
Analytic gradient: fast (:), exact (:), error-prone :()

In practice: Derive analytic gradient, check your implementation with numerical gradient

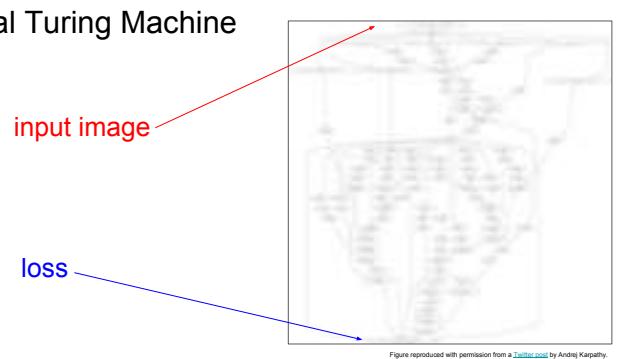
Computational graphs



Convolutional network (AlexNet)



Neural Turing Machine



Fei-Fei Li & Justin Johnson & Serena Yeung

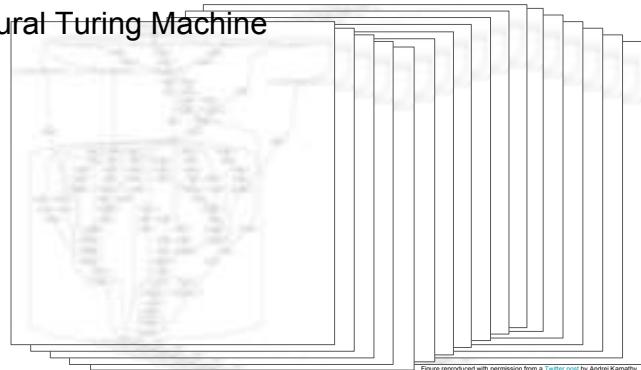
Lecture 4 - 8

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 9

April 12, 2018

Neural Turing Machine



Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 -

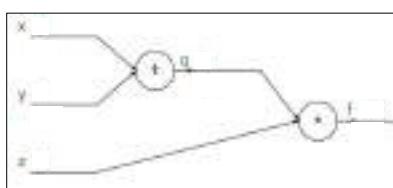
April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 11

April 12, 2018

Backpropagation: a simple example

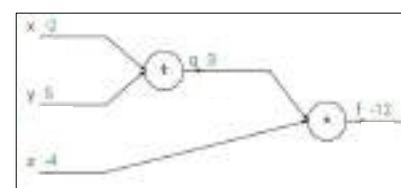
$$f(x, y, z) = (x + y)z$$



Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 12

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 13

April 13, 2018

Backpropagation: a simple example

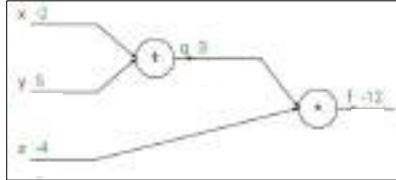
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

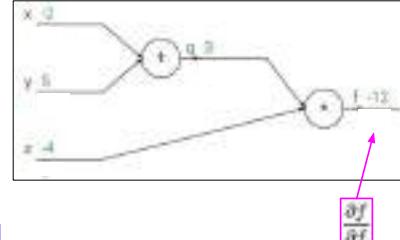
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

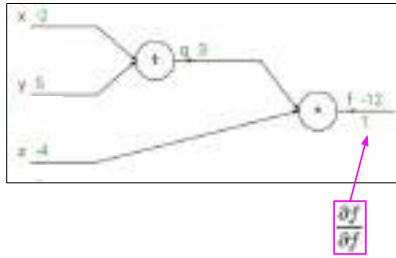
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

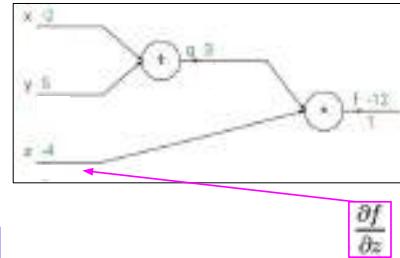
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

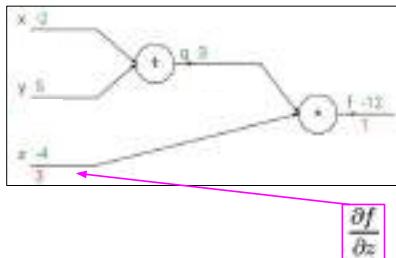
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

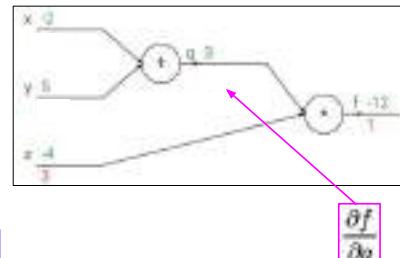
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

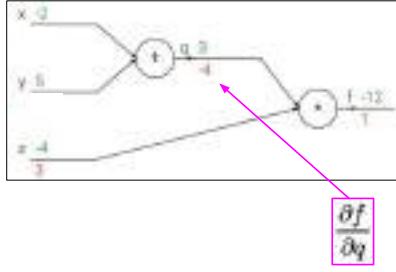
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

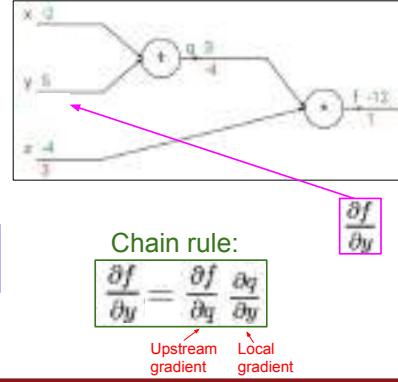
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

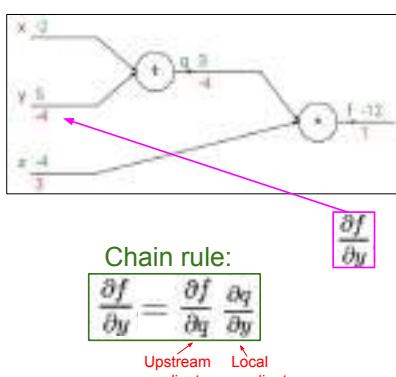
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

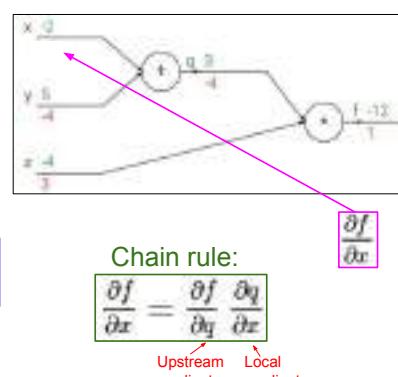
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: a simple example

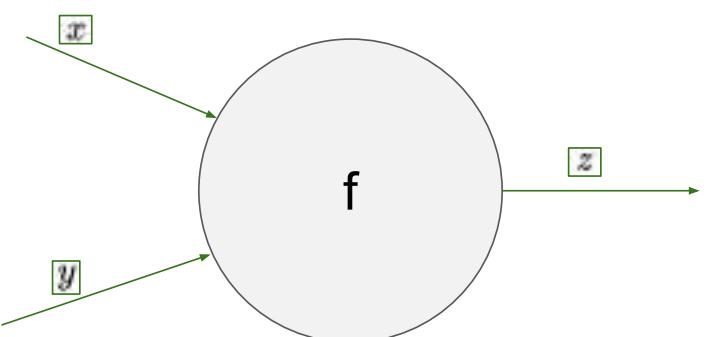
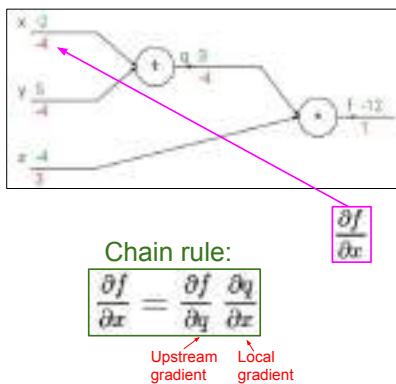
$$f(x, y, z) = (x + y)z$$

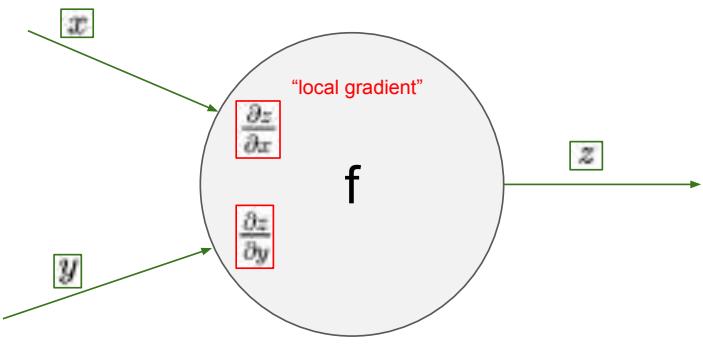
e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$





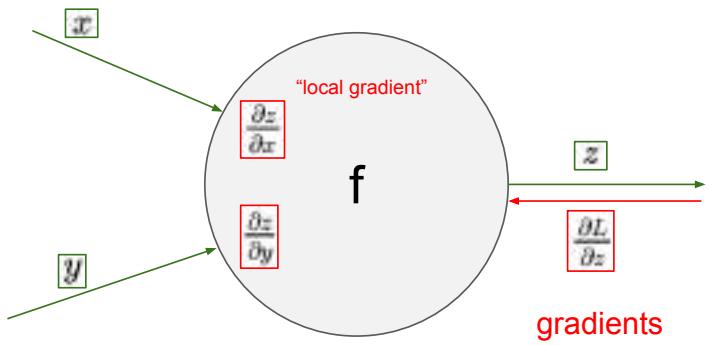
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 26

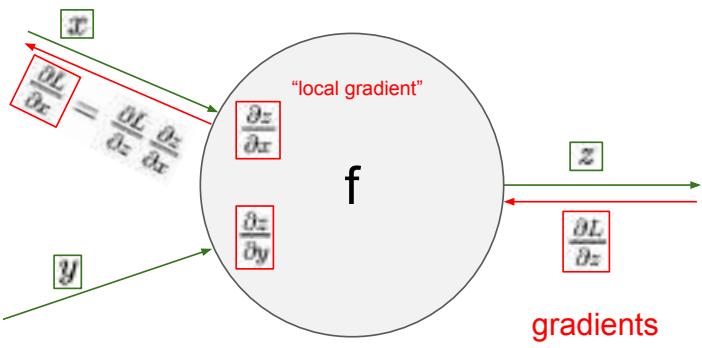
April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 27

April 12, 2018



gradients



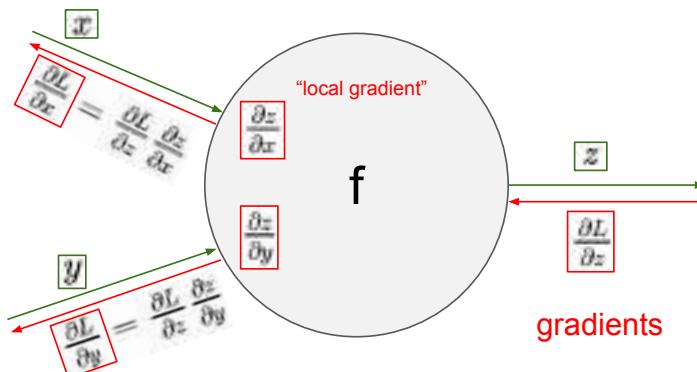
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 28

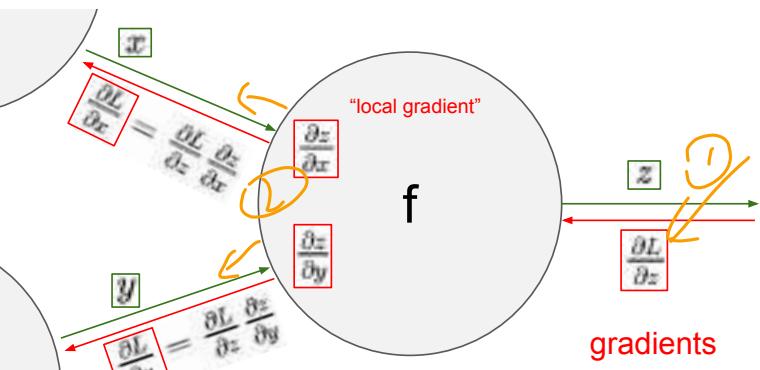
April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 29

April 12, 2018



gradients



Fei Li & Justin Johnson & Serena Yeung

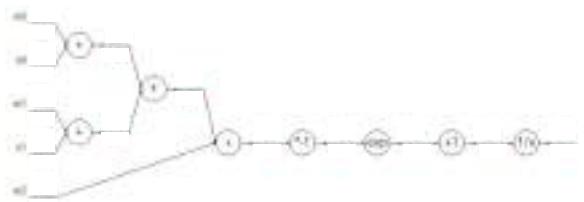
Lecture 4 - 30

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

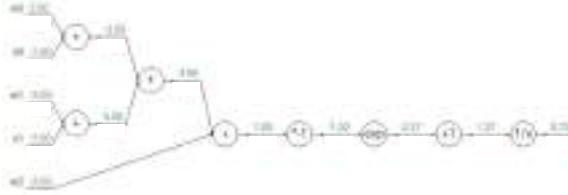
Lecture 4 - 31

April 12, 2018

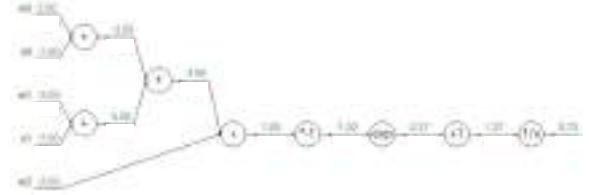
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

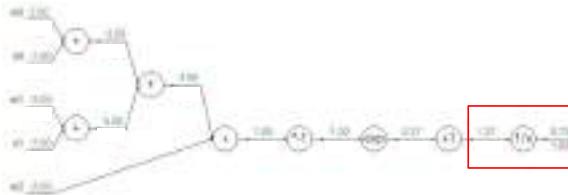


Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$\begin{array}{c|c} f(x) = e^x & \rightarrow \\ f_x(x) = ux & \rightarrow \\ \hline \end{array} \quad \begin{array}{c|c} \frac{df}{dx} = e^x & \rightarrow \\ \frac{df}{dx} = u & \rightarrow \\ \hline \end{array} \quad \begin{array}{c|c} f(x) = \frac{1}{e^x} & \rightarrow \\ f_x(x) = c + x & \rightarrow \\ \hline \end{array} \quad \begin{array}{c|c} \frac{df}{dx} = -1/e^x & \rightarrow \\ \frac{df}{dx} = 1 & \rightarrow \\ \hline \end{array}$$

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



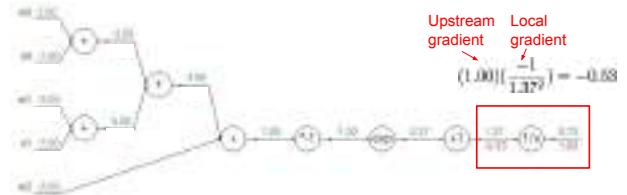
$$\begin{array}{c|c} f(x) = e^x & \rightarrow \\ f_x(x) = ux & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} \frac{df}{dx} = e^x & \rightarrow \\ \frac{df}{dx} = u & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} f(x) = \frac{1}{e^x} & \rightarrow \\ f_x(x) = c + x & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} \frac{df}{dx} = -1/e^x & \rightarrow \\ \frac{df}{dx} = 1 & \rightarrow \\ \hline \end{array}$$

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



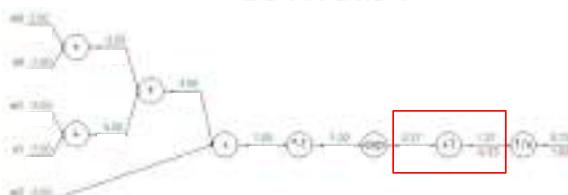
$$\begin{array}{c|c} f(x) = e^x & \rightarrow \\ f_x(x) = ux & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} \frac{df}{dx} = e^x & \rightarrow \\ \frac{df}{dx} = u & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} f(x) = \frac{1}{e^x} & \rightarrow \\ f_x(x) = c + x & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} \frac{df}{dx} = -1/e^x & \rightarrow \\ \frac{df}{dx} = 1 & \rightarrow \\ \hline \end{array}$$

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



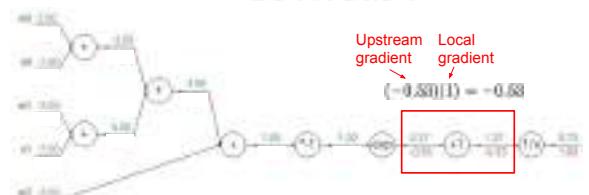
$$\begin{array}{c|c} f(x) = e^x & \rightarrow \\ f_x(x) = ux & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} \frac{df}{dx} = e^x & \rightarrow \\ \frac{df}{dx} = u & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} f(x) = \frac{1}{e^x} & \rightarrow \\ f_x(x) = c + x & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} \frac{df}{dx} = -1/e^x & \rightarrow \\ \frac{df}{dx} = 1 & \rightarrow \\ \hline \end{array}$$

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



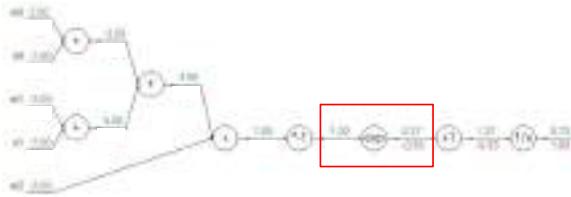
$$\begin{array}{c|c} f(x) = e^x & \rightarrow \\ f_x(x) = ux & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} \frac{df}{dx} = e^x & \rightarrow \\ \frac{df}{dx} = u & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} f(x) = \frac{1}{e^x} & \rightarrow \\ f_x(x) = c + x & \rightarrow \\ \hline \end{array}$$

$$\begin{array}{c|c} \frac{df}{dx} = -1/e^x & \rightarrow \\ \frac{df}{dx} = 1 & \rightarrow \\ \hline \end{array}$$

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

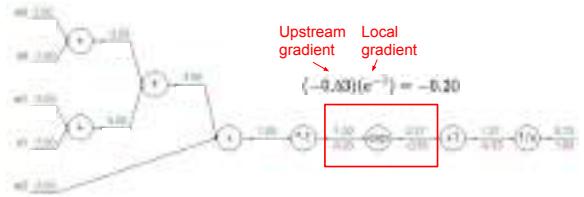


$$\begin{array}{c} f(x) = e^x \\ f_s(x) = ux \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = e^x \\ \frac{df}{dx} = u \end{array}$$

$$\begin{array}{c} f(x) = \frac{1}{e} \\ f_t(x) = c + x \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = -1/e^x \\ \frac{df}{dx} = 1 \end{array}$$

$$\begin{array}{c} f(x) = e^x \\ f_s(x) = ux \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = e^x \\ \frac{df}{dx} = u \end{array}$$

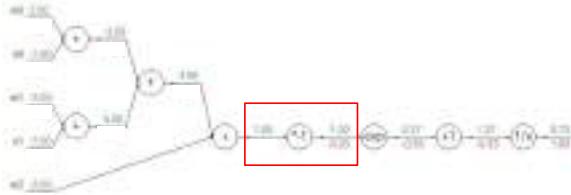
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$\begin{array}{c} f(x) = e^x \\ f_s(x) = ux \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = e^x \\ \frac{df}{dx} = u \end{array}$$

$$\begin{array}{c} f(x) = \frac{1}{e} \\ f_t(x) = c + x \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = -1/e^x \\ \frac{df}{dx} = 1 \end{array}$$

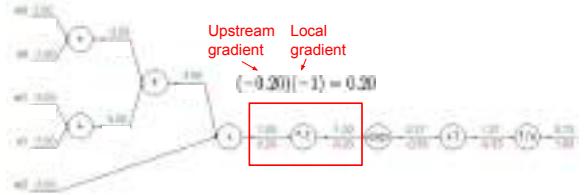
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$\begin{array}{c} f(x) = e^x \\ f_s(x) = ux \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = e^x \\ \frac{df}{dx} = u \end{array}$$

$$\begin{array}{c} f(x) = \frac{1}{e} \\ f_t(x) = c + x \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = -1/e^x \\ \frac{df}{dx} = 1 \end{array}$$

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$\begin{array}{c} f(x) = e^x \\ f_s(x) = ux \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = e^x \\ \frac{df}{dx} = u \end{array}$$

$$\begin{array}{c} f(x) = \frac{1}{e} \\ f_t(x) = c + x \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = -1/e^x \\ \frac{df}{dx} = 1 \end{array}$$

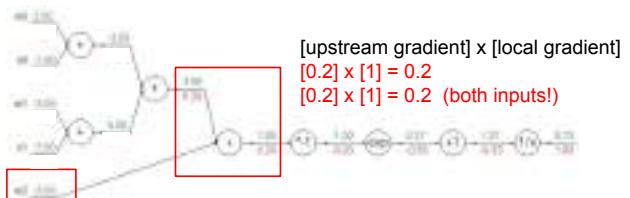
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$\begin{array}{c} f(x) = e^x \\ f_s(x) = ux \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = e^x \\ \frac{df}{dx} = u \end{array}$$

$$\begin{array}{c} f(x) = \frac{1}{e} \\ f_t(x) = c + x \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = -1/e^x \\ \frac{df}{dx} = 1 \end{array}$$

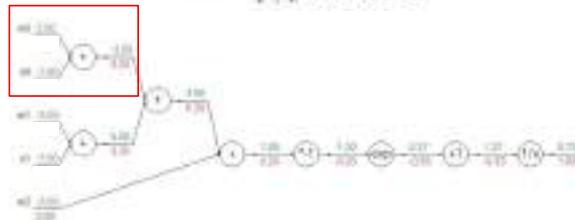
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$\begin{array}{c} f(x) = e^x \\ f_s(x) = ux \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = e^x \\ \frac{df}{dx} = u \end{array}$$

$$\begin{array}{c} f(x) = \frac{1}{e} \\ f_t(x) = c + x \end{array} \rightarrow \begin{array}{c} \frac{df}{dx} = -1/e^x \\ \frac{df}{dx} = 1 \end{array}$$

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



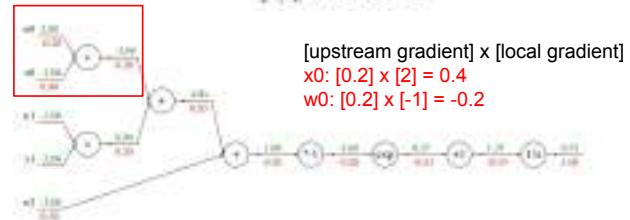
$$\begin{aligned} f(x) &= e^x \\ f_a(x) &= ux \end{aligned}$$

$$\begin{aligned} \frac{df}{dx} &= e^x \\ \frac{df}{dx} &= u \end{aligned}$$

$$\begin{aligned} f(x) &= \frac{1}{e} \\ f_a(x) &= c + x \end{aligned}$$

$$\begin{aligned} \frac{df}{dx} &= -1/e^x \\ \frac{df}{dx} &= 1 \end{aligned}$$

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$\begin{aligned} [\text{upstream gradient}] \times [\text{local gradient}] \\ x0: [0.2] \times [2] = 0.4 \\ w0: [0.2] \times [-1] = -0.2 \end{aligned}$$

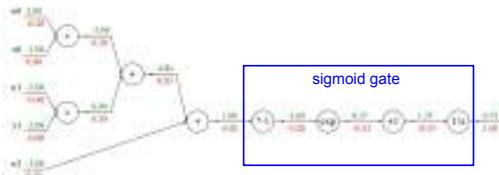
$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

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

sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\sigma(x)$$



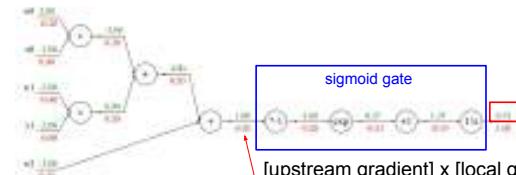
$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

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

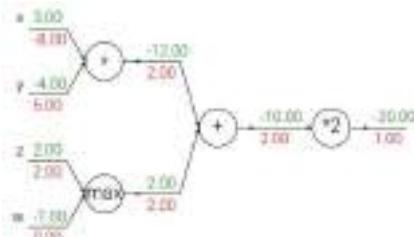
sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\sigma(x)$$



Patterns in backward flow

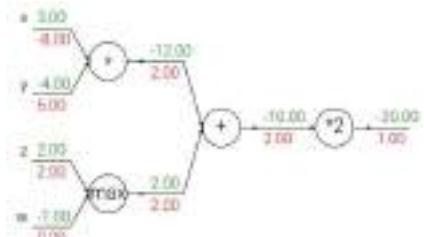
add gate: gradient distributor



Patterns in backward flow

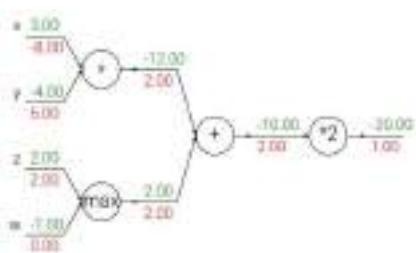
add gate: gradient distributor

Q: What is a **max gate**?



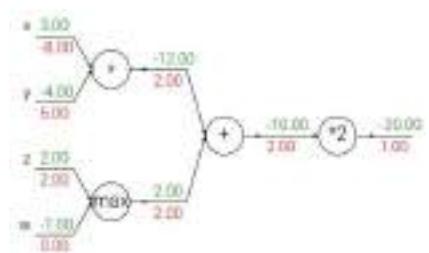
Patterns in backward flow

add gate: gradient distributor
max gate: gradient router



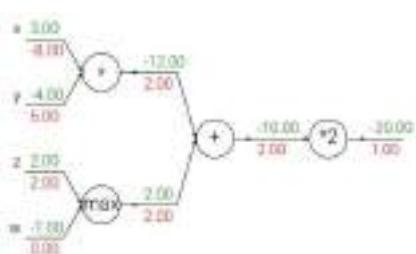
Patterns in backward flow

add gate: gradient distributor
max gate: gradient router
Q: What is a mul gate?

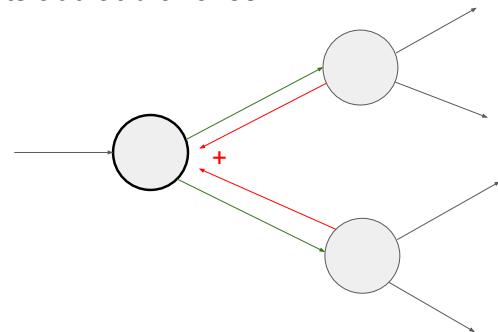


Patterns in backward flow

add gate: gradient distributor
max gate: gradient router
mul gate: gradient switcher



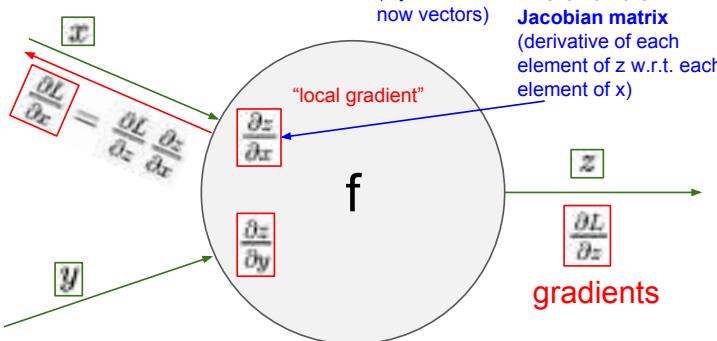
Gradients add at branches



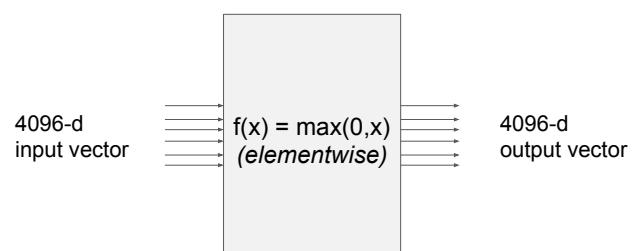
Gradients for vectorized code

(x, y, z are now vectors)

This is now the **Jacobian matrix** (derivative of each element of z w.r.t. each element of x)

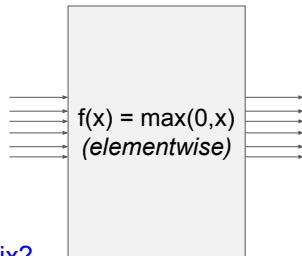


Vectorized operations



Vectorized operations

4096-d input vector



Q: what is the size of the Jacobian matrix?

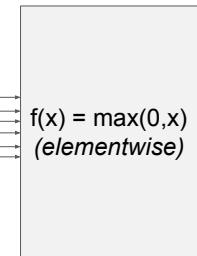
$$\frac{\partial L}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial L}{\partial f} \end{bmatrix}$$

Jacobian matrix

4096-d output vector

Vectorized operations

4096-d input vector



Q: what is the size of the Jacobian matrix?
[4096 x 4096!]

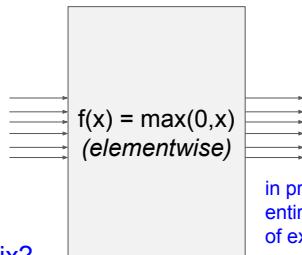
$$\frac{\partial L}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial L}{\partial f} \end{bmatrix}$$

Jacobian matrix

4096-d output vector

Vectorized operations

4096-d input vector



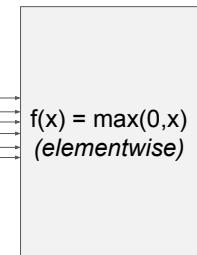
Q: what is the size of the Jacobian matrix?
[4096 x 4096!]

4096-d output vector

in practice we process an entire minibatch (e.g. 100) of examples at one time:
i.e. Jacobian would technically be a [409,600 x 409,600] matrix :)

Vectorized operations

4096-d input vector



Q: what is the size of the Jacobian matrix?
[4096 x 4096!]

$$\frac{\partial L}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial L}{\partial f} \end{bmatrix}$$

Jacobian matrix

4096-d output vector

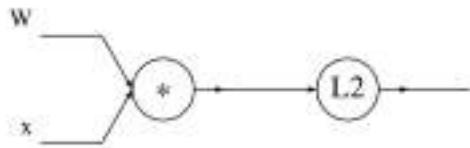
Q2: what does it look like?

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$

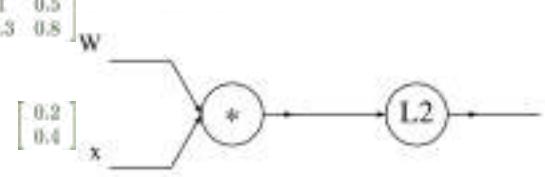
A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$

$$\in \mathbb{R}^n \in \mathbb{R}^{n \times n}$$

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



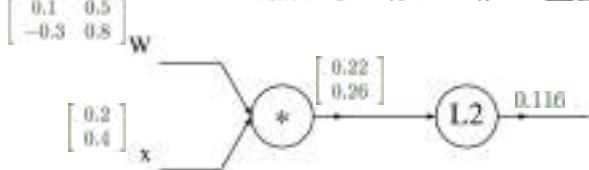
A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



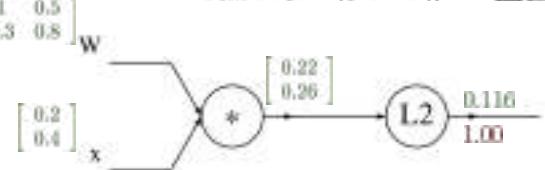
$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



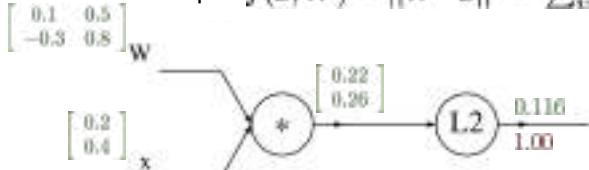
A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



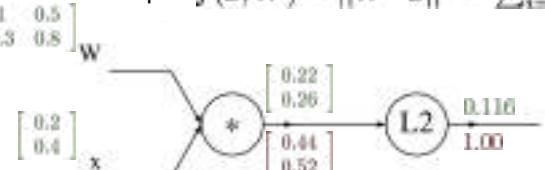
$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



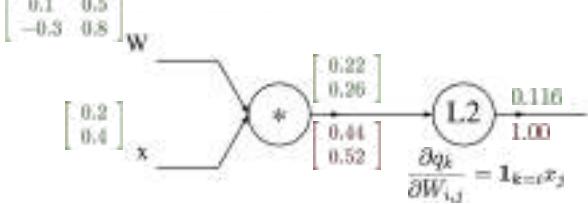
$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial f}{\partial q_i} = 2q_i$$

$$\nabla_q f = 2q$$

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W \quad \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} x \quad \begin{bmatrix} 0.22 \\ 0.26 \\ 0.44 \\ 0.52 \end{bmatrix} \quad \text{L2} \quad \begin{bmatrix} 0.116 \\ 1.00 \end{bmatrix}$$

$$\frac{\partial q_k}{\partial W_{i,j}} = \mathbf{1}_{k=i} x_j$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 68

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 69

April 12, 2018

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$

Diagram illustrating the calculation of the squared norm of the weighted input:

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W \quad \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} x \quad \begin{bmatrix} 0.22 \\ 0.26 \\ 0.44 \\ 0.52 \end{bmatrix} \quad \text{L2} \quad \begin{bmatrix} 0.116 \\ 1.00 \end{bmatrix}$$

$$\frac{\partial q_k}{\partial W_{i,j}} = \mathbf{1}_{k=i} x_j$$

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j)$$

$$= 2q_i x_j$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 68

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 69

April 12, 2018

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$

Diagram illustrating the calculation of the squared norm of the weighted input:

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 70

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 71

April 12, 2018

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$

Diagram illustrating the calculation of the squared norm of the weighted input:

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W \quad \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} x \quad \begin{bmatrix} 0.22 \\ 0.26 \\ 0.44 \\ 0.52 \end{bmatrix} \quad \text{L2} \quad \begin{bmatrix} 0.116 \\ 1.00 \end{bmatrix}$$

$$\frac{\partial q_k}{\partial W_{i,j}} = \mathbf{1}_{k=i} x_j$$

$$\frac{\partial f}{\partial W_{i,j}} = \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial W_{i,j}}$$

$$= \sum_k (2q_k)(\mathbf{1}_{k=i} x_j)$$

$$= 2q_i x_j$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 70

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 71

April 12, 2018

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$

Diagram illustrating the calculation of the squared norm of the weighted input:

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 72

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 73

April 12, 2018

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$

Diagram illustrating the calculation of the squared norm of the weighted input:

$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

$$\begin{bmatrix} 0.1 & 0.5 \\ -0.3 & 0.8 \end{bmatrix} W \quad \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix} x \quad \begin{bmatrix} 0.22 \\ 0.26 \\ 0.44 \\ 0.52 \end{bmatrix} \quad \text{L2} \quad \begin{bmatrix} 0.116 \\ 1.00 \end{bmatrix}$$

$$\frac{\partial q_k}{\partial x_i} = W_{k,i}$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 72

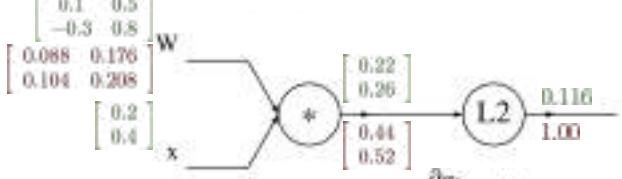
April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 73

April 12, 2018

Always check: The gradient with respect to a variable should have the same shape as the variable

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



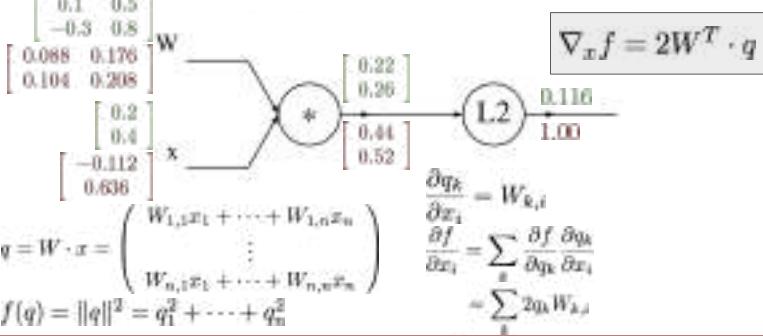
$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

Lecture 4 - 74

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



Lecture 4 - 75

April 12, 2018

In discussion section: A matrix example...

$$z_1 = XW_1$$

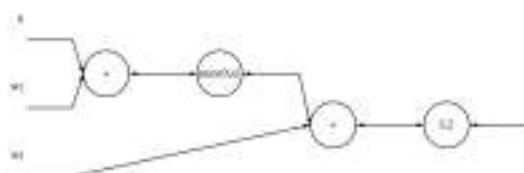
$$h_1 = \text{ReLU}(z_1)$$

$$\hat{y} = h_1 W_2$$

$$L = \|\hat{y}\|_2^2$$

$$\frac{\partial L}{\partial W_2} = ?$$

$$\frac{\partial L}{\partial W_1} = ?$$



Modularized implementation: forward / backward API

Graph (or Net) object (rough pseudo code)

```
class GraphNetObject:
    def forward(self, inputs):
        # 1. pass inputs to input ports ...
        # 2. forward the computational graph
        # 3. get a list graph nodes, topologically sorted
        #     (graph.reorder)
        # 4. find head # the final node in the graph outputs the last
        #     backward()
        # 5. get a reversed list graph.nodes_topologically_reversed()
        #     graph.backward()
        # 6. backward()

    def backward(self, upstream_gradients):
        # compute gradients
        # ...
```

Fei-Fei Li & Justin Johnson & Serena Yeung

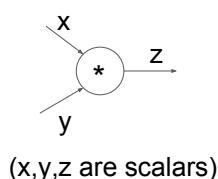
Lecture 4 - 76

April 13, 2017 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 77

April 12, 2018

Modularized implementation: forward / backward API



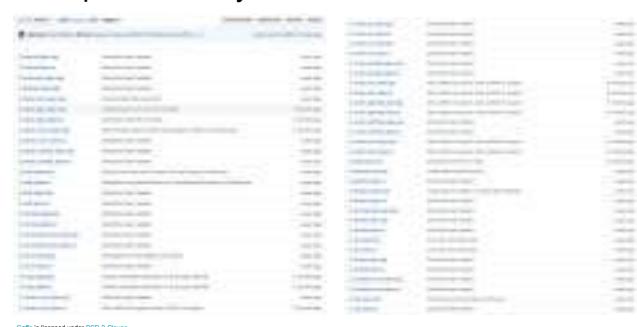
(x,y,z are scalars)

```
class MultiOpProtocolObject:
    def forward(self, p):
        z = p[0]*p[1]
        mat, r, s = p[2].mat, p[2].row, p[2].sum
        mat[r, s] = z
        p[2] = z
        return p

    def backward(self, d):
        dx = self[0].p + d * (self[1].p + d * self[2].p)
        dy = self[1].p + d * (self[0].p + d * self[2].p)
        d[2] = (dx, dy)
```

Local gradient Upstream gradient variable

Example: Caffe layers



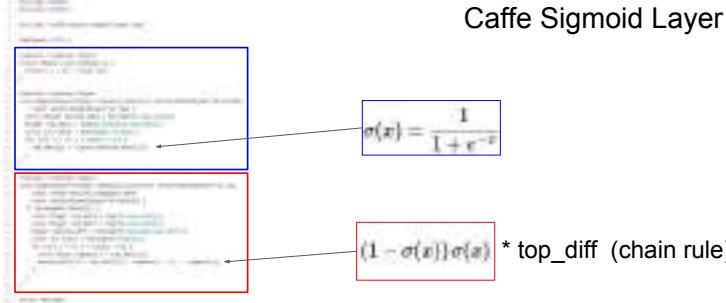
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 78

April 12, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 4 - 79

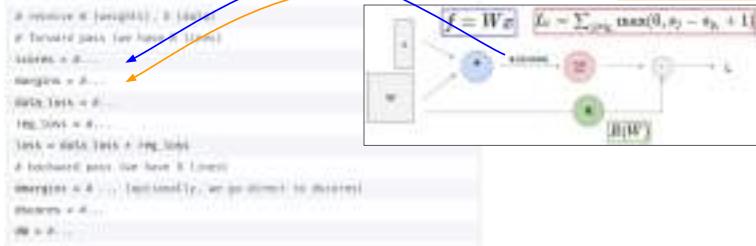
April 12, 2018



In Assignment 1: Writing SVM / Softmax

Stage your forward/backward computation!

E.g. for the SVM:



Summary so far...

- neural nets will be very large: impractical to write down gradient formula by hand for all parameters
 - **backpropagation** = recursive application of the chain rule along a computational graph to compute the gradients of all inputs/parameters/intermediates
 - implementations maintain a graph structure, where the nodes implement the **forward()** / **backward()** API
 - **forward**: compute result of an operation and save any intermediates needed for gradient computation in memory
 - **backward**: apply the chain rule to compute the gradient of the loss function with respect to the inputs

Next: Neural Networks

Neural networks: without the brain stuff

Neural networks: without the brain stuff

(Before) Linear score function: $f = Wx$

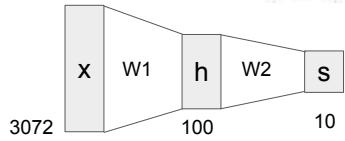
(Before) Linear score function: $f = Wx$

(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$

Neural networks: without the brain stuff

(Before) Linear score function: $f = Wx$

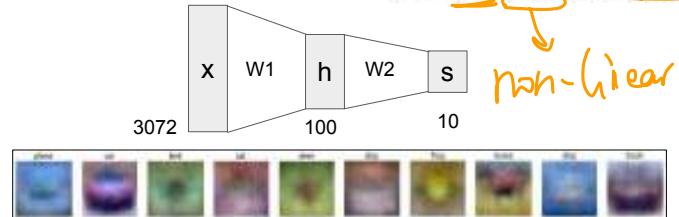
(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$



Neural networks: without the brain stuff

(Before) Linear score function: $f = Wx$

(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$



Neural networks: without the brain stuff

(Before) Linear score function: $f = Wx$

(Now) 2-layer Neural Network
or 3-layer Neural Network

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

Full implementation of training a 2-layer Neural Network needs ~20 lines:

```

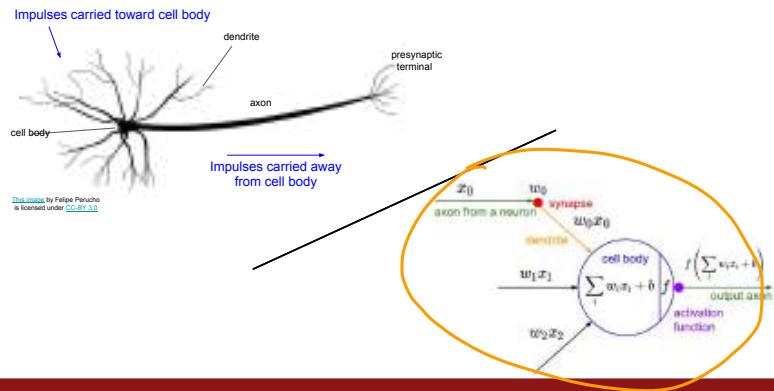
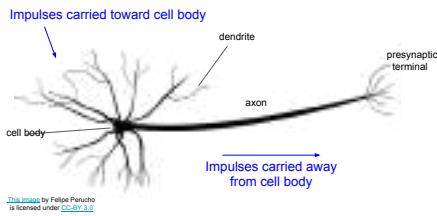
1 import numpy as np
2 from scipy.random import randn
3
4 H0, H1, H2, H3, loss = np.zeros((100, 100, 100, 10))
5 H0[0] = randn(10, 100); H1[0] = randn(100, 100)
6 W1, W2 = randn(100, 100); W3 = randn(100, 10)
7
8 for t in range(10000):
9     q = -1 / t + 10 * np.exp(-t * np.pi / 10)
10    a0 = np.zeros(100); a0[0] = q
11    a0[1] = np.exp(q) / sum(a0)
12    a1 = a0.T.dot(H0) - q * H1[0]
13    a1[1] = np.exp(a1) / sum(a1)
14    a2 = a1.T.dot(W1) - q * H2[0]
15    a2[1] = np.exp(a2) / sum(a2)
16    a3 = a2.T.dot(W2) - q * H3[0]
17    a3[1] = np.exp(a3) / sum(a3)
18    a4 = a3.T.dot(W3)
19    a4[1] = np.exp(a4) / sum(a4)
20
21    grad_w3 = 2 * (a4 - y) * (a4 * (1 - a4) * p1)
22    grad_a3 = (a4.T).dot(grad_w3 * p1)
23    grad_h3 = grad_w3 * a3 * (1 - a3) * p1
24    grad_w2 = a3.T.dot(grad_h3) + b2 * 1e-06
25
26    grad_w1 = a2.T.dot(grad_h2) + b1 * 1e-06
27    grad_a2 = (a2.T).dot(grad_w1) + b2 * 1e-06
28    grad_h2 = grad_w1 * a2 * (1 - a2) * p1
29
30    grad_w3 -= 0.01 * grad_w3
31    grad_w1 -= 0.01 * grad_w1
  
```

In HW: Writing a 2-layer net

```

# receive N0,N1,N2,losses), # inputs
# forward pass:
b1 = ...; # bias of h1,W1,b1
scores = ...; # function of h1,W1,b1
loss = ...; # several lines of code to evaluate softmax loss
# backward pass:
dcores = ...;
dh1,dW1,db1 = ...;
dW1,dh1 = ...
  
```



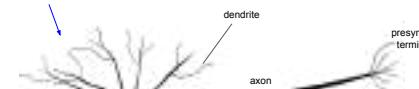


Impulses carried toward cell body

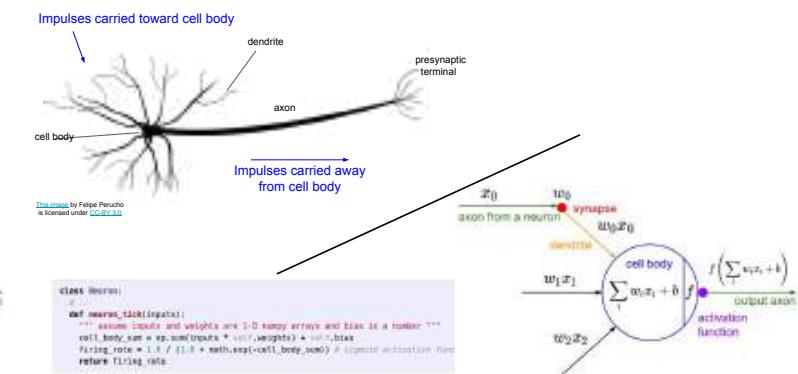


Impulses carried away from cell body

Impulses carried toward cell body



Impulses carried away from cell body



Be very careful with your brain analogies!

Biological Neurons:

- Many different types
- Dendrites can perform complex non-linear computations
- Synapses are not a single weight but a complex non-linear dynamical system
- Rate code may not be adequate

[Dendritic Computation, London and Häusser]

Activation functions

Sigmoid

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

tanh

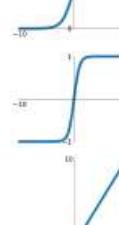
$$\tanh(x)$$

ReLU

$$\max(0, x)$$

Leaky ReLU

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

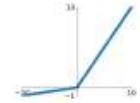


Maxout

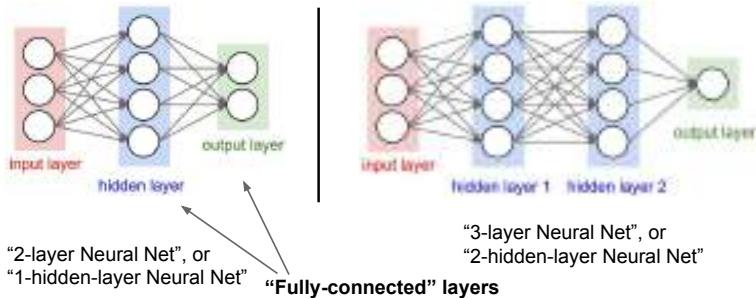
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Neural networks: Architectures

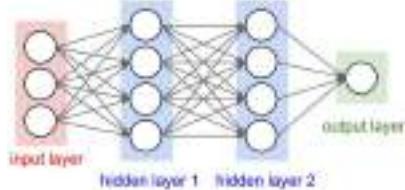


Example feed-forward computation of a neural network

```
class Neuron:  
    def neuron_tick(inputs):  
        """ assume inputs and weights are 1-D numpy arrays and bias is a number """  
        cell_body_sum = np.sum(inputs * self.weights) + self.bias  
        firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation function  
        return firing_rate
```

We can efficiently evaluate an entire layer of neurons.

Example feed-forward computation of a neural network



Summary

- We arrange neurons into fully-connected layers
- The abstraction of a **layer** has the nice property that it allows us to use efficient vectorized code (e.g. matrix multiplies)
- Neural networks are not really *neural*
- Next time: Convolutional Neural Networks

```
# Function class of a 2-layer neural net:  
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function: error sigmoid  
x = np.random.randn(3, 10) # random input vector, of shape (num nodes, num)  
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activation (size 10)  
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activation (size 5)  
out = np.dot(W3, h2) + b3 # random numbers r/200
```

Administrative

Lecture 5: Convolutional Neural Networks

Assignment 1 due **Wednesday April 18**, 11:59pm
Assignment 2 will also be released Wednesday

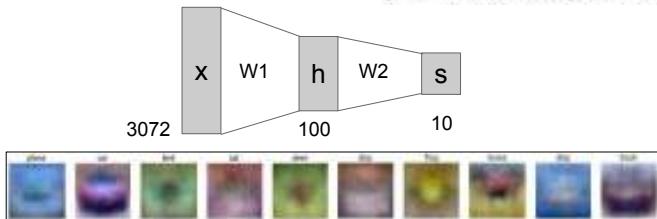
Last time: Neural Networks

Linear score function:

$$f = Wx$$

2-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 3

April 17, 2018,

Next: Convolutional Neural Networks

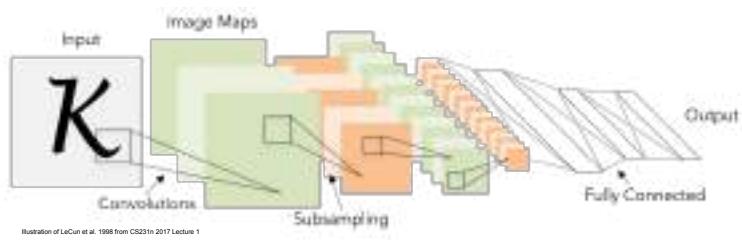


Illustration of LeCun et al. 1998 from CS231n Lecture 1

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 4

April 17, 2018,

A bit of history...

The **Mark I Perceptron** machine was the first implementation of the perceptron algorithm.

The machine was connected to a camera that used 20×20 cadmium sulfide photocells to produce a 400-pixel image.

recognized letters of the alphabet

$$f(x) = \begin{cases} 1 & \text{if } \sum w_i x_i + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

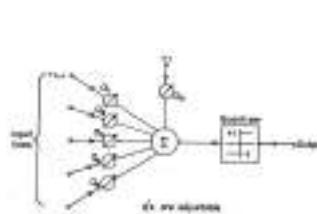
From a neuron's point of view:
 x_1, x_2, \dots, x_n : net input
 w_1, w_2, \dots, w_n : weights
 b : bias
 $\sigma(\cdot)$: activation function

Frank Rosenblatt, ~1957: Perceptron



[This image](#) by Rocky Acosta is licensed under [CC-BY 3.0](#)

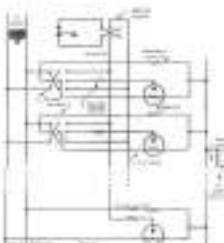
A bit of history...



Widrow and Hoff, ~1960: Adaline/Madaline



These figures are reproduced from [Widrow, 1960, Stanford Electronics Laboratories Technical Report](#) with permission from [Stanford University Special Collections](#).



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 5

April 17, 2018,

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 6

April 17, 2018,

A bit of history...

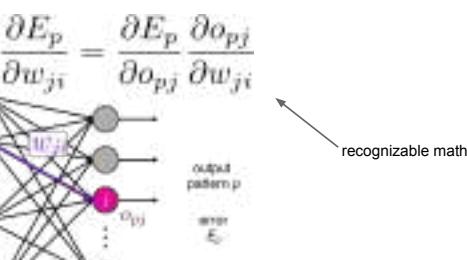


Illustration of Rumelhart et al., 1986 by Lane Mcintosh, copyright CS231n 2017

Rumelhart et al., 1986: First time back-propagation became popular

A bit of history...

[Hinton and Salakhutdinov 2006]

Reinvigorated research in Deep Learning

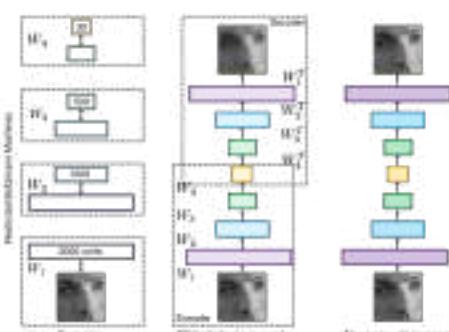


Illustration of Hinton and Salakhutdinov 2006 by Lane McIntosh, copyright CS231n 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 7

April 17, 2018,

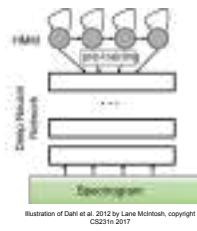
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 8

April 17, 2018,

First strong results

Acoustic Modeling using Deep Belief Networks
 Abdel-rahman Mohamed, George Dahl, Geoffrey Hinton, 2010
Context-Dependent Pre-trained Deep Neural Networks for Large Vocabulary Speech Recognition
 George Dahl, Dong Yu, Li Deng, Alex Acero, 2012



Imagenet classification with deep convolutional neural networks
 Alex Krizhevsky, Ilya Sutskever, Geoffrey E Hinton, 2012



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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 9 April 17, 2018

A bit of history:

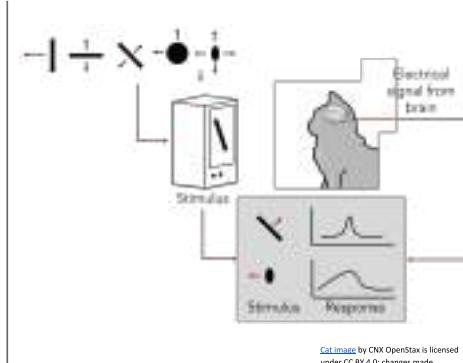
Hubel & Wiesel, 1959

RECEPTIVE FIELDS OF SINGLE NEURONES IN THE CAT'S STRIATE CORTEX

1962

RECEPTIVE FIELDS, BINOCULAR INTERACTION AND FUNCTIONAL ARCHITECTURE IN THE CAT'S VISUAL CORTEX

1968...

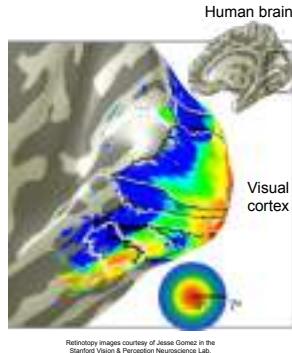
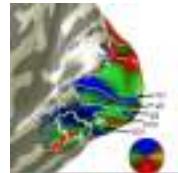


Lecture 5 - 10 April 17, 2018

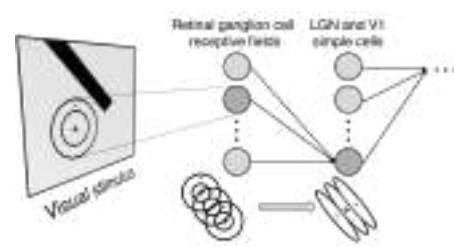
Fei-Fei Li & Justin Johnson & Serena Yeung

A bit of history

Topographical mapping in the cortex:
 nearby cells in cortex represent nearby regions in the visual field



Hierarchical organization



Simple cells:
Response to light orientation

Complex cells:
Response to light orientation and movement

Hypercomplex cells:
Response to movement with an end point



Fei-Fei Li & Justin Johnson & Serena Yeung

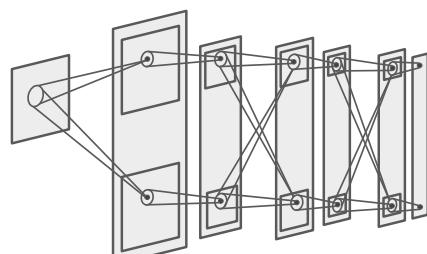
Lecture 5 - 11 April 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 12 April 17, 2018

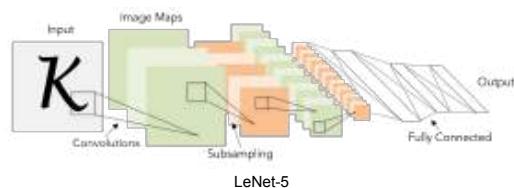
A bit of history:

Neocognitron [Fukushima 1980]



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

A bit of history:
Gradient-based learning applied to document recognition
[LeCun, Bottou, Bengio, Haffner 1998]



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 13 April 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 14 April 17, 2018

A bit of history:

ImageNet Classification with Deep Convolutional Neural Networks

[Krizhevsky, Sutskever, Hinton, 2012]

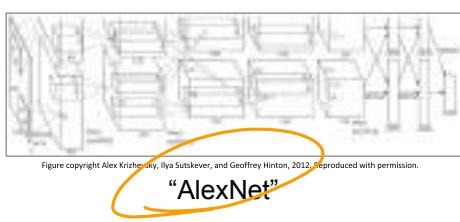


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Fast-forward to today: ConvNets are everywhere



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

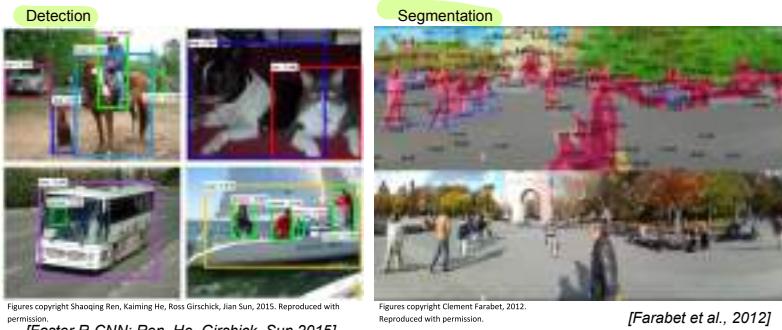
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 15 April 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 16 April 17, 2018

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]

[Farabet et al., 2012]

Fast-forward to today: ConvNets are everywhere



self-driving cars



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

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

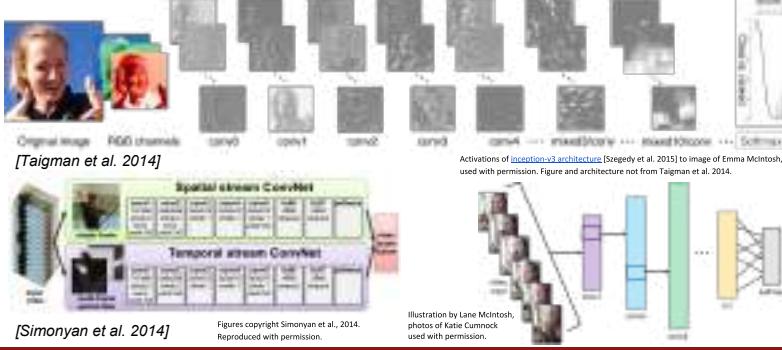
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 17 April 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 18 April 17, 2018

Fast-forward to today: ConvNets are everywhere



[Taigman et al. 2014]

Activations of *Inception-v3 architecture* [Szegedy et al. 2015] to image of Emma McIntosh, used with permission. Figure and architecture not from Taigman et al. 2014.

[Simonyan et al. 2014]

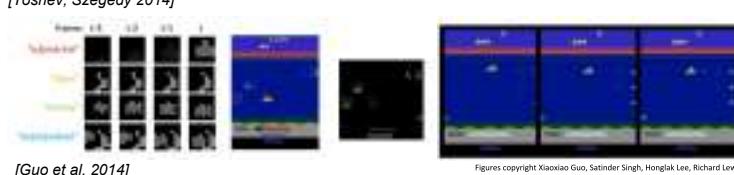
Figures copyright Simonyan et al., 2014. Reproduced with permission.

Fast-forward to today: ConvNets are everywhere



[Toshev, Szegedy 2014]

Images are examples of pose estimation, not actually from Toshev & Szegedy 2014. Copyright Lane McIntosh.



[Guo et al. 2014]

Figures copyright Xiaoxiao Guo, Satinder Singh, Honglak Lee, Richard Lewis, and Xiaohi Wang, 2014. Reproduced with permission.

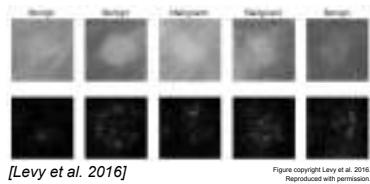
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 19 April 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 20 April 17, 2018

Fast-forward to today: ConvNets are everywhere



[Levy et al. 2016]

Figure copyright Levy et al. 2016.
Reproduced with permission.



[Serianet et al. 2011]
[Ciresan et al.]

This image by Christin Khan is in the public domain
and originally came from the U.S. NOAA.



Whale recognition, Kaggle Challenge



Mnih and Hinton, 2010

Photo and figure by Lane Mcintosh; not actual example from Mnih and Hinton, 2010 paper.



[Dieleman et al. 2014]

From left to right: photo courtesy of NASA, image (c) reproduced by
ESA/Hubble, photo courtesy of NASA, and photo courtesy of

Lecture 5 - 21 April 17, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 22 April 17, 2018

Lecture 5 - 22 April 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 5 - 21 April 17, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 5 - 22 April 17, 2018

No errors

Minor errors

Somewhat related

Image Captioning

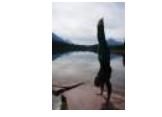
[Vinyals et al., 2015]
[Karpathy and Fei-Fei, 2015]



A white teddy bear sitting in
the grass

A man in a baseball
uniform throwing a ball

A woman is holding a
cat in her hand



A man riding a wave on
top of a surfboard

A cat sitting on a
suitcase on the floor

A woman standing on a
beach holding a surfboard

All images are CC0 Public domain:
<https://creativecommons.org/publicdomain/zero/1.0/>
<https://creativecommons.org/publicdomain/zero/1.0/licenses.html?lang=en>
<https://creativecommons.org/licenses/by-sa/2.0/>
<https://creativecommons.org/licenses/by-nd/2.0/>
<https://creativecommons.org/licenses/by-nc/2.0/>
<https://creativecommons.org/licenses/by-nd-nc/2.0/>

Captions generated by Justin Johnson using Neuraltalk2

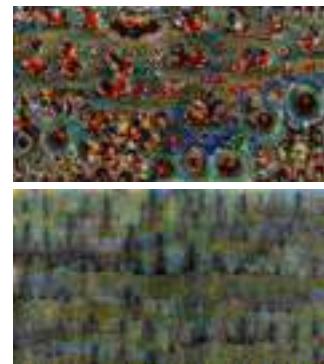


Figure copyright Justin Johnson, 2015. Reproduced with permission. Generated using the Inceptionism approach
from a <https://code.google.com/p/neuraltalk2/> by Google Research.

Original images in CC0 public domain:
https://commons.wikimedia.org/wiki/File:Van_Gogh_Saint_Remy_1889.jpg

Stylized images copyright Justin Johnson, 2015; reproduced with permission

Gatys et al., "Image Style Transfer using Convolutional Neural Networks", CVPR 2016

Gatys et al., "Controlling Perceptual Factors in Neural Style Transfer", CVPR 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 23 April 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 24 April 17, 2018

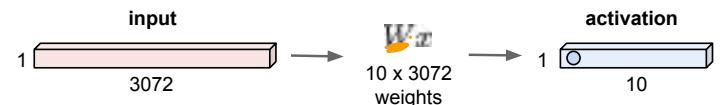
Lecture 5 - 24 April 17, 2018

Convolutional Neural Networks

(First without the brain stuff)

Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 25 April 17, 2018

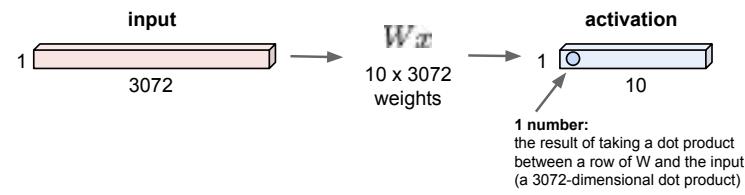
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 26 April 17, 2018

Lecture 5 - 26 April 17, 2018

Fully Connected Layer

32x32x3 image \rightarrow stretch to 3072 x 1

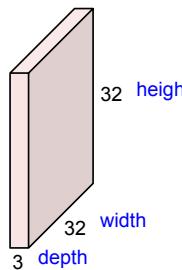


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 27 April 17, 2018

Convolution Layer

32x32x3 image \rightarrow preserve spatial structure

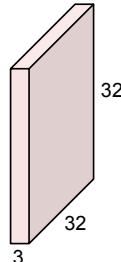


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 28 April 17, 2018

Convolution Layer

32x32x3 image

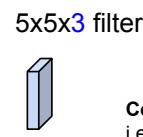
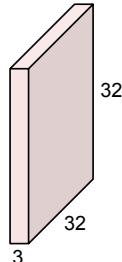


Convolve the filter with the image
i.e. "slide over the image spatially,
computing dot products"

slide

Convolution Layer

32x32x3 image



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"

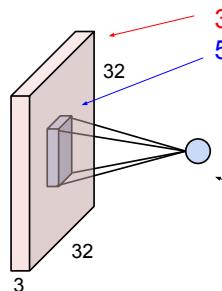
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 29 April 17, 2018

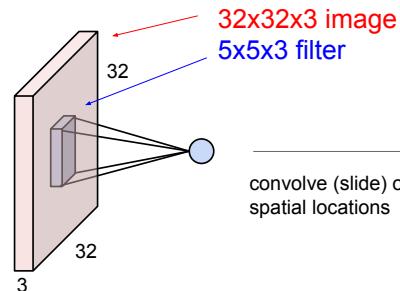
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 30 April 17, 2018

Convolution Layer



Convolution Layer



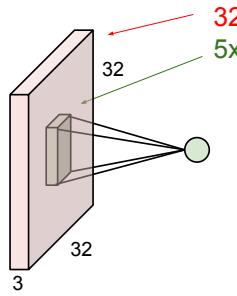
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 31 April 17, 2018

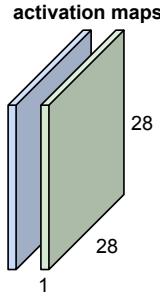
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 32 April 17, 2018

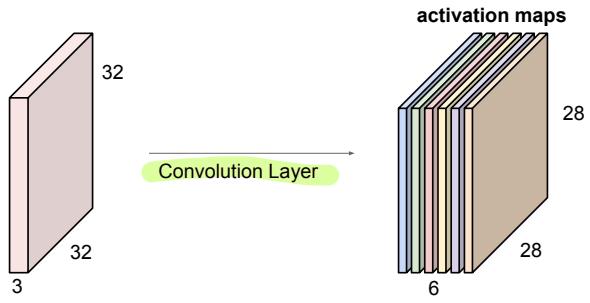
Convolution Layer



consider a second, green filter

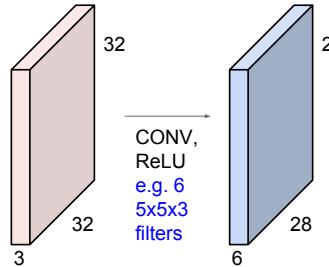


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

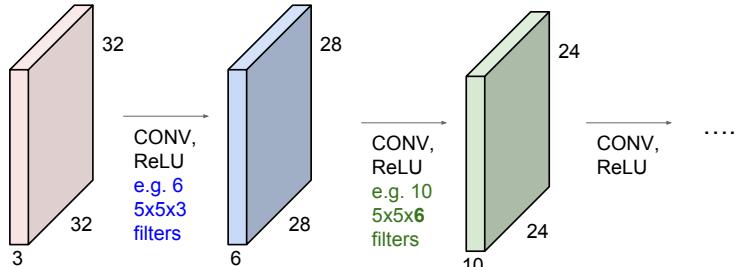


We stack these up to get a “new image” of size 28x28x6!

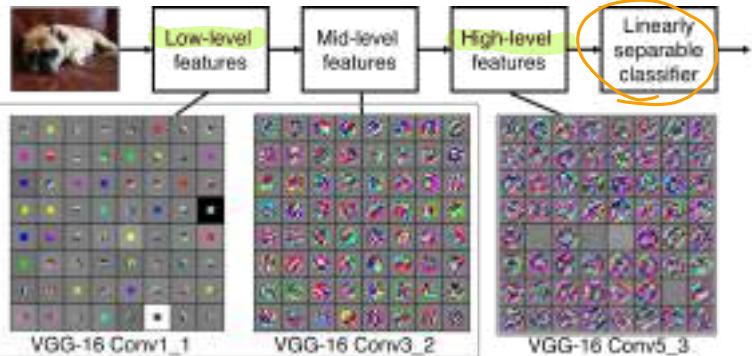
Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



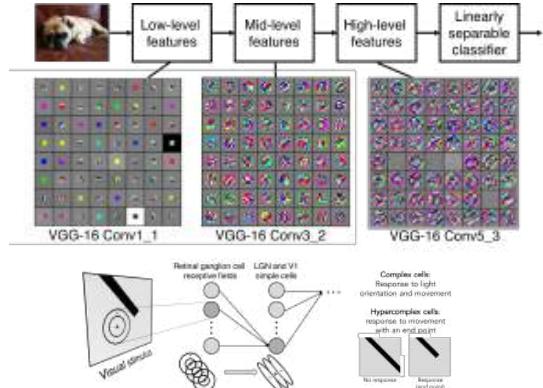
Preview



[Zeiler and Fergus 2013]

Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].

Preview



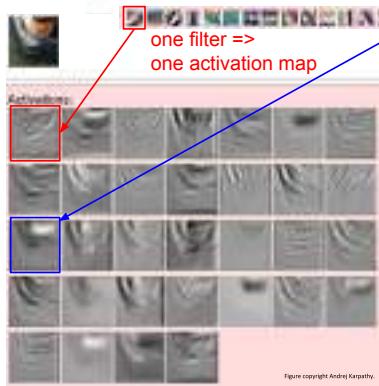


Figure copyright Andrej Karpathy.

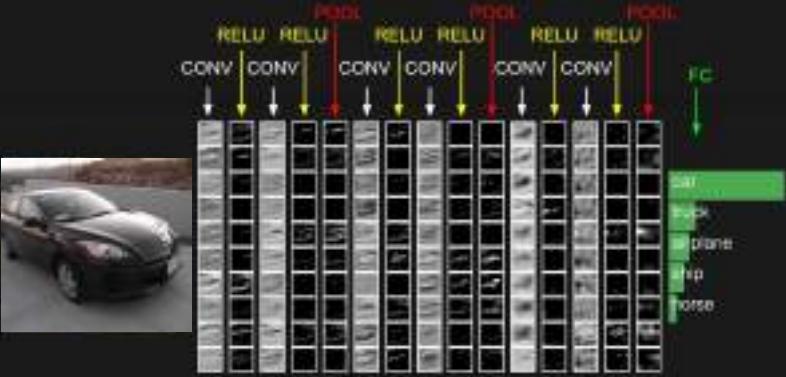
example 5x5 filters
(32 total)

We call the layer convolutional because it is related to convolution of two signals:

$$f(x, y) * g(x, y) = \sum_{i_1} \sum_{i_2} f[i_1, i_2] \cdot g[x - i_1, y - i_2]$$

↑
elementwise multiplication and sum of a filter and the signal (image)

preview:



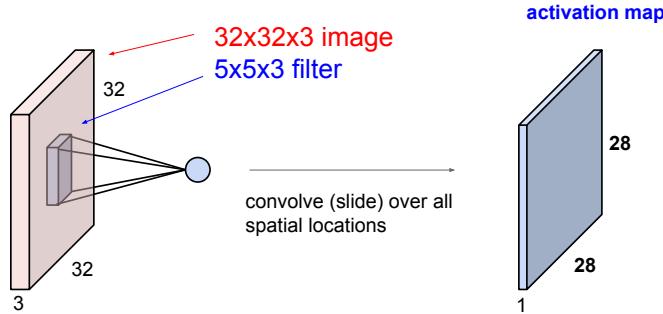
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 39 April 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 40 April 17, 2018

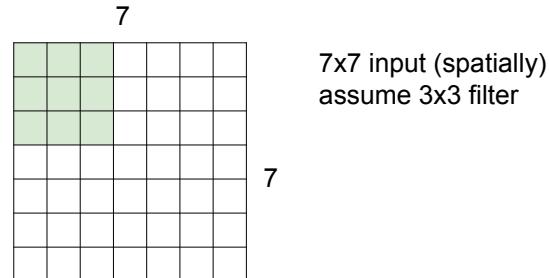
A closer look at spatial dimensions:



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 41 April 17, 2018

A closer look at spatial dimensions:



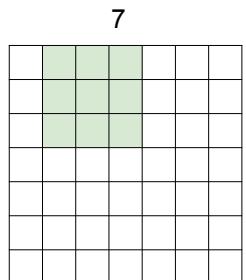
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 41 April 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

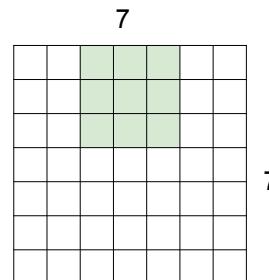
Lecture 5 - 42 April 17, 2018

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

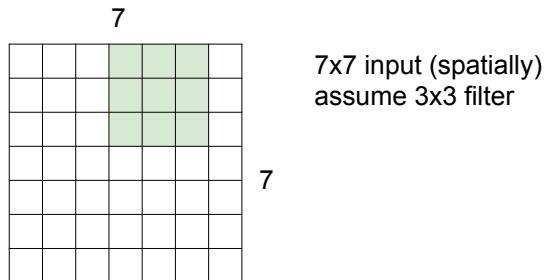
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 43 April 17, 2018

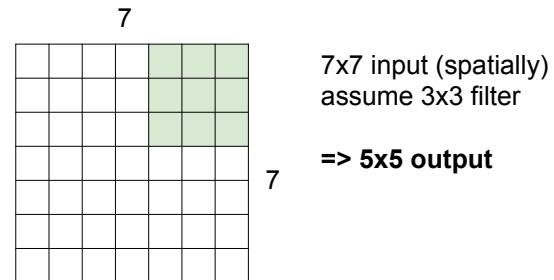
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 44 April 17, 2018

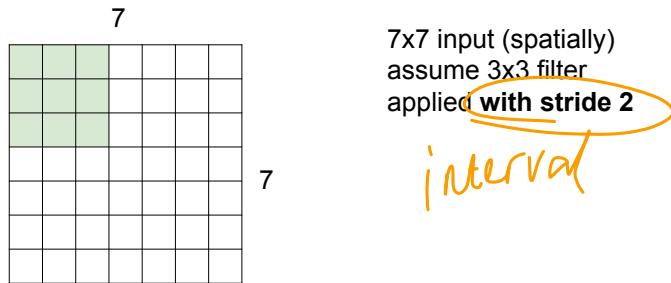
A closer look at spatial dimensions:



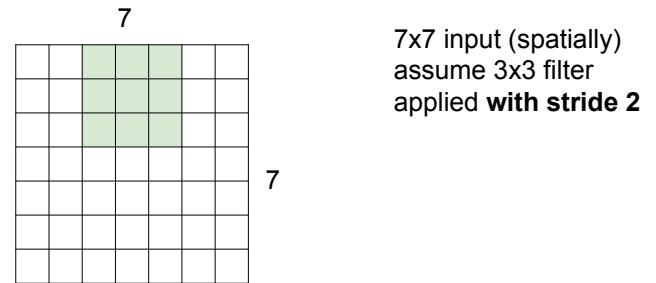
A closer look at spatial dimensions:



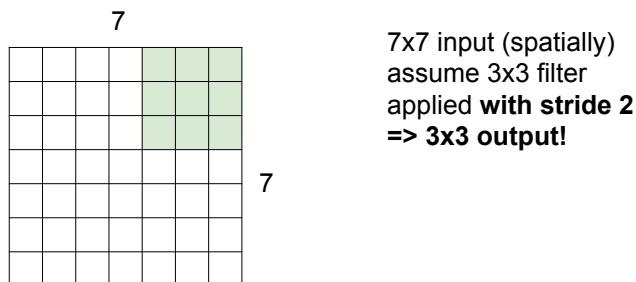
A closer look at spatial dimensions:



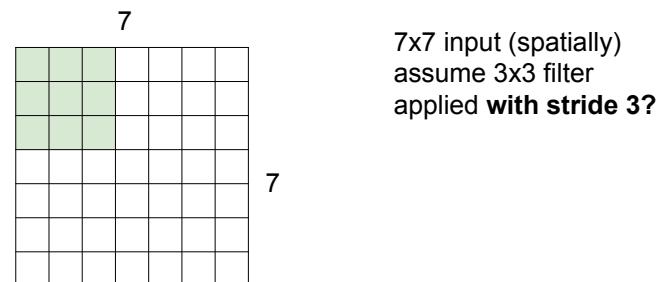
A closer look at spatial dimensions:



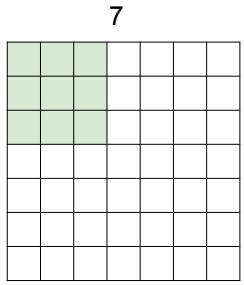
A closer look at spatial dimensions:



A closer look at spatial dimensions:

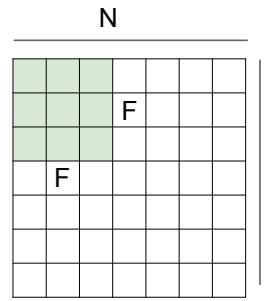


A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied with stride 3?

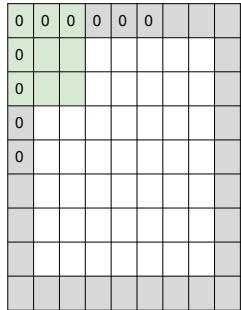
doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.



Output size:
 $(N - F) / \text{stride} + 1$

e.g. N = 7, F = 3:
stride 1 => $(7 - 3)/1 + 1 = 5$
stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33 \dots$

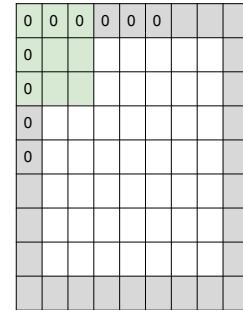
In practice: Common to zero pad the border



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

(recall):
 $(N - F) / \text{stride} + 1$

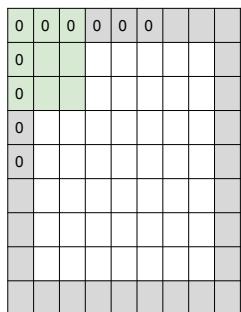
In practice: Common to zero pad the border



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

7x7 output!

In practice: Common to zero pad the border

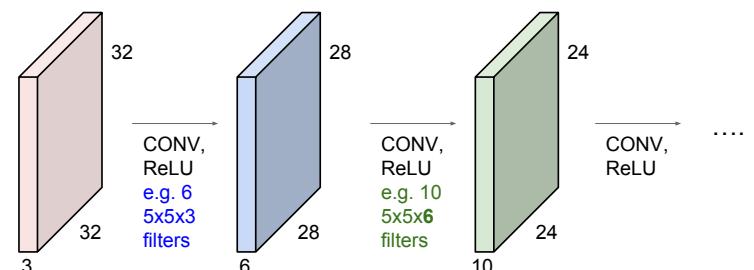


e.g. input 7x7
3x3 filter, applied with stride 1
pad with 1 pixel border => what is the output?

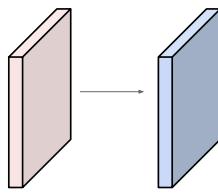
7x7 output!
in general, common to see CONV layers with
stride 1, filters of size FxF, and zero-padding with
(F-1)/2. (will preserve size spatially)
e.g. F = 3 => zero pad with 1
F = 5 => zero pad with 2
F = 7 => zero pad with 3

Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!
(32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.



Examples time:

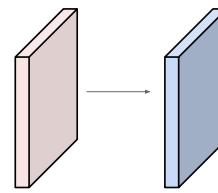


Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size: ?

Examples time:



Input volume: **32x32x3**

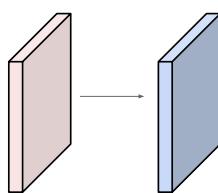
10 5x5 filters with stride 1, pad 2

Output volume size:

$$(32+2 \cdot 2 - 5)/1 + 1 = 32 \text{ spatially, so } \mathbf{32x32x10}$$



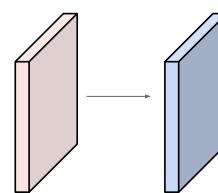
Examples time:



Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Examples time:



Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?

Number of parameters in this layer?
each filter has $5 \cdot 5 \cdot 3 + 1 = 76$ params
 $\Rightarrow 76 \cdot 10 = 760$

(+1 for bias)

Summary: To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K .
 - their spatial extent F .
 - the stride S .
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed usually by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d-th depth slice (if size $W_2 \times H_2$) is the result of performing a valid convolution of the d-th filter over the input volume with a stride of S , and then offset by d-th bias.

Summary: To summarize, the Conv Layer:

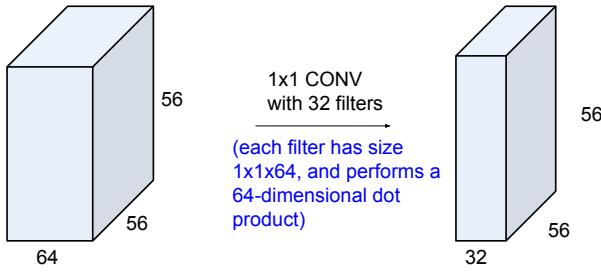
- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K .
 - their spatial extent F .
 - the stride S .
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed usually by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d-th depth slice (if size $W_2 \times H_2$) is the result of performing a valid convolution of the d-th filter over the input volume with a stride of S , and then offset by d-th bias.

Common settings:

$K = (\text{powers of 2, e.g. } 32, 64, 128, 512)$

- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$ (whatever fits)
- $F = 1, S = 1, P = 0$

(btw, 1x1 convolution layers make perfect sense)



Example: CONV layer in Torch

Spatial Convolution

https://github.com/torch/torch7/blob/master/doc/api/torch7/nn/docs/SpatialConvolution.md

SPATIAL CONVOLUTIONS

SPATIAL CONVOLUTIONS ARE 2D CONVOLUTIONS OVER STRIDES OF IMAGE PATCHES. THE IMAGE SIZE IN THE CHANNEL DIMENSION IS PRESERVED AS FOR A 2D CONVOLUTION (IN CONTRAST TO A 1D CONVOLUTION).

The convolution has the following:

- `in_channels`: The number of input channels to the image given to the function.
- `out_channels`: The number of output planes of the convolution layer or feature map.
- `kernel_size`: The kernel size of the convolution.
- `stride`: The stride of the convolution in the spatial dimensions. Default is 1.
- `dilation`: The dilation factor applied to the receptive field. Default is 1. Exponentiate it to get the stride.
- `ceil_mode`: The additional stride added (as well as the padding) to the result planes. Default is 0. (applies to `ceil` mode only).

Additional parameters are available for the convolutional layer, such as the padding, the bias, and the activation function.

Copyright © 2010-2012 The Torch7 Authors. All rights reserved. This code is distributed under the terms of the MIT license.

Example: CONV layer in Caffe

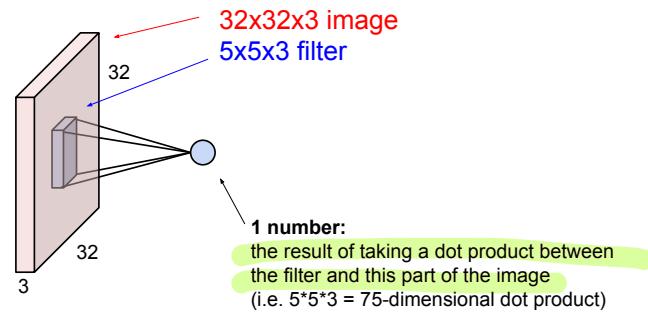
Summary: To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K .
 - their spatial extent F .
 - the stride S .
 - the amount of zero padding P .

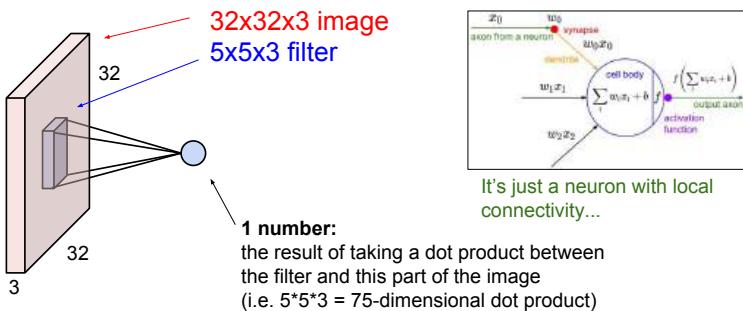
```
layer {
  type: "Convolution"
  bottom: "data"
  top: "conv1"
  # learning rate and decay multipliers for the filters
  param_lr: 1. decay_mult: 3.3
  # learning rate and decay multipliers for the biases
  param_lr: 2. decay_mult: 0.0
  convolution_param {
    num_output: 64
    kernel_size: 3
    stride: 1
    pad: 1
    weight_filler {
      type: "gaussian"
      std: 0.01
      distribution: "uniform"
    }
    bias_filler {
      type: "constant"
      value: 0
    }
  }
}
```

Caffe is licensed under [BSD 2-Clause](#).

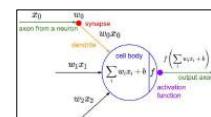
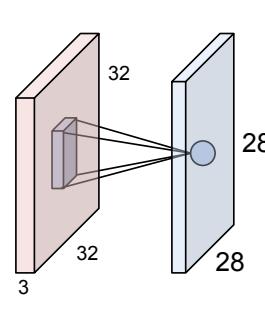
The brain/neuron view of CONV Layer



The brain/neuron view of CONV Layer



The brain/neuron view of CONV Layer

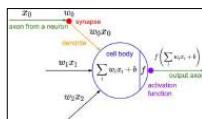
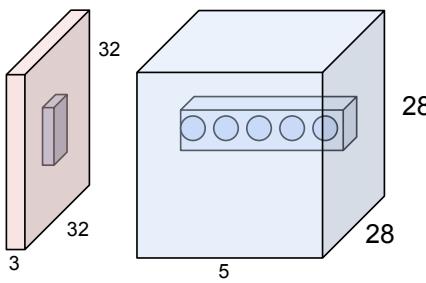


An activation map is a 28×28 sheet of neuron outputs:

1. Each is connected to a small region in the input
2. All of them share parameters

"5x5 filter" \rightarrow "5x5 receptive field for each neuron"

The brain/neuron view of CONV Layer



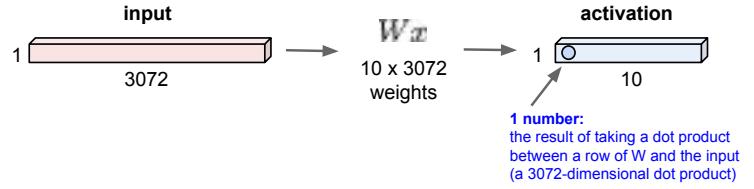
E.g. with 5 filters,
CONV layer consists of
neurons arranged in a 3D grid
($28 \times 28 \times 5$)

There will be 5 different
neurons all looking at the same
region in the input volume

Reminder: Fully Connected Layer

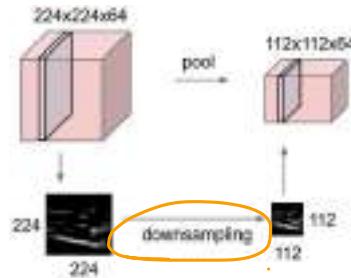
$32 \times 32 \times 3$ image \rightarrow stretch to 3072×1

Each neuron
looks at the full
input volume



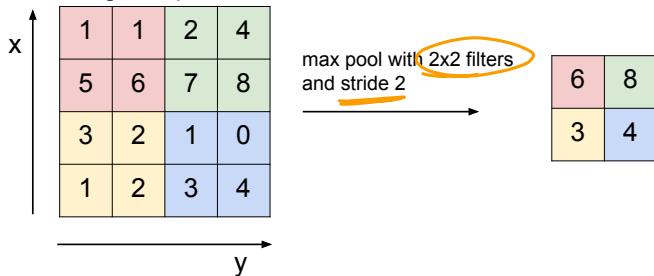
Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:



MAX POOLING

Single depth slice



- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
 - the spatial extent F ,
 - the stride S .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F)/S + 1$
 - $H_2 = (H_1 - F)/S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero padding for Pooling layers

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
 - filter spatial extent F ,
 - the stride S .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F)/S + 1$
 - $H_2 = (H_1 - F)/S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input.
- Note that it is not common to use max pooling for Pooling layers.

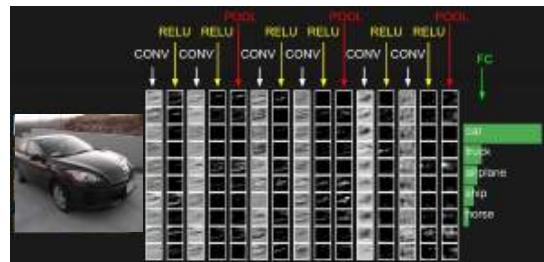
Common settings:

$$F = 2, S = 2$$

$$F = 3, S = 2$$

Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



[ConvNetJS demo: training on CIFAR-10]



<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

Summary

- ConvNets stack CONV,POOL,FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like

$$[(CONV-RELU)^*N-POOL?]^*M-(FC-RELU)^*K,SOFTMAX$$
 where N is usually up to ~5, M is large, $0 \leq K \leq 2$.
 - but recent advances such as ResNet/GoogLeNet challenge this paradigm

Administrative

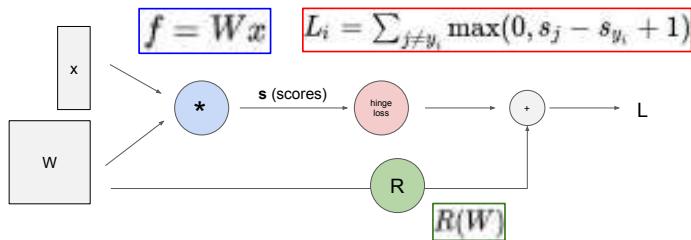
Assignment 1 was due yesterday.

Assignment 2 is out, due Wed May 2. Q5 will be released in a few days, keep an eye out for announcements.

Project proposal due Wed April 25.

Where we are now...

Computational graphs



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 3

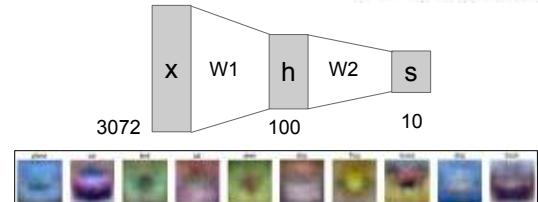
April 19, 2018

Where we are now...

Neural Networks

Linear score function:

2-layer Neural Network

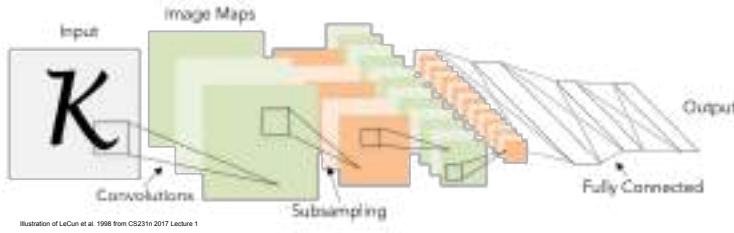


Lecture 6 - 4

April 19, 2018

Where we are now...

Convolutional Neural Networks



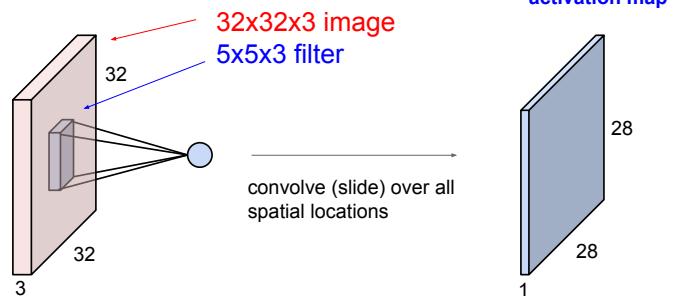
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 5

April 19, 2018

Where we are now...

Convolutional Layer

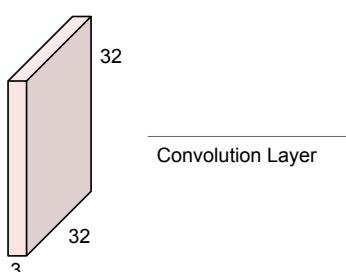


Lecture 6 - 6

April 19, 2018

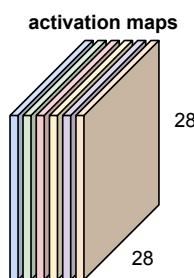
Where we are now...

Convolutional Layer



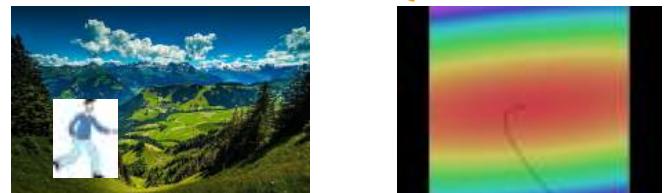
We stack these up to get a "new image" of size 28x28x6!

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



Where we are now...

Learning network parameters through optimization



Landscape image CC0 1.0 public domain
Walking man image CC0 1.0 public domain

```
# Vanilla Gradient Descent
while True:
    weights.grad = evaluate_gradient(loss_fn, data, weights)
    weights += -step_size * weights.grad # perform parameter update
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 7

April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 8

April 19, 2018

Where we are now...

Mini-batch SGD

Loop:

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

Next: Training Neural Networks

Overview

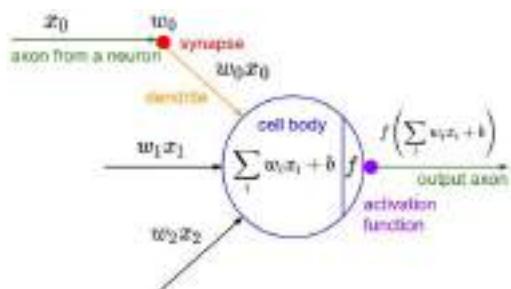
1. **One time setup**
activation functions, preprocessing, weight initialization, regularization, gradient checking
2. **Training dynamics**
babysitting the learning process, parameter updates, hyperparameter optimization
3. **Evaluation**
model ensembles

Part 1

- Activation Functions
- Data Preprocessing
- Weight Initialization
- Batch Normalization
- Babysitting the Learning Process
- Hyperparameter Optimization

Activation Functions

Activation Functions



Activation Functions

Sigmoid



tanh
 $\tanh(x)$



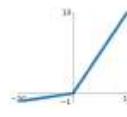
ReLU

$\max(0, x)$



Leaky ReLU

$\max(0.1x, x)$



Maxout

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

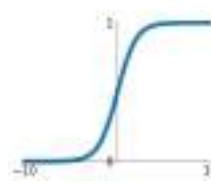
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

Activation Functions

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

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron



Sigmoid

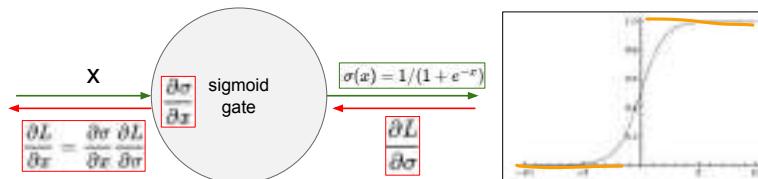
Activation Functions

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

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

3 problems:

1. Saturated neurons "kill" the gradients



What happens when $x = -10$?
What happens when $x = 0$?
What happens when $x = 10$?

Activation Functions

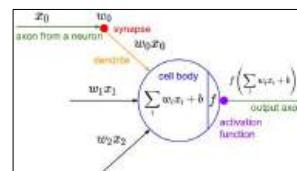
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

3 problems:

1. Saturated neurons "kill" the gradients
2. Sigmoid outputs are not zero-centered

Consider what happens when the input to a neuron (x) is always positive.

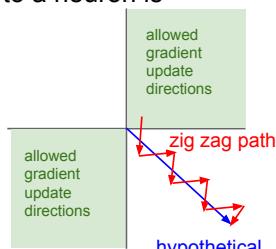


$$f\left(\sum_i w_i x_i + b\right)$$

What can we say about the gradients on w ?

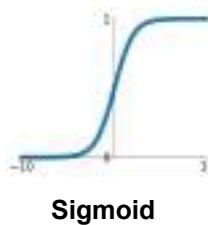
Consider what happens when the input to a neuron is always positive...

$$f\left(\sum_i w_i x_i + b\right)$$



What can we say about the gradients on w ?
Always all positive or all negative :(
(this is also why you want zero-mean data!)

Activation Functions



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

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

3 problems:

1. Saturated neurons "kill" the gradients
2. Sigmoid outputs are not zero-centered
3. `exp()` is a bit compute expensive

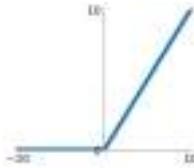
Activation Functions



$\tanh(x)$

- Squashes numbers to range [-1,1]
- zero centered (nice)
- still kills gradients when saturated :(

Activation Functions



ReLU
(Rectified Linear Unit)

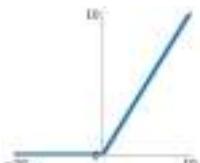
[LeCun et al., 1991]

- Computes $f(x) = \max(0,x)$

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid

[Krizhevsky et al., 2012]

Activation Functions

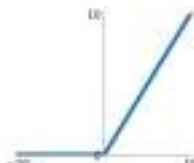


ReLU
(Rectified Linear Unit)

- Computes $f(x) = \max(0,x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid

- Not zero-centered output

Activation Functions



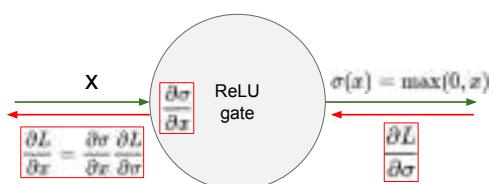
ReLU
(Rectified Linear Unit)

- Computes $f(x) = \max(0,x)$

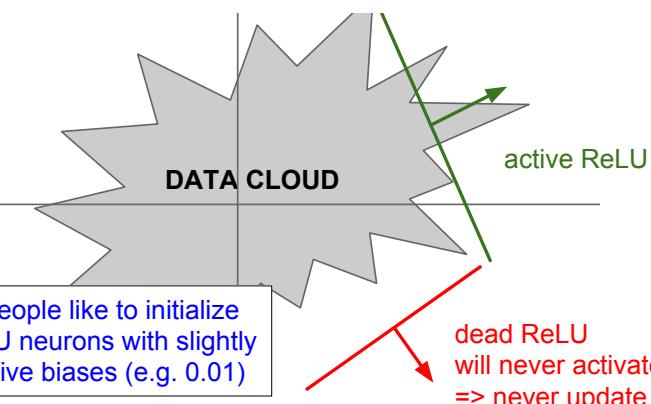
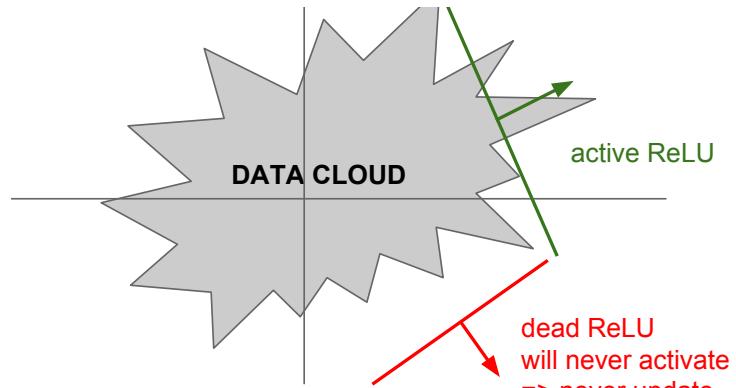
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid

- Not zero-centered output
- An annoyance:

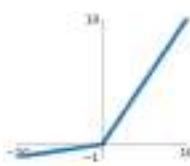
hint: what is the gradient when $x < 0$?



What happens when $x = -10$?
 What happens when $x = 0$?
 What happens when $x = 10$?



Activation Functions



Leaky ReLU

$$f(x) = \max(0.01x, x)$$

[Mass et al., 2013]
 [He et al., 2015]

- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- will not “die”.

Activation Functions

[Mass et al., 2013]
 [He et al., 2015]

- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- will not “die”.

Parametric Rectifier (PReLU)

$$f(x) = \max(\alpha x, x)$$

backprop into \alpha (parameter)

Leaky ReLU

$$f(x) = \max(0.01x, x)$$

Activation Functions

[Clevert et al., 2015]

Exponential Linear Units (ELU)



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU adds some robustness to noise
- Computation requires exp()

[Goodfellow et al., 2013]

Maxout “Neuron”

- Does not have the basic form of dot product -> nonlinearity
- Generalizes ReLU and Leaky ReLU
- Linear Regime! Does not saturate! Does not die!

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

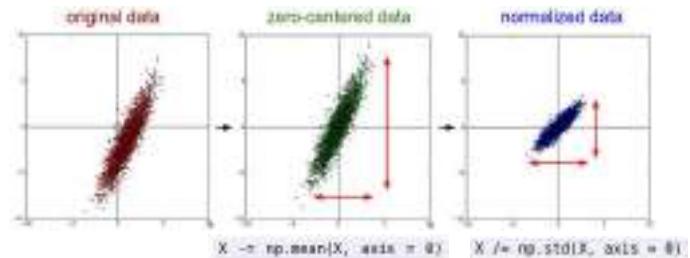
Problem: doubles the number of parameters/neuron :(

TLDR: In practice:

- Use ReLU. Be careful with your learning rates
- Try out Leaky ReLU / Maxout / ELU
- Try out tanh but don't expect much
- Don't use sigmoid

Data Preprocessing

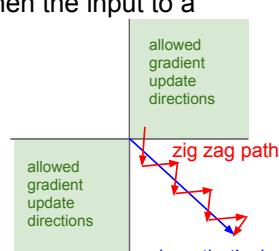
Step 1: Preprocess the data



(Assume X [NxD] is data matrix,
each example in a row)

Remember: Consider what happens when the input to a neuron is always positive...

$$f\left(\sum_i w_i x_i + b\right)$$

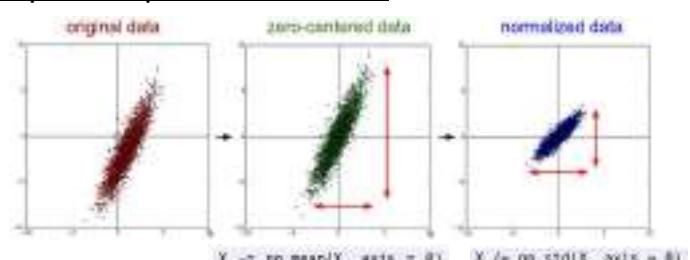


What can we say about the gradients on w ?

Always all positive or all negative :(

(this is also why you want zero-mean data!)

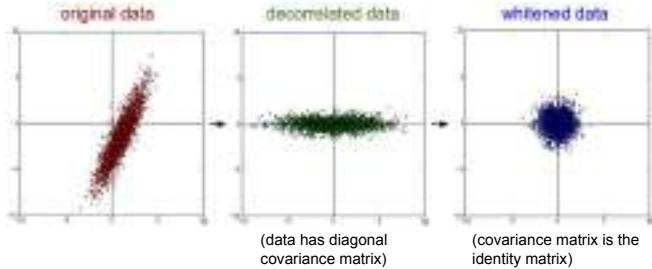
Step 1: Preprocess the data



(Assume X [NxD] is data matrix,
each example in a row)

Step 1: Preprocess the data

In practice, you may also see **PCA** and **Whitening** of the data



TLDR: In practice for Images: center only

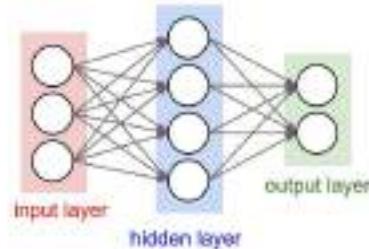
e.g. consider CIFAR-10 example with [32,32,3] images

- Subtract the mean image (e.g. AlexNet)
(mean image = [32,32,3] array)
- Subtract per-channel mean (e.g. VGGNet)
(mean along each channel = 3 numbers)

Not common to normalize variance, to do PCA or whitening

Weight Initialization

- Q: what happens when W=constant init is used?



- First idea: **Small random numbers**
(gaussian with zero mean and 1e-2 standard deviation)

```
W = 0.01* np.random.randn(D,H)
```

- First idea: **Small random numbers**
(gaussian with zero mean and 1e-2 standard deviation)

```
W = 0.01* np.random.randn(D,H)
```

Works ~okay for small networks, but problems with deeper networks.

Lets look at
some
activation
statistics

E.g. 10-layer net with 500 neurons on each layer, using tanh non-linearities, and initializing as described in last slide

Fei-Fei Li & Justin Johnson & Serena Yeung

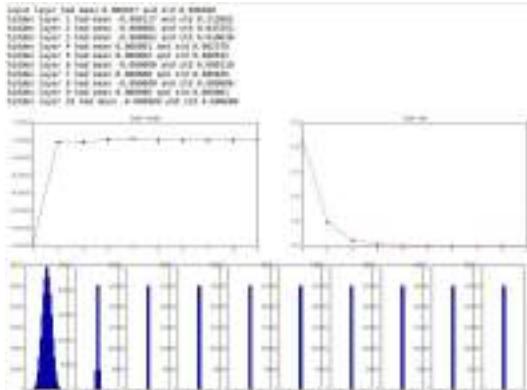
Lecture 6 - 45

April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 46

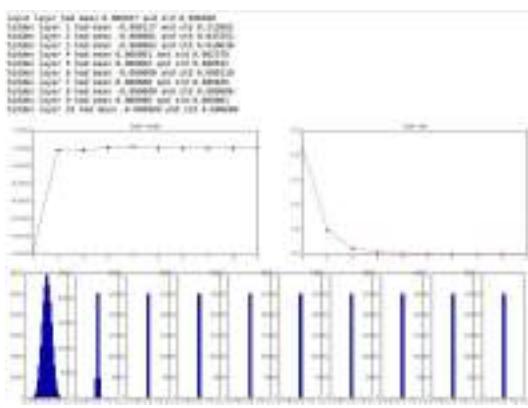
April 19, 2018



All activations
become zero!

Q: think about the backward pass.
What do the gradients look like?

Hint: think about backward pass for a W^*X gate.



Fei-Fei Li & Justin Johnson & Serena Yeung

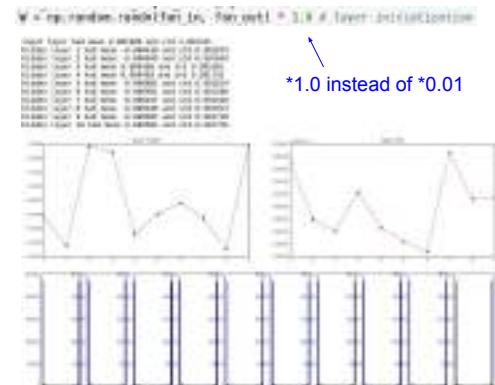
Lecture 6 - 47

April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 48

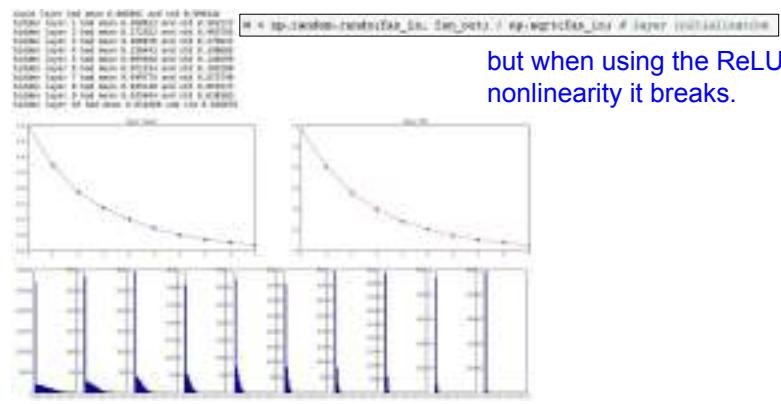
April 19, 2018



Almost all neurons completely saturated, either -1 and 1. Gradients will be all zero.

“Xavier initialization” [Glorot et al., 2010]

Reasonable initialization. (Mathematical derivation assumes linear activations)



but when using the ReLU nonlinearity it breaks.

Fei-Fei Li & Justin Johnson & Serena Yeung

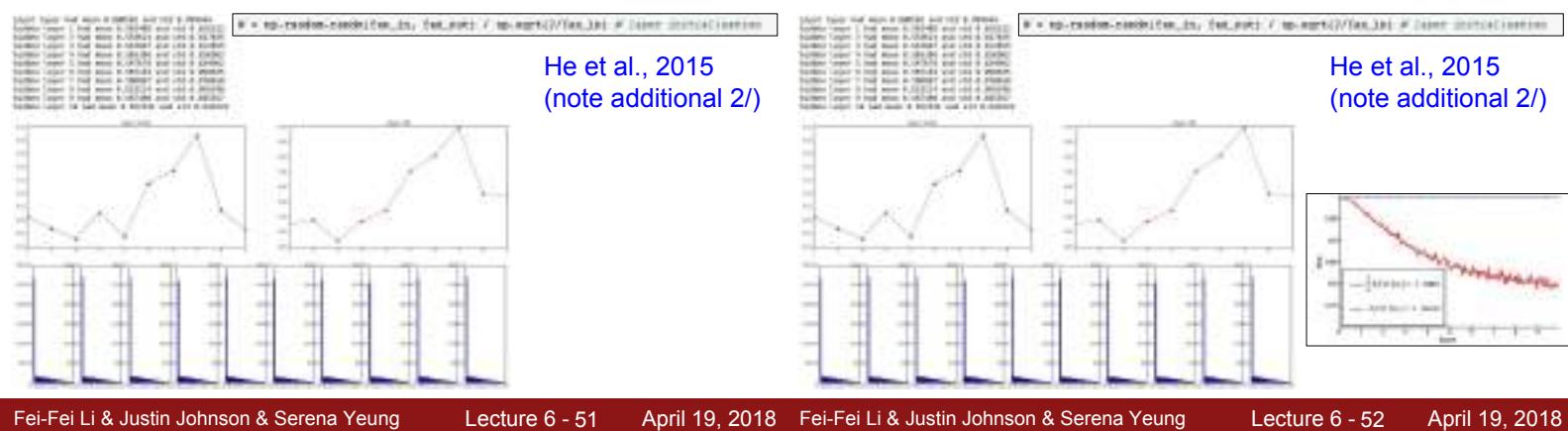
Lecture 6 - 49

April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 50

April 19, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 51

April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 52

April 19, 2018

Proper initialization is an active area of research...

Understanding the difficulty of training deep feedforward neural networks
by Glorot and Bengio, 2010

Exact solutions to the nonlinear dynamics of learning in deep linear neural networks
by Saxe et al, 2013

Random walk initialization for training very deep feedforward networks
by Sussillo and Abbott, 2014

Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification
by He et al., 2015

Data-dependent Initializations of Convolutional Neural Networks
by Krähenbühl et al., 2015

All you need is a good init, Mishkin and Matas, 2015

...

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 53

April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 54

April 19, 2018

Batch Normalization

Batch Normalization

[Ioffe and Szegedy, 2015]

“you want zero-mean unit-variance activations? just make them so.”

consider a batch of activations at some layer. To make each dimension zero-mean unit-variance, apply:

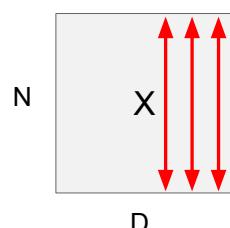
$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{Var[x^{(k)}]}}$$

this is a vanilla
differentiable function...

Batch Normalization

[Ioffe and Szegedy, 2015]

“you want zero-mean unit-variance activations? just make them so.”



1. compute the empirical mean and variance independently for each dimension.

2. Normalize

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{Var[x^{(k)}]}}$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 55

April 19, 2018

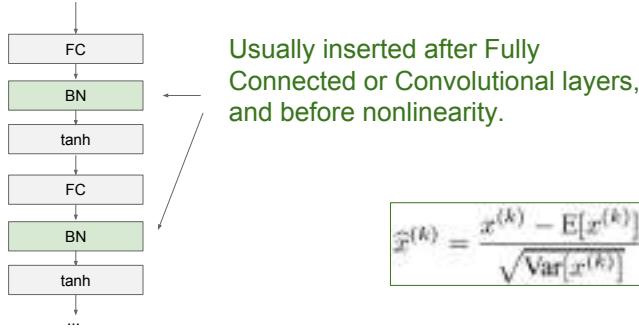
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 56

April 19, 2018

Batch Normalization

[Ioffe and Szegedy, 2015]



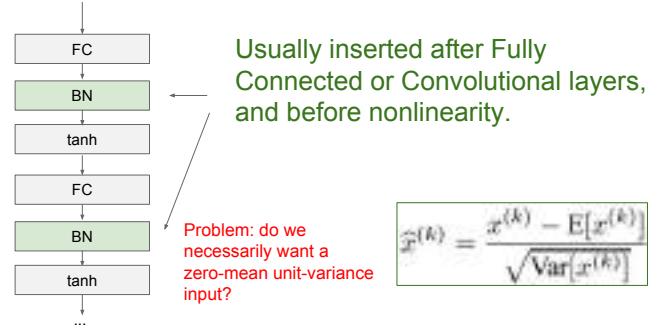
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 57

April 19, 2018

Batch Normalization

[Ioffe and Szegedy, 2015]



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 58

April 19, 2018

Batch Normalization

[Ioffe and Szegedy, 2015]

Normalize:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

And then allow the network to squash the range if it wants to:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

Note, the network can learn:
 $\gamma^{(k)} = \sqrt{\text{Var}[x^{(k)}]}$
 $\beta^{(k)} = \mathbb{E}[x^{(k)}]$
 to recover the identity mapping.

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: Values of x over a mini-batch: $B = \{x_1, \dots, x_m\}$
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

```

 $\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$  // mini-batch mean
 $\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$  // mini-batch variance
 $\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$  // normalize
 $y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$  // scale and shift
    
```

- Improves gradient flow through the network
- Allows higher learning rates
- Reduces the strong dependence on initialization
- Acts as a form of regularization in a funny way, and slightly reduces the need for dropout, maybe

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 59

April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 60

April 19, 2018

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: Values of x over a mini-batch: $B = \{x_1, \dots, x_m\}$
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

```

 $\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$  // mini-batch mean
 $\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$  // mini-batch variance
 $\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$  // normalize
 $y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$  // scale and shift
    
```

Note: at test time BatchNorm layer functions differently:

The mean/std are not computed based on the batch. Instead, a single fixed empirical mean of activations during training is used.

(e.g. can be estimated during training with running averages)

Babysitting the Learning Process

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 61

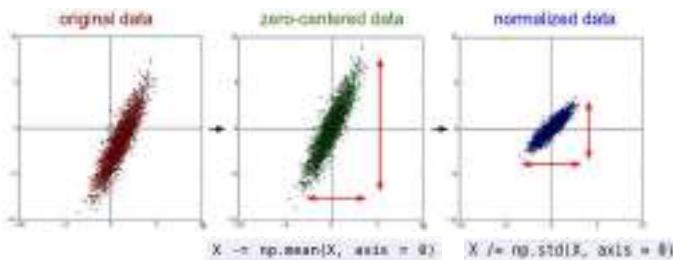
April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 62

April 19, 2018

Step 1: Preprocess the data



(Assume X [NxD] is data matrix,
each example in a row)

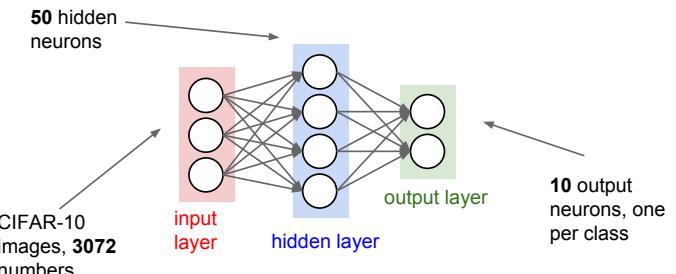
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 63

April 19, 2018

April 19, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 6 - 64 April 19, 2018

Step 2: Choose the architecture:
say we start with one hidden layer of 50 neurons:



Double check that the loss is reasonable:

Double check that the loss is reasonable:

```

def init_two_layer_model(input_size, hidden_size, output_size):
    np.random.seed(42)
    model = {}
    model["W1"] = 0.0001 * np.random.randn(input_size, hidden_size)
    model["b1"] = np.zeros(hidden_size)
    model["W2"] = 0.0001 * np.random.randn(hidden_size, output_size)
    model["b2"] = np.zeros(output_size)
    return model

```

```

def init_two_layer_model(input_size, hidden_size, output_size):
    np.random.seed(0)
    model = {}
    model['W1'] = np.random.randn(input_size, hidden_size)
    model['b1'] = np.zeros(hidden_size)
    model['W2'] = np.random.randn(hidden_size, output_size)
    model['b2'] = np.zeros(output_size)
    return model

```

```
model = init two-layer neural(1000000, 10, 30 + len(classnames), 1000000, num_of_classes)
loss, grad = two_layer_net(X_train, model, y_train, 0.01)    crank up regularization
print loss
3.86859726483


loss went up, good, (sanity check)


```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 65

April 19, 2018

April 19, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 6 - 66 April 19, 2018

Lets try to train now...

Tip: Make sure that you can overfit very small portion of the training data

The above code:

- take the first 20 examples from CIFAR-10
 - turn off regularization ($\text{reg} = 0.0$)
 - use simple vanilla 'sgd'

Lets try to train now...

Tip: Make sure that you can overfit very small portion of the training data

Very small loss,
train accuracy 1.00,
nice!

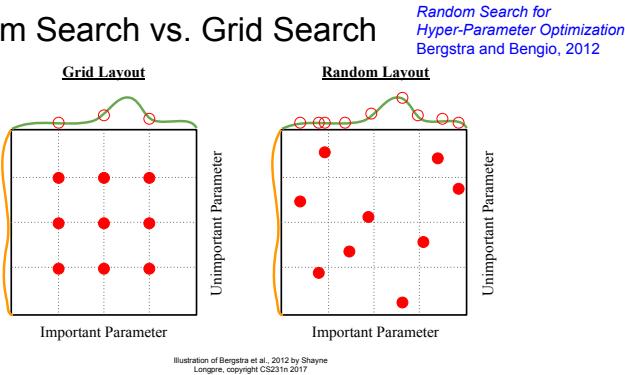
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 67

April 19, 2018

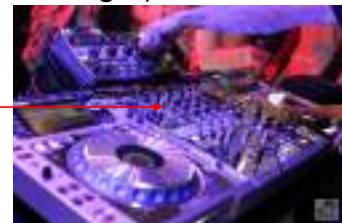
Lecture 6 - 68 April 19, 2018

Random Search vs. Grid Search



Hyperparameters to play with:

- network architecture
- learning rate, its decay schedule, update type
- regularization (L2/Dropout strength)

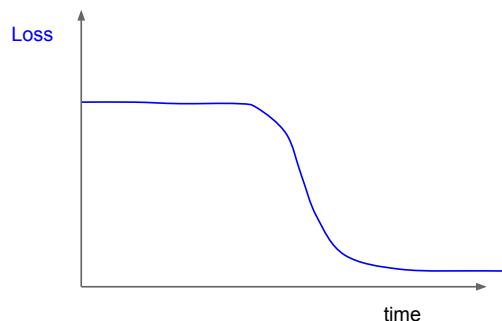
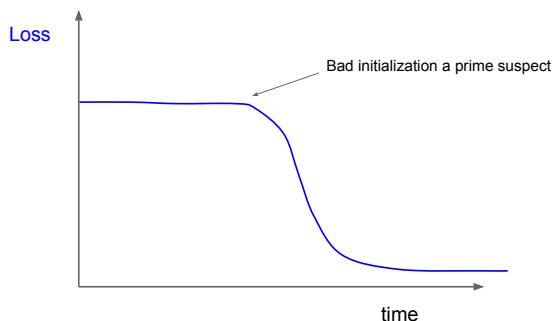
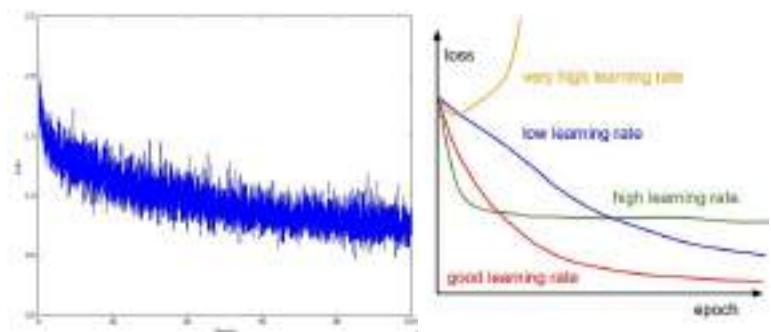


This image by Paolo Guerri is licensed under CC BY 2.0

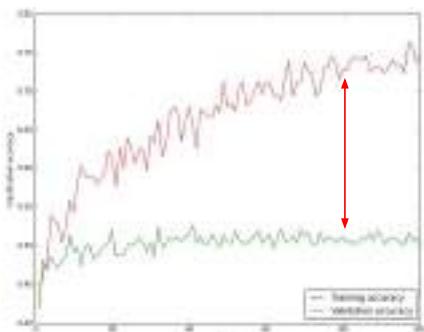
Cross-validation
“command center”



Monitor and visualize the loss curve



Monitor and visualize the accuracy:



big gap = overfitting
=> increase regularization strength?

no gap
=> increase model capacity?

Track the ratio of weight updates / weight magnitudes:

```
# assume parameter vector W and its gradient vector dW
param_scale = np.linalg.norm(W.ravel())
update = -learning_rate*dW # simple SGD update
update_scale = np.linalg.norm(update.ravel())
W += update # the actual update
print update_scale / param_scale # want ~1e-3
```

ratio between the updates and values: $\sim 0.0002 / 0.02 = 0.01$ (about okay)
want this to be somewhere around 0.001 or so

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 87

April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 88

April 19, 2018

Summary

We looked in detail at:

- Activation Functions ([use ReLU](#))
- Data Preprocessing ([images: subtract mean](#))
- Weight Initialization ([use Xavier/He init](#))
- Batch Normalization ([use](#))
- Babysitting the Learning process
- Hyperparameter Optimization
([random sample hyperparams, in log space when appropriate](#))

TLDRs

Next time:

Training Neural Networks, Part 2

- Parameter update schemes
- Learning rate schedules
- Gradient checking
- Regularization (Dropout etc.)
- Evaluation (Ensembles etc.)
- Transfer learning / fine-tuning

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 89

April 19, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 6 - 90

April 19, 2018

Administrative

- Assignment 1 is being graded, stay tuned
- Project proposals due tomorrow by 11:59pm on Gradescope
- Assignment 2 is out, due Wednesday 5/2 11:59pm

Lecture 7: Training Neural Networks, Part 2

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 1

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

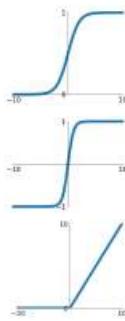
Lecture 7 - 2

April 24, 2018

Last time: Activation Functions

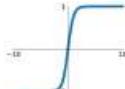
Sigmoid

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



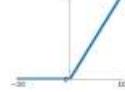
tanh

$$\tanh(x)$$



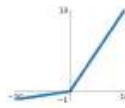
ReLU

$$\max(0, x)$$



Leaky ReLU

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



Maxout

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

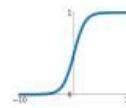
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

Last time: Activation Functions

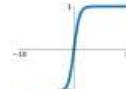
Sigmoid

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



tanh

$$\tanh(x)$$



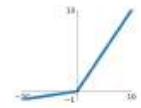
ReLU

$$\max(0, x)$$

Good default choice

Leaky ReLU

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



Maxout

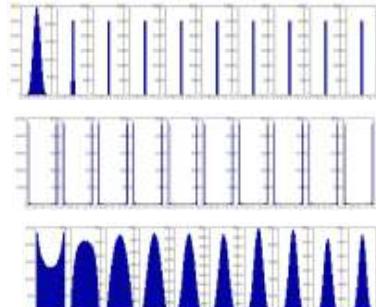
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Last time: Weight Initialization



Initialization too small:

Activations go to zero, gradients also zero,
No learning

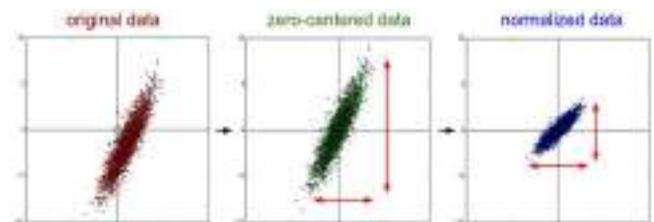
Initialization too big:

Activations saturate (for tanh),
Gradients zero, no learning

Initialization just right:

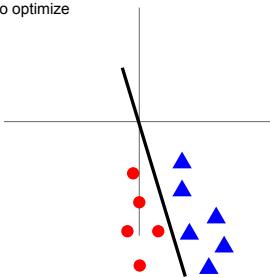
Nice distribution of activations at all layers,
Learning proceeds nicely

Last time: Data Preprocessing

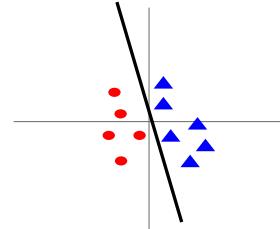


Last time: Data Preprocessing

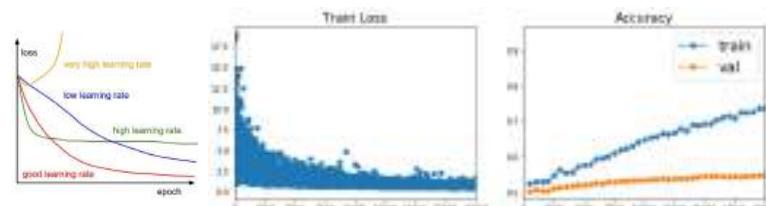
Before normalization: classification loss
very sensitive to changes in weight matrix;
hard to optimize



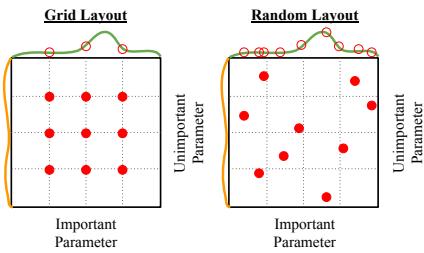
After normalization: less sensitive to small
changes in weights; easier to optimize



Last time: Babysitting Learning



Last time: Hyperparameter Search



Coarse to fine search

```

wsl_001: 0.003069, lr: 1.49720e-05, reg: 0.79794e-03, 13 / 3801
wsl_002: 0.003069, lr: 2.10087e-05, reg: 0.81510e-03, 13 / 3801
wsl_003: 0.003069, lr: 3.39113e-05, reg: 4.17890e-05, 14 / 3801
wsl_004: 0.003069, lr: 4.25039e-05, reg: 0.80110e-03, 14 / 3801
wsl_005: 0.003069, lr: 5.00000e-05, reg: 0.80110e-03, 14 / 3801
wsl_006: 0.003069, lr: 1.99751e-05, reg: 0.28283e-03, 15 / 3801
wsl_007: 0.003069, lr: 4.42180e-05, reg: 1.09820e-04, 16 / 3801
wsl_008: 0.003069, lr: 3.65090e-05, reg: 4.51920e-03, 15 / 3801

```

Today

- More normalization
- Fancier optimization
- Regularization
- Transfer Learning

Last time: Batch Normalization

Input: $x : N \times D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Learnable params:

$$\gamma, \beta : D$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

Intermediates: $\mu, \sigma : D$
 $\hat{x} : N \times D$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Output: $y : N \times D$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Last time: Batch Normalization

Input: $x : N \times D$

Learnable params:

$$\gamma, \beta : D$$

Intermediates: $\mu, \sigma : D$
 $\hat{x} : N \times D$

Output: $y : N \times D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Batch Normalization: Test Time

Input: $x : N \times D$

μ_j = (Running) average of values seen during training

Learnable params:

$$\gamma, \beta : D$$

σ_j^2 = (Running) average of values seen during training

Intermediates: $\mu, \sigma : D$
 $\hat{x} : N \times D$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Output: $y : N \times D$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Batch Normalization for ConvNets

Batch Normalization for
fully-connected networks

x: $N \times D$

Normalize

$$\mu, \sigma: 1 \times D$$

$$\gamma, \beta: 1 \times D$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

Batch Normalization for
convolutional networks
(Spatial Batchnorm, BatchNorm2D)

x: $N \times C \times H \times W$

Normalize

$$\mu, \sigma: 1 \times C \times 1 \times 1$$

$$\gamma, \beta: 1 \times C \times 1 \times 1$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

Layer Normalization

Batch Normalization for fully-connected networks

$$\begin{array}{l} \mathbf{x}: N \times D \\ \text{Normalize} \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma}: 1 \times D \\ \mathbf{y}, \beta: 1 \times D \\ \mathbf{y} = \mathbf{y}(\mathbf{x}-\boldsymbol{\mu}) / \boldsymbol{\sigma} + \beta \end{array}$$

Layer Normalization for fully-connected networks
Same behavior at train and test!
Can be used in recurrent networks

$$\begin{array}{l} \mathbf{x}: N \times D \\ \text{Normalize} \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma}: N \times 1 \\ \mathbf{y}, \beta: 1 \times D \\ \mathbf{y} = \mathbf{y}(\mathbf{x}-\boldsymbol{\mu}) / \boldsymbol{\sigma} + \beta \end{array}$$

Instance Normalization

Batch Normalization for convolutional networks

$$\begin{array}{l} \mathbf{x}: N \times C \times H \times W \\ \text{Normalize} \quad \downarrow \quad \downarrow \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma}: 1 \times C \times 1 \times 1 \\ \mathbf{y}, \beta: 1 \times C \times 1 \times 1 \\ \mathbf{y} = \mathbf{y}(\mathbf{x}-\boldsymbol{\mu}) / \boldsymbol{\sigma} + \beta \end{array}$$

Instance Normalization for convolutional networks
Same behavior at train / test!

$$\begin{array}{l} \mathbf{x}: N \times C \times H \times W \\ \text{Normalize} \quad \downarrow \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma}: N \times C \times 1 \times 1 \\ \mathbf{y}, \beta: 1 \times C \times 1 \times 1 \\ \mathbf{y} = \mathbf{y}(\mathbf{x}-\boldsymbol{\mu}) / \boldsymbol{\sigma} + \beta \end{array}$$

Ba, Kiros, and Hinton, "Layer Normalization", arXiv 2016

Lecture 7 - 15

April 24, 2018

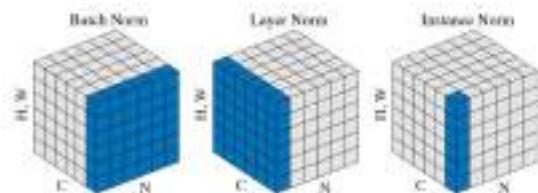
Ulyanov et al, Improved Texture Networks: Maximizing Quality and Diversity in Feed-forward Stylization and Texture Synthesis, CVPR 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 16

April 24, 2018

Comparison of Normalization Layers

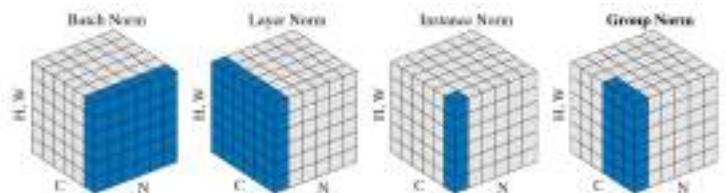


Wu and He, "Group Normalization", arXiv 2018

Lecture 7 - 17

April 24, 2018

Group Normalization



Wu and He, "Group Normalization", arXiv 2018 (Appeared 3/22/2018)

Lecture 7 - 17

April 24, 2018

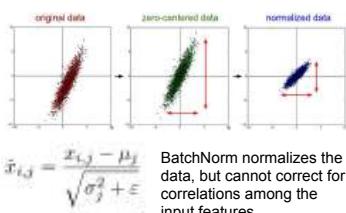
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 18

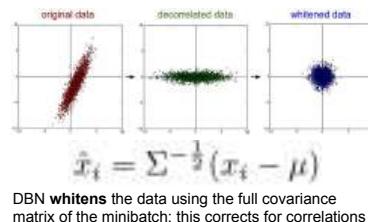
April 24, 2018

Decorrelated Batch Normalization

Batch Normalization



Decorrelated Batch Normalization



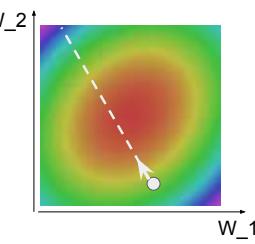
Huang et al, "Decorrelated Batch Normalization", arXiv 2018 (Appeared 4/23/2018)

Lecture 7 -

April 24, 2018

Optimization

```
# Vanilla Gradient-Descent
while True:
    weights_grad = evaluate_gradient(loss_fn, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```



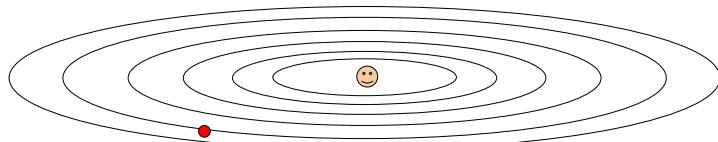
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 20

April 24, 2018

Optimization: Problems with SGD

What if loss changes quickly in one direction and slowly in another?
What does gradient descent do?



Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

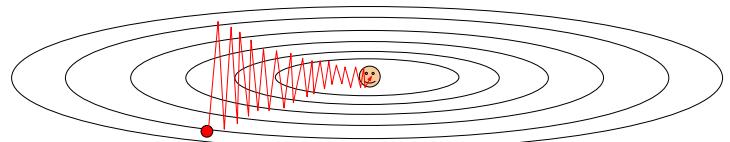
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 21

April 24, 2018

Optimization: Problems with SGD

What if loss changes quickly in one direction and slowly in another?
What does gradient descent do?
Very slow progress along shallow dimension, jitter along steep direction



Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

Lecture 7 - 22

April 24, 2018

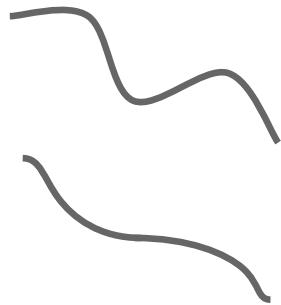
Fei-Fei Li & Justin Johnson & Serena Yeung

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 21 April 24, 2018

Optimization: Problems with SGD

What if the loss function has a **local minima or saddle point**?



Fei-Fei Li & Justin Johnson & Serena Yeung

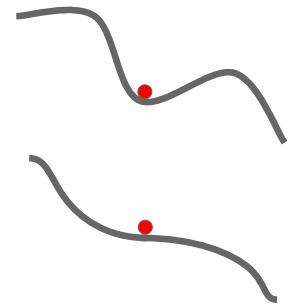
Lecture 7 - 23

April 24, 2018

Optimization: Problems with SGD

What if the loss function has a **local minima or saddle point**?

Zero gradient, gradient descent gets stuck



Lecture 7 - 24

April 24, 2018

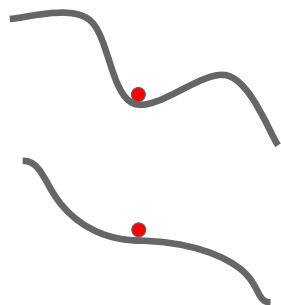
Fei-Fei Li & Justin Johnson & Serena Yeung

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 23 April 24, 2018

Optimization: Problems with SGD

What if the loss function has a **local minima or saddle point**?



Saddle points much more common in high dimension

Dauphin et al., "Identifying and attacking the saddle point problem in high-dimensional non-convex optimization", NIPS 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 25

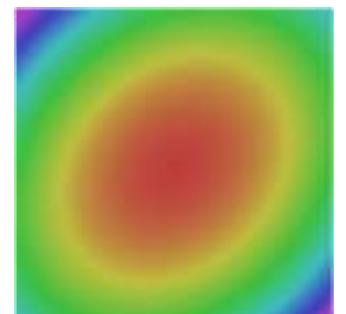
April 24, 2018

Optimization: Problems with SGD

Our gradients come from minibatches so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W)$$



Lecture 7 - 26

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 25 April 24, 2018

SGD + Momentum

SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
while True:
    dx = compute_gradient(x)
    x -= learning_rate * dx
```

- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically rho=0.9 or 0.99

SGD+Momentum

$$\begin{aligned} v_{t+1} &= \rho v_t + \nabla f(x_t) \\ x_{t+1} &= x_t - \alpha v_{t+1} \end{aligned}$$

```
vx = B
while True:
    dx = compute_gradient(x)
    vx = rho * vx + dx
    x = learning_rate * vx + vx
```

SGD + Momentum

SGD+Momentum

$$\begin{aligned} v_{t+1} &= \rho v_t - \alpha \nabla f(x_t) \\ x_{t+1} &= x_t + v_{t+1} \end{aligned}$$

```
vx = B
while True:
    dx = compute_gradient(x)
    vx = rho * vx - learning_rate * dx
    x += vx
```

- You may see SGD+Momentum formulated different ways, but they are equivalent - give same sequence of x

SGD+Momentum

$$\begin{aligned} v_{t+1} &= \rho v_t + \nabla f(x_t) \\ x_{t+1} &= x_t - \alpha v_{t+1} \end{aligned}$$

```
vx = B
while True:
    dx = compute_gradient(x)
    vx = rho * vx + dx
    x = learning_rate * vx + vx
```

Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013

Lecture 7 - 27

April 24, 2018

Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013

Lecture 7 - 28

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Fei-Fei Li & Justin Johnson & Serena Yeung

April 24, 2018

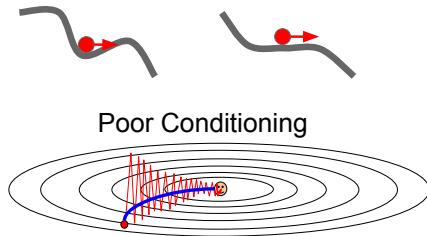
Fei-Fei Li & Justin Johnson & Serena Yeung

Fei-Fei Li & Justin Johnson & Serena Yeung

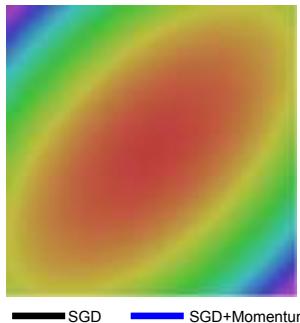
April 24, 2018

SGD + Momentum

Local Minima Saddle points

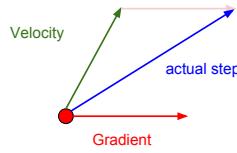


Gradient Noise



SGD+Momentum

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

Nesterov, “A method of solving a convex programming problem with convergence rate $O(1/k^2)$ ”, 1983

Nesterov, “Introductory lectures on convex optimization: a basic course”, 2004

Sutskever et al, “On the importance of initialization and momentum in deep learning”, ICML 2013

Fei-Fei Li & Justin Johnson & Serena Yeung

Fei-Fei Li & Justin Johnson & Serena Yeung

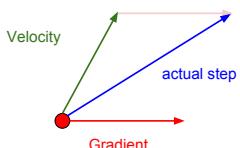
April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

April 24, 2018

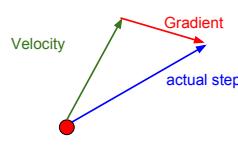
Nesterov Momentum

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

Nesterov Momentum

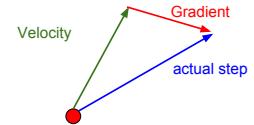


“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$



“Look ahead” to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Fei-Fei Li & Justin Johnson & Serena Yeung

Fei-Fei Li & Justin Johnson & Serena Yeung

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Fei-Fei Li & Justin Johnson & Serena Yeung

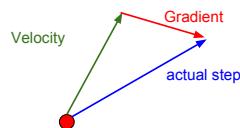
April 24, 2018

Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

Annoying, usually we want update in terms of $x_t, \nabla f(x_t)$



"Look ahead" to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

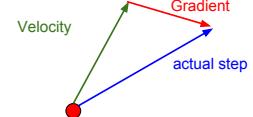
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 33

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Annoying, usually we want update in terms of $x_t, \nabla f(x_t)$



Change of variables $\tilde{x}_t = x_t + \rho v_t$ and rearrange:

$$v_{t+1} = \rho v_t - \alpha \nabla f(\tilde{x}_t)$$

$$\tilde{x}_{t+1} = \tilde{x}_t - \rho v_t + (1 + \rho)v_{t+1}$$

$$= \tilde{x}_t + v_{t+1} + \rho(v_{t+1} - v_t)$$

"Look ahead" to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Lecture 7 - 34 April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 33

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 34 April 24, 2018

Nesterov Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

Annoying, usually we want update in terms of $x_t, \nabla f(x_t)$

Change of variables $\tilde{x}_t = x_t + \rho v_t$ and rearrange:

$$v_{t+1} = \rho v_t - \alpha \nabla f(\tilde{x}_t)$$

$$\tilde{x}_{t+1} = \tilde{x}_t - \rho v_t + (1 + \rho)v_{t+1}$$

$$= \tilde{x}_t + v_{t+1} + \rho(v_{t+1} - v_t)$$

`dx = compute_gradient(x)`

`old_v = v`

`v = rho * v - learning_rate * dx`

`x = -rho * old_v + (1 + rho) * v`

Fei-Fei Li & Justin Johnson & Serena Yeung

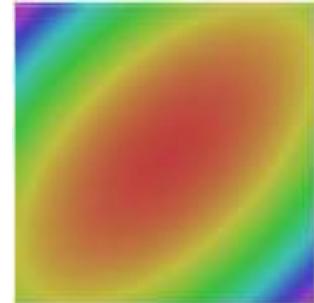
Lecture 7 - 35

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 36 April 24, 2018

Nesterov Momentum



— SGD
— SGD+Momentum
— Nesterov

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 35

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 36 April 24, 2018

AdaGrad

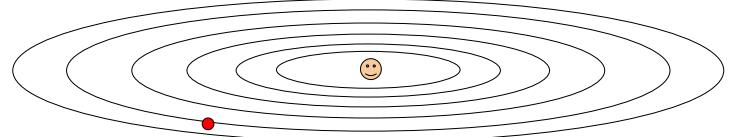
```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

Added element-wise scaling of the gradient based on the historical sum of squares in each dimension

"Per-parameter learning rates" or "adaptive learning rates"

AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q: What happens with AdaGrad?

Duchi et al., "Adaptive subgradient methods for online learning and stochastic optimization", JMLR 2011

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 37

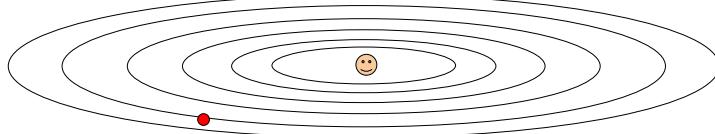
April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 38 April 24, 2018

AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q: What happens with AdaGrad?

Progress along "steep" directions is damped;
progress along "flat" directions is accelerated

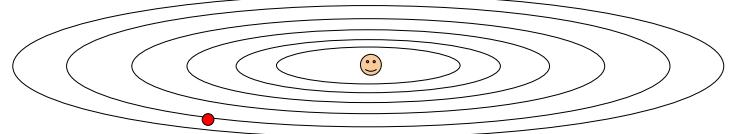
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 39

April 24, 2018

AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q2: What happens to the step size over long time?

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 40

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 39

April 24, 2018

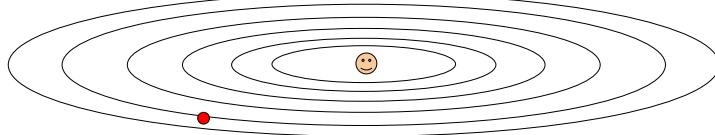
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 40

April 24, 2018

AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q2: What happens to the step size over long time?

Decays to zero

RMSProp

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



RMSProp

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 41

April 24, 2018

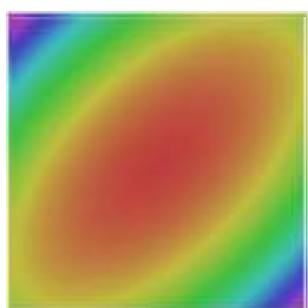
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 42

April 24, 2018

Tieleman and Hinton, 2012

RMSProp



- SGD
- SGD+Momentum
- RMSProp

Adam (almost)

```
first_moment = 0
second_moment = 0
while True:
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    x -= learning_rate * first_moment / (np.sqrt(second_moment) + 1e-7)
```

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 43

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 44

April 24, 2018

Adam (almost)

```

first_moment = 0
second_moment = 0
while Train:
    dx = compute_gradients()
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    x -= learning_rate * first_moment / (np.sqrt(second_moment) + 1e-7)

```

Sort of like RMSProp with momentum

Momentum
AdaGrad / RMSProp

Q: What happens at first timestep?

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 45

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 46

April 24, 2018

Adam (full form)

```

first_moment = 0
second_moment = 0
for t in range(1, num_iterations):
    dx = compute_gradients()
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_meanvar = first_moment / (1 - beta1 ** t)
    second_meanvar = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_meanvar / (np.sqrt(second_meanvar) + 1e-7)

```

Momentum

Bias correction

AdaGrad / RMSProp

Bias correction for the fact that
first and second moment
estimates start at zero

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 45

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 46

April 24, 2018

Adam (full form)

```

first_moment = 0
second_moment = 0
for t in range(1, num_iterations):
    dx = compute_gradients()
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_meanvar = first_moment / (1 - beta1 ** t)
    second_meanvar = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_meanvar / (np.sqrt(second_meanvar) + 1e-7)

```

Momentum
Bias correction
AdaGrad / RMSProp

Bias correction for the fact that
first and second moment
estimates start at zero

Adam with beta1 = 0.9,
beta2 = 0.999, and learning_rate = 1e-3 or 5e-4
is a great starting point for many models!

Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 47

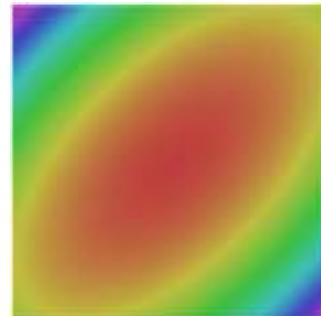
April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 48

April 24, 2018

Adam



- SGD
- SGD+Momentum
- RMSProp
- Adam

SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



Q: Which one of these
learning rates is best to use?

SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



=> Learning rate decay over time!

step decay:
e.g. decay learning rate by half every few epochs.

exponential decay:

$$\alpha = \alpha_0 e^{-kt}$$

1/t decay:

$$\alpha = \alpha_0 / (1 + kt)$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 49

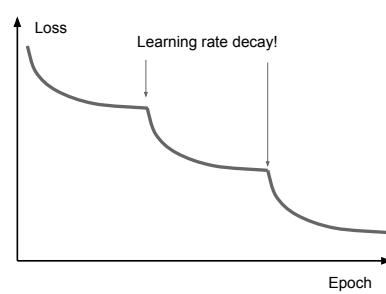
April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

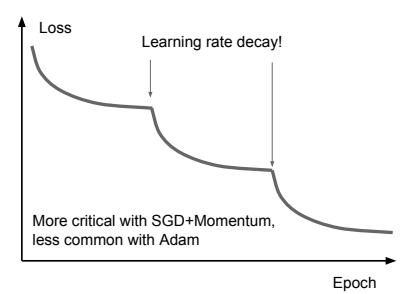
Lecture 7 - 50

April 24, 2018

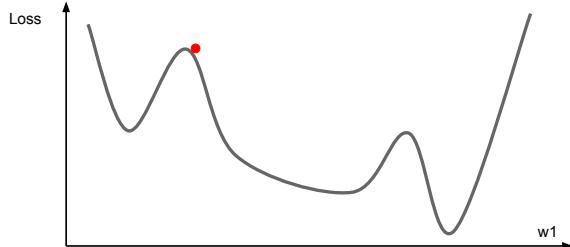
SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.

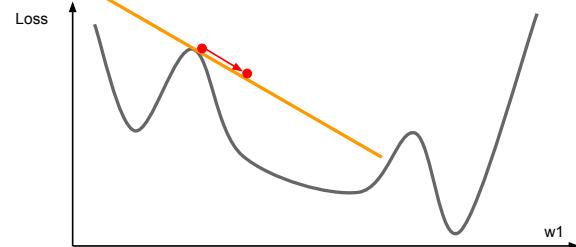


First-Order Optimization

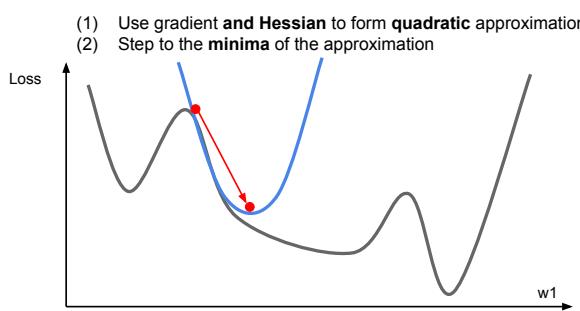


First-Order Optimization

- (1) Use gradient form linear approximation
- (2) Step to minimize the approximation



Second-Order Optimization



Second-Order Optimization

second-order Taylor expansion:

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^T \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^T H(\theta - \theta_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

Q: What is nice about this update?

Second-Order Optimization

second-order Taylor expansion:

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^\top \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^\top H(\theta - \theta_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

No hyperparameters!
No learning rate!
(Though you might use one in practice)

Q: What is nice about this update?

Second-Order Optimization

second-order Taylor expansion:

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^\top \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^\top H(\theta - \theta_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

Q2: Why is this bad for deep learning?

Second-Order Optimization

second-order Taylor expansion:

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^\top \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^\top H(\theta - \theta_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

Hessian has $O(N^2)$ elements
Inverting takes $O(N^3)$
 $N =$ (Tens or Hundreds of Millions)

Q2: Why is this bad for deep learning?

Second-Order Optimization

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

- Quasi-Newton methods (**BGFS** most popular):
instead of inverting the Hessian ($O(n^3)$), approximate inverse Hessian with rank 1 updates over time ($O(n^2)$ each).
- **L-BFGS** (Limited memory BFGS):
Does not form/store the full inverse Hessian.

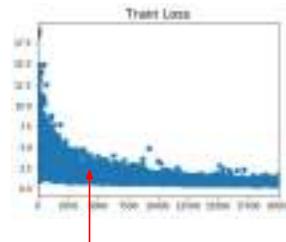
L-BFGS

- **Usually works very well in full batch, deterministic mode**
i.e. if you have a single, deterministic $f(x)$ then L-BFGS will probably work very nicely
- **Does not transfer very well to mini-batch setting.** Gives bad results. Adapting second-order methods to large-scale, stochastic setting is an active area of research.

In practice:

- **Adam** is a good default choice in many cases
- **SGD+Momentum** with learning rate decay often outperforms Adam by a bit, but requires more tuning
- If you can afford to do full batch updates then try out **L-BFGS** (and don't forget to disable all sources of noise)

Beyond Training Error

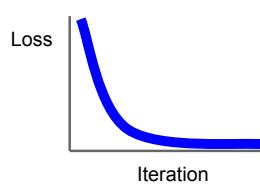


Better optimization algorithms help reduce training loss

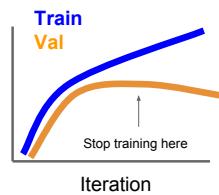


But we really care about error on new data - how to reduce the gap?

Early Stopping



Stop training the model when accuracy on the validation set decreases
Or train for a long time, but always keep track of the model snapshot that worked best on val



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 63

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 64

April 24, 2018

Model Ensembles

1. Train multiple independent models
2. At test time average their results

(Take average of predicted probability distributions, then choose argmax)

Enjoy 2% extra performance

Model Ensembles: Tips and Tricks

Instead of training independent models, use multiple snapshots of a single model during training!



Loshchilov and Hutter, "SGDR: Stochastic gradient descent with restarts", arXiv 2016
Huang et al., "Snapshot ensembles: train 1, get M for free", ICLR 2017
Figures copyright Yuxuan Li and Geoff Pleiss, 2017. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 65

April 24, 2018

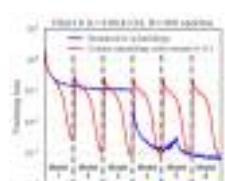
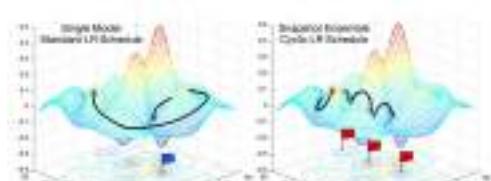
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 66

April 24, 2018

Model Ensembles: Tips and Tricks

Instead of training independent models, use multiple snapshots of a single model during training!



Cyclic learning rate schedules can make this work even better!

Model Ensembles: Tips and Tricks

Instead of using actual parameter vector, keep a moving average of the parameter vector and use that at test time (Polyak averaging)

```
while True:  
    data_batch = dataset.sample_data_batch()  
    loss = network.forward(data_batch)  
    dx = network.backward()  
    x += learning_rate * dx  
    x_test = 0.995*x_test + 0.005*x # incorporate test data
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 67

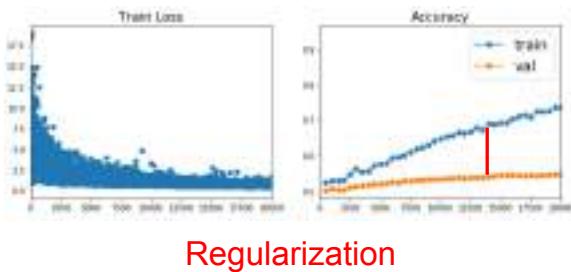
April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 68

April 24, 2018

How to improve single-model performance?



Regularization: Add term to loss

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

In common use:

L2 regularization

L1 regularization

Elastic net (L1 + L2)

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (\text{Weight decay})$$

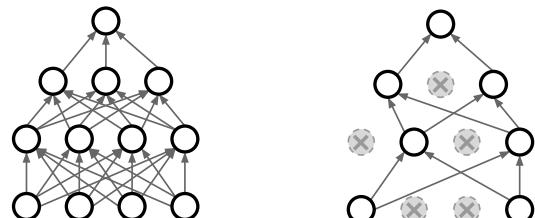
$$R(W) = \sum_k \sum_l |W_{k,l}|$$

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

Regularization: Dropout

In each forward pass, randomly set some neurons to zero

Probability of dropping is a hyperparameter; 0.5 is common



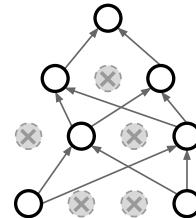
Srivastava et al., "Dropout: A simple way to prevent neural networks from overfitting", JMLR 2014

Regularization: Dropout

```
# 0.5 = probability of keeping a unit active; higher = less robust
def train_step(X):
    """ It contains the data """
    # forward pass: first example 3-layer neural network
    H0 = np.maximum(0, np.dot(W0, X) + b0)
    H1 = np.random.rand(H0.shape) * p # first dropout mask
    H2 = np.maximum(0, np.dot(W1, H1) + b1)
    H3 = np.random.rand(H2.shape) * p # second dropout mask
    out = H3.dot(H2, H1) + b0

    # backward pass: compute gradients... (not shown)
    # update parameter update... (not shown)
```

Example forward pass with a 3-layer network using dropout



Regularization: Dropout

How can this possibly be a good idea?

Forces the network to have a redundant representation; Prevents co-adaptation of features



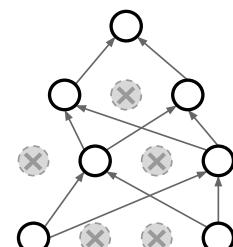
Regularization: Dropout

How can this possibly be a good idea?

Another interpretation:

Dropout is training a large **ensemble** of models (that share parameters).

Each binary mask is one model



An FC layer with 4096 units has $2^{4096} \sim 10^{1233}$ possible masks! Only $\sim 10^{82}$ atoms in the universe...

Dropout: Test time

Dropout makes our output random!

$$\text{Output (label)} \quad \text{Input (image)}$$

$$y = f_W(x|z)$$

Random mask

Want to “average out” the randomness at test-time

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

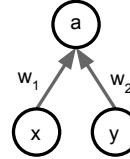
But this integral seems hard ...

Dropout: Test time

Want to approximate the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.



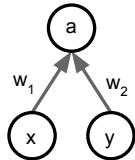
Dropout: Test time

Want to approximate the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.

At test time we have: $E[a] = w_1x + w_2y$

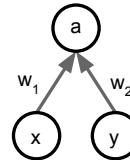


Dropout: Test time

Want to approximate the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.



At test time we have: $E[a] = w_1x + w_2y$

$$\begin{aligned} \text{During training we have: } E[a] &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) \\ &\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y) \end{aligned}$$

Dropout: Test time

Want to approximate the integral

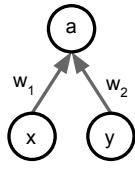
$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.

At test time we have: $E[a] = w_1x + w_2y$

During training we have: $E[a] = \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y)$

$$\begin{aligned} &\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y) \end{aligned}$$



Dropout: Test time

```

def predict(X):
    # forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
    out = np.dot(W3, H2) + b3
  
```

At test time all neurons are active always
 \Rightarrow We must scale the activations so that for each neuron:
output at test time = expected output at training time

Dropout Summary

More common: "Inverted dropout"

drop in forward pass

scale at test time

test time is unchanged!

Regularization: A common pattern

Training: Add some kind of randomness

$$y = f_W(x, z)$$

Testing: Average out randomness (sometimes approximate)

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Regularization: A common pattern

Training: Add some kind of randomness

$$y = f_W(x, z)$$

Testing: Average out randomness (sometimes approximate)

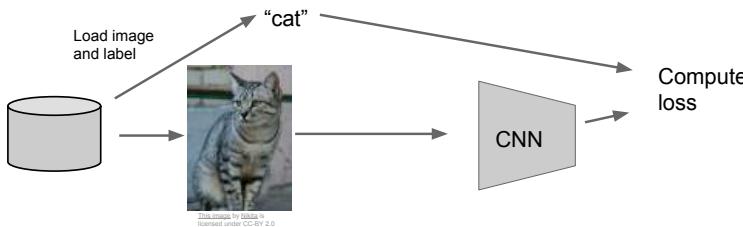
$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Example: Batch Normalization

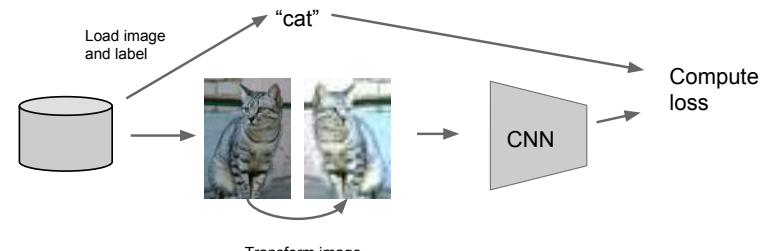
Training: Normalize using stats from random minibatches

Testing: Use fixed stats to normalize

Regularization: Data Augmentation

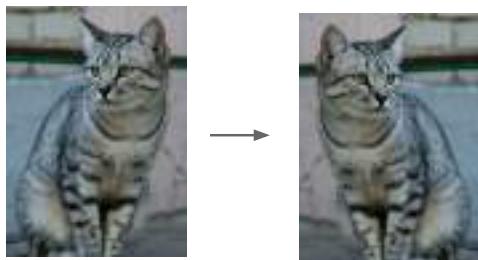


Regularization: Data Augmentation



Data Augmentation

Horizontal Flips



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 87

April 24, 2018

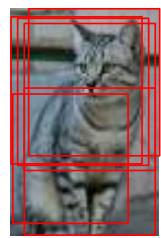
Data Augmentation

Random crops and scales

Training: sample random crops / scales

ResNet:

1. Pick random L in range [256, 480]
2. Resize training image, short side = L
3. Sample random 224 x 224 patch



Lecture 7 - 88

April 24, 2018

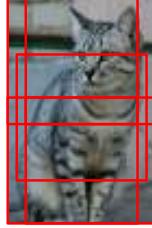
Data Augmentation

Random crops and scales

Training: sample random crops / scales

ResNet:

1. Pick random L in range [256, 480]
2. Resize training image, short side = L
3. Sample random 224 x 224 patch



Testing: average a fixed set of crops

ResNet:

1. Resize image at 5 scales: {224, 256, 384, 480, 640}
2. For each size, use 10 224 x 224 crops: 4 corners + center, + flips

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 89

April 24, 2018

Data Augmentation

Color Jitter

Simple: Randomize contrast and brightness



Lecture 7 - 90

April 24, 2018

Data Augmentation

Color Jitter

Simple: Randomize contrast and brightness



More Complex:

1. Apply PCA to all [R, G, B] pixels in training set
2. Sample a "color offset" along principal component directions
3. Add offset to all pixels of a training image

(As seen in [Krizhevsky et al. 2012], ResNet, etc)

Data Augmentation

Get creative for your problem!

Random mix/combinations of :

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 91

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 92

April 24, 2018

Regularization: A common pattern

Training: Add random noise
Testing: Marginalize over the noise

Examples:

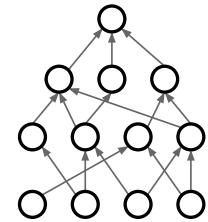
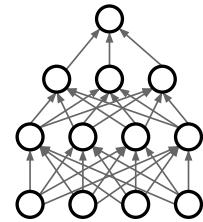
Dropout
Batch Normalization
Data Augmentation

Regularization: A common pattern

Training: Add random noise
Testing: Marginalize over the noise

Examples:

Dropout
Batch Normalization
Data Augmentation
DropConnect



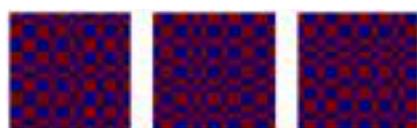
Wan et al., "Regularization of Neural Networks using DropConnect", ICML 2013

Regularization: A common pattern

Training: Add random noise
Testing: Marginalize over the noise

Examples:

Dropout
Batch Normalization
Data Augmentation
DropConnect
Fractional Max Pooling



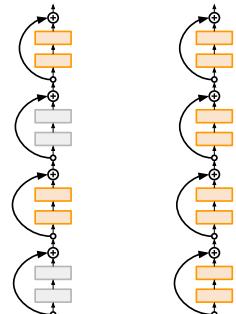
Graham, "Fractional Max Pooling", arXiv 2014

Regularization: A common pattern

Training: Add random noise
Testing: Marginalize over the noise

Examples:

Dropout
Batch Normalization
Data Augmentation
DropConnect
Fractional Max Pooling
Stochastic Depth



Huang et al., "Deep Networks with Stochastic Depth", ECCV 2016

Transfer Learning

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

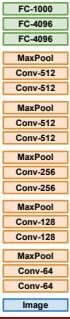
Transfer Learning

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

BUSTED

Transfer Learning with CNNs

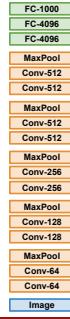
1. Train on Imagenet



Donahue et al. "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al. "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

Transfer Learning with CNNs

1. Train on Imagenet



Donahue et al. "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al. "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

2. Small Dataset (C classes)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 99

April 24, 2018

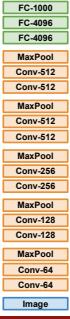
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 10
0

April 24, 2018

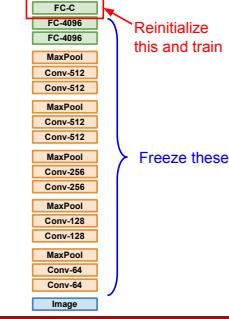
Transfer Learning with CNNs

1. Train on Imagenet

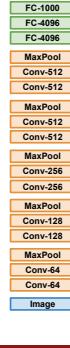
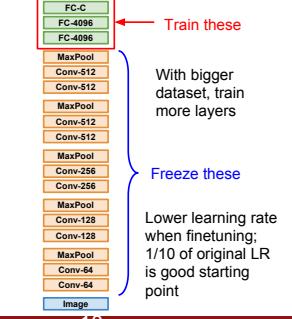


Donahue et al. "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al. "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

2. Small Dataset (C classes)



3. Bigger dataset



More specific
More generic

	very similar dataset	very different dataset
very little data	?	?
quite a lot of data	?	?

Fei-Fei Li & Justin Johnson & Serena Yeung

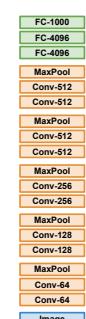
Lecture 7 - 10
1

April 24, 2018

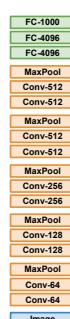
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 10
2

April 24, 2018



	very similar dataset	very different dataset
very little data	Use Linear Classifier on top layer	?
quite a lot of data	Finetune a few layers	?



	very similar dataset	very different dataset
very little data	Use Linear Classifier on top layer	You're in trouble... Try linear classifier from different stages
quite a lot of data	Finetune a few layers	Finetune a larger number of layers

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 10
3

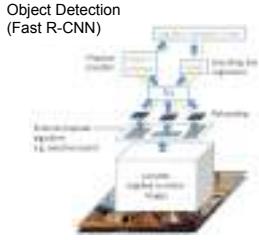
April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 10
4

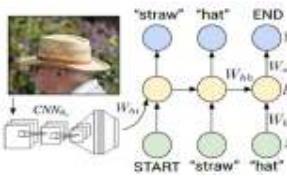
April 24, 2018

Transfer learning with CNNs is pervasive...
(it's the norm, not an exception)



Girshick, "Fast R-CNN", ICCV 2015
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Image Captioning: CNN + RNN

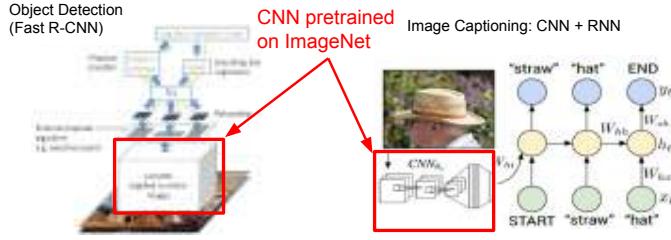


Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015
Figure copyright IEEE, 2015. Reproduced for educational purposes.

Lecture 7 - 10
5

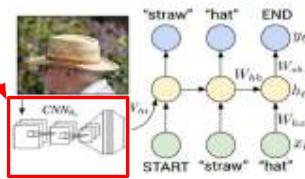
April 24, 2018

Transfer learning with CNNs is pervasive...
(it's the norm, not an exception)



Girshick, "Fast R-CNN", ICCV 2015
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Image Captioning: CNN + RNN



Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015
Figure copyright IEEE, 2015. Reproduced for educational purposes.

Lecture 7 - 10
6

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

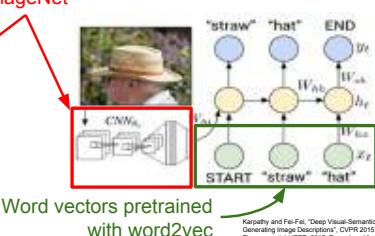
Transfer learning with CNNs is pervasive...
(it's the norm, not an exception)

Object Detection
(Fast R-CNN)



CNN pretrained on ImageNet

Image Captioning: CNN + RNN



Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015
Figure copyright IEEE, 2015. Reproduced for educational purposes.

Word vectors pretrained with word2vec

Lecture 7 - 10
7

April 24, 2018

Takeaway for your projects and beyond:
Have some dataset of interest but it has $< \sim 1M$ images?

1. Find a very large dataset that has similar data, train a big ConvNet there
2. Transfer learn to your dataset

Deep learning frameworks provide a "Model Zoo" of pre-trained models so you don't need to train your own

Caffe: <https://github.com/BVLC/caffe/wiki/Model-Zoo>

TensorFlow: <https://github.com/tensorflow/models>

PyTorch: <https://github.com/pytorch/vision>

Girshick, "Fast R-CNN", ICCV 2015
Figure copyright Ross Girshick, 2015. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 10
8

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Summary

- Lots of Batch Normalization variants
- Optimization
 - Momentum, RMSProp, Adam, etc
- Regularization
 - Dropout, etc
- Transfer learning
 - Use this for your projects!

Next time: Deep Learning Software!

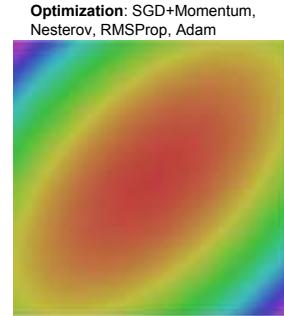
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 11
0

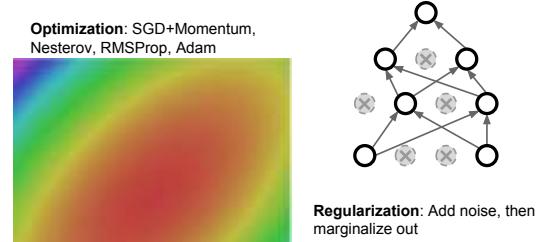
April 24, 2018

Lecture 8: Hardware and Software

Last time



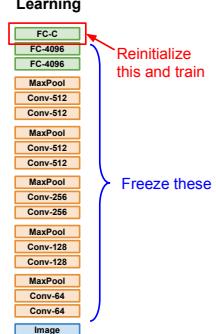
Regularization: Dropout



Train $y = f_W(x, z)$

Test $y = f(x) = E_z[f(x, z)]$

Transfer Learning



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 11 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 22 April 26, 2018

Deep learning hardware

- Deep learning hardware
 - CPU, GPU, TPU
- Deep learning software
 - PyTorch and TensorFlow
 - Static vs Dynamic computation graphs

Deep Learning Hardware

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 33 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 44 April 26, 2018

My computer



Spot the CPU!
(central processing unit)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 55 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 66 April 26, 2018

Spot the GPUs! (graphics processing unit)



NVIDIA vs AMD

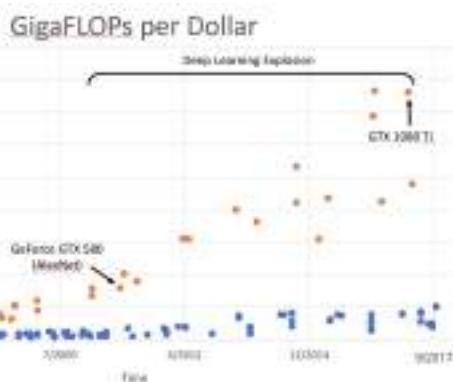
NVIDIA vs AMD

CPU vs GPU

	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

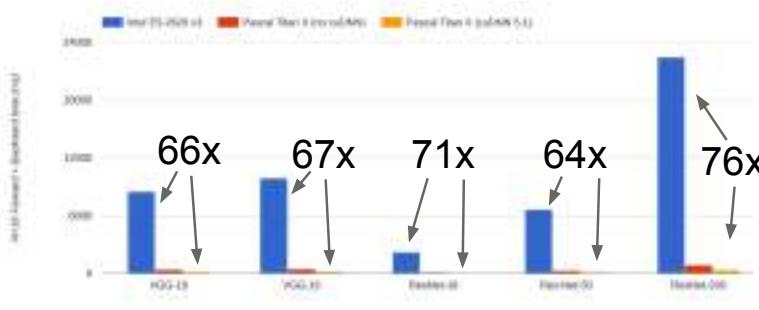
CPU: Fewer cores, but each core is much faster and much more capable; great at sequential tasks

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



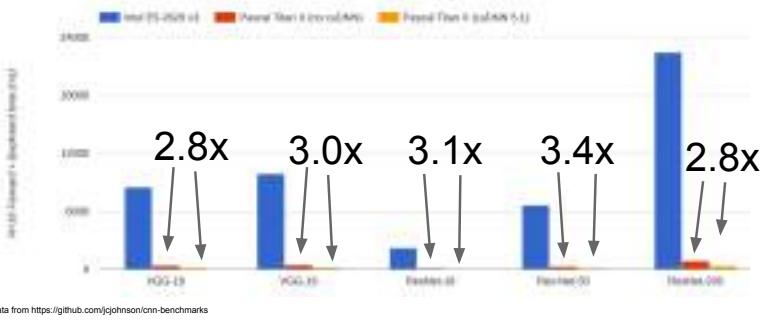
CPU vs GPU in practice

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



CPU vs GPU in practice

cuDNN much faster than "unoptimized" CUDA



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 13 April 26, 2018

CPU vs GPU

	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

Fei-Fei Li & Justin Johnson & Serena Yeung

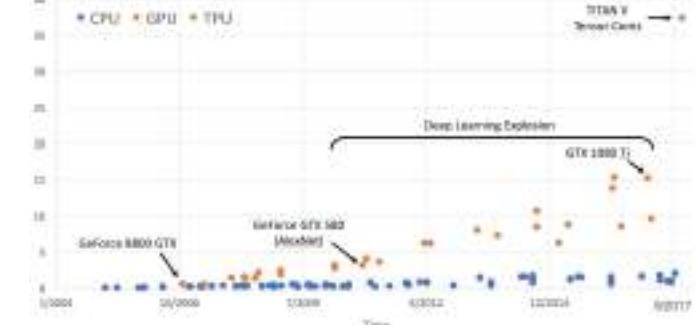
Lecture 8 - 14 April 26, 2018

CPU vs GPU

	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

NOTE: TITAN V isn't technically a "TPU" since that's a Google term, but both have hardware specialized for deep learning

GigaFLOPs per Dollar



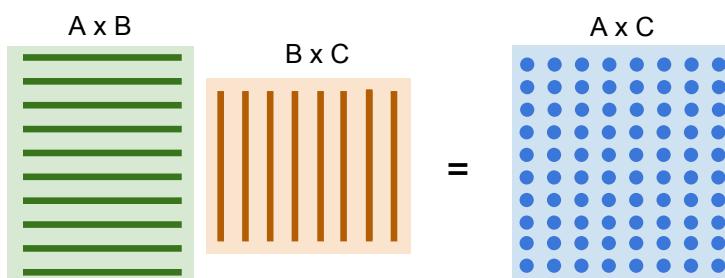
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 15 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 16 April 26, 2018

Example: Matrix Multiplication



Programming GPUs

- CUDA (NVIDIA only)
 - Write C-like code that runs directly on the GPU
 - Optimized APIs: cuBLAS, cuFFT, cuDNN, etc
- OpenCL
 - Similar to CUDA, but runs on anything
 - Usually slower on NVIDIA hardware
- HIP <https://github.com/ROCm-Developer-Tools/HIP>
 - New project that automatically converts CUDA code to something that can run on AMD GPUs
- Udacity: Intro to Parallel Programming <https://www.udacity.com/course/cs344>
 - For deep learning just use existing libraries

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 17 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 18 April 26, 2018

CPU / GPU Communication



Model
is here

Data is here

CPU / GPU Communication



Model
is here

Data is here

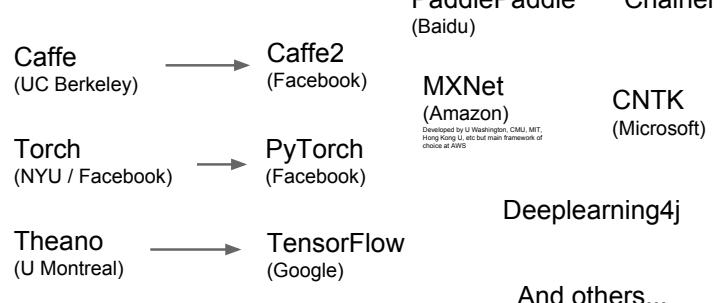
If you aren't careful, training can bottleneck on reading data and transferring to GPU!

Solutions:

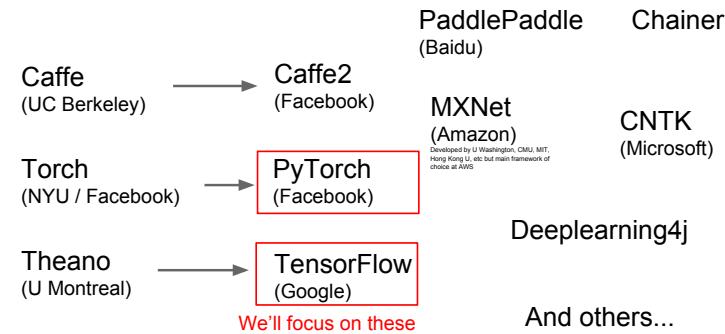
- Read all data into RAM
- Use SSD instead of HDD
- Use multiple CPU threads to prefetch data

Deep Learning Software

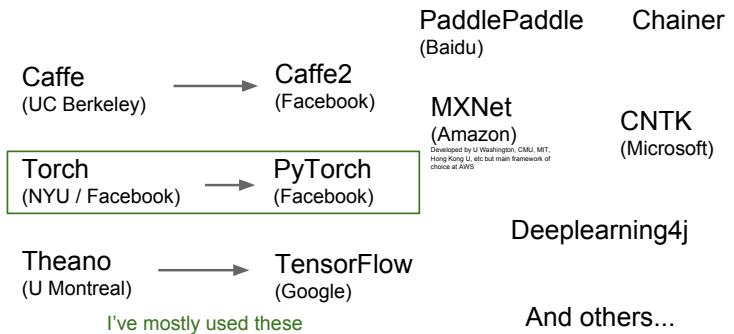
A zoo of frameworks!



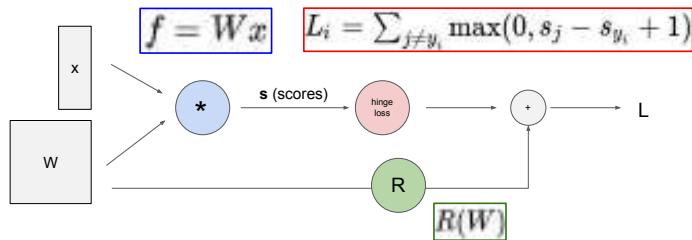
A zoo of frameworks!



A zoo of frameworks!



Recall: Computational Graphs



Recall: Computational Graphs

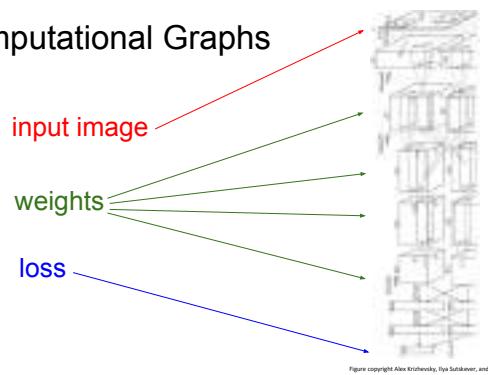
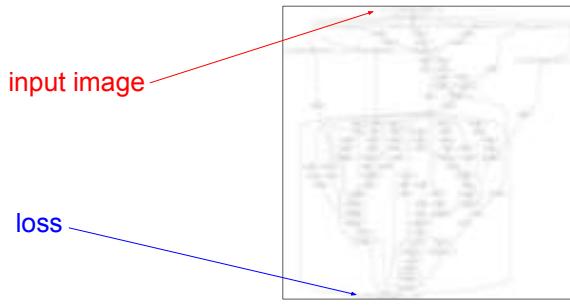


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Recall: Computational Graphs



The point of deep learning frameworks

- (1) Quick to develop and test new ideas
- (2) Automatically compute gradients
- (3) Run it all efficiently on GPU (wrap cuDNN, cuBLAS, etc)

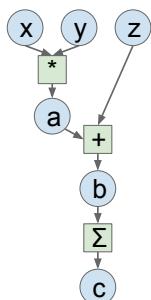
Computational Graphs

Numpy

```
import numpy as np
np.random.seed(0)

x, y = np.random.randint(10, size=(2, 3))
z = np.random.randint(10, size=3)
a = np.multiply(x, y)
b = np.add(a, z)
c = np.sum(b)

print(c)
```



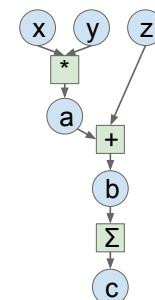
Computational Graphs

Numpy

```
import numpy as np
np.random.seed(0)

x, y = np.random.randint(10, size=(2, 3))
z = np.random.randint(10, size=3)
a = np.multiply(x, y)
b = np.add(a, z)
c = np.sum(b)

grad_a_0 = 1.0
grad_a_1 = grad_a_0 * np.ones((3, 3))
grad_a_2 = grad_a_1 * np.ones(3)
grad_a_3 = grad_a_2 * np.ones(3)
grad_a_4 = grad_a_3 * np.ones(3)
grad_a_5 = grad_a_4 * np.ones(3)
```



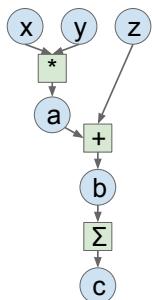
Computational Graphs

Numpy

```
import numpy as np
np.random.seed(0)
x, y = np.zeros(2), 0
a = np.random.randn(2, 2)
y = np.random.randn(2, 2)
z = np.random.randn(2, 2)

a = a + y
b = a + z
c = np.sum(b)

grad_a = 1.0
grad_b = grad_a * np.ones((2, 2))
grad_c = grad_b * np.ones(2)
grad_a = grad_a * y
grad_c = grad_a * c
```



Good:

- Clean API, easy to write numeric code

Bad:

- Have to compute our own gradients
- Can't run on GPU

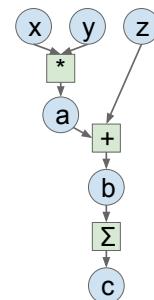
Computational Graphs

Numpy

```
import numpy as np
np.random.seed(0)
x, y = np.zeros(2), 0
a = np.random.randn(2, 2)
y = np.random.randn(2, 2)
z = np.random.randn(2, 2)

a = a + y
b = a + z
c = np.sum(b)

grad_a = 1.0
grad_b = grad_a * np.ones((2, 2))
grad_c = grad_b * np.ones(2)
grad_a = grad_a * y
grad_c = grad_a * c
```



PyTorch

```
import torch
x, y = torch.zeros(2), 0
a = torch.randn(2, 2)
y = torch.randn(2, 2)
z = torch.randn(2, 2)

a = a + y
b = a + z
c = torch.sum(b)

grad_a = 1.0
grad_b = grad_a * torch.ones((2, 2))
grad_c = grad_b * torch.ones(2)
grad_a = grad_a * y
grad_c = grad_a * c
```

Looks exactly like numpy!

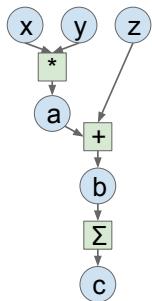
Computational Graphs

Numpy

```
import numpy as np
np.random.seed(0)
x, y = np.zeros(2), 0
a = np.random.randn(2, 2)
y = np.random.randn(2, 2)
z = np.random.randn(2, 2)

a = a + y
b = a + z
c = np.sum(b)

grad_a = 1.0
grad_b = grad_a * np.ones((2, 2))
grad_c = grad_b * np.ones(2)
grad_a = grad_a * y
grad_c = grad_a * c
```



PyTorch

```
import torch
x, y = torch.zeros(2), 0
a = torch.randn(2, 2)
y = torch.randn(2, 2)
z = torch.randn(2, 2)

a = a + y
b = a + z
c = torch.sum(b)

grad_a = 1.0
grad_b = grad_a * torch.ones((2, 2))
grad_c = grad_b * torch.ones(2)
grad_a = grad_a * y
grad_c = grad_a * c
```

PyTorch handles gradients for us!

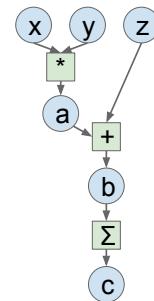
Computational Graphs

Numpy

```
import numpy as np
np.random.seed(0)
x, y = np.zeros(2), 0
a = np.random.randn(2, 2)
y = np.random.randn(2, 2)
z = np.random.randn(2, 2)

a = a + y
b = a + z
c = np.sum(b)

grad_a = 1.0
grad_b = grad_a * np.ones((2, 2))
grad_c = grad_b * np.ones(2)
grad_a = grad_a * y
grad_c = grad_a * c
```



PyTorch

```
import torch
x, y = torch.zeros(2), 0
a = torch.randn(2, 2, requires_grad=True)
y = torch.randn(2, 2, requires_grad=True)
z = torch.randn(2, 2, requires_grad=True)

a = a + y
b = a + z
c = torch.sum(b)

grad_a = 1.0
grad_b = grad_a * torch.ones((2, 2))
grad_c = grad_b * torch.ones(2)
grad_a = grad_a * y
grad_c = grad_a * c
```

Trivial to run on GPU - just construct arrays on a different device!

PyTorch: Fundamental Concepts

Tensor: Like a numpy array, but can run on GPU

Autograd: Package for building computational graphs out of Tensors, and automatically computing gradients

Module: A neural network layer; may store state or learnable weights

PyTorch
(More detail)

PyTorch: Versions

For this class we are using **PyTorch version 0.4** which was released Tuesday 4/24

This version makes a lot of changes to some of the core APIs around autograd, Tensor construction, Tensor datatypes / devices, etc

Be careful if you are looking at older PyTorch code!

PyTorch: Tensors

Running example: Train a two-layer ReLU network on random data with L2 loss

```
import torch
device = torch.device('cpu')

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in, device=device)
y = torch.randn(N, D_out, device=device)
w1 = torch.rand(H, D_in, device=device)
w2 = torch.rand(H, D_out, device=device)

learning_rate = 1e-6
for t in range(1000):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2.t())
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_y_relu = grad_y_pred.masked_fill_(h < 0, 0)
    grad_h = grad_y_relu.mm(w2.t())
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch: Tensors

PyTorch Tensors are just like numpy arrays, but they can run on GPU.

PyTorch Tensor API looks almost exactly like numpy!

Here we fit a two-layer net using PyTorch Tensors:

```
import torch
device = torch.device('cpu')

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in, device=device)
y = torch.randn(N, D_out, device=device)
w1 = torch.rand(H, D_in, device=device)
w2 = torch.rand(H, D_out, device=device)

learning_rate = 1e-6
for t in range(1000):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2.t())
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_y_relu = grad_y_pred.masked_fill_(h < 0, 0)
    grad_h = grad_y_relu.mm(w2.t())
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch: Tensors

Create random tensors for data and weights

```
import torch
device = torch.device('cpu')

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in, device=device)
y = torch.randn(N, D_out, device=device)
w1 = torch.rand(H, D_in, device=device)
w2 = torch.rand(H, D_out, device=device)

learning_rate = 1e-6
for t in range(1000):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2.t())
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_y_relu = grad_y_pred.masked_fill_(h < 0, 0)
    grad_h = grad_y_relu.mm(w2.t())
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch: Tensors

Forward pass: compute predictions and loss

```
import torch
device = torch.device('cpu')

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in, device=device)
y = torch.randn(N, D_out, device=device)
w1 = torch.rand(H, D_in, device=device)
w2 = torch.rand(H, D_out, device=device)

learning_rate = 1e-6
for t in range(1000):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2.t())
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_y_relu = grad_y_pred.masked_fill_(h < 0, 0)
    grad_h = grad_y_relu.mm(w2.t())
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch: Tensors

Backward pass:
manually compute gradients

```
import torch
device = torch.device('cpu')

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in, device=device)
y = torch.randn(N, D_out, device=device)
w1 = torch.rand(H, D_in, device=device)
w2 = torch.rand(H, D_out, device=device)

learning_rate = 1e-6
for t in range(1000):
    h = x.mm(w1)
    h_relu = h.clamp(min=0)
    y_pred = h_relu.mm(w2.t())
    loss = (y_pred - y).pow(2).sum()

    grad_y_pred = 2.0 * (y_pred - y)
    grad_w2 = h_relu.t().mm(grad_y_pred)
    grad_y_relu = grad_y_pred.masked_fill_(h < 0, 0)
    grad_h = grad_y_relu.mm(w2.t())
    grad_w1 = x.t().mm(grad_h)

    w1 -= learning_rate * grad_w1
    w2 -= learning_rate * grad_w2
```

PyTorch: Tensors

```
import torch  
device = torch.device('cpu')  
  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in, device=device)  
y = torch.randn(H, D_out, device=device)  
w1 = torch.randn(D_in, H, device=device)  
w2 = torch.randn(H, D_out, device=device)  
  
learning_rate = 1e-6  
for t in range(500):  
    x = x.mm(w1)  
    x_min = x.clamp(min=0)  
    y_pred = x_min.mm(w2)  
    loss = (y_pred - y).pow(2).mean()  
  
    grad_y_pred = 2.0 * (y_pred - y)  
    grad_w2 = x_min.t().mm(y_pred.grad)  
    grad_x_min = grad_y_pred.mm(w2.t())  
    grad_b1 = grad_x_min.clamp_(0)  
    grad_b2 = y - y_pred  
    grad_w1 = x.t().mm(grad_b1)  
  
    w1 -= learning_rate * grad_w1  
    w2 -= learning_rate * grad_w2
```

Gradient descent step on weights

PyTorch: Tensors

To run on GPU, just use a different device!

```
import torch  
device = torch.device('cuda:0')  
  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in, device=device)  
y = torch.randn(H, D_out, device=device)  
w1 = torch.randn(D_in, H, device=device)  
w2 = torch.randn(H, D_out, device=device)  
  
learning_rate = 1e-6  
for t in range(500):  
    x = x.mm(w1)  
    x_min = x.clamp(min=0)  
    y_pred = x_min.mm(w2)  
    loss = (y_pred - y).pow(2).mean()  
  
    grad_y_pred = 2.0 * (y_pred - y)  
    grad_w2 = x_min.t().mm(y_pred.grad)  
    grad_x_min = grad_y_pred.mm(w2.t())  
    grad_b1 = grad_x_min.clamp_(0)  
    grad_b2 = y - y_pred  
    grad_w1 = x.t().mm(grad_b1)  
  
    w1 -= learning_rate * grad_w1  
    w2 -= learning_rate * grad_w2
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 43 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 44 April 26, 2018

PyTorch: Autograd

Creating Tensors with requires_grad=True enables autograd

Operations on Tensors with requires_grad=True cause PyTorch to build a computational graph

```
import torch  
  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = 1000.0 * torch.randn(H, D_out)  
w1 = torch.randn(D_in, H, requires_grad=True)  
w2 = torch.randn(H, D_out, requires_grad=True)  
  
learning_rate = 1e-6  
for t in range(500):  
    y_pred = x.mm(w1).clamp(min=0).mm(w2)  
    loss = (y_pred - y).pow(2).mean()  
  
    loss.backward()  
  
    with torch.no_grad():  
        w1 += -learning_rate * w1.grad  
        w2 += -learning_rate * w2.grad  
        w1.grad.zero_()  
        w2.grad.zero_()
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 45 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 46 April 26, 2018

PyTorch: Autograd

We will not want gradients (of loss) with respect to data

Do want gradients with respect to weights

```
import torch  
  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = 1000.0 * torch.randn(H, D_out)  
w1 = torch.randn(D_in, H, requires_grad=True)  
w2 = torch.randn(H, D_out, requires_grad=True)  
  
learning_rate = 1e-6  
for t in range(500):  
    y_pred = x.mm(w1).clamp(min=0).mm(w2)  
    loss = (y_pred - y).pow(2).mean()  
  
    loss.backward()  
  
    with torch.no_grad():  
        w1 += -learning_rate * w1.grad  
        w2 += -learning_rate * w2.grad  
        w1.grad.zero_()  
        w2.grad.zero_()
```

PyTorch: Autograd

Forward pass looks exactly the same as before, but we don't need to track intermediate values - PyTorch keeps track of them for us in the graph

```
import torch  
  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = 1000.0 * torch.randn(H, D_out)  
w1 = torch.randn(D_in, H, requires_grad=True)  
w2 = torch.randn(H, D_out, requires_grad=True)  
  
learning_rate = 1e-6  
for t in range(500):  
    y_pred = x.mm(w1).clamp(min=0).mm(w2)  
    loss = (y_pred - y).pow(2).mean()  
  
    loss.backward()  
  
    with torch.no_grad():  
        w1 += -learning_rate * w1.grad  
        w2 += -learning_rate * w2.grad  
        w1.grad.zero_()  
        w2.grad.zero_()
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 47 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 48 April 26, 2018

PyTorch: Autograd

Compute gradient of loss with respect to w1 and w2

```
import torch  
  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = 1000.0 * torch.randn(H, D_out)  
w1 = torch.randn(D_in, H, requires_grad=True)  
w2 = torch.randn(H, D_out, requires_grad=True)  
  
learning_rate = 1e-6  
for t in range(500):  
    y_pred = x.mm(w1).clamp(min=0).mm(w2)  
    loss = (y_pred - y).pow(2).mean()  
  
    loss.backward()  
  
    with torch.no_grad():  
        w1 += -learning_rate * w1.grad  
        w2 += -learning_rate * w2.grad  
        w1.grad.zero_()  
        w2.grad.zero_()
```

Lecture 8 - 48 April 26, 2018

PyTorch: Autograd

```
import torch
N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    loss.backward()
    with torch.no_grad():
        w1 += learning_rate * w1.grad
        w2 += learning_rate * w2.grad
        w1.grad.zero_()
        w2.grad.zero_()
```

Make gradient step on weights, then zero them. Torch.no_grad means "don't build a computational graph for this part"

PyTorch: Autograd

```
import torch
N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    loss.backward()
    with torch.no_grad():
        w1 += learning_rate * w1.grad
        w2 += learning_rate * w2.grad
        w1.grad.zero_()
        w2.grad.zero_()
```

PyTorch methods that end in underscore modify the Tensor in-place; methods that don't return a new Tensor

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 49 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 50 April 26, 2018

PyTorch: New Autograd Functions

Define your own autograd functions by writing forward and backward functions for Tensors

Very similar to modular layers in A2! Use ctx object to "cache" values for the backward pass, just like cache objects from A2

```
class MyReLU(torch.autograd.Function):
    @staticmethod
    def forward(ctx, x):
        ctx.save_for_backward(x)
        return x.clamp(min=0)

    @staticmethod
    def backward(ctx, grad_y):
        x, = ctx.saved_tensors
        grad_input = grad_y.clone()
        grad_input[x < 0] = 0
        return grad_input
```

PyTorch: New Autograd Functions

Define your own autograd functions by writing forward and backward functions for Tensors

Very similar to modular layers in A2! Use ctx object to "cache" values for the backward pass, just like cache objects from A2

Define a helper function to make it easy to use the new function

```
class MyReLU(torch.autograd.Function):
    @staticmethod
    def forward(ctx, x):
        ctx.save_for_backward(x)
        return x.clamp(min=0)

    @staticmethod
    def backward(ctx, grad_y):
        x, = ctx.saved_tensors
        grad_input = grad_y.clone()
        grad_input[x < 0] = 0
        return grad_input

    def my_relu(x):
        return MyReLU.apply(x)
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 51 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 52 April 26, 2018

PyTorch: New Autograd Functions

```
class MyReLU(torch.autograd.Function):
    @staticmethod
    def forward(ctx, x):
        ctx.save_for_backward(x)
        return x.clamp(min=0)

    @staticmethod
    def backward(ctx, grad_y):
        x, = ctx.saved_tensors
        grad_input = grad_y.clone()
        grad_input[x < 0] = 0
        return grad_input

def my_relu(x):
    return MyReLU.apply(x)
```

Can use our new autograd function in the forward pass

```
N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = my_relu(x.mm(w1)).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    loss.backward()
    with torch.no_grad():
        w1 += learning_rate * w1.grad
        w2 += learning_rate * w2.grad
        w1.grad.zero_()
        w2.grad.zero_()
```

PyTorch: New Autograd Functions

```
def my_relu(x):
    return x.clamp(min=0)
```

In practice you almost never need to define new autograd functions! Only do it when you need custom backward. In this case we can just use a normal Python function

```
N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = my_relu(x.mm(w1)).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    loss.backward()
    with torch.no_grad():
        w1 += learning_rate * w1.grad
        w2 += learning_rate * w2.grad
        w1.grad.zero_()
        w2.grad.zero_()
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 53 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 54 April 26, 2018

PyTorch: nn

Higher-level wrapper for working with neural nets

Use this! It will make your life easier

```
import torch  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = torch.randn(N, D_out)  
  
model = torch.nn.Sequential(  
    torch.nn.Linear(D_in, H),  
    torch.nn.ReLU(),  
    torch.nn.Linear(H, D_out))  
  
learning_rate = 1e-3  
for t in range(500):  
    y_pred = model(x)  
    loss = torch.nn.functional.mse_loss(y_pred, y)  
  
    loss.backward()  
  
    with torch.no_grad():  
        for param in model.parameters():  
            param -= learning_rate * param.grad  
    model.zero_grad()
```

PyTorch: nn

Define our model as a sequence of layers; each layer is an object that holds learnable weights

```
import torch  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = torch.randn(N, D_out)  
  
model = torch.nn.Sequential(  
    torch.nn.Linear(D_in, H),  
    torch.nn.ReLU(),  
    torch.nn.Linear(H, D_out))  
  
learning_rate = 1e-3  
for t in range(500):  
    y_pred = model(x)  
    loss = torch.nn.functional.mse_loss(y_pred, y)  
  
    loss.backward()  
  
    with torch.no_grad():  
        for param in model.parameters():  
            param -= learning_rate * param.grad  
    model.zero_grad()
```

PyTorch: nn

Forward pass: feed data to model, and compute loss

```
import torch  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = torch.randn(N, D_out)  
  
model = torch.nn.Sequential(  
    torch.nn.Linear(D_in, H),  
    torch.nn.ReLU(),  
    torch.nn.Linear(H, D_out))  
  
learning_rate = 1e-3  
for t in range(500):  
    y_pred = model(x)  
    loss = torch.nn.functional.mse_loss(y_pred, y)  
  
    loss.backward()  
  
    with torch.no_grad():  
        for param in model.parameters():  
            param -= learning_rate * param.grad  
    model.zero_grad()
```

PyTorch: nn

Forward pass: feed data to model, and compute loss

```
import torch  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = torch.randn(N, D_out)  
  
model = torch.nn.Sequential(  
    torch.nn.Linear(D_in, H),  
    torch.nn.ReLU(),  
    torch.nn.Linear(H, D_out))  
  
learning_rate = 1e-3  
for t in range(500):  
    y_pred = model(x)  
    loss = torch.nn.functional.mse_loss(y_pred, y)  
  
    loss.backward()  
  
    with torch.no_grad():  
        for param in model.parameters():  
            param -= learning_rate * param.grad  
    model.zero_grad()
```

PyTorch: nn

Backward pass: compute gradient with respect to all model weights (they have `requires_grad=True`)

```
import torch  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = torch.randn(N, D_out)  
  
model = torch.nn.Sequential(  
    torch.nn.Linear(D_in, H),  
    torch.nn.ReLU(),  
    torch.nn.Linear(H, D_out))  
  
learning_rate = 1e-3  
for t in range(500):  
    y_pred = model(x)  
    loss = torch.nn.functional.mse_loss(y_pred, y)  
  
    loss.backward()  
  
    with torch.no_grad():  
        for param in model.parameters():  
            param -= learning_rate * param.grad  
    model.zero_grad()
```

PyTorch: nn

Make gradient step on each model parameter (with gradients disabled)

```
import torch  
N, D_in, H, D_out = 64, 1000, 100, 10  
x = torch.randn(N, D_in)  
y = torch.randn(N, D_out)  
  
model = torch.nn.Sequential(  
    torch.nn.Linear(D_in, H),  
    torch.nn.ReLU(),  
    torch.nn.Linear(H, D_out))  
  
learning_rate = 1e-3  
for t in range(500):  
    y_pred = model(x)  
    loss = torch.nn.functional.mse_loss(y_pred, y)  
  
    loss.backward()  
  
    with torch.no_grad():  
        for param in model.parameters():  
            param -= learning_rate * param.grad  
    model.zero_grad()
```


PyTorch: nn Define new Modules

Define forward pass using child modules

No need to define backward - autograd will handle it

```
import torch
from torch import nn, optim
from torch.nn import functional as F
from torch.autograd import Variable

def forward(x):
    h_relu = nn.ReLU()(x).clamp(min=0)
    y_pred = nn.Linear(h_relu, 10)
    return y_pred

H, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(H, D_in))
y = Variable(torch.randn(H, D_out))

model = nn.Sequential(nn.Linear(D_in, H), nn.ReLU(),
                      nn.Linear(H, D_out))

optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
for t in range(100):
    y_pred = model(x)
    loss = F.cross_entropy(y_pred, y)

    loss.backward()
    optimizer.step()
    optimizer.zero_grad()
```

PyTorch: nn Define new Modules

Construct and train an instance of our model

```
import torch
from torch import nn, optim
from torch.nn import functional as F
from torch.autograd import Variable

def forward(x):
    h_relu = nn.ReLU()(x).clamp(min=0)
    y_pred = nn.Linear(h_relu, 10)
    return y_pred

H, D_in, H, D_out = 64, 1000, 100, 10
x = Variable(torch.randn(H, D_in))
y = Variable(torch.randn(H, D_out))

model = nn.Sequential(nn.Linear(D_in, H), nn.ReLU(),
                      nn.Linear(H, D_out))

optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
for t in range(100):
    y_pred = model(x)
    loss = F.cross_entropy(y_pred, y)

    loss.backward()
    optimizer.step()
    optimizer.zero_grad()
```

PyTorch: nn Define new Modules

Very common to mix and match custom Module subclasses and Sequential containers

```
import torch
from torch import nn, optim
from torch.nn import functional as F
from torch.autograd import Variable

def forward(x):
    h1, h2, D_out = 64, 1000, 100, 10
    x = nn.Linear(1000, D_in)(x)
    p = nn.Linear(1000, D_out)

    model = nn.Sequential(nn.Linear(D_in, 100),
                          nn.ReLU(),
                          nn.Linear(100, D_out))

    optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
    for t in range(100):
        y_pred = model(x)
        loss = F.cross_entropy(y_pred, y)
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()
```

PyTorch: nn Define new Modules

Define network component as a Module subclass

```
import torch
from torch import nn, optim
from torch.nn import functional as F
from torch.autograd import Variable

def forward(x):
    h1, h2, D_out = 64, 1000, 100, 10
    x = nn.Linear(1000, D_in)(x)
    p = nn.Linear(1000, D_out)

    model = nn.Sequential(nn.Linear(D_in, 100),
                          nn.ReLU(),
                          nn.Linear(100, D_out))

    optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
    for t in range(100):
        y_pred = model(x)
        loss = F.cross_entropy(y_pred, y)
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()
```

PyTorch: nn Define new Modules

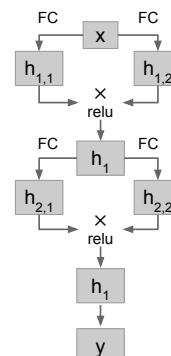
Stack multiple instances of the component in a sequential

```
import torch
from torch import nn, optim
from torch.nn import functional as F
from torch.autograd import Variable

def forward(x):
    h1, h2, D_out = 64, 1000, 100, 10
    x = nn.Linear(1000, D_in)(x)
    p = nn.Linear(1000, D_out)

    model = nn.Sequential(nn.Linear(D_in, 100),
                          nn.ReLU(),
                          nn.Linear(100, D_out))

    optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
    for t in range(100):
        y_pred = model(x)
        loss = F.cross_entropy(y_pred, y)
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()
```



```
import torch
from torch import nn, optim
from torch.nn import functional as F
from torch.autograd import Variable

def forward(x):
    h1, h2, D_out = 64, 1000, 100, 10
    x = nn.Linear(1000, D_in)(x)
    p = nn.Linear(1000, D_out)

    model = nn.Sequential(nn.Linear(D_in, 100),
                          nn.ReLU(),
                          nn.Linear(100, D_out))

    optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
    for t in range(100):
        y_pred = model(x)
        loss = F.cross_entropy(y_pred, y)
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()
```

PyTorch: DataLoaders

A **DataLoader** wraps a **Dataset** and provides minibatching, shuffling, multithreading, for you

When you need to load custom data, just write your own **Dataset** class

```
import torch
from torch.utils.data import TensorDataset, DataLoader

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)

loader = DataLoader(TensorDataset(x, y), batch_size=64)
model = TwoLayerNet(D_in, H, D_out)

optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
for epoch in range(25):
    for i, (x, y) in enumerate(loader):
        y_pred = model(x)
        loss = torch.nn.functional.nll_loss(y_pred, y)
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()
```

PyTorch: DataLoaders

Iterate over loader to form minibatches

```
import torch
from torch.utils.data import TensorDataset, DataLoader

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)

loader = DataLoader(TensorDataset(x, y), batch_size=64)
model = TwoLayerNet(D_in, H, D_out)

optimizer = torch.optim.SGD(model.parameters(), lr=1e-2)
for epoch in range(25):
    for i, (x, y) in enumerate(loader):
        y_pred = model(x)
        loss = torch.nn.functional.nll_loss(y_pred, y)
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()
```

PyTorch: Pretrained Models

Super easy to use pretrained models with torchvision
<https://github.com/pytorch/vision>

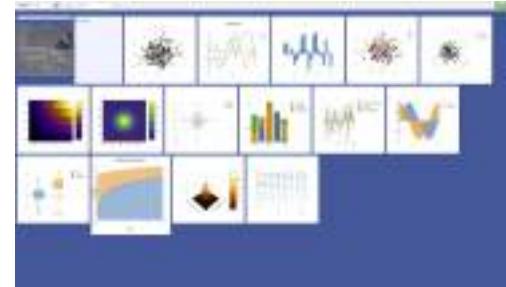
```
import torch
import torchvision

alexnet = torchvision.models.alexnet(pretrained=True)
vgg16 = torchvision.models.vgg16(pretrained=True)
resnet101 = torchvision.models.resnet101(pretrained=True)
```

PyTorch: Visdom

Visualization tool: add logging to your code, then visualize in a browser

Can't visualize computational graph structure (yet?)



PyTorch: Dynamic Computation Graphs

```
import torch

N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    loss.backward()
```

PyTorch: Dynamic Computation Graphs

x w1 w2 y

```
import torch

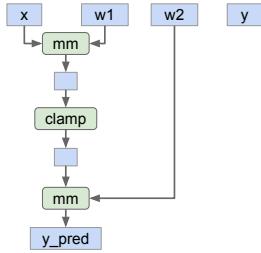
N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

    loss.backward()
```

Create Tensor objects

PyTorch: Dynamic Computation Graphs



```

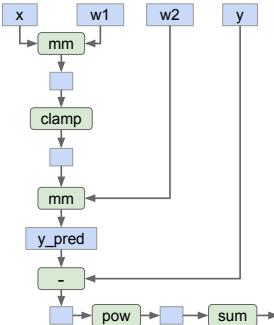
import torch
N, D_in, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, D, requires_grad=True)
w2 = torch.randn(D, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()
    loss.backward()

loss.backward()
  
```

Build graph data structure AND perform computation

PyTorch: Dynamic Computation Graphs



```

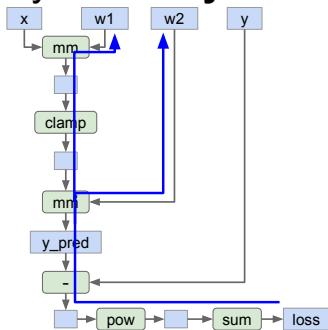
import torch
N, D_in, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, D, requires_grad=True)
w2 = torch.randn(D, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()
    loss.backward()

loss.backward()
  
```

Build graph data structure AND perform computation

PyTorch: Dynamic Computation Graphs



```

import torch
N, D_in, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, D, requires_grad=True)
w2 = torch.randn(D, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()
    loss.backward()

loss.backward()
  
```

Search for path between loss and w1, w2 (for backprop) AND perform computation

PyTorch: Dynamic Computation Graphs



```

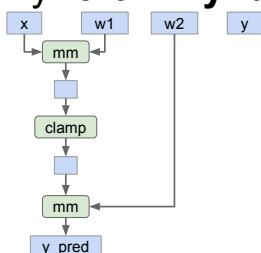
import torch
N, D_in, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, D, requires_grad=True)
w2 = torch.randn(D, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()
    loss.backward()

loss.backward()
  
```

Throw away the graph, backprop path, and rebuild it from scratch on every iteration

PyTorch: Dynamic Computation Graphs



```

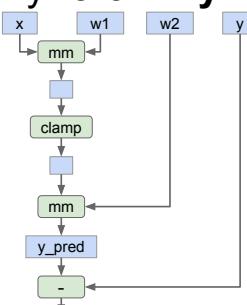
import torch
N, D_in, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, D, requires_grad=True)
w2 = torch.randn(D, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()
    loss.backward()

loss.backward()
  
```

Build graph data structure AND perform computation

PyTorch: Dynamic Computation Graphs



```

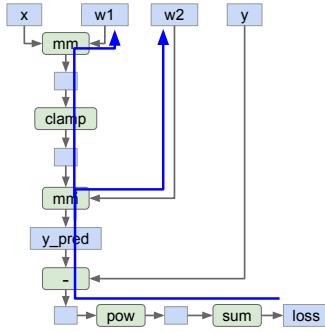
import torch
N, D_in, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, D, requires_grad=True)
w2 = torch.randn(D, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()
    loss.backward()

loss.backward()
  
```

Build graph data structure AND perform computation

PyTorch: Dynamic Computation Graphs



```

import torch
N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(H, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()
    loss.backward()

```

Search for path between loss and w1, w2
(for backprop) AND perform computation

PyTorch: Dynamic Computation Graphs

Building the graph and **computing** the graph happen at the same time.

Seems inefficient, especially if we are building the same graph over and over again...

```

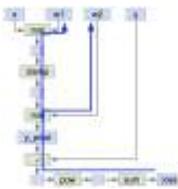
import torch
N, D_in, H, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(H, D_out)
w1 = torch.randn(D_in, H, requires_grad=True)
w2 = torch.randn(H, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()
    loss.backward()

```

loss.backward()

Static Computation Graphs



Alternative: **Static** graphs

Step 1: Build computational graph describing our computation (including finding paths for backprop)

Step 2: Reuse the same graph on every iteration

```

graph = build_graph()
for x_batch, y_batch in loader:
    run_graph(graph, x=x_batch, y=y_batch)

```

TensorFlow

TensorFlow

TensorFlow: Neural Net

```

import numpy as np
import tensorflow as tf

# Assume imports at the top of each snippet

N, D_in, H, D_out = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D_in])
y = tf.placeholder(tf.float32, shape=[H, D_out])
w1 = tf.placeholder(tf.float32, shape=[D_in, H])
w2 = tf.placeholder(tf.float32, shape=[H, D_out])

z = tf.matmul(tf.matmul(x, w1), w2)
y_pred = tf.nn.relu(z)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.square(diff ** 2, axis=-1)))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = {x: np.random.randn(N, D_in),
              w1: np.random.randn(D_in, H),
              w2: np.random.randn(H, D_out),
              y: np.random.randn(H, D_out)}
    out = sess.run([loss, grad_w1, grad_w2],
                  feed_dict=values)
    loss_val, grad_w1_val, grad_w2_val = out

```

TensorFlow: Neural Net

First **define** computational graph

Then **run** the graph many times

```

N, D_in, H, D_out = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D_in])
y = tf.placeholder(tf.float32, shape=[H, D_out])
w1 = tf.placeholder(tf.float32, shape=[D_in, H])
w2 = tf.placeholder(tf.float32, shape=[H, D_out])

z = tf.matmul(tf.matmul(x, w1), w2)
y_pred = tf.nn.relu(z)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.square(diff ** 2, axis=-1)))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

```

```

with tf.Session() as sess:
    values = {x: np.random.randn(N, D_in),
              w1: np.random.randn(D_in, H),
              w2: np.random.randn(H, D_out),
              y: np.random.randn(H, D_out)}
    out = sess.run([loss, grad_w1, grad_w2],
                  feed_dict=values)
    loss_val, grad_w1_val, grad_w2_val = out

```

TensorFlow: Neural Net

Create placeholders for input x, weights w1 and w2, and targets y

```

x, D, N = 88, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.placeholder(tf.float32, shape=[D, 10])
w2 = tf.placeholder(tf.float32, shape=[10, D])

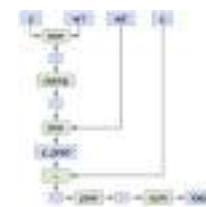
h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff) ** 2, axis=1))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = (x, np.random.randn(N, D),
              w1, np.random.randn(D, 10),
              w2, np.random.randn(10, D),
              y, np.random.randn(N, D))
    out = sess.run([loss, grad_w1, grad_w2],
                  feed_dict=values)
    loss_val, grad_w1_val, grad_w2_val = out

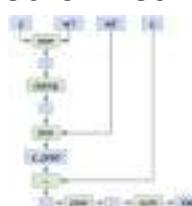
```

TensorFlow: Neural Net



Forward pass: compute prediction for y and loss. No computation - just building graph

TensorFlow: Neural Net



```

x, D, N = 88, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.placeholder(tf.float32, shape=[D, 10])
w2 = tf.placeholder(tf.float32, shape=[10, D])

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff) ** 2, axis=1))

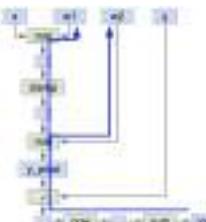
grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = (x, np.random.randn(N, D),
              w1, np.random.randn(D, 10),
              w2, np.random.randn(10, D),
              y, np.random.randn(N, D))
    out = sess.run([loss, grad_w1, grad_w2],
                  feed_dict=values)
    loss_val, grad_w1_val, grad_w2_val = out

```

Tell TensorFlow to compute loss of gradient with respect to w1 and w2. No compute - just building the graph

TensorFlow: Neural Net



Find paths between loss and w1, w2

```

x, D, N = 88, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.placeholder(tf.float32, shape=[D, 10])
w2 = tf.placeholder(tf.float32, shape=[10, D])

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff) ** 2, axis=1))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = (x, np.random.randn(N, D),
              w1, np.random.randn(D, 10),
              w2, np.random.randn(10, D),
              y, np.random.randn(N, D))
    out = sess.run([loss, grad_w1, grad_w2],
                  feed_dict=values)
    loss_val, grad_w1_val, grad_w2_val = out

```

TensorFlow: Neural Net



```

x, D, N = 88, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.placeholder(tf.float32, shape=[D, 10])
w2 = tf.placeholder(tf.float32, shape=[10, D])

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff) ** 2, axis=1))

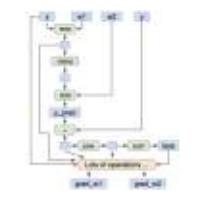
grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = (x, np.random.randn(N, D),
              w1, np.random.randn(D, 10),
              w2, np.random.randn(10, D),
              y, np.random.randn(N, D))
    out = sess.run([loss, grad_w1, grad_w2],
                  feed_dict=values)
    loss_val, grad_w1_val, grad_w2_val = out

```

Add new operators to the graph which compute grad_w1 and grad_w2

TensorFlow: Neural Net



Now done building our graph, so we enter a session so we can actually run the graph

```

x, D, N = 88, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.placeholder(tf.float32, shape=[D, 10])
w2 = tf.placeholder(tf.float32, shape=[10, D])

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff) ** 2, axis=1))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = (x, np.random.randn(N, D),
              w1, np.random.randn(D, 10),
              w2, np.random.randn(10, D),
              y, np.random.randn(N, D))
    out = sess.run([loss, grad_w1, grad_w2],
                  feed_dict=values)
    loss_val, grad_w1_val, grad_w2_val = out

```

TensorFlow: Neural Net



Create numpy arrays that will fill in the placeholders above

```

x, b, N = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.placeholder(tf.float32, shape=[D, 5])
w2 = tf.placeholder(tf.float32, shape=[5, D])

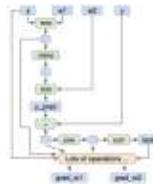
h = tf.matmul(tf.matmul(x, w1), b)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff), axis=1))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = (x, np.random.randn(N, D),
              w1, np.random.randn(D, 5),
              w2, np.random.randn(5, D))
    out = sess.run([loss, grad_w1, grad_w2],
                  feed_dict=values)
    loss_val, grad_w1_val, grad_w2_val = out

```

TensorFlow: Neural Net



Run the graph: feed in the numpy arrays for x, y, w1, and w2; get numpy arrays for loss, grad_w1, and grad_w2

```

x, b, N = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.placeholder(tf.float32, shape=[D, 5])
w2 = tf.placeholder(tf.float32, shape=[5, D])

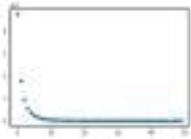
h = tf.matmul(tf.matmul(x, w1), b)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff), axis=1))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = (x, np.random.randn(N, D),
              w1, np.random.randn(D, 5),
              w2, np.random.randn(5, D))
    out = sess.run([loss, grad_w1, grad_w2],
                  feed_dict=values)
    loss_val, grad_w1_val, grad_w2_val = out

```

TensorFlow: Neural Net



Train the network: Run the graph over and over, use gradient to update weights

```

x, b, N = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.placeholder(tf.float32, shape=[D, 5])
w2 = tf.placeholder(tf.float32, shape=[5, D])

h = tf.matmul(tf.matmul(x, w1), b)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff), axis=1))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = (x, np.random.randn(N, D),
              w1, np.random.randn(D, 5),
              w2, np.random.randn(5, D),
              y: np.random.randn(N, D))
    learning_rate = 1e-2
    for t in range(50):
        sess.run([loss, grad_w1, grad_w2],
                feed_dict=values)
        loss_val, grad_w1_val, grad_w2_val = sess.run([loss, grad_w1, grad_w2],
                                                       feed_dict=values)
        values[w1] = learning_rate * grad_w1_val
        values[w2] = learning_rate * grad_w2_val

```

TensorFlow: Neural Net

Problem: copying weights between CPU / GPU each step

Train the network: Run the graph over and over, use gradient to update weights

```

x, b, N = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.placeholder(tf.float32, shape=[D, 5])
w2 = tf.placeholder(tf.float32, shape=[5, D])

h = tf.matmul(tf.matmul(x, w1), b)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff), axis=1))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

with tf.Session() as sess:
    values = (x, np.random.randn(N, D),
              w1, np.random.randn(D, 5),
              w2, np.random.randn(5, D),
              y: np.random.randn(N, D))
    learning_rate = 1e-2
    for t in range(50):
        sess.run([loss, grad_w1, grad_w2],
                feed_dict=values)
        loss_val, grad_w1_val, grad_w2_val = sess.run([loss, grad_w1, grad_w2],
                                                       feed_dict=values)
        values[w1] = learning_rate * grad_w1_val
        values[w2] = learning_rate * grad_w2_val

```

TensorFlow: Neural Net

Change w1 and w2 from **placeholder** (fed on each call) to **Variable** (persists in the graph between calls)

```

x, b, N = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.Variable(tf.random_normal([D, 5]))
w2 = tf.Variable(tf.random_normal([5, D]))

h = tf.matmul(tf.matmul(x, w1), b)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff), axis=1))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

learning_rate = 1e-2
new_w1 = w1.assign(w1 - learning_rate * grad_w1)
new_w2 = w2.assign(w2 - learning_rate * grad_w2)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = (x, np.random.randn(N, D),
              y: np.random.randn(N, D))
    for t in range(50):
        loss_val, _ = sess.run([loss, new_w1], feed_dict=values)

```

TensorFlow: Neural Net

Add **assign** operations to update w1 and w2 as part of the graph!

```

x, b, N = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.Variable(tf.random_normal([D, 5]))
w2 = tf.Variable(tf.random_normal([5, D]))

h = tf.matmul(tf.matmul(x, w1), b)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.abs(diff), axis=1))

grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

learning_rate = 1e-2
new_w1 = w1.assign(w1 - learning_rate * grad_w1)
new_w2 = w2.assign(w2 - learning_rate * grad_w2)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = (x, np.random.randn(N, D),
              y: np.random.randn(N, D))
    for t in range(50):
        loss_val, _ = sess.run([loss, new_w1, new_w2], feed_dict=values)

```

TensorFlow: Neural Net

```

N, D, M = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.Variable(tf.random_normal([D, M]))
w2 = tf.Variable(tf.random_normal([M, D]))

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.square(diff), axis=1))
grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

learning_rate = 1e-1
new_w1 = w1.assign(w1 - learning_rate * grad_w1)
new_w2 = w2.assign(w2 - learning_rate * grad_w2)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = [x: np.random.randn(N, D),
              y: np.random.randn(N, D)]
    for t in range(50):
        loss_val, _ = sess.run([loss, new_w1, new_w2], feed_dict=values)

```

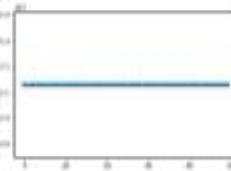
Run graph once to initialize w1 and w2

Run many times to train

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 103 April 26, 2018

TensorFlow: Neural Net



Problem: loss not going down! Assign calls not actually being executed!

```

N, D, M = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.Variable(tf.random_normal([D, M]))
w2 = tf.Variable(tf.random_normal([M, D]))

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.square(diff), axis=1))
grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])

learning_rate = 1e-1
new_w1 = w1.assign(w1 - learning_rate * grad_w1)
new_w2 = w2.assign(w2 - learning_rate * grad_w2)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = [x: np.random.randn(N, D),
              y: np.random.randn(N, D)]
    for t in range(50):
        loss_val, _ = sess.run([loss, new_w1, new_w2], feed_dict=values)

```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 104 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

TensorFlow: Neural Net

```

N, D, M = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.Variable(tf.random_normal([D, M]))
w2 = tf.Variable(tf.random_normal([M, D]))

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
diff = y_pred - y
loss = tf.reduce_mean(tf.reduce_mean(tf.square(diff), axis=1))
grad_w1, grad_w2 = tf.gradients(loss, [w1, w2])
updates = tf.group(w1, new_w2)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = [x: np.random.randn(N, D),
              y: np.random.randn(N, D)]
    losses = []
    for t in range(50):
        loss_val, _ = sess.run([loss, updates], feed_dict=values)

```

Add dummy graph node that depends on updates

Tell TensorFlow to compute dummy node

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 105 April 26, 2018

TensorFlow: Optimizer

Can use an **optimizer** to compute gradients and update weights

Remember to execute the output of the optimizer!

```

optimizer = tf.train.GradientDescentOptimizer(1e-1)
updates = optimizer.minimize(loss)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = [x: np.random.randn(N, D),
              y: np.random.randn(N, D)]
    losses = []
    for t in range(50):
        loss_val, _ = sess.run([loss, updates], feed_dict=values)

```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 106 April 26, 2018

TensorFlow: Loss

Use predefined common losses

```

N, D, M = 64, 1000, 100
x = tf.placeholder(tf.float32, shape=[N, D])
y = tf.placeholder(tf.float32, shape=[N, D])
w1 = tf.Variable(tf.random_normal([D, M]))
w2 = tf.Variable(tf.random_normal([M, D]))

h = tf.maximum(tf.matmul(x, w1), 0)
y_pred = tf.matmul(h, w2)
loss = tf.losses.mean_squared_error(y_pred, y)

optimizer = tf.train.GradientDescentOptimizer(1e-1)
updates = optimizer.minimize(loss)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = [x: np.random.randn(N, D),
              y: np.random.randn(N, D)]
    for t in range(50):
        loss_val, _ = sess.run([loss, updates], feed_dict=values)

```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 107 April 26, 2018

TensorFlow: Layers

Use He initializer

```

init = tf.variance_scaling_initializer(2.0)
h = tf.layers.dense(inputs=x, units=M,
                    activation=tf.nn.relu, kernel_initializer=init)
y_pred = tf.layers.dense(inputs=h, units=D,
                        kernel_initializer=init)

loss = tf.losses.mean_squared_error(y_pred, y)

optimizer = tf.train.GradientDescentOptimizer(1e-1)
updates = optimizer.minimize(loss)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = [x: np.random.randn(N, D),
              y: np.random.randn(N, D)]
    for t in range(50):
        loss_val, _ = sess.run([loss, updates], feed_dict=values)

```

tf.layers automatically sets up weight and (and bias) for us!

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 108 April 26, 2018

Keras: High-Level Wrapper

Keras is a layer on top of TensorFlow, makes common things easy to do

(Used to be third-party, now merged into TensorFlow)

```
# x, y = np.random.randn(100, 100)
x = tf.placeholder(tf.float32, shape=[None, D])
y = tf.placeholder(tf.float32, shape=[None, D])

model = tf.keras.Sequential()
model.add(tf.keras.layers.Dense(10, input_shape=[D], activation=tf.nn.relu))
model.add(tf.keras.layers.Dense(10, activation=tf.nn.relu))
x_pured = model(x)

loss = tf.losses.mean_squared_error(y, y_pured)

optimizer = tf.train.GradientDescentOptimizer(1.0)
updates = optimizer.minimize(loss)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = {x: np.random.randn(N, D),
              y: np.random.randn(N, D)}
    for t in range(50):
        loss_val, _ = sess.run([loss, updates],
                             feed_dict=values)
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 109 April 26, 2018

Keras: High-Level Wrapper

Define model as a sequence of layers

Get output by calling the model

```
# x, y = np.random.randn(100, 100)
x = tf.placeholder(tf.float32, shape=[None, D])
y = tf.placeholder(tf.float32, shape=[None, D])

model = tf.keras.Sequential()
model.add(tf.keras.layers.Dense(10, input_shape=[D], activation=tf.nn.relu))
model.add(tf.keras.layers.Dense(10, activation=tf.nn.relu))
x_pured = model(x)

loss = tf.losses.mean_squared_error(y, y_pured)

optimizer = tf.train.GradientDescentOptimizer(1.0)
updates = optimizer.minimize(loss)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = {x: np.random.randn(N, D),
              y: np.random.randn(N, D)}
    for t in range(50):
        loss_val, _ = sess.run([loss, updates],
                             feed_dict=values)
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 110 April 26, 2018

Keras: High-Level Wrapper

Keras can handle the training loop for you!
No sessions or
feed_dict

```
# x, y = np.random.randn(100, 100)
model = tf.keras.Sequential()
model.add(tf.keras.layers.Dense(10, input_shape=[D], activation=tf.nn.relu))
model.add(tf.keras.layers.Dense(10))

model.compile(loss=tf.keras.losses.mean_squared_error,
              optimizer=tf.keras.optimizers.SGD(lr=1e-6))

x = np.random.randn(N, D)
y = np.random.randn(N, D)

history = model.fit(x, y, epochs=50, batch_size=10)
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 111 April 26, 2018

TensorFlow: High-Level Wrappers

tf.keras (https://www.tensorflow.org/api_docs/python/tf/keras)

tf.layers (https://www.tensorflow.org/api_docs/python/tf/layers)

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 111 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 112 April 26, 2018

TensorFlow: High-Level Wrappers

Keras (<https://keras.io/>)

tf.keras (https://www.tensorflow.org/api_docs/python/tf/keras)

tf.layers (https://www.tensorflow.org/api_docs/python/tf/layers)

TensorFlow: High-Level Wrappers

Keras (<https://keras.io/>)

tf.keras (https://www.tensorflow.org/api_docs/python/tf/keras)

tf.layers (https://www.tensorflow.org/api_docs/python/tf/layers)

tf.estimator (https://www.tensorflow.org/api_docs/python/tf/estimator)

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 113 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 114 April 26, 2018

TensorFlow: High-Level Wrappers

Keras (<https://keras.io/>)

tf.keras (https://www.tensorflow.org/api_docs/python/tf/keras)

tf.layers (https://www.tensorflow.org/api_docs/python/tf/layers)

tf.estimator (https://www.tensorflow.org/api_docs/python/tf/estimator)

tf.contrib.estimator (https://www.tensorflow.org/api_docs/python/tf/contrib/estimator)

tf.contrib.layers (https://www.tensorflow.org/api_docs/python/tf/contrib/layers)

tf.contrib.slim (<https://github.com/tensorflow/tree/master/tensorflow/contrib/slim>)

tf.contrib.learn (https://www.tensorflow.org/api_docs/python/tf/contrib/learn)

Ships with TensorFlow

TensorFlow: High-Level Wrappers

Keras (<https://keras.io/>)

tf.keras (https://www.tensorflow.org/api_docs/python/tf/keras)

tf.layers (https://www.tensorflow.org/api_docs/python/tf/layers)

tf.estimator (https://www.tensorflow.org/api_docs/python/tf/estimator)

tf.contrib.estimator (https://www.tensorflow.org/api_docs/python/tf/contrib/estimator)

tf.contrib.layers (https://www.tensorflow.org/api_docs/python/tf/contrib/layers)

tf.contrib.slim (<https://github.com/tensorflow/tree/master/tensorflow/contrib/slim>)

tf.contrib.learn (https://www.tensorflow.org/api_docs/python/tf/contrib/learn) DEPRECATED

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 1:5 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 1:6 April 26, 2018

TensorFlow: High-Level Wrappers

Keras (<https://keras.io/>)

Ships with TensorFlow

tf.keras (https://www.tensorflow.org/api_docs/python/tf/keras)

tf.layers (https://www.tensorflow.org/api_docs/python/tf/layers)

tf.estimator (https://www.tensorflow.org/api_docs/python/tf/estimator)

tf.contrib.estimator (https://www.tensorflow.org/api_docs/python/tf/contrib/estimator)

tf.contrib.layers (https://www.tensorflow.org/api_docs/python/tf/contrib/layers)

tf.contrib.slim (<https://github.com/tensorflow/tree/master/tensorflow/contrib/slim>)

tf.contrib.learn (https://www.tensorflow.org/api_docs/python/tf/contrib/learn) DEPRECATED

Sonnet (<https://github.com/deepmind/sonnet>)

By DeepMind

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 1:7 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 1:8 April 26, 2018

TensorFlow: High-Level Wrappers

Keras (<https://keras.io/>)

Ships with TensorFlow

tf.keras (https://www.tensorflow.org/api_docs/python/tf/keras)

tf.layers (https://www.tensorflow.org/api_docs/python/tf/layers)

tf.estimator (https://www.tensorflow.org/api_docs/python/tf/estimator)

tf.contrib.estimator (https://www.tensorflow.org/api_docs/python/tf/contrib/estimator)

tf.contrib.layers (https://www.tensorflow.org/api_docs/python/tf/contrib/layers)

tf.contrib.slim (<https://github.com/tensorflow/tree/master/tensorflow/contrib/slim>)

tf.contrib.learn (https://www.tensorflow.org/api_docs/python/tf/contrib/learn) DEPRECATED

Sonnet (<https://github.com/deepmind/sonnet>)

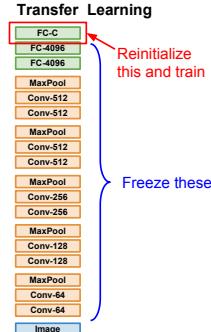
By DeepMind

TFLearn (<http://tflearn.org/>)

TensorLayer (<http://tensorlayer.readthedocs.io/en/latest/>)

Third-Party

TensorFlow: Pretrained Models



tf.keras: (https://www.tensorflow.org/api_docs/python/tf/keras/applications)

TF-Slim: (<https://github.com/tensorflow/models/tree/master/slim/nets>)

TensorFlow: Tensorboard

Add logging to code to record loss, stats, etc

Run server and get pretty graphs!



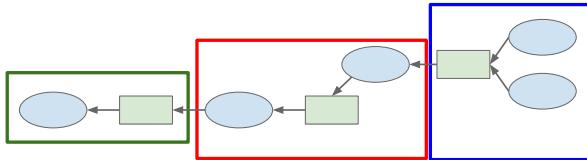
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 1:9 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 120 April 26, 2018

TensorFlow: Distributed Version

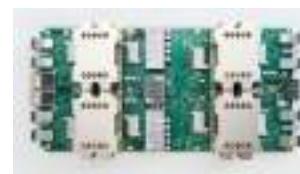


Split one graph
over multiple
machines!



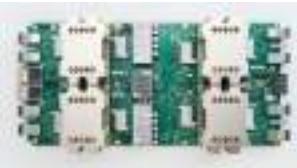
<https://www.tensorflow.org/deploy/distributed>

TensorFlow: Tensor Processing Units



Google Cloud TPU
= 180 TFLOPs of compute!

TensorFlow: Tensor Processing Units



Google Cloud TPU
= 180 TFLOPs of compute!



NVIDIA Tesla V100
= 125 TFLOPs of compute

TensorFlow: Tensor Processing Units



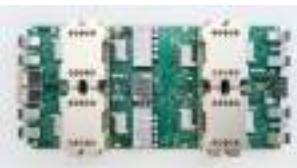
Google Cloud TPU
= 180 TFLOPs of compute!



NVIDIA Tesla V100
= 125 TFLOPs of compute

NVIDIA Tesla P100 = 11 TFLOPs of compute
GTX 580 = 0.2 TFLOPs

TensorFlow: Tensor Processing Units



Google Cloud TPU
= 180 TFLOPs of compute!



Google Cloud TPU Pod
= 64 Cloud TPUs
= 11.5 PFLOPs of compute!

https://www.tensorflow.org/versions/master/programmers_guide/using_tpu

Static vs Dynamic Graphs

TensorFlow: Build graph once, then run many times (**static**)

```
# x, y, z = tf.placeholder(tf.float32, shape=[None, None])
# w = tf.Variable(tf.random_normal([None, None]))
# b = tf.Variable(tf.random_normal([None, 1]))
# z = tf.matmul(x, w) + b
# loss = tf.reduce_mean(tf.reduce_mean(tf.square(z - y), axis=1)))
# grad_w, grad_b = tf.gradients(loss, [w, b])
# updates = tf.assign(w, w - learning_rate * grad_w)
# updates = tf.assign(b, b - learning_rate * grad_b)

learning_rate = 0.01
grad_w = w1.assign(w1 - learning_rate * grad_w)
grad_b = w2.assign(w2 - learning_rate * grad_b)
updates = tf.group(grad_w, grad_b)

with tf.Session() as sess:
    sess.run(tf.global_variables_initializer())
    values = (x1, np.random.randn(100, 1),
              y1, np.random.randn(100, 1))
    losses = []
    for t in range(100):
        loss, _ = sess.run([loss, updates], feed_dict=values)
        losses.append(loss)
```

PyTorch: Each forward pass defines a new graph (**dynamic**)

```
import torch
N, D_in, N, D_out = 64, 1000, 100, 10
x = torch.randn(N, D_in)
y = torch.randn(N, D_out)
w1 = torch.randn(D_in, N, requires_grad=True)
w2 = torch.randn(N, D_out, requires_grad=True)

learning_rate = 1e-6
for t in range(500):
    y_pred = x.mm(w1).clamp(min=0).mm(w2)
    loss = (y_pred - y).pow(2).sum()

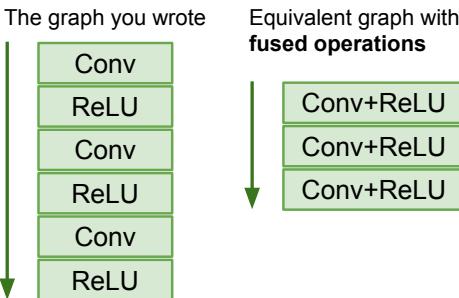
    loss.backward()
```

Build graph

Run each iteration

Static vs Dynamic: Optimization

With static graphs, framework can **optimize** the graph for you before it runs!



Static vs Dynamic: Serialization

Static

Once graph is built, can **serialize** it and run it without the code that built the graph!

Dynamic

Graph building and execution are intertwined, so always need to keep code around

Static vs Dynamic: Conditional

$$y = \begin{cases} w1 * x & \text{if } z > 0 \\ w2 * x & \text{otherwise} \end{cases}$$

Static vs Dynamic: Conditional

$$y = \begin{cases} w1 * x & \text{if } z > 0 \\ w2 * x & \text{otherwise} \end{cases}$$

PyTorch: Normal Python

```

x1, x2, y = 5, 4, 3
x = torch.tensor(x1, requires_grad=True)
y1 = torch.tensor(x1, requires_grad=True)
w1 = torch.tensor(1.0, requires_grad=True)
x2 = torch.tensor(x2, requires_grad=True)
w2 = torch.tensor(1.0, requires_grad=True)

a = 10
if y < 0:
    y = w1*x1
else:
    y = w2*x2
print(y)
  
```

Static vs Dynamic: Conditional

$$y = \begin{cases} w1 * x & \text{if } z > 0 \\ w2 * x & \text{otherwise} \end{cases}$$

PyTorch: Normal Python

```

x1, x2, y = 5, 4, 3
x = torch.tensor(x1, requires_grad=True)
y1 = torch.tensor(x1, requires_grad=True)
w1 = torch.tensor(1.0, requires_grad=True)
x2 = torch.tensor(x2, requires_grad=True)
w2 = torch.tensor(1.0, requires_grad=True)

a = 10
if y < 0:
    y = w1*x1
else:
    y = w2*x2
print(y)
  
```

TensorFlow: Special TF control flow operator!

```

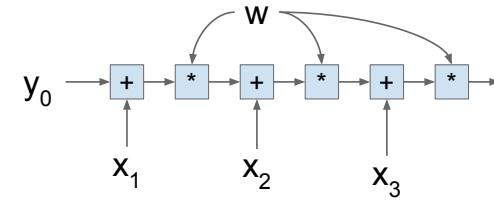
x1, x2, y = 5, 4, 3
x = tf.placeholder(tf.float32, shape=[1])
y = tf.placeholder(tf.float32, shape=[1])
w1 = tf.placeholder(tf.float32, shape=[1])
x1 = tf.placeholder(tf.float32, shape=[1])
w2 = tf.placeholder(tf.float32, shape=[1])

def tf_if():
    with tf.control_dependencies([tf.assert_equal(y, 3)]):
        val1 = tf.cond(tf.less(y, 0), lambda: w1*x1, lambda: w2*x2)
        val2 = tf.cond(tf.greater(y, 0), lambda: w1*x1, lambda: w2*x2)
        y_val = tf.where(tf.equal(y, 3), val1, val2)

with tf.Session() as sess:
    sess.run(tf_if())
    print(sess.run(y))
  
```

Static vs Dynamic: Loops

$$y_t = (y_{t-1} + x_t) * w$$



Static vs Dynamic: Loops

$$y_t = (y_{t-1} + x_t) * w$$

PyTorch: Normal Python

```

T, D = 3, 8
y0 = torch.randn(D, requires_grad=True)
x = torch.randn(T, D)
w = torch.randn(D)

y = ly0]
for t in range(T):
    prev_x = x[-1]
    next_x = (prev_x + x[t]) * w
    y.append(next_x)

```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 133 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

```

TensorFlow: Special TF control flow

T, N, D = 5, 4, 2
x = tf.placeholder(tf.float32, shape=[T, N, D])
y = tf.placeholder(tf.float32, shape=[D])
w = tf.placeholder(tf.float32, shape=[D])

def f(prev_y, cur_x):
    return [prev_y + cur_x] * w

y = tf.foldl(f, x, y0)

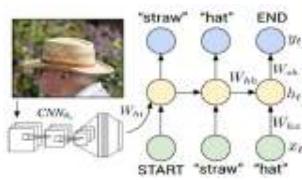
with tf.Session() as sess:
    values = [
        x1=np.random.randn(T, D),
        y0=np.random.randn(D),
        w1=np.random.randn(D),
    ]
    y_val = sess.run(y, feed_dict=values)

```

Lecture 8 - 134 April 26, 2018

Dynamic Graph Applications

- ### - Recurrent networks



Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015

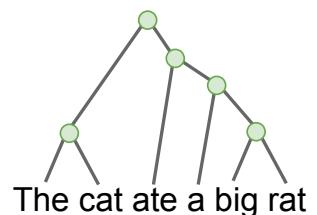
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 135 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 136 April 26, 2018

Dynamic Graph Applications

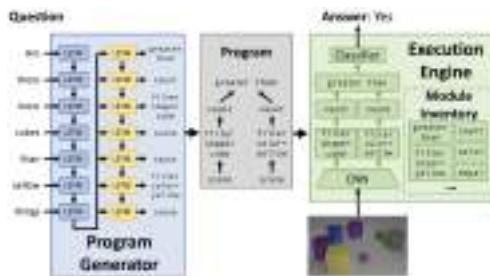
- Recurrent networks
 - Recursive networks



The cat ate a big rat

Dynamic Graph Applications

- Recurrent networks
 - Recursive networks
 - Modular Networks



Andreas et al, "Neural Module Networks", CVPR 2016
Andreas et al, "Learning to Compose Neural Networks for Question Answering", NAACL 2016
Johnson et al, "Inferring and Executing Programs for Visual Reasoning", ICCV 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 137 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 138 April 26, 2018

Dynamic Graph Applications

- Recurrent networks
 - Recursive networks
 - Modular Networks
 - (Your creative idea here)

PyTorch vs TensorFlow, Static vs Dynamic

PyTorch
Dynamic Graphs

TensorFlow
Static Graphs

PyTorch vs TensorFlow, Static vs Dynamic

PyTorch
Dynamic Graphs

TensorFlow
Static Graphs

Lines are blurring! PyTorch is adding static features, and TensorFlow is adding dynamic features.

Dynamic TensorFlow: Dynamic Batching

TensorFlow Fold make dynamic graphs easier in TensorFlow through **dynamic batching**

Dynamic TensorFlow: Eager Execution

TensorFlow 1.7 added **eager execution** which allows dynamic graphs!

```
import tensorflow as tf
import tensorflow.contrib.eager as tfe
tfe.enable_eager_execution()

N, D = 3, 4
x = tfe.Variable(tf.random_normal([N, D]))
y = tfe.Variable(tf.random_normal([N, D]))
z = tfe.Variable(tf.random_normal([N, D]))
with tfe.GradientTape() as tape:
    a = x + y
    b = a + z
    c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tape.gradient(c, [x, y, z])
print(grad_x)
```

Dynamic TensorFlow: Eager Execution

Enable eager mode at the start of the program: it's a global switch

```
import tensorflow as tf
import tensorflow.contrib.eager as tfe
tfe.enable_eager_execution()

N, D = 3, 4
x = tfe.Variable(tf.random_normal([N, D]))
y = tfe.Variable(tf.random_normal([N, D]))
z = tfe.Variable(tf.random_normal([N, D]))
with tfe.GradientTape() as tape:
    a = x + y
    b = a + z
    c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tape.gradient(c, [x, y, z])
print(grad_x)
```

Dynamic TensorFlow: Eager Execution

These calls to `tf.random_normal` produce concrete values! No need for placeholders / sessions

Wrap values in a `tfe.Variable` if we might want to compute grads for them

```
import tensorflow as tf
import tensorflow.contrib.eager as tfe
tfe.enable_eager_execution()

N, D = 3, 4
x = tfe.Variable(tf.random_normal([N, D]))
y = tfe.Variable(tf.random_normal([N, D]))
z = tfe.Variable(tf.random_normal([N, D]))
with tfe.GradientTape() as tape:
    a = x + y
    b = a + z
    c = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tape.gradient(c, [x, y, z])
print(grad_x)
```

Dynamic TensorFlow: Eager Execution

```
import tensorflow as tf
import tensorflow.contrib.eager as tfe

tf.enable_eager_execution()

N, D = 3, 4
x = tfe.Variable(tf.random_normal([N, D]))
y = tfe.Variable(tf.random_normal([N, D]))
z = tfe.Variable(tf.random_normal([N, D]))
with tfe.GradientTape() as tape:
    a = x * z
    b = a + y
    o = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tape.gradient(o, [x, y, z])
print(grad_x)
```

Operations scoped under a GradientTape will build a dynamic graph, similar to PyTorch

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 145 April 26, 2018

Dynamic TensorFlow: Eager Execution

```
import tensorflow as tf
import tensorflow.contrib.eager as tfe

tf.enable_eager_execution()

N, D = 3, 4
x = tfe.Variable(tf.random_normal([N, D]))
y = tfe.Variable(tf.random_normal([N, D]))
z = tfe.Variable(tf.random_normal([N, D]))
with tfe.GradientTape() as tape:
    a = x * z
    b = a + y
    o = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tape.gradient(o, [x, y, z])
print(grad_x)
```

Use the tape to compute gradients, like .backward() in PyTorch. The print statement works!

Lecture 8 - 146 April 26, 2018

Dynamic TensorFlow: Eager Execution

Eager execution still pretty new, not fully supported in all TensorFlow APIs

Try it out!

```
import tensorflow as tf
import tensorflow.contrib.eager as tfe

tf.enable_eager_execution()

N, D = 3, 4
x = tfe.Variable(tf.random_normal([N, D]))
y = tfe.Variable(tf.random_normal([N, D]))
z = tfe.Variable(tf.random_normal([N, D]))
with tfe.GradientTape() as tape:
    a = x * z
    b = a + y
    o = tf.reduce_sum(b)

grad_x, grad_y, grad_z = tape.gradient(o, [x, y, z])
print(grad_x)
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 147 April 26, 2018

Static PyTorch: Caffe2 <https://caffe2.ai/>

- Deep learning framework developed by Facebook
- Static graphs, somewhat similar to TensorFlow
- Core written in C++
- Nice Python interface
- Can train model in Python, then serialize and deploy without Python
- Works on iOS / Android, etc

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 148 April 26, 2018

Static PyTorch: ONNX Support

ONNX is an open-source standard for neural network models

Goal: Make it easy to train a network in one framework, then run it in another framework

Supported by PyTorch, Caffe2, Microsoft CNTK, Apache MXNet

<https://github.com/onnx/onnx>

Static PyTorch: ONNX Support

You can export a PyTorch model to ONNX

Run the graph on a dummy input, and save the graph to a file

Will only work if your model doesn't actually make use of dynamic graph - must build same graph on every forward pass, no loops / conditionals

```
import torch

N, D_in, H, D_out = 64, 1024, 100, 10
model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out))

dummy_input = torch.randn(N, D_in)
torch.onnx.export(model, dummy_input,
                  'model1.onnx',
                  verbose=True)
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 149 April 26, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 150 April 26, 2018

Static PyTorch: ONNX Support

```
graph(%0 : Float(64, 1000)
    %1 : Float(100, 1000)
    %2 : Float(100)
    %3 : Float(10, 100)
    %4 : Float(10) {
        %5 : Float(64, 100) =
        onnx::Gemm(alpha=1, beta=1, broadcast=1,
        transB=1) (%0, %1, %2), scope:
        Sequential/Linear[0]
        %6 : Float(64, 100) = onnx::Relu(%5),
        scope: Sequential/ReLU[1]
        %7 : Float(64, 10) = onnx::Gemm(alpha=1,
        beta=1, broadcast=1, transB=1) (%6, %3,
        %4), scope: Sequential/Linear[2]
        return (%7);
    }
```

```
import torch

N, D_in, H, D_out = 64, 1000, 100, 10
model = torch.nn.Sequential(
    torch.nn.Linear(D_in, H),
    torch.nn.ReLU(),
    torch.nn.Linear(H, D_out))

dummy_input = torch.randn(N, D_in)
torch.onnx.export(model, dummy_input,
                  'model1.onnx',
                  verbose=True)
```

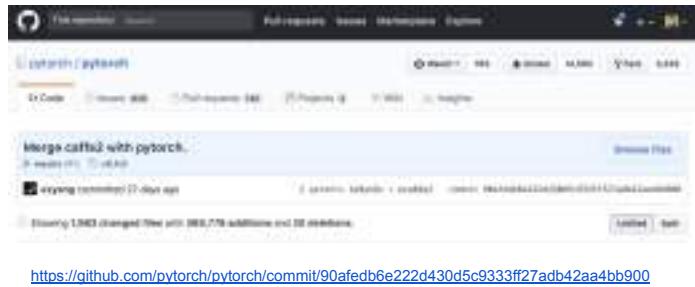
After exporting to ONNX, can run the PyTorch model in Caffe2

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 151 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 152 April 26, 2018

Static PyTorch: Future???



PyTorch vs TensorFlow, Static vs Dynamic

PyTorch
Dynamic Graphs
Static: ONNX, Caffe2

TensorFlow
Static Graphs
Dynamic: Eager

My Advice:

PyTorch is my personal favorite. Clean API, dynamic graphs make it very easy to develop and debug. Can build model in PyTorch then export to Caffe2 with ONNX for production / mobile

TensorFlow is a safe bet for most projects. Not perfect but has huge community, wide usage. Can use same framework for research and production. Probably use a high-level framework. Only choice if you want to run on TPUs.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 153 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 154 April 26, 2018

Next Time:
CNN Architecture Case Studies

Lecture 9: CNN Architectures

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 8 - 155 April 26, 2018 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 1 May 1, 2018

Administrative

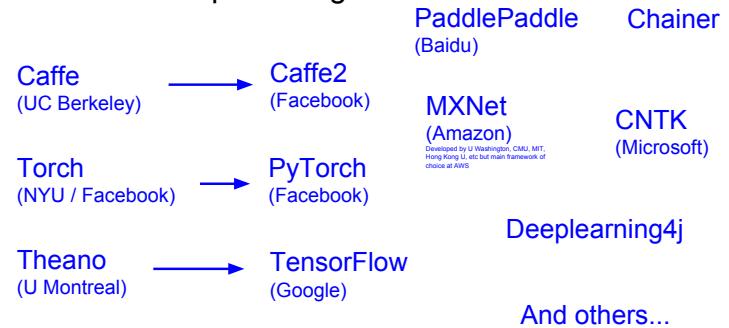
A2 due Wed May 2

Midterm: In-class Tue May 8. Covers material through Lecture 10 (Thu May 3).

Sample midterm released on piazza.

Midterm review session: Fri May 4 discussion section

Last time: Deep learning frameworks



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 2

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 3

May 1, 2018

Last time: Deep learning frameworks

- (1) Easily build big computational graphs
- (2) Easily compute gradients in computational graphs
- (3) Run it all efficiently on GPU (wrap cuDNN, cuBLAS, etc)

Today: CNN Architectures

Case Studies

- AlexNet
- VGG
- GoogLeNet
- ResNet

Also....

- NiN (Network in Network)
- Wide ResNet
- ResNeXT
- Stochastic Depth
- Squeeze-and-Excitation Network
- DenseNet
- FractalNet
- SqueezeNet
- NASNet

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 4

May 1, 2018

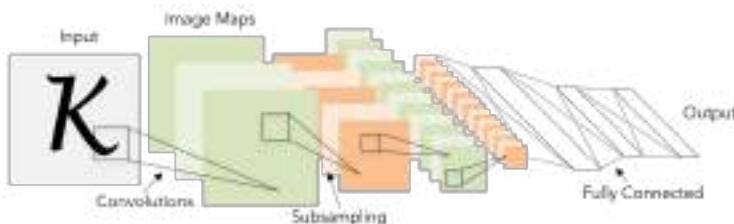
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 5

May 1, 2018

Review: LeNet-5

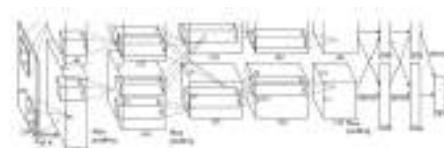
[LeCun et al., 1998]



Conv filters were 5x5, applied at stride 1
Subsampling (Pooling) layers were 2x2 applied at stride 2
i.e. architecture is [CONV-POOL-CONV-POOL-FC-FC]

Case Study: AlexNet

[Krizhevsky et al. 2012]



Architecture:

- CONV1
- MAX POOL1
- NORM1
- CONV2
- MAX POOL2
- NORM2
- CONV3
- CONV4
- CONV5
- Max POOL3
- FC6
- FC7
- FC8

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 6

May 1, 2018

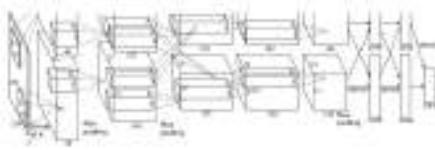
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 7

May 1, 2018

Case Study: AlexNet

[Krizhevsky et al. 2012]



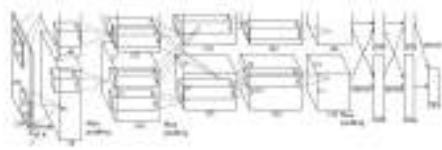
Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4

=>
Q: what is the output volume size? Hint: $(227-11)/4+1 = 55$

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4

=>
Output volume **[55x55x96]**

Q: What is the total number of parameters in this layer?

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

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

Fei-Fei Li & Justin Johnson & Serena Yeung

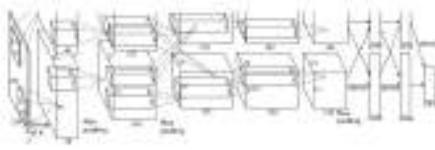
Lecture 9 - 8 May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 9 May 1, 2018

Case Study: AlexNet

[Krizhevsky et al. 2012]



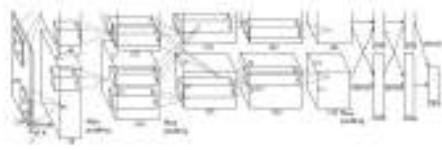
Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4

=>
Output volume **[55x55x96]**
Parameters: $(11 \times 11 \times 3) \times 96 = 35K$

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images
After CONV1: 55x55x96

Second layer (POOL1): 3x3 filters applied at stride 2

Q: what is the output volume size? Hint: $(55-3)/2+1 = 27$

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

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

Fei-Fei Li & Justin Johnson & Serena Yeung

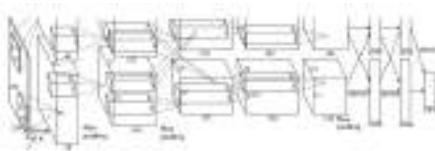
Lecture 9 - 10 May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 11 May 1, 2018

Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images
After CONV1: 55x55x96

Second layer (POOL1): 3x3 filters applied at stride 2
Output volume: 27x27x96

Q: what is the number of parameters in this layer?

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

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

Fei-Fei Li & Justin Johnson & Serena Yeung

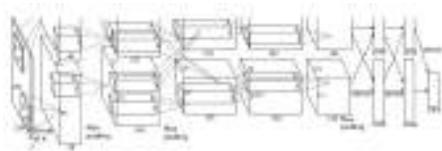
Lecture 9 - 12 May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 13 May 1, 2018

Case Study: AlexNet

[Krizhevsky et al. 2012]

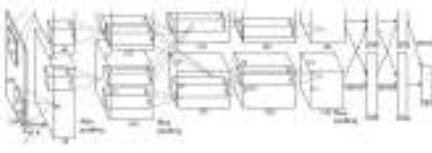


Input: 227x227x3 images
After CONV1: 55x55x96

Second layer (POOL1): 3x3 filters applied at stride 2
Output volume: 27x27x96
Parameters: 0!

Case Study: AlexNet

[Krizhevsky et al. 2012]



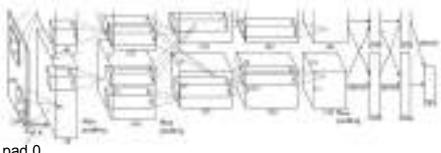
Input: 227x227x3 images
After CONV1: 55x55x96
After POOL1: 27x27x96

...

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

Case Study: AlexNet

[Krizhevsky et al. 2012]



Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
[1000] FC8: 1000 neurons (class scores)

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 14

May 1, 2018

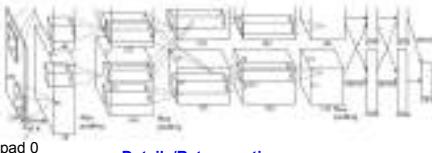
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 15

May 1, 2018

Case Study: AlexNet

[Krizhevsky et al. 2012]



Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
[1000] FC8: 1000 neurons (class scores)

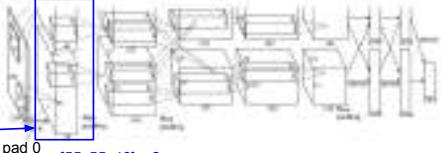
Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

Case Study: AlexNet

[Krizhevsky et al. 2012]



Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
[1000] FC8: 1000 neurons (class scores)

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

Historical note: Trained on GTX 580 GPU with only 3 GB of memory. Network spread across 2 GPUs, half the neurons (feature maps) on each GPU.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 16

May 1, 2018

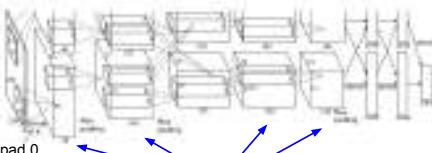
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 17

May 1, 2018

Case Study: AlexNet

[Krizhevsky et al. 2012]



Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
[1000] FC8: 1000 neurons (class scores)

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

CONV1, CONV2, CONV4, CONV5:
Connections only with feature maps
on same GPU

Case Study: AlexNet

[Krizhevsky et al. 2012]



Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
[1000] FC8: 1000 neurons (class scores)

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

CONV3, FC6, FC7, FC8:
Connections with all feature maps in
preceding layer, communication
across GPUs

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 18

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

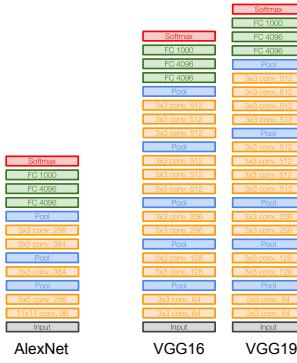
Lecture 9 - 19

May 1, 2018

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 26

May 1, 2018

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers
has same **effective receptive field** as
one 7x7 conv layer



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 26

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 27

May 1, 2018

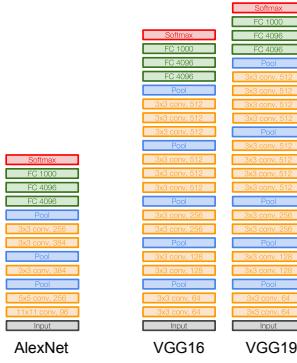
Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers
has same **effective receptive field** as
one 7x7 conv layer

[7x7]



Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers
has same **effective receptive field** as
one 7x7 conv layer

But deeper, more non-linearities

And fewer parameters: $3 * (3^2 C^2)$ vs.
 $7^2 C^2$ for C channels per layer



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 28

May 1, 2018

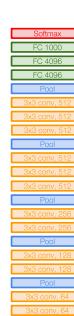
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 29

May 1, 2018

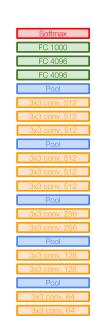
```

INPUT: [224x224x3] memory: 224*224*3=150K params: 0 (not counting biases)
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*64)*64 = 36,864
POOL2: [112x112x64] memory: 112*112*64=800K params: 0
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*128)*128 = 147,456
POOL2: [56x56x128] memory: 56*56*128=400K params: 0
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*128)*256 = 294,912
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824
POOL2: [28x28x256] memory: 28*28*256=200K params: 0
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*256)*512 = 1,179,648
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512] memory: 14*14*512=100K params: 0
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512] memory: 77*512=25K params: 0
FC: [1x1x4096] memory: 4096 params: 77*512*4096 = 102,760,448
FC: [1x1x4096] memory: 4096 params: 4096*4096 = 16,777,216
FC: [1x1x1000] memory: 1000 params: 4096*1000 = 4,096,000
  
```



```

INPUT: [224x224x3] memory: 224*224*3=150K params: 0 (not counting biases)
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*64)*64 = 36,864
POOL2: [112x112x64] memory: 112*112*64=800K params: 0
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*128)*128 = 147,456
POOL2: [56x56x128] memory: 56*56*128=400K params: 0
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*128)*256 = 294,912
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824
POOL2: [28x28x256] memory: 28*28*256=200K params: 0
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*256)*512 = 1,179,648
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512] memory: 14*14*512=100K params: 0
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512] memory: 77*512=25K params: 0
FC: [1x1x4096] memory: 4096 params: 77*512*4096 = 102,760,448
FC: [1x1x4096] memory: 4096 params: 4096*4096 = 16,777,216
FC: [1x1x1000] memory: 1000 params: 4096*1000 = 4,096,000
  
```



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 30

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

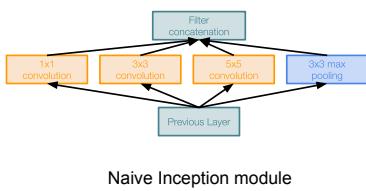
Lecture 9 - 31

May 1, 2018

TOTAL memory: 24M * 4 bytes ~ 96MB / image (for a forward pass)
TOTAL params: 138M parameters

Case Study: GoogLeNet

[Szegedy et al., 2014]



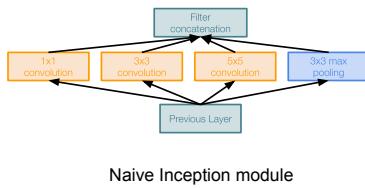
Apply parallel filter operations on the input from previous layer:

- Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)
- Pooling operation (3x3)

Concatenate all filter outputs together depth-wise

Case Study: GoogLeNet

[Szegedy et al., 2014]



Apply parallel filter operations on the input from previous layer:

- Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)
- Pooling operation (3x3)

Concatenate all filter outputs together depth-wise

Q: What is the problem with this?
[Hint: Computational complexity]

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 38

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

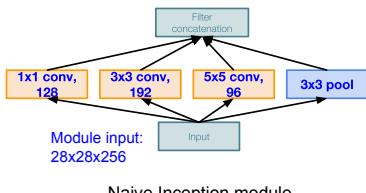
Lecture 9 - 39

May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

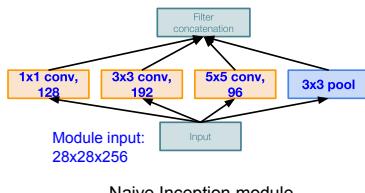


Q: What is the problem with this?
[Hint: Computational complexity]

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example: Q1: What is the output size of the 1x1 conv, with 128 filters?



Q: What is the problem with this?
[Hint: Computational complexity]

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 40

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 41

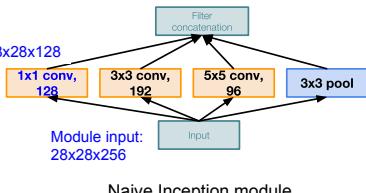
May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q1: What is the output size of the 1x1 conv, with 128 filters?

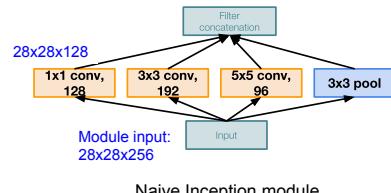


Q: What is the problem with this?
[Hint: Computational complexity]

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example: Q2: What are the output sizes of all different filter operations?



Q: What is the problem with this?
[Hint: Computational complexity]

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 42

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 43

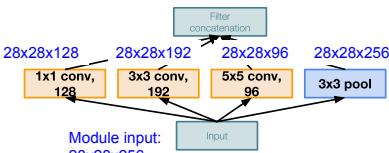
May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q2: What are the output sizes of all different filter operations?



Naive Inception module

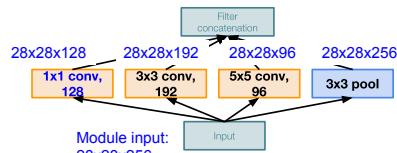
Q: What is the problem with this?
[Hint: Computational complexity]

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q3: What is output size after filter concatenation?



Naive Inception module

Q: What is the problem with this?
[Hint: Computational complexity]

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 44

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 45

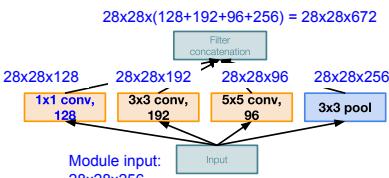
May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q3: What is output size after filter concatenation?



Naive Inception module

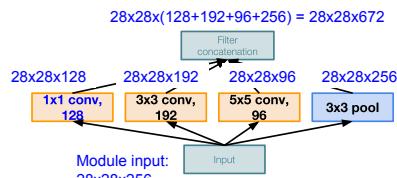
Q: What is the problem with this?
[Hint: Computational complexity]

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q3: What is output size after filter concatenation?



Naive Inception module

Q: What is the problem with this?
[Hint: Computational complexity]

Conv Ops:

[1x1 conv, 128] $28 \times 28 \times 128 \times 1 \times 1 \times 256$
 [3x3 conv, 192] $28 \times 28 \times 192 \times 3 \times 3 \times 256$
 [5x5 conv, 96] $28 \times 28 \times 96 \times 5 \times 5 \times 256$

Total: 854M ops

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 46

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 47

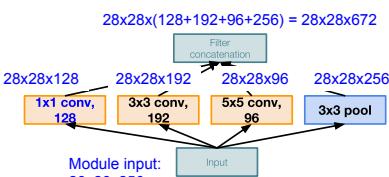
May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q3: What is output size after filter concatenation?



Naive Inception module

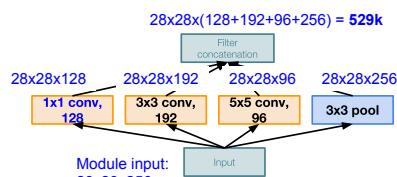
Q: What is the problem with this?
[Hint: Computational complexity]

Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

Q3: What is output size after filter concatenation?



Naive Inception module

Q: What is the problem with this?
[Hint: Computational complexity]

Solution: "bottleneck" layers that use 1x1 convolutions to reduce feature depth

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 48

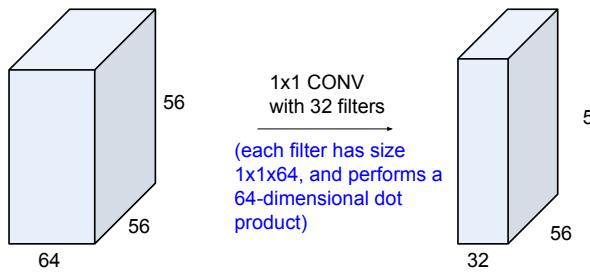
May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

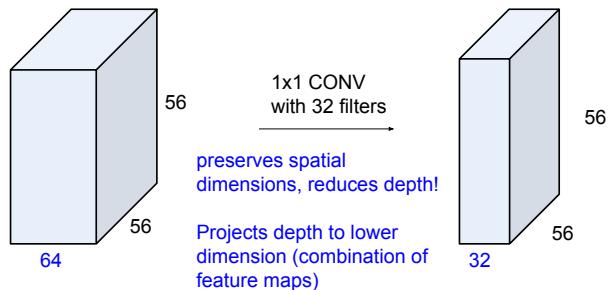
Lecture 9 - 49

May 1, 2018

Reminder: 1x1 convolutions



Reminder: 1x1 convolutions



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 50

May 1, 2018

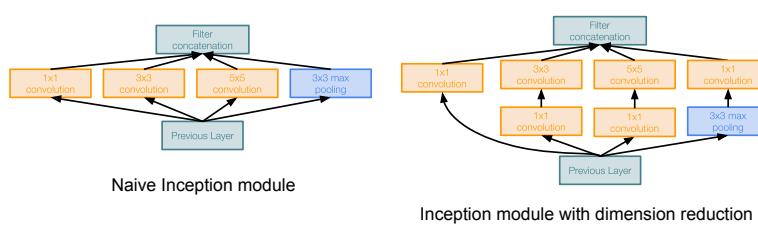
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 51

May 1, 2018

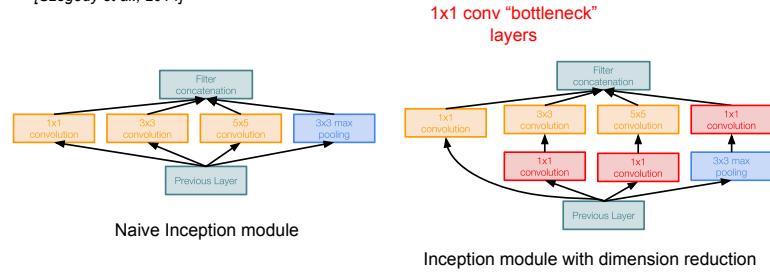
Case Study: GoogLeNet

[Szegedy et al., 2014]



Case Study: GoogLeNet

[Szegedy et al., 2014]



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 52

May 1, 2018

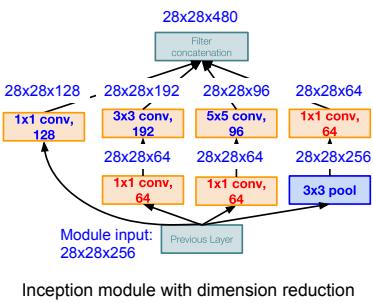
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 53

May 1, 2018

Case Study: GoogLeNet

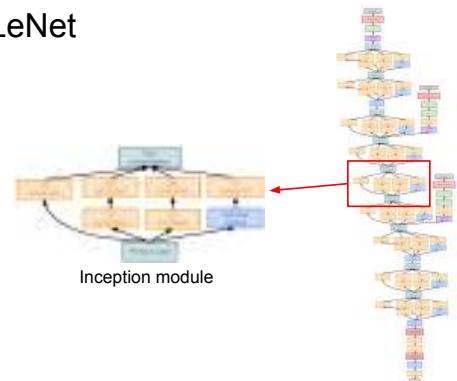
[Szegedy et al., 2014]



Case Study: GoogLeNet

[Szegedy et al., 2014]

Stack Inception modules with dimension reduction on top of each other



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 54

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

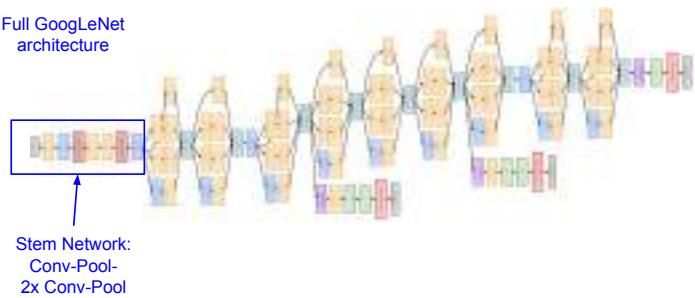
Lecture 9 - 55

May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet architecture



Stem Network:
Conv-Pool-
2x Conv-Pool

Fei-Fei Li & Justin Johnson & Serena Yeung

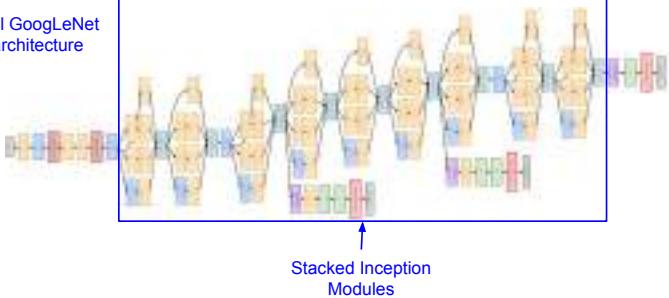
Lecture 9 - 56

May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet architecture



Stacked Inception
Modules

Fei-Fei Li & Justin Johnson & Serena Yeung

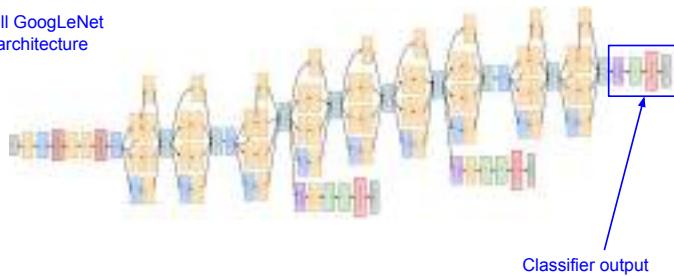
Lecture 9 - 57

May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet architecture



Classifier output

Fei-Fei Li & Justin Johnson & Serena Yeung

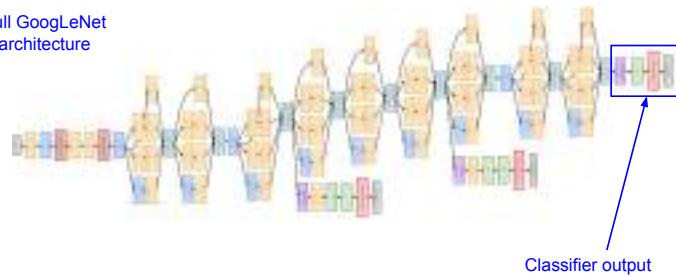
Lecture 9 - 58

May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet architecture



Classifier output
(removed expensive FC layers!)

Fei-Fei Li & Justin Johnson & Serena Yeung

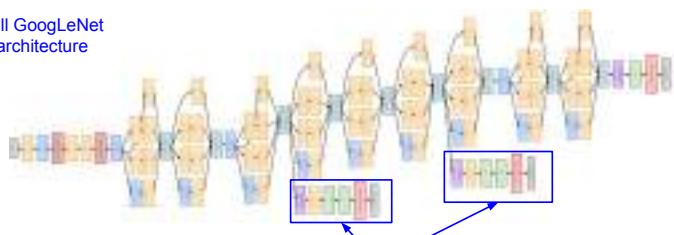
Lecture 9 - 59

May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet architecture



Auxiliary classification outputs to inject additional gradient at lower layers
(AvgPool-1x1Conv-FC-FC-Softmax)

Fei-Fei Li & Justin Johnson & Serena Yeung

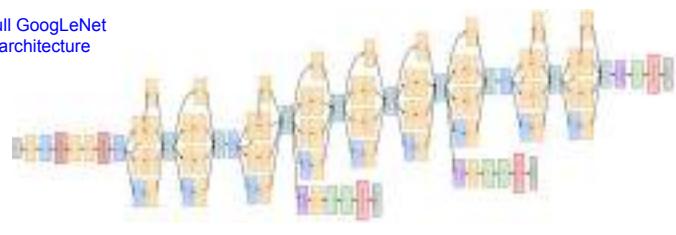
Lecture 9 - 60

May 1, 2018

Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet architecture



22 total layers with weights
(parallel layers count as 1 layer => 2 layers per Inception module. Don't count auxiliary output layers)

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 61

May 1, 2018

Case Study: ResNet

[He et al., 2015]

Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

Case Study: ResNet

[He et al., 2015]

Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

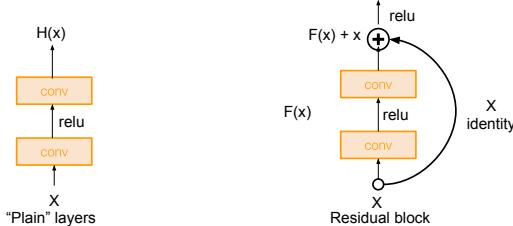
The deeper model should be able to perform at least as well as the shallower model.

A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping.

Case Study: ResNet

[He et al., 2015]

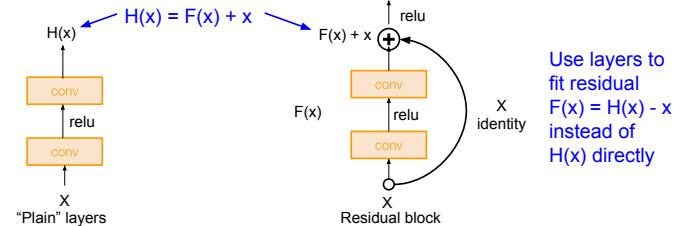
Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



Case Study: ResNet

[He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

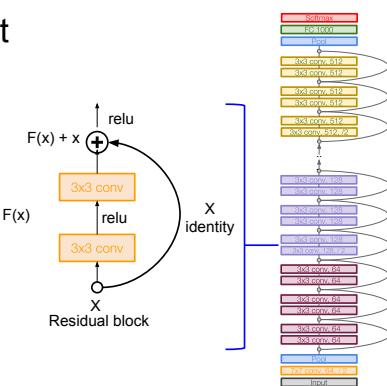


Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers

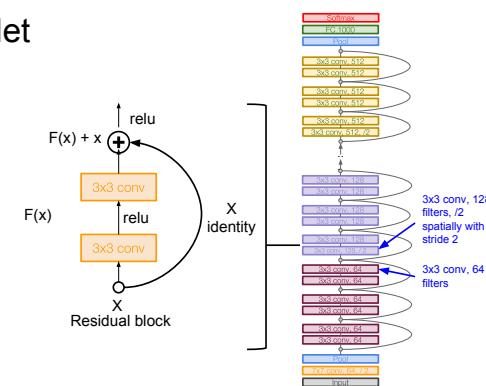


Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)

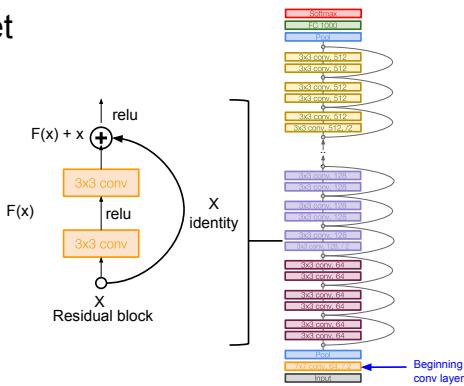


Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
- Additional conv layer at the beginning



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 74

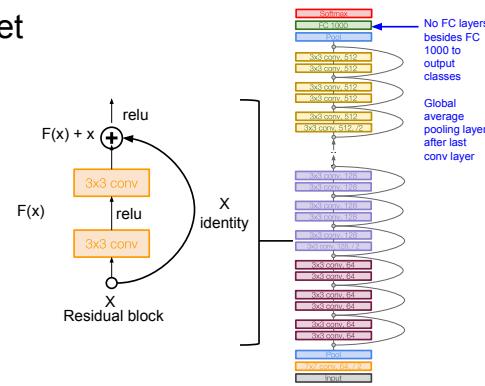
May 1, 2018

Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
- Additional conv layer at the beginning
- No FC layers at the end (only FC 1000 to output classes)



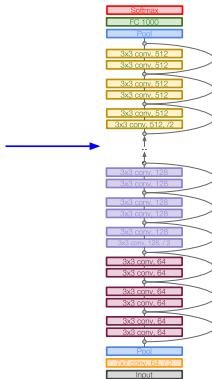
Lecture 9 - 75

May 1, 2018

Case Study: ResNet

[He et al., 2015]

Total depths of 34, 50, 101, or 152 layers for ImageNet



Fei-Fei Li & Justin Johnson & Serena Yeung

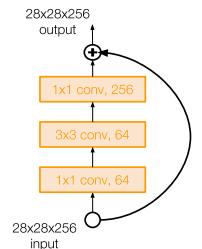
Lecture 9 - 76

May 1, 2018

Case Study: ResNet

[He et al., 2015]

For deeper networks (ResNet-50+), use "bottleneck" layer to improve efficiency (similar to GoogLeNet)



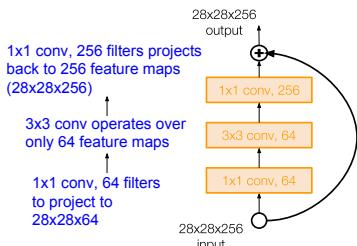
Lecture 9 - 77

May 1, 2018

Case Study: ResNet

[He et al., 2015]

For deeper networks (ResNet-50+), use "bottleneck" layer to improve efficiency (similar to GoogLeNet)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 78

May 1, 2018

Case Study: ResNet

[He et al., 2015]

Training ResNet in practice:

- Batch Normalization after every CONV layer
- Xavier 2/ initialization from He et al.
- SGD + Momentum (0.9)
- Learning rate: 0.1, divided by 10 when validation error plateaus
- Mini-batch size 256
- Weight decay of 1e-5
- No dropout used

Lecture 9 - 79

May 1, 2018

Case Study: ResNet

[He et al., 2015]

Experimental Results

- Able to train very deep networks without degrading (152 layers on ImageNet, 1202 on Cifar)
- Deeper networks now achieve lower training error as expected
- Swept 1st place in all ILSVRC and COCO 2015 competitions

MSRA @ ILSVRC & COCO 2015 Competitions

- + top places in all five main tracks
- ImageNet Classification: "After being" -- better than 100k-layer nets
- ImageNet Detection: 10% better than 2nd
- ImageNet Localization: 37% better than 2nd
- COCO Detection: 13% better than 2nd
- COCO Localization: 13% better than 2nd

Case Study: ResNet

[He et al., 2015]

Experimental Results

- Able to train very deep networks without degrading (152 layers on ImageNet, 1202 on Cifar)
- Deeper networks now achieve lower training error as expected
- Swept 1st place in all ILSVRC and COCO 2015 competitions

MSRA @ ILSVRC & COCO 2015 Competitions

- + top places in all five main tracks
- ImageNet Classification: "After being" -- better than 100k-layer nets
- ImageNet Detection: 10% better than 2nd
- ImageNet Localization: 37% better than 2nd
- COCO Detection: 13% better than 2nd
- COCO Localization: 13% better than 2nd

ILSVRC 2015 classification winner (3.6% top 5 error) -- better than "human performance"! (Russakovsky 2014)

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 80

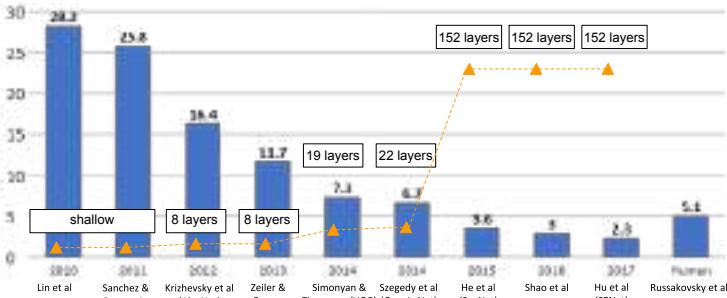
May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

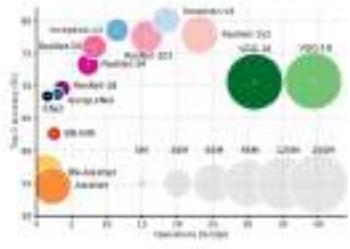
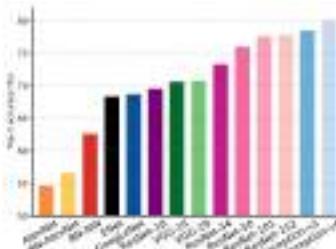
Lecture 9 - 81

May 1, 2018

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Comparing complexity...



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 82

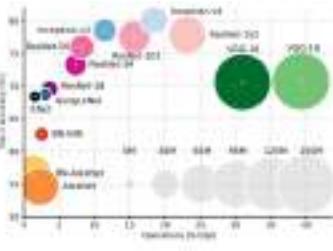
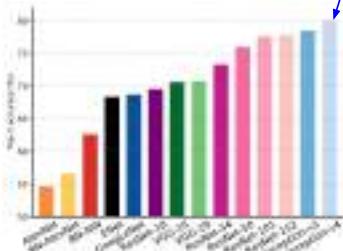
May 1, 2018

Lecture 9 - 83

May 1, 2018

Comparing complexity...

Inception-v4: Resnet + Inception!

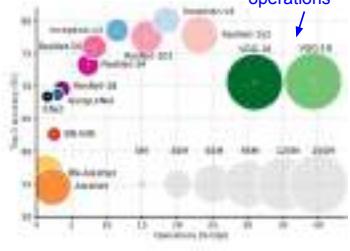
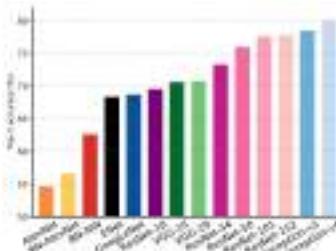


An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

Comparing complexity...

VGG: Highest memory, most operations



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

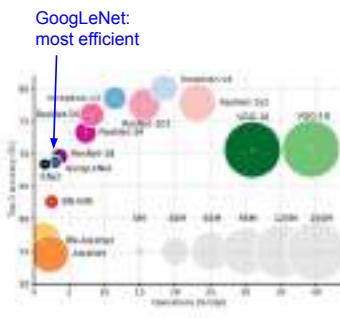
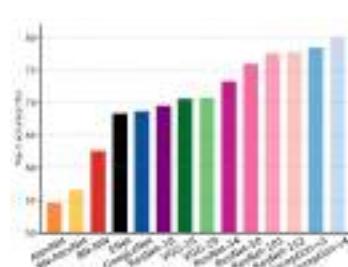
Lecture 9 - 84

May 1, 2018

Lecture 9 - 85

May 1, 2018

Comparing complexity...



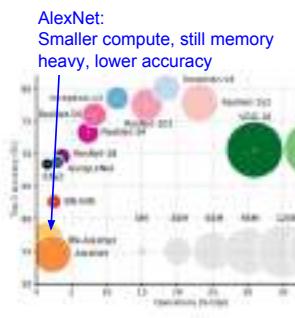
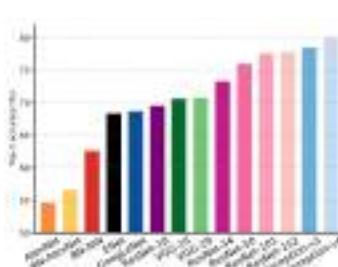
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

Lecture 9 - 86

May 1, 2018

Comparing complexity...



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

Lecture 9 - 86

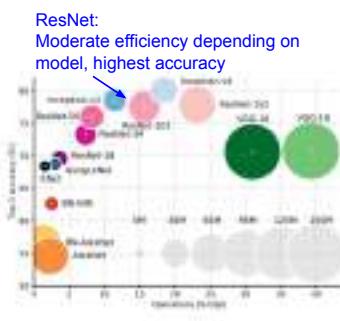
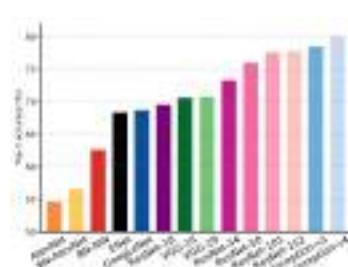
May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 87

May 1, 2018

Comparing complexity...



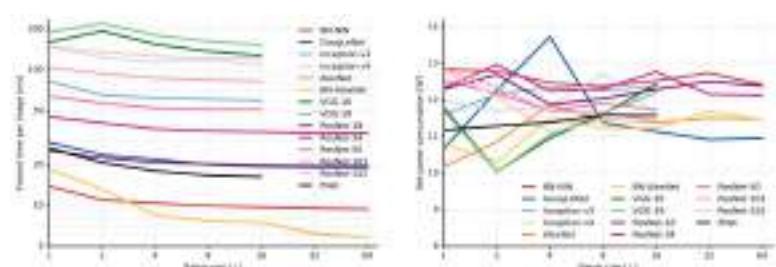
An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

Lecture 9 - 88

May 1, 2018

Forward pass time and power consumption



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

Lecture 9 - 88

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 89

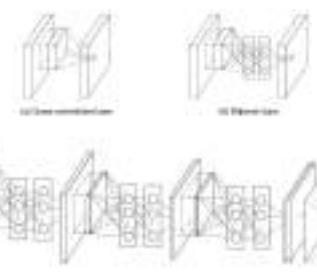
May 1, 2018

Other architectures to know...

Network in Network (NiN)

[Lin et al. 2014]

- Mlpconv layer with "micronetwork" within each conv layer to compute more abstract features for local patches
- Micronetwork uses multilayer perceptron (FC, i.e. 1x1 conv layers)
- Precursor to GoogLeNet and ResNet "bottleneck" layers
- Philosophical inspiration for GoogLeNet



Figures copyright Lin et al., 2014. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 90

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 91

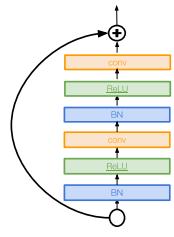
May 1, 2018

Improving ResNets...

Identity Mappings in Deep Residual Networks

[He et al. 2016]

- Improved ResNet block design from creators of ResNet
- Creates a more direct path for propagating information throughout network (moves activation to residual mapping pathway)
- Gives better performance

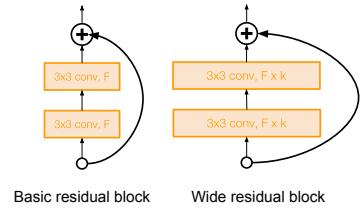


Improving ResNets...

Wide Residual Networks

[Zagoruyko et al. 2016]

- Argues that residuals are the important factor, not depth
- Use wider residual blocks ($F \times k$ filters instead of F filters in each layer)
- 50-layer wide ResNet outperforms 152-layer original ResNet
- Increasing width instead of depth more computationally efficient (parallelizable)

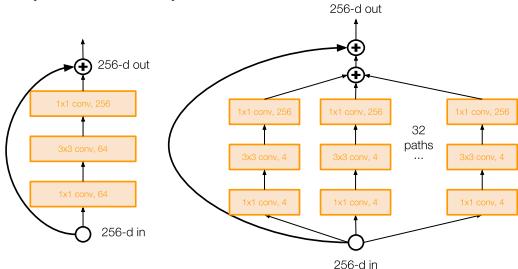


Improving ResNets...

Aggregated Residual Transformations for Deep Neural Networks (ResNeXt)

[Xie et al. 2016]

- Also from creators of ResNet
- Increases width of residual block through multiple parallel pathways ("cardinality")
- Parallel pathways similar in spirit to Inception module

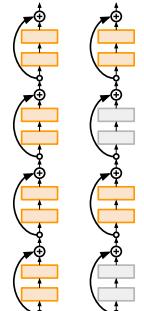


Improving ResNets...

Deep Networks with Stochastic Depth

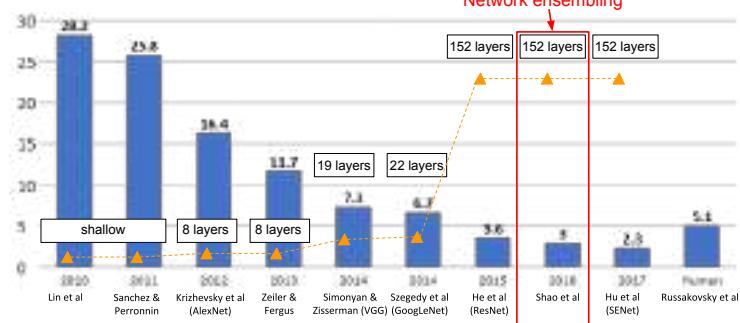
[Huang et al. 2016]

- Motivation: reduce vanishing gradients and training time through short networks during training
- Randomly drop a subset of layers during each training pass
- Bypass with identity function
- Use full deep network at test time



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

Network ensembling



Improving ResNets...

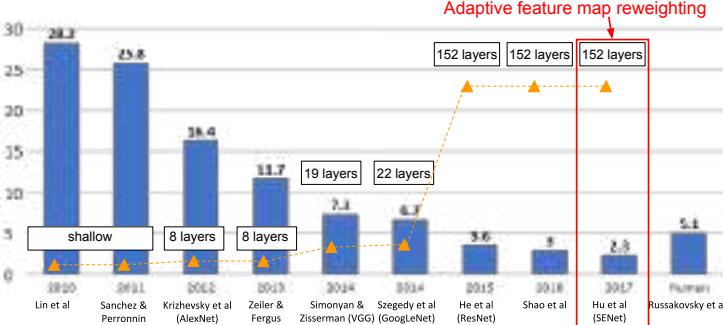
"Good Practices for Deep Feature Fusion"

[Shao et al. 2016]

- Multi-scale ensembling of Inception, Inception-Resnet, Resnet, Wide Resnet models
- ILSVRC'16 classification winner



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 98

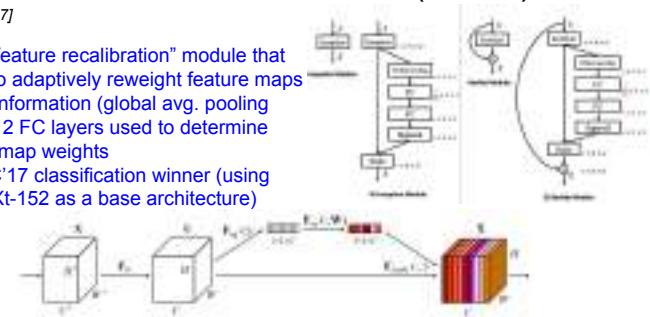
May 1, 2018

Improving ResNets...

Squeeze-and-Excitation Networks (SENet)

[Hu et al. 2017]

- Add a “feature recalibration” module that learns to adaptively reweight feature maps
- Global information (global avg. pooling layer) + 2 FC layers used to determine feature map weights
- ILSVRC’17 classification winner (using ResNeXt-152 as a base architecture)



Lecture 9 - 99

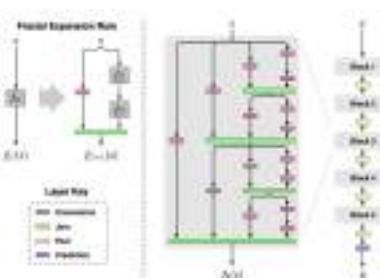
May 1, 2018

Beyond ResNets...

FractalNet: Ultra-Deep Neural Networks without Residuals

[Larsson et al. 2017]

- Argues that key is transitioning effectively from shallow to deep and residual representations are not necessary
- Fractal architecture with both shallow and deep paths to output
- Trained with dropping out sub-paths
- Full network at test time



Figures copyright Larsson et al., 2017. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 10

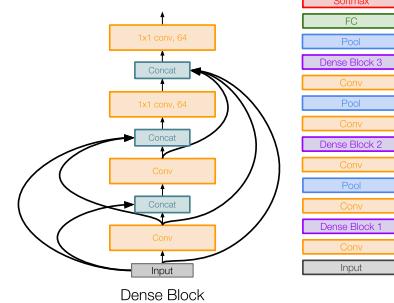
May 1, 2018

Beyond ResNets...

Densely Connected Convolutional Networks

[Huang et al. 2017]

- Dense blocks where each layer is connected to every other layer in feedforward fashion
- Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse



Lecture 9 -

May 1, 2018

Efficient networks...

SqueezeNet: AlexNet-level Accuracy With 50x Fewer Parameters and <0.5Mb Model Size

[Iandola et al. 2017]

- Fire modules consisting of a ‘squeeze’ layer with 1x1 filters feeding an ‘expand’ layer with 1x1 and 3x3 filters
- AlexNet level accuracy on ImageNet with 50x fewer parameters
- Can compress to 510x smaller than AlexNet (0.5Mb)

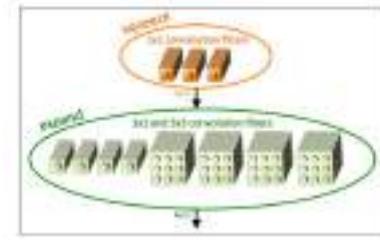


Figure copyright Iandola, Han, Moskewicz, Ashraf, Dally, Keutzer, 2017. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 -

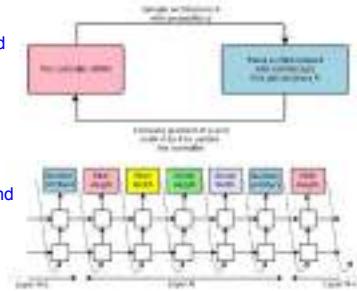
May 1, 2018

Meta-learning: Learning to learn network architectures...

Neural Architecture Search with Reinforcement Learning (NAS)

[Zoph et al. 2016]

- “Controller” network that learns to design a good network architecture (output a string corresponding to network design)
- Iterate:
 - 1) Sample an architecture from search space
 - 2) Train the architecture to get a “reward” R corresponding to accuracy
 - 3) Compute gradient of sample probability, and scale by R to perform controller parameter update (i.e. increase likelihood of good architecture being sampled, decrease likelihood of bad architecture)



Lecture 9 - 10
3

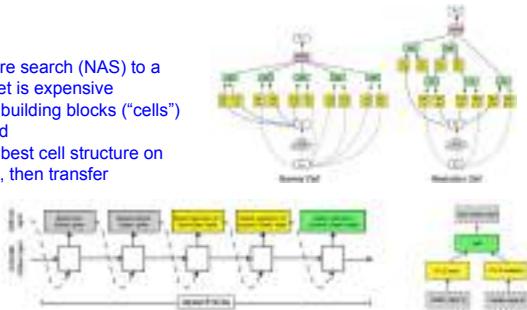
May 1, 2018

Meta-learning: Learning to learn network architectures...

Learning Transferable Architectures for Scalable Image Recognition

[Zoph et al. 2017]

- Applying neural architecture search (NAS) to a large dataset like ImageNet is expensive
- Design a search space of building blocks ("cells") that can be flexibly stacked
- NASNet: Use NAS to find best cell structure on smaller CIFAR-10 dataset, then transfer architecture to ImageNet



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 10
4

May 1, 2018

Summary: CNN Architectures

Case Studies

- AlexNet
- VGG
- GoogLeNet
- ResNet

Also....

- NiN (Network in Network)
- Wide ResNet
- ResNeXT
- Stochastic Depth
- Squeeze-and-Excitation Network
- DenseNet
- FractalNet
- SqueezeNet
- NASNet

Fei-Fei Li & Justin Johnson & Serena Yeung

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 10
5

May 1, 2018

Summary: CNN Architectures

- VGG, GoogLeNet, ResNet all in wide use, available in model zoos
- ResNet current best default, also consider SENet when available
- Trend towards extremely deep networks
- Significant research centers around design of layer / skip connections and improving gradient flow
- Efforts to investigate necessity of depth vs. width and residual connections
- Even more recent trend towards meta-learning
- Next time: Recurrent neural networks

Lecture 10: Recurrent Neural Networks

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 9 - 106

May 1, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 1

May 3, 2018

Administrative

A1 regrade deadline is tonight

A2 due yesterday

Redeem your Google Cloud coupons by Sunday 5/6

Administrative: Midterm

In-class midterm on Tuesday!
Details on Piazza

If you need an alternate exam time then let us know by tomorrow

Midterm review session tomorrow

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 2

May 3, 2018

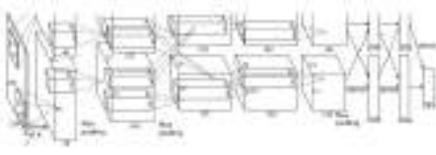
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 3

May 3, 2018

Last Time: CNN Architectures

AlexNet



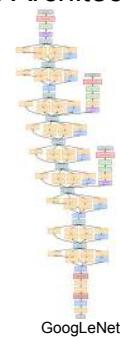
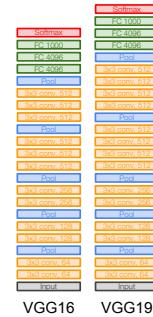
Resolution of Depth

Figure copyright Kaiming He, 2016. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 4 May 3, 2018

Last Time: CNN Architectures

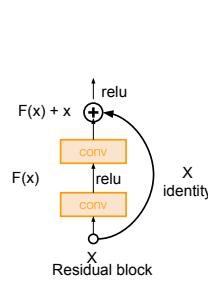


Resolution of Depth

Figure copyright Kaiming He, 2016. Reproduced with permission.

Lecture 10 - 5 May 3, 2018

Last Time: CNN Architectures



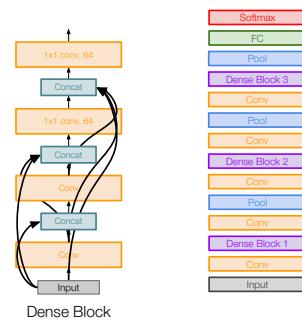
Resolution of Depth

Figure copyright Kaiming He, 2016. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

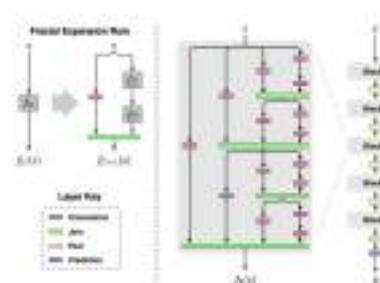
Lecture 10 - 6 May 3, 2018

DenseNet



Fei-Fei Li & Justin Johnson & Serena Yeung

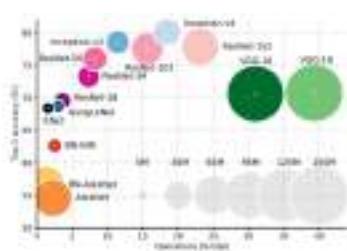
FractalNet



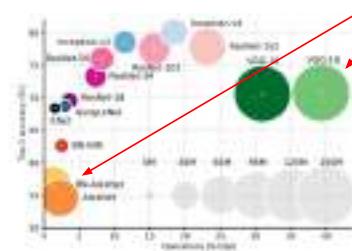
Figures copyright Larsson et al., 2017. Reproduced with permission.

Lecture 10 - 7 May 3, 2018

Last Time: CNN Architectures



Last Time: CNN Architectures



AlexNet and VGG have tons of parameters in the fully connected layers

AlexNet: ~62M parameters

FC6: 256x6x6 -> 4096: 38M params
FC7: 4096 -> 4096: 17M params
FC8: 4096 -> 1000: 4M params
~59M params in FC layers!

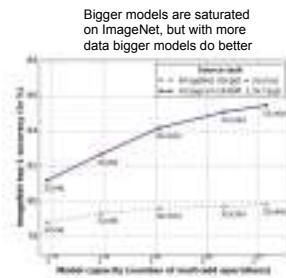
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 8 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 9 May 3, 2018

ImageNet pretraining -> Instagram pretraining

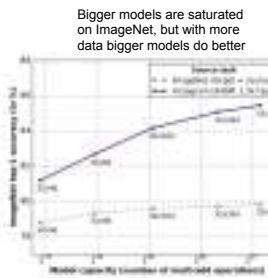


Biggest network was pretrained on 3.5B Instagram images

Trained on 336 GPUs for 22 days

Mahajan et al., "Exploring the Limits of Weakly Supervised Pretraining", arXiv 2018

ImageNet pretraining -> Instagram pretraining



Biggest network was pretrained on 3.5B Instagram images

Trained on 336 GPUs for 22 days
≈ \$129,000 on Google Cloud

Mahajan et al., "Exploring the Limits of Weakly Supervised Pretraining", arXiv 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

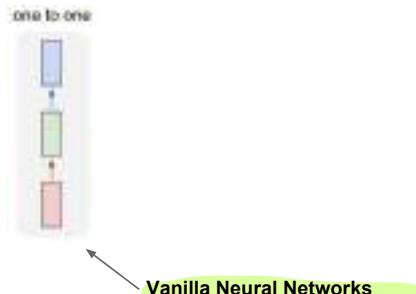
Lecture 10 - 10 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 11 May 3, 2018

Today: Recurrent Neural Networks

"Vanilla" Neural Network



Vanilla Neural Networks

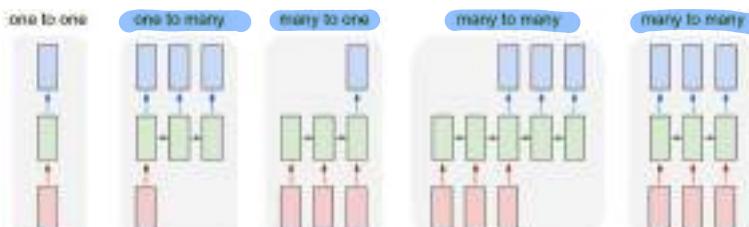
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 12 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

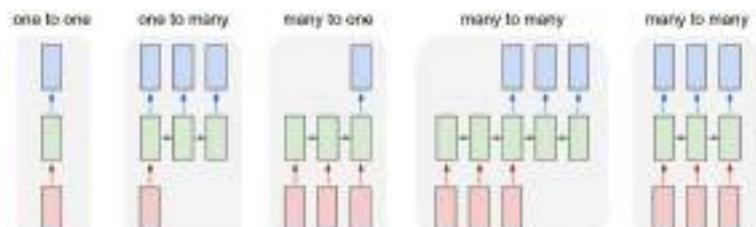
Lecture 10 - 13 May 3, 2018

Recurrent Neural Networks: Process Sequences



e.g. **Image Captioning**
image -> sequence of words

Recurrent Neural Networks: Process Sequences



e.g. **Sentiment Classification**
sequence of words -> sentiment

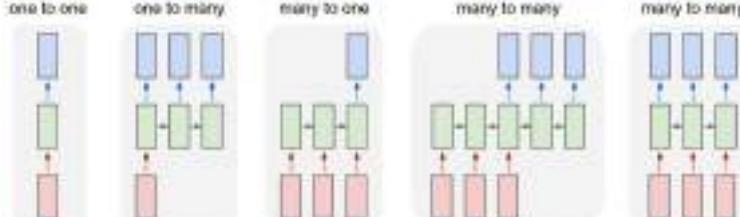
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 14 May 3, 2018

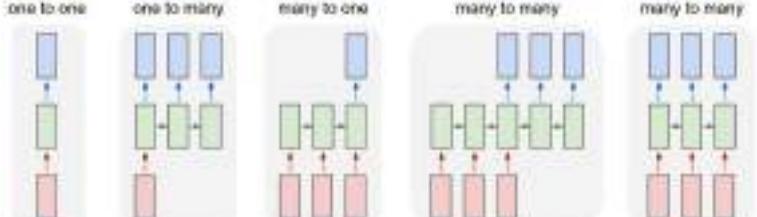
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 15 May 3, 2018

Recurrent Neural Networks: Process Sequences



Recurrent Neural Networks: Process Sequences



e.g. Machine Translation
seq of words -> seq of words

e.g. Video classification on frame level

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 16 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 17 May 3, 2018

Sequential Processing of Non-Sequence Data



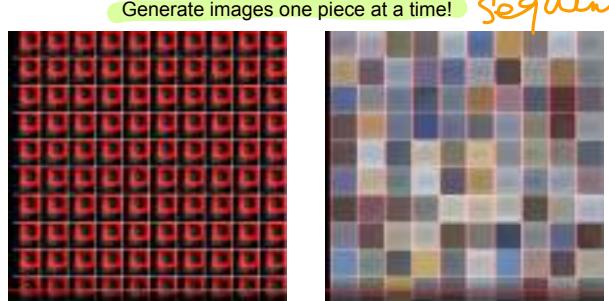
Classify images by taking a series of "glimpses"

Ba, Mihai, and Karol Gregor. "Multiple Object Recognition with Visual Attention". ICLR 2015.
Gregor et al., "DRAW: A Recurrent Neural Network For Image Generation", ICML 2015
Figure copyright Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra, 2015. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 18 May 3, 2018

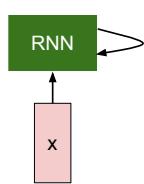
Sequential Processing of Non-Sequence Data



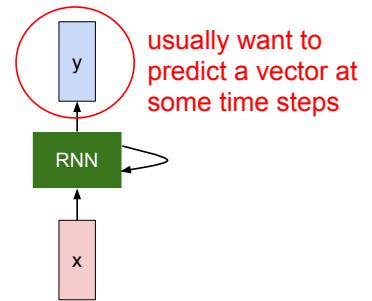
Gregor et al., "DRAW: A Recurrent neural network for image generation", ICLR, 2015
Figure copyright Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra, 2015. Reproduced with permission.

Lecture 10 - 19 May 3, 2018

Recurrent Neural Network



Recurrent Neural Network



Fei-Fei Li & Justin Johnson & Serena Yeung

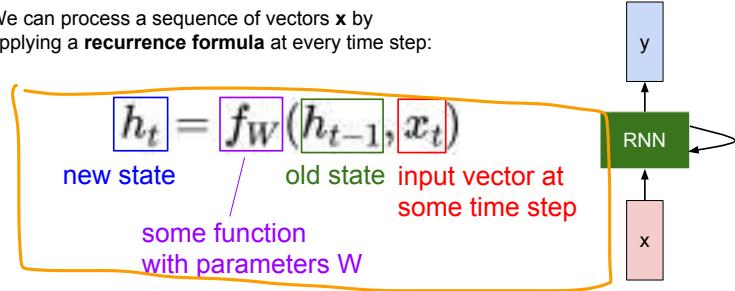
Lecture 10 - 20 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 21 May 3, 2018

Recurrent Neural Network

We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

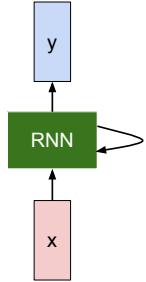


Recurrent Neural Network

We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

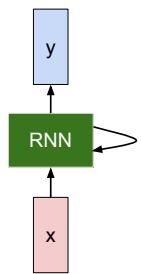
$$h_t = f_W(h_{t-1}, x_t)$$

Notice: the same function and the same set of parameters are used at every time step.



(Simple) Recurrent Neural Network

The state consists of a single "hidden" vector \mathbf{h} :

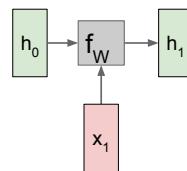


$$h_t = f_W(h_{t-1}, x_t)$$

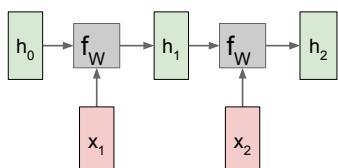
$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ y_t &= W_{hy}h_t \end{aligned}$$

Sometimes called a "Vanilla RNN" or an "Elman RNN" after Prof. Jeffrey Elman

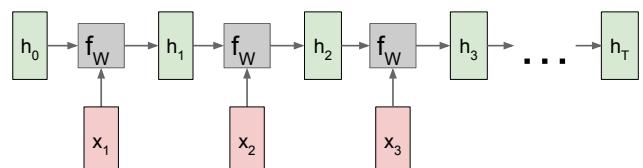
RNN: Computational Graph



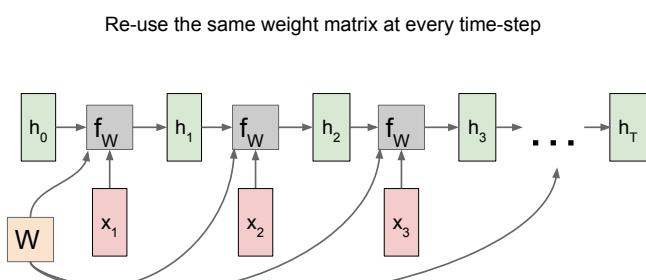
RNN: Computational Graph



RNN: Computational Graph



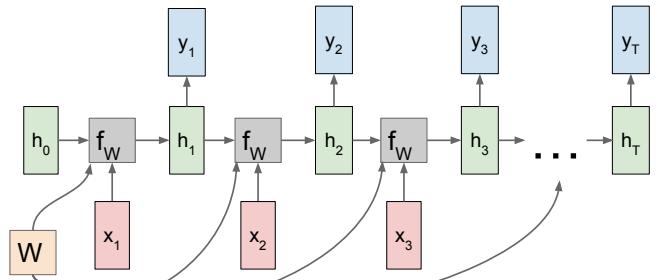
RNN: Computational Graph



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 28 May 3, 2018

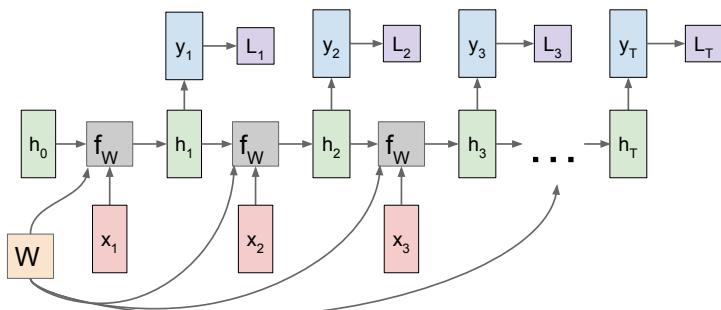
RNN: Computational Graph: Many to Many



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 29 May 3, 2018

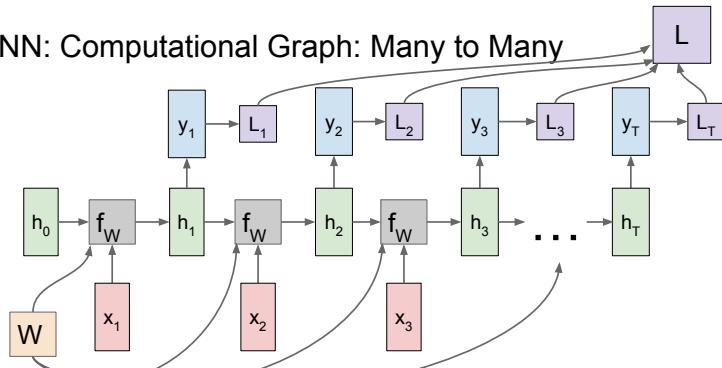
RNN: Computational Graph: Many to Many



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 30 May 3, 2018

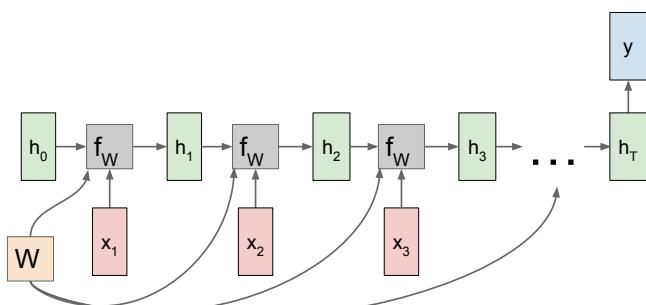
RNN: Computational Graph: Many to Many



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 31 May 3, 2018

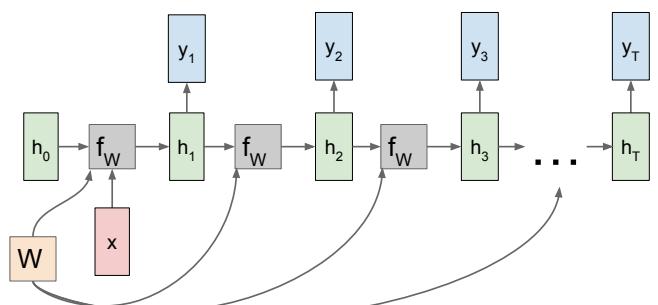
RNN: Computational Graph: Many to One



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 32 May 3, 2018

RNN: Computational Graph: One to Many

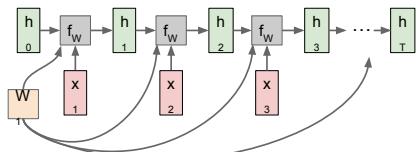


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 33 May 3, 2018

Sequence to Sequence: Many-to-one + one-to-many

Many to one: Encode input sequence in a single vector



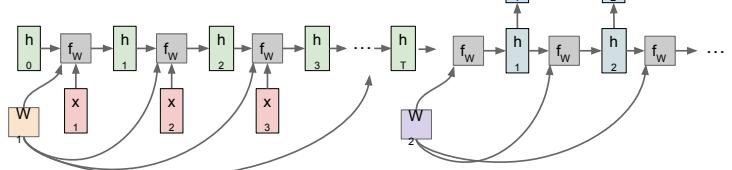
Sutskever et al., "Sequence to Sequence Learning with Neural Networks", NIPS 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 34 May 3, 2018

Sequence to Sequence: Many-to-one + one-to-many

Many to one: Encode input sequence in a single vector



One to many: Produce output sequence from single input vector

Sutskever et al., "Sequence to Sequence Learning with Neural Networks", NIPS 2014

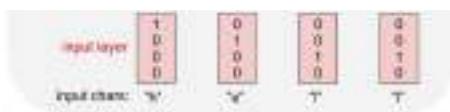
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 35 May 3, 2018

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
"hello"



Fei-Fei Li & Justin Johnson & Serena Yeung

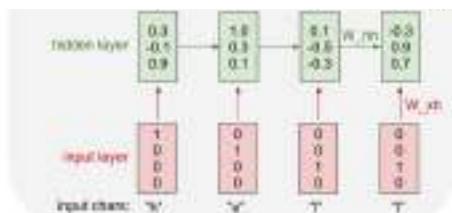
Lecture 10 - 36 May 3, 2018

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

Example training
sequence:
"hello"

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$



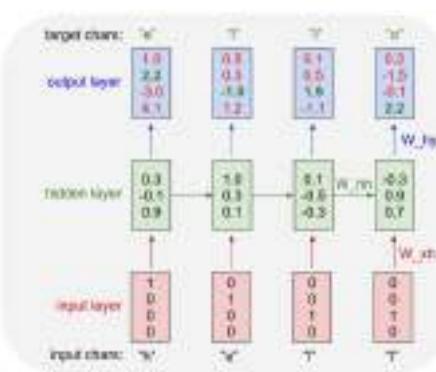
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 37 May 3, 2018

Example: Character-level Language Model

Vocabulary:
[h,e,l,o]

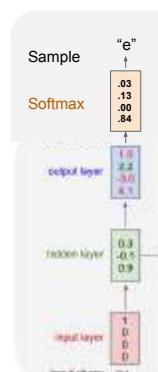
Example training
sequence:
"hello"



Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a time,
feed back to model



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 38 May 3, 2018

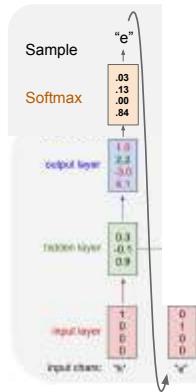
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 39 May 3, 2018

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a time,
feed back to model



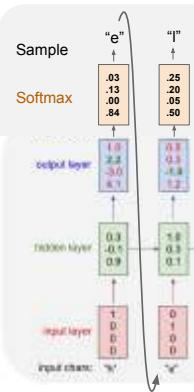
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 40 May 3, 2018

Example: Character-level Language Model Sampling

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a time,
feed back to model



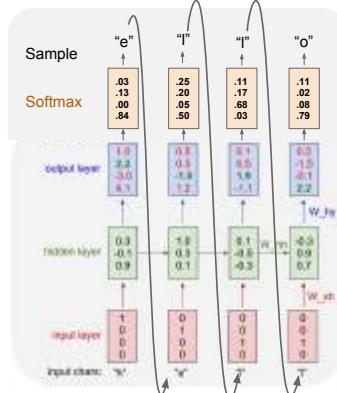
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 41 May 3, 2018

Example: Character-level Language Model Sampling

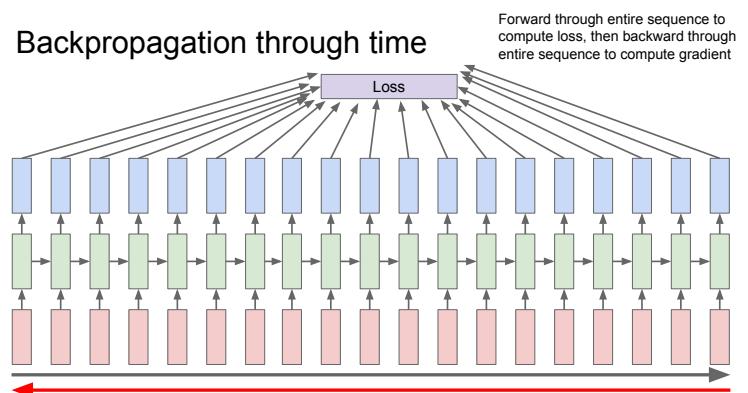
Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a time,
feed back to model



Lecture 10 - 42 May 3, 2018

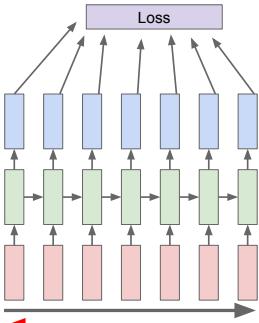
Backpropagation through time



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 43 May 3, 2018

Truncated Backpropagation through time

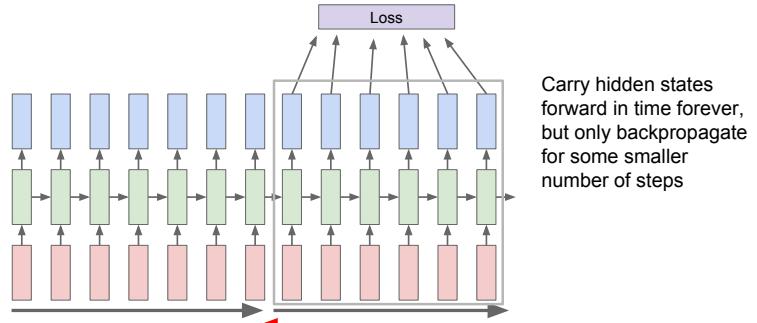


Run forward and backward
through chunks of the
sequence instead of whole
sequence

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 44 May 3, 2018

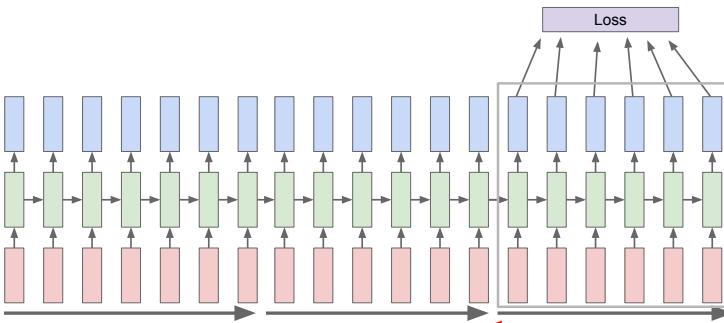
Truncated Backpropagation through time



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 45 May 3, 2018

Truncated Backpropagation through time



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 46 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 47 May 3, 2018

[min-char-rnn.py](https://gist.github.com/karpathy/d4dee566867f8291f086) gist: 112 lines of Python

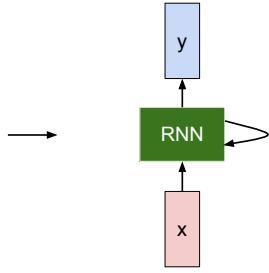


(<https://gist.github.com/karpathy/d4dee566867f8291f086>)

THE SONNETS

by William Shakespeare

Shall I compare thee to a summer's day?
Thou art more lovely and more temperate:
Rough winds do shake the darling buds of May,
And thy sweet sunne doth endevour them
To overset with winde and wat're.



at first:

"yon sonnete" (sonnet) (sonnet)
Shall I compare thee to a summer's day?
Thou art more lovely and more temperate:
Rough winds do shake the darling buds of May,

↓ train more

"Your tractay" (sonnet) (sonnet)
Shall I compare thee to a summer's day?
Thou art more lovely and more temperate:
Rough winds do shake the darling buds of May,

↓ train more

"Atkyt folc unscut that the batt tis Prince Willyam's that we st.
Her heire, and heis to be alwaies fyring were to tht belay, þow my fylling master.
Now, and Sogiton as so peccyal and other."

↓ train more

"Why do what that day?" replied WILLYAM, and waiting to answer the
princess, Princess MARY was easier, fed him and rymed him.
Pierre eating his soul, came to the meads and croon as his father-in-law wove.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 48 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 49 May 3, 2018

SONNETS
Shall I compare thee to a summer's day?
Thou art more lovely and more temperate:
Rough winds do shake the darling buds of May,
And thy sweet sunne doth endevour them
To overset with winde and wat're.

SECOND SONNET
They say unto me, "Thou art a good man, precentor upon my soul."
WILLIAM AND MARYLLA shold be ALLYED, when I precent
the earth and universall of mannes.

THIRD SONNET
Well, you say, as in the name of wife and child...

FOURTH SONNET
They make me eylet when this shold, and
ME TELL unto begin set of the path, to be changed.
When adde sonne, I'll have the breath of mannes.

FIFTH SONNET
Well, when I will make old bound your swetly.

SIXTH SONNET

VII
Shall I compare thee to a summer's day?
Thou art more lovely and more temperate:
Rough winds do shake the darling buds of May,
And thy sweet sunne doth endevour them
To overset with winde and wat're.

VIII
Shall I compare thee to a summer's day?
Thou art more lovely and more temperate:
Rough winds do shake the darling buds of May,
And thy sweet sunne doth endevour them
To overset with winde and wat're.

IX
Shall I compare thee to a summer's day?
Thou art more lovely and more temperate:
Rough winds do shake the darling buds of May,
And thy sweet sunne doth endevour them
To overset with winde and wat're.

X
Shall I compare thee to a summer's day?
Thou art more lovely and more temperate:
Rough winds do shake the darling buds of May,
And thy sweet sunne doth endevour them
To overset with winde and wat're.

The Stacks Project: open source algebraic geometry textbook

Part	Chapter	titles	Topics	View pdf
Preliminaries	1. Introduction	titles	topics	View pdf
	2. Commutative	titles	topics	View pdf
	3. Set Theory	titles	topics	View pdf
	4. Categories	titles	topics	View pdf
	5. Topology	titles	topics	View pdf
	6. Sheaves on Spaces	titles	topics	View pdf
	7. Sheaves and Sheaf Cohomology	titles	topics	View pdf
	8. Schemes	titles	topics	View pdf
	9. Fibrations	titles	topics	View pdf
	10. Commutative Algebra	titles	topics	View pdf

Latex source

<http://stacks.math.columbia.edu/>
The stacks project is licensed under the [GNU Free Documentation License](#)

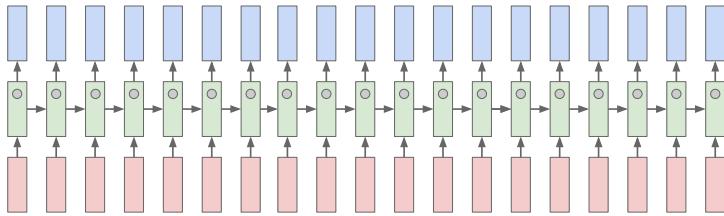
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 50 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 51 May 3, 2018

Searching for interpretable cells



Searching for interpretable cells



Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

Lecture 10 - 58 May 3, 2018

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 59 May 3, 2018

Searching for interpretable cells



quote detection cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

Searching for interpretable cells



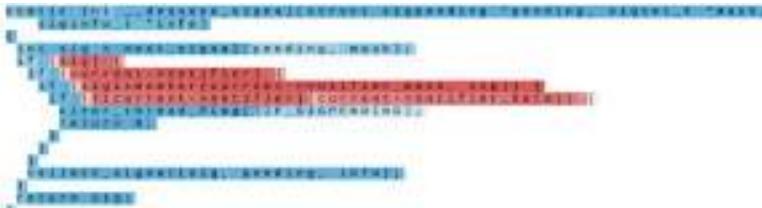
line length tracking cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 61 May 3, 2018

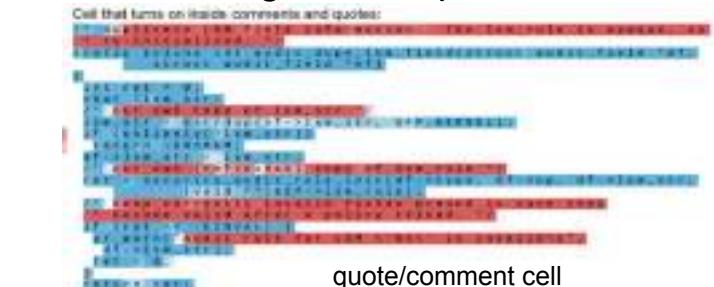
Searching for interpretable cells



if statement cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

Searching for interpretable cells



quote/comment cell

Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

Figures copyright Karpathy, Johnson, and Fei-Fei, 2016; reproduced with permission

Fei-Fei Li & Justin Johnson & Serena Yeung

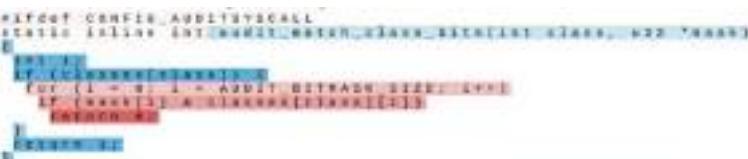
Lecture 10 - 62 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 63 May 3, 2018

Searching for interpretable cells

code depth cell



Karpathy, Johnson, and Fei-Fei: Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016
Premier Copyright Karpathy, Johnson, and Fei-Fei, 2016. Reprinted with permission.

Image Captioning

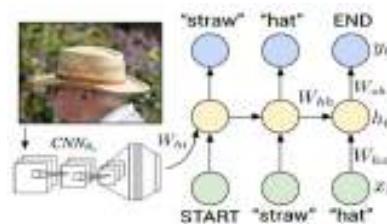


Figure from Karpathy et al., "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015. ©2015 IEEE. Reproduced for educational purposes.

Explain Images with Multimodal Recurrent Neural Networks, Mao et al.
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei
Show and Tell: A Neural Image Caption Generator, Vinyals et al.
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

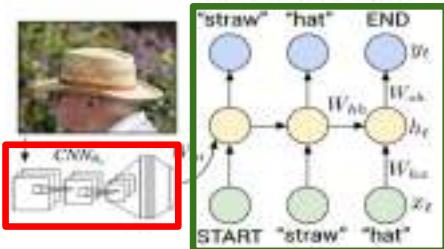
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 64 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 65 May 3, 2018

Recurrent Neural Network



test image

Convolutional Neural Network

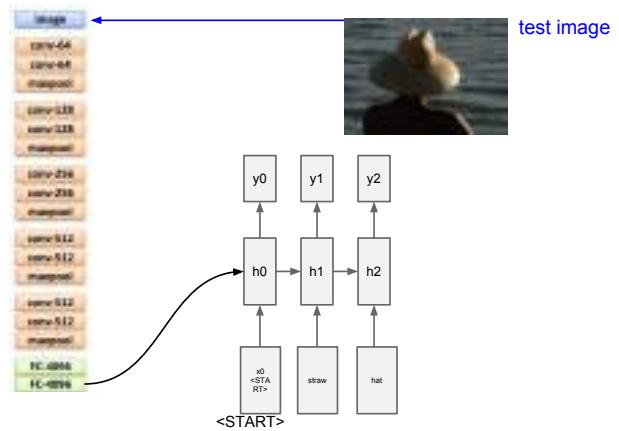
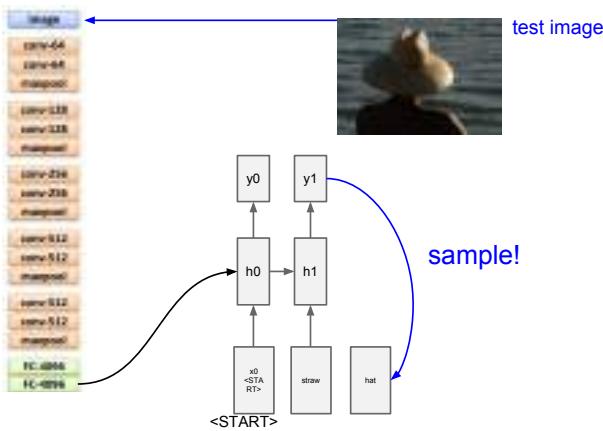
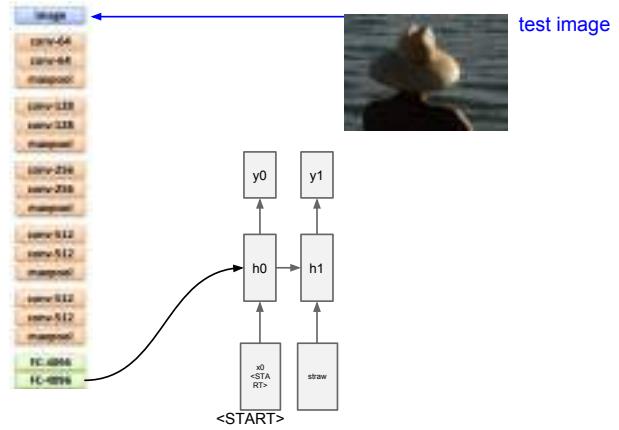
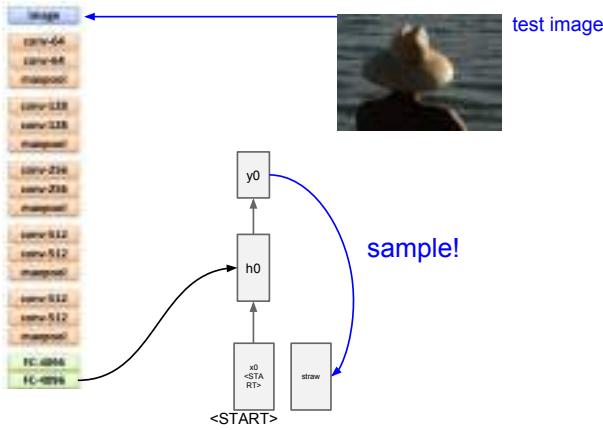
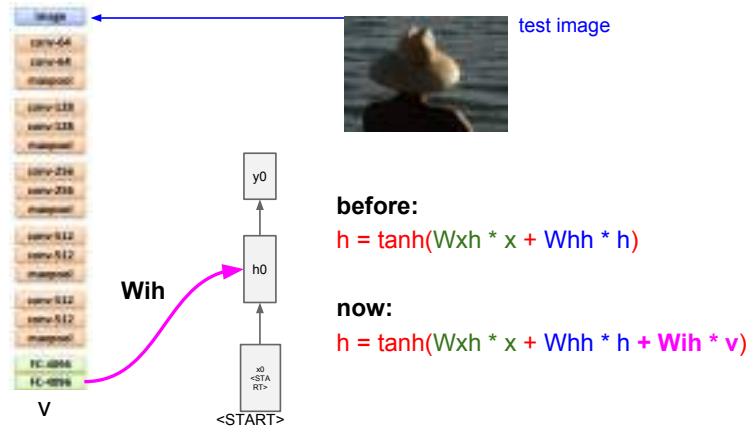
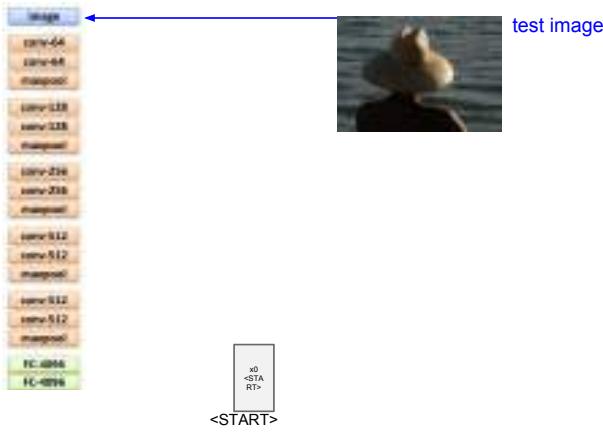
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 66 May 3, 2018



test image





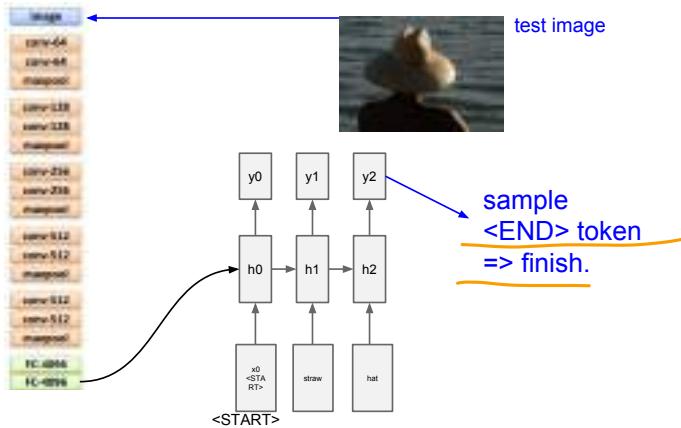


Image Captioning: Example Results



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 77 May 3, 2018

Image Captioning: Failure Cases

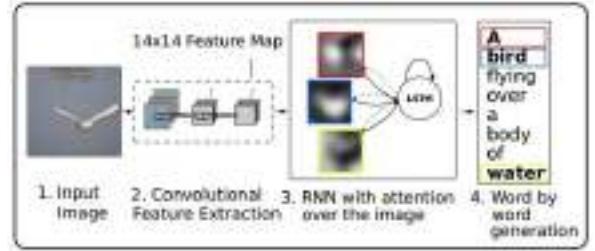


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 78 May 3, 2018

Image Captioning with Attention

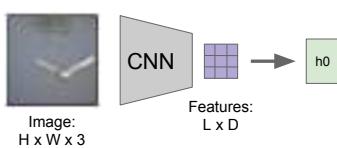
RNN focuses its attention at a different spatial location when generating each word



Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015
Figure copyright Kelvin Xu, Jimmy Lin, Jiajun Xiao, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard E. Zemel, and Yoshua Bengio, 2015. Reproduced with permission.

Lecture 10 - 79 May 3, 2018

Image Captioning with Attention

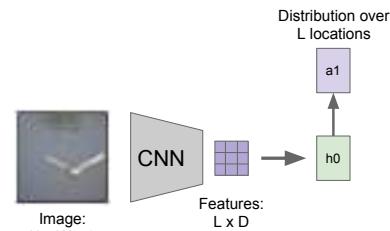


Xu et al., "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 80 May 3, 2018

Image Captioning with Attention

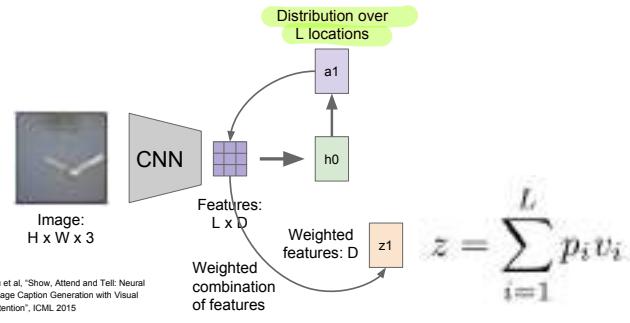


Xu et al., "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 81 May 3, 2018

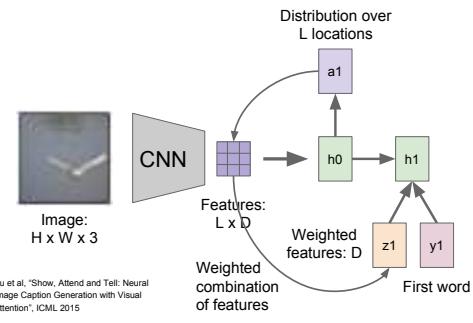
Image Captioning with Attention



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 82 May 3, 2018

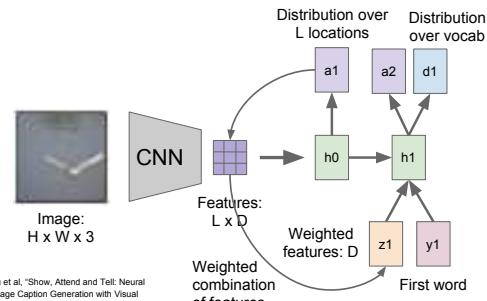
Image Captioning with Attention



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 83 May 3, 2018

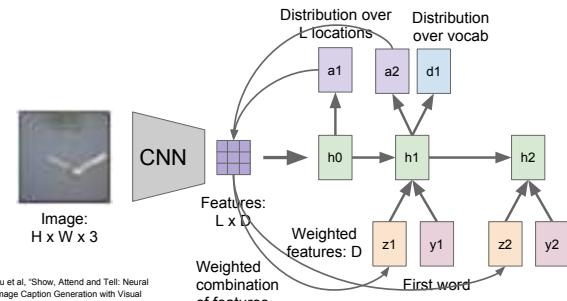
Image Captioning with Attention



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 84 May 3, 2018

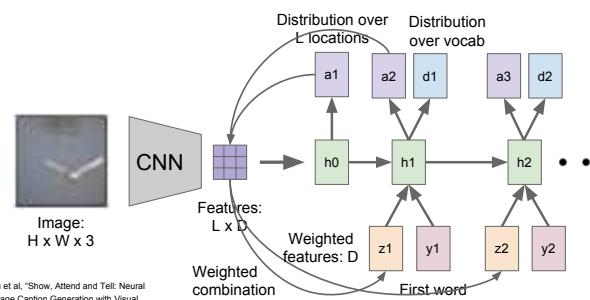
Image Captioning with Attention



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 85 May 3, 2018

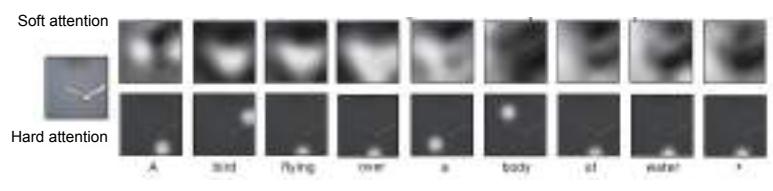
Image Captioning with Attention



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 86 May 3, 2018

Image Captioning with Attention



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 87 May 3, 2018

Image Captioning with Attention



Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Figure copyright Kelvin Xu, Jimmy Lin, Jamie Koh, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio, 2015. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 88 May 3, 2018

Visual Question Answering



Agarwal et al., "VQA: Visual Question Answering", ICCV 2015

Zhu et al., "Visual7W: Grounded Question Answering in Images", CVPR 2016

Figure from Zhu et al., copyright IEEE 2016. Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 89 May 3, 2018

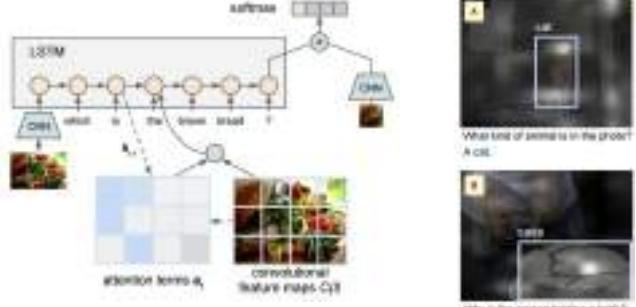
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 88 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 89 May 3, 2018

Visual Question Answering: RNNs with Attention



Zhu et al., "Visual7W: Grounded Question Answering in Images", CVPR 2016

Figure from Zhu et al., copyright IEEE 2016. Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 90 May 3, 2018

Multilayer RNNs

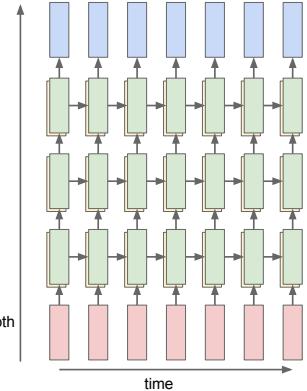
$$h_t^l = \tanh W^l \begin{pmatrix} h_{t-1}^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

LSTM:

$$\begin{pmatrix} i \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_{t-1}^{l-1} \\ h_t^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$



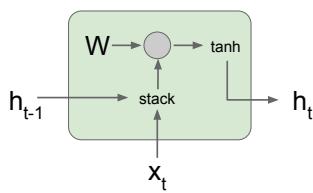
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 91 May 3, 2018

Vanilla RNN Gradient Flow

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994

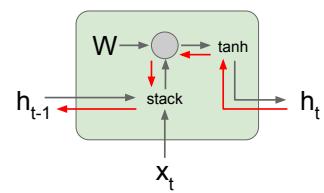
Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



$$\begin{aligned} h_t &= \tanh(W_{hh} h_{t-1} + W_{xh} x_t) \\ &= \tanh \left((W_{hh} \quad W_{xh}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

Vanilla RNN Gradient Flow

Backpropagation from h_t to h_{t-1} multiplies by W (actually W_{hh}^T)



$$\begin{aligned} h_t &= \tanh(W_{hh} h_{t-1} + W_{xh} x_t) \\ &= \tanh \left((W_{hh} \quad W_{xh}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

Fei-Fei Li & Justin Johnson & Serena Yeung

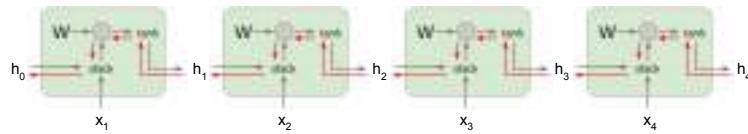
Lecture 10 - 92 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 93 May 3, 2018

Vanilla RNN Gradient Flow

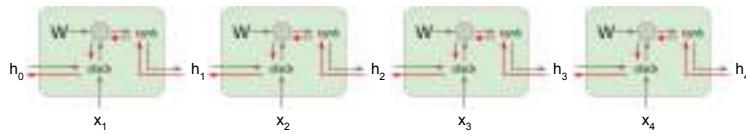
Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Vanilla RNN Gradient Flow

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



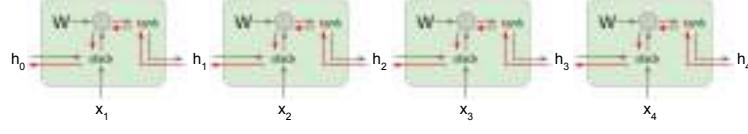
Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1:
Exploding gradients

Largest singular value < 1:
Vanishing gradients

Vanilla RNN Gradient Flow

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



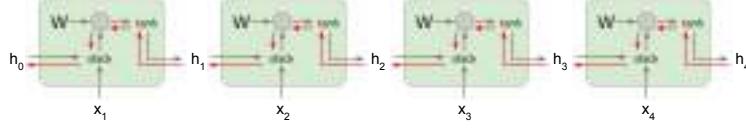
Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1:
Exploding gradients

Largest singular value < 1:
Vanishing gradients

Vanilla RNN Gradient Flow

Bengio et al., "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al., "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1:
Exploding gradients

Largest singular value < 1:
Vanishing gradients

→ Change RNN architecture

Long Short Term Memory (LSTM)

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

$$\begin{pmatrix} f \\ \sigma \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

call State

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

Long Short Term Memory (LSTM) [Hochreiter et al., 1997]

vector from below (x)

vector from before (h)

$4h \times 2h$

i

f

o

g

$tanh$

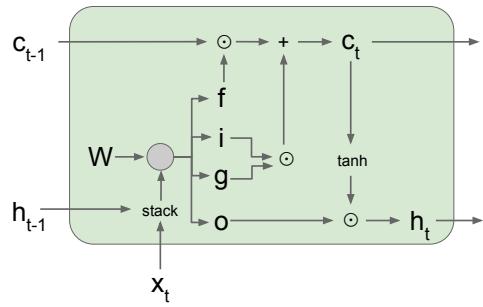
i: Input gate, whether to write to cell
f: Forget gate, Whether to erase cell
o: Output gate, How much to reveal cell
g: Gate gate (?), How much to write to cell

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Long Short Term Memory (LSTM) [Hochreiter et al., 1997]



$$\begin{pmatrix} i \\ f \\ \sigma \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

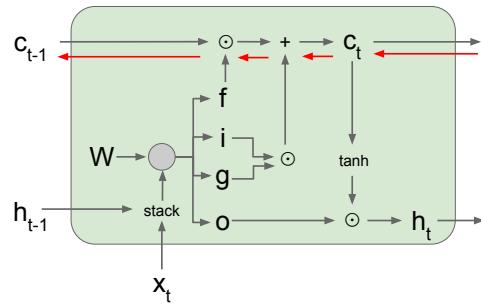
$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 100 May 3, 2018

Long Short Term Memory (LSTM): Gradient Flow [Hochreiter et al., 1997]



Backpropagation from c_t to c_{t-1} only elementwise multiplication by f , no matrix multiply by W

$$\begin{pmatrix} i \\ f \\ \sigma \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

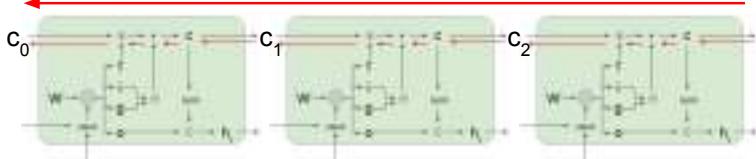
$$h_t = o \odot \tanh(c_t)$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 101 May 3, 2018

Long Short Term Memory (LSTM): Gradient Flow [Hochreiter et al., 1997]

Uninterrupted gradient flow!

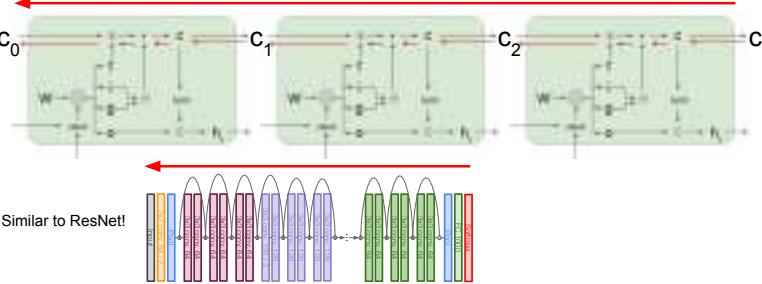


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 102 May 3, 2018

Long Short Term Memory (LSTM): Gradient Flow [Hochreiter et al., 1997]

Uninterrupted gradient flow!

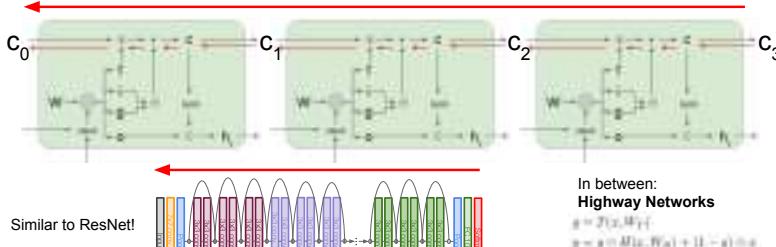


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 103 May 3, 2018

Long Short Term Memory (LSTM): Gradient Flow [Hochreiter et al., 1997]

Uninterrupted gradient flow!



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 104 May 3, 2018

Other RNN Variants

[An Empirical Exploration of Recurrent Network Architectures, Jozefowicz et al., 2015]

GRU [Learning phrase representations using rnn encoder-decoder for statistical machine translation, Cho et al. 2014]

$$r_t = \sigma(W_{sr}x_t + W_{hr}h_{t-1} + b_r)$$

$$z_t = \sigma(W_{sz}x_t + W_{hz}h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W_{sh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h)$$

$$h_t = z_t \odot \tilde{h}_t + (1 - z_t) \odot \tilde{h}_t$$

[LSTM: A Search Space Odyssey, Greff et al., 2015]

MUT1

$$i_t = \text{sigmoid}(W_{xi}x_t + W_{hi}h_{t-1} + b_i)$$

$$r_t = \text{sigmoid}(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

$$h_{t-1} = \tanh(W_{xh}x_t \odot i_t + W_{hh}r_t \odot h_{t-1} + b_h)$$

$$h_t = h_{t-1} \odot (1 - i_t)$$

MUT2

$$i_t = \text{sigmoid}(W_{xi}x_t + W_{hi}h_{t-1} + b_i)$$

$$r_t = \text{sigmoid}(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

$$h_{t-1} = \tanh(W_{xh}x_t \odot i_t + W_{hh}r_t \odot h_{t-1} + b_h)$$

$$h_t = h_t \odot (1 - i_t)$$

MUT3

$$i_t = \text{sigmoid}(W_{xi}x_t + W_{hi}h_{t-1} + b_i)$$

$$r_t = \text{sigmoid}(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

$$h_{t-1} = \tanh(W_{xh}x_t \odot i_t + W_{hh}r_t \odot h_{t-1} + b_h)$$

$$h_t = h_t \odot (1 - i_t)$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 105 May 3, 2018

Summary

- RNNs allow a lot of flexibility in architecture design
 - Vanilla RNNs are simple but don't work very well
 - Common to use LSTM or GRU: their additive interactions improve gradient flow
 - Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
 - Better/simpler architectures are a hot topic of current research
 - Better understanding (both theoretical and empirical) is needed.

Next time: Midterm!

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 106 May 3, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 107 May 3, 2018

Administrative

Lecture 11:

Detection and Segmentation

1. Midterms being graded
 - **Please don't** discuss midterms until next week - some students not yet taken
 2. A2 being graded
 3. Project milestones due Wed 5/16

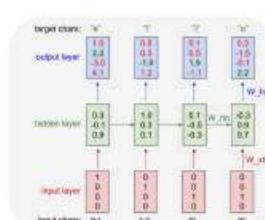
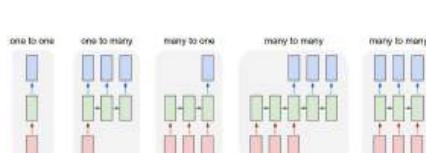
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 1 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 2 May 10, 2018

Last Time: Recurrent Networks



Last Time: Recurrent Networks

For example, when $\alpha_0 = 0$ we have that $\alpha_1 = \alpha_2 = \alpha_3 = 0$ in (2) and (3). In this case, $\beta_0 = \beta_1 = \beta_2 = \beta_3 = 0$ and $\gamma_0 = \gamma_1 = \gamma_2 = \gamma_3 = 0$. Thus, (4) is a consequence of (3).

Proof. First of all, it is clear that $\beta_0 = 0$.

Let $\beta_0 \neq 0$ in (3). By (3), $\beta_1 = \beta_2 = \beta_3 = 0$ and, by (4), $\alpha_1 = \alpha_2 = \alpha_3 = 0$. Then, by (2), $\gamma_1 = \gamma_2 = \gamma_3 = 0$. Let $\alpha_0 = 0$ in (2). Then, by (3), $\beta_0 = \beta_1 = \beta_2 = \beta_3 = 0$ and, by (4), $\gamma_0 = \gamma_1 = \gamma_2 = \gamma_3 = 0$. Thus, $\alpha_0 = \beta_0 = \gamma_0 = 0$ and any solution $B = (B_0, B_1, B_2, B_3)$ of (3) is a solution of (4).

be should be more implemented and the day
should be available to attend to living issues and
not a child's absence of class.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 3 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 4 May 10, 2018

Last Time: Recurrent Networks

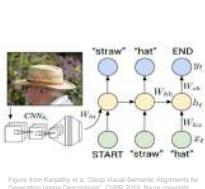
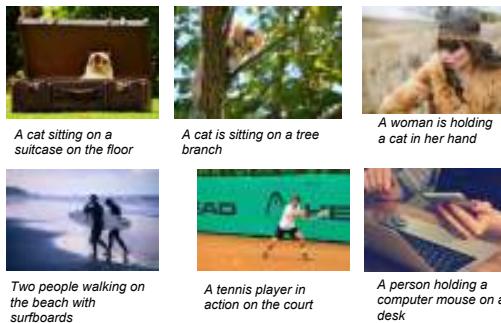
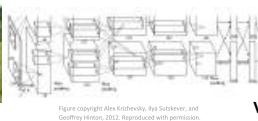
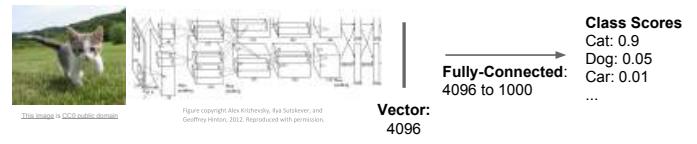


Figure from Karpathy et al., "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015; figure copyright © 2015 IEEE. All rights reserved. Reproduced for educational purposes.



So far: Image Classification



Fully-Connected:
4096 to 1000

Vector:
4096

Class Scores
Cat: 0.9
Dog: 0.05
Car: 0.01
...

Fei-Fei Li & Justin Johnson & Serena Yeung

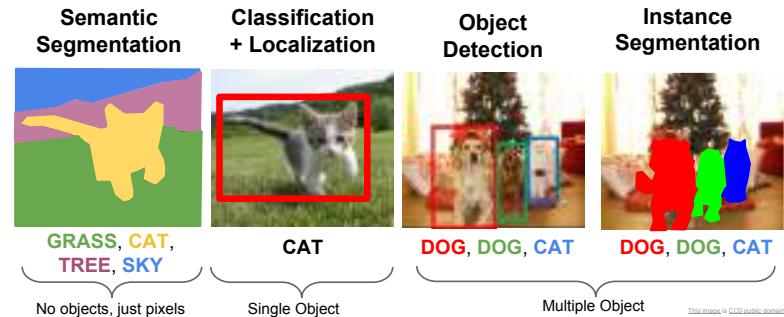
Lecture 11 - 5 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 6 May 10, 2018

Today: Detection, Segmentation

Other Computer Vision Tasks



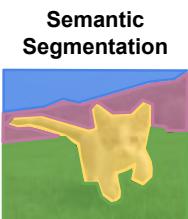
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 7 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 8 May 10, 2018

Other Computer Vision Tasks



No objects, just pixels

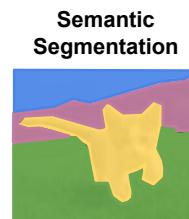


Object categories +
2D bounding boxes



Object categories +
3D bounding boxes

Semantic Segmentation



No objects, just pixels



Object categories +
2D bounding boxes



Object categories +
3D bounding boxes

Fei-Fei Li & Justin Johnson & Serena Yeung

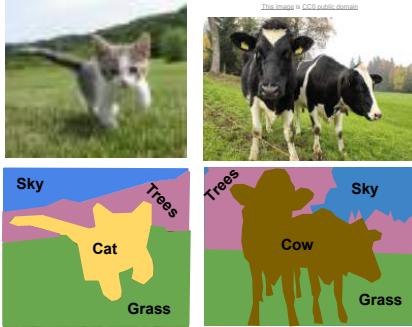
Lecture 11 - 9 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 10 May 10, 2018

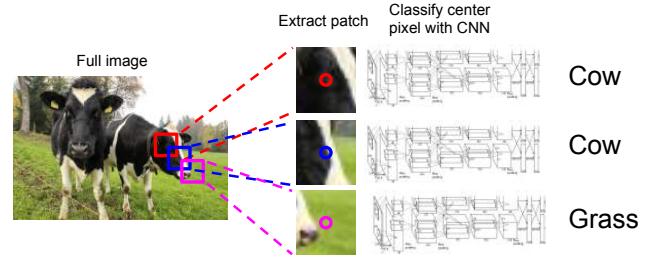
Semantic Segmentation

Label each pixel in the image with a category label



Don't differentiate instances, only care about pixels

Semantic Segmentation Idea: Sliding Window



Farabet et al., "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

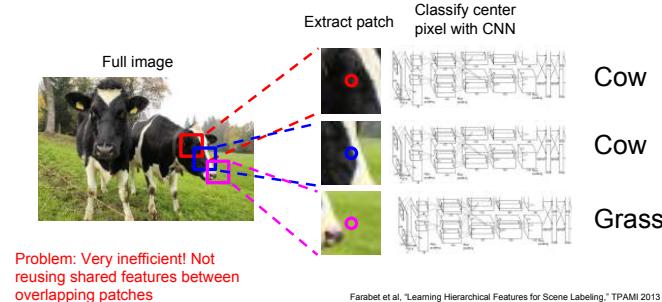
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 11 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 12 May 10, 2018

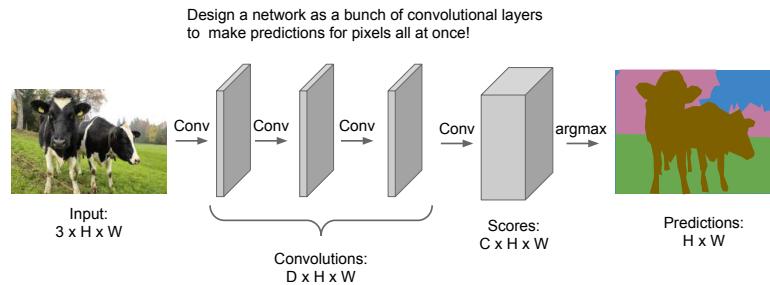
Semantic Segmentation Idea: Sliding Window



Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al., "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation Idea: Fully Convolutional



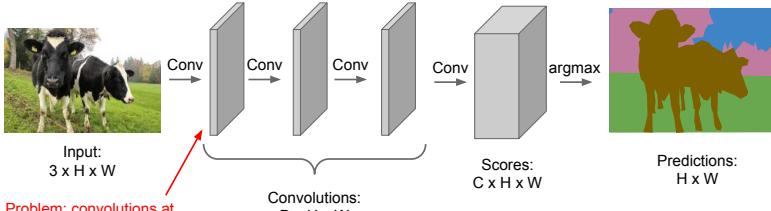
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 13 May 10, 2018

Lecture 11 - 14 May 10, 2018

Semantic Segmentation Idea: Fully Convolutional

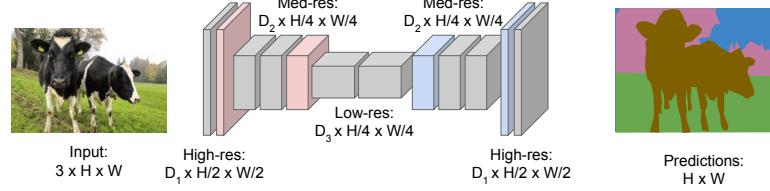
Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Problem: convolutions at original image resolution will be very expensive ...

Semantic Segmentation Idea: Fully Convolutional

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al., "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 15 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 16 May 10, 2018

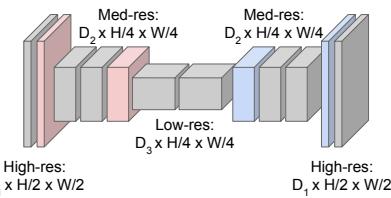
Semantic Segmentation Idea: Fully Convolutional

Downsampling:
Pooling, stride
convolution

Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!



Input:
 $3 \times H \times W$



Upsampling:
???



Predictions:
 $H \times W$

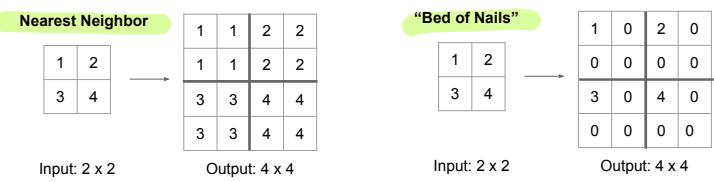
Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Noh et al., "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 17 May 10, 2018

In-Network upsampling: "Unpooling"



Input: 2×2

Output: 4×4

Input: 2×2

Output: 4×4

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 18 May 10, 2018

In-Network upsampling: "Max Unpooling"

Max Pooling
Remember which element was max.

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input: 4×4

5	6
7	8

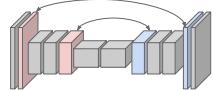
Output: 2×2

Max Unpooling
Use positions from
pooling layer

0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

Output: 4×4

Corresponding pairs of
downsampling and
upsampling layers



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 19 May 10, 2018

Learnable Upsampling: Transpose Convolution

Recall: Typical 3×3 convolution, stride 1 pad 1

Input: 4×4

Output: 4×4

Fei-Fei Li & Justin Johnson & Serena Yeung

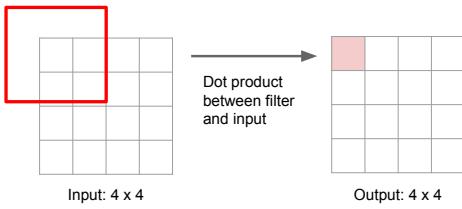
Lecture 11 - 19 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 20 May 10, 2018

Learnable Upsampling: Transpose Convolution

Recall: Normal 3×3 convolution, stride 1 pad 1

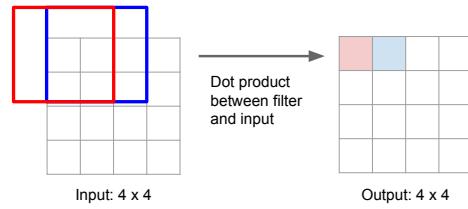


Input: 4×4

Output: 4×4

Learnable Upsampling: Transpose Convolution

Recall: Normal 3×3 convolution, stride 1 pad 1



Input: 4×4

Output: 4×4

Fei-Fei Li & Justin Johnson & Serena Yeung

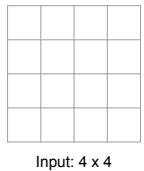
Lecture 11 - 21 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 22 May 10, 2018

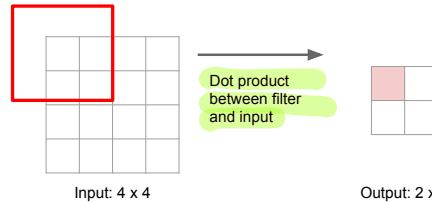
Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1

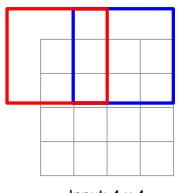


Input: 4 x 4

Output: 2 x 2

Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1

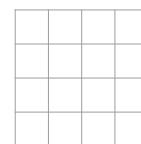


Filter moves 2 pixels in the input for every one pixel in the output

Stride gives ratio between movement in input and output

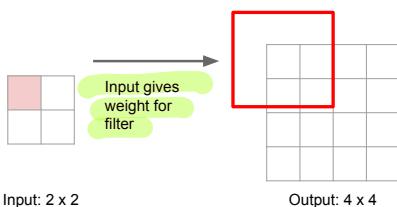
Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1



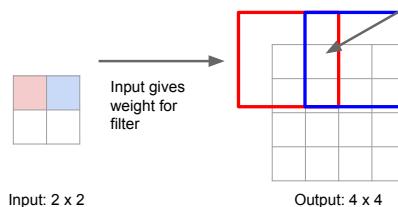
Learnable Upsampling: Transpose Convolution

3 x 3 transpose convolution, stride 2 pad 1

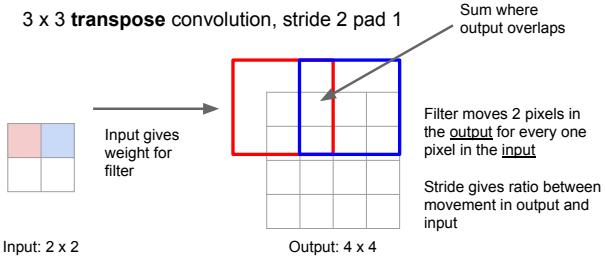


Learnable Upsampling: Transpose Convolution

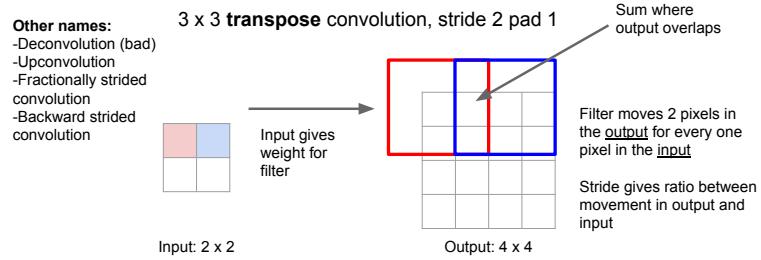
3 x 3 transpose convolution, stride 2 pad 1



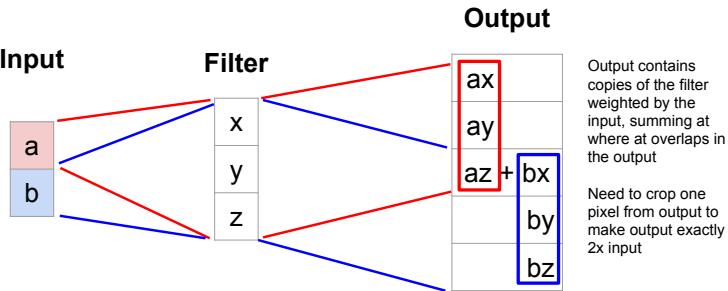
Learnable Upsampling: Transpose Convolution



Learnable Upsampling: Transpose Convolution



Learnable Upsampling: 1D Example



Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & x & y & x & 0 & 0 \\ 0 & 0 & x & y & x & 0 \\ 0 & 0 & 0 & x & y & x \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax + \delta z \\ ay + bz + cz \\ az + bx + cx \\ by \\ bz + dy \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & x & y & x & 0 & 0 \\ 0 & 0 & x & y & x & 0 \\ 0 & 0 & 0 & x & y & x \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax + \delta z \\ ay + bz + cz \\ az + bx + cx \\ by \\ bz + dy \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ ay + bx \\ az + by + cx \\ bz + cy + dx \\ cz + dy \\ dz \end{bmatrix}$$

When stride=1, convolution transpose is just a regular convolution (with different padding rules)

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & 0 & x & y & x & 0 \\ 0 & 0 & 0 & x & y & x \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ ay + bz \\ az + bx + cx \\ by \\ bz + dy \\ dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Convolution as Matrix Multiplication (1D Example)

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & 0 & x & y & z & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax + bz \\ ay + cz \\ az + bx \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=2, padding=1

Convolution transpose multiplies by the transpose of the same matrix:

$$\vec{x}^T * \vec{a} = X^T \vec{a}$$

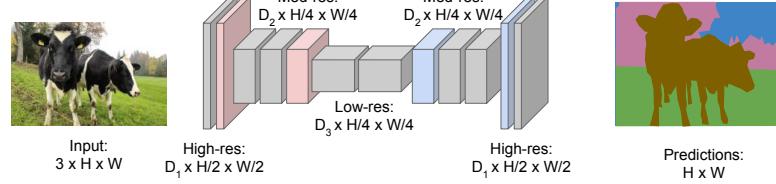
$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

When stride>1, convolution transpose is no longer a normal convolution!

Semantic Segmentation Idea: Fully Convolutional

Downsampling:
Pooling, strided convolution

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Upsampling:
Unpooling or strided transpose convolution



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 35 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 36 May 10, 2018

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Aside: Multi-view 3D Reconstruction



Chey, C. B., Xu, D., Gwak, J., Chen, K., & Savarese, S. (2016, October). 3dr2n2: A unified approach for single and multi-view 3d object reconstruction. In European Conference on Computer Vision (pp. 628-644). Springer, Cham.

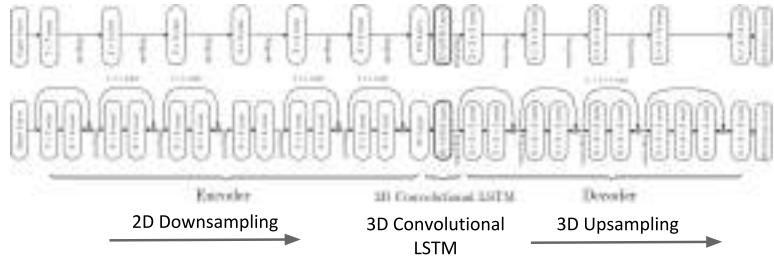
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 37 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

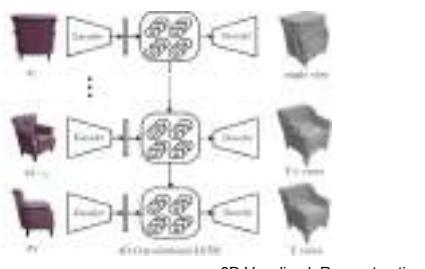
Lecture 11 - 38 May 10, 2018

Aside: Multi-view 3D Reconstruction



Chey, C. B., Xu, D., Gwak, J., Chen, K., & Savarese, S. (2016, October). 3dr2n2: A unified approach for single and multi-view 3d object reconstruction. In European Conference on Computer Vision (pp. 628-644). Springer, Cham.

Aside: Multi-view 3D Reconstruction



Chey, C. B., Xu, D., Gwak, J., Chen, K., & Savarese, S. (2016, October). 3dr2n2: A unified approach for single and multi-view 3d object reconstruction. In European Conference on Computer Vision (pp. 628-644). Springer, Cham.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 39 May 10, 2018

2D Object Detection

Semantic Segmentation



No objects, just pixels

2D Object Detection



Object categories + 2D bounding boxes

3D Object Detection



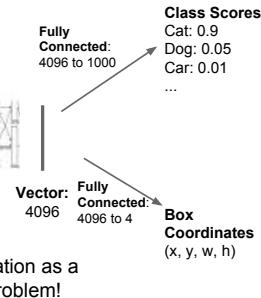
Object categories + 3D bounding boxes

This image is CC-BY public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 40 May 10, 2018

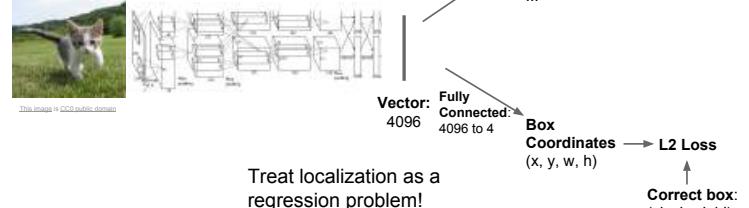
Classification + Localization



were 11-15

Classification + Localization

identify object



Correct label: Cat
↓
Softmax Loss

Correct box: (x', y', w', h')
↑
L2 Loss

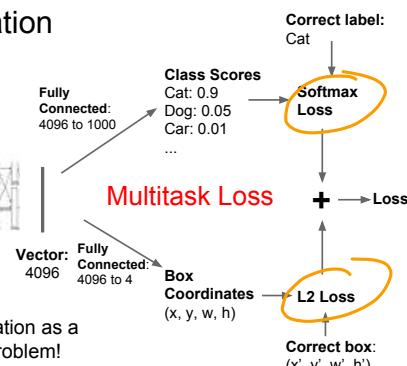
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 41 May 10, 2018

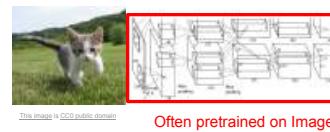
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 42 May 10, 2018

Classification + Localization



Classification + Localization



Correct label: Cat
↓
Softmax Loss

Correct box: (x', y', w', h')
↑
L2 Loss

Often pretrained on ImageNet (Transfer learning)

+

Loss

Treat localization as a regression problem!

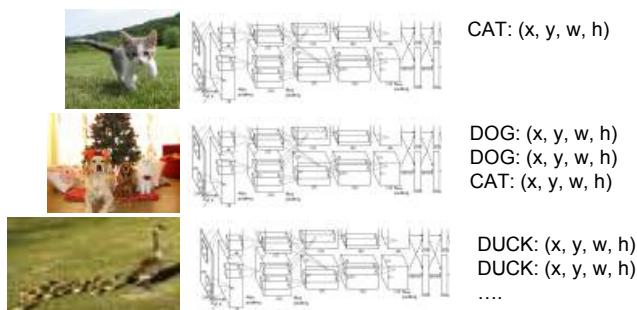
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 43 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

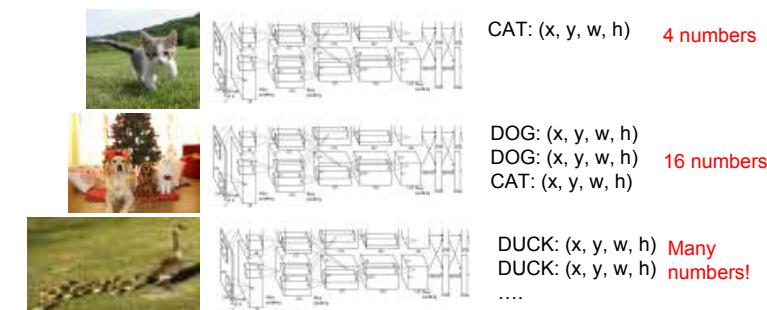
Lecture 11 - 44 May 10, 2018

Object Detection as Regression?



Object Detection as Regression?

Each image needs a different number of outputs!



Fei-Fei Li & Justin Johnson & Serena Yeung

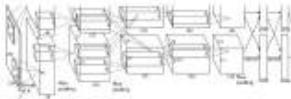
Lecture 11 - 45 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 46 May 10, 2018

Object Detection as Classification: Sliding Window

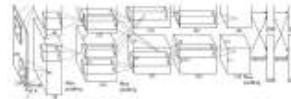
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? NO
Background? YES

Object Detection as Classification: Sliding Window

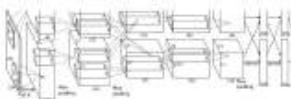
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection as Classification: Sliding Window

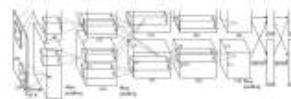
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? YES
Cat? NO
Background? NO

Object Detection as Classification: Sliding Window

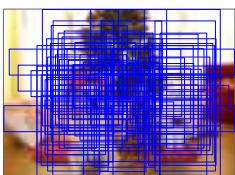
Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO
Cat? YES
Background? NO

Object Detection as Classification: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

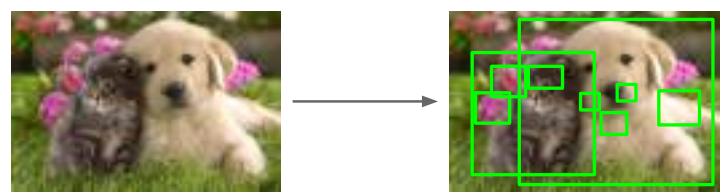


Dog? NO
Cat? YES
Background? NO

Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!

Region Proposals / Selective Search

- Find "blobby" image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



Alexe et al., "Measuring the objectness of image windows", TPAMI 2012
Uijlings et al., "Selective Search for Object Recognition", IJCV 2013
Chen et al., "Efficient multi-scale region proposal generation at 300fps", CVPR 2014
Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014

R-CNN



Input image

Fei-Fei Li & Justin Johnson & Serena Yeung

R-CNN



Regions of Interest (RoI) from a proposal method (~2k)

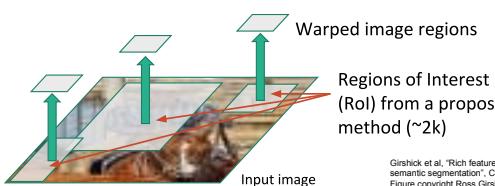
Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015. [source](#). Reproduced with permission.

Lecture 11 - 53 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 54 May 10, 2018

R-CNN



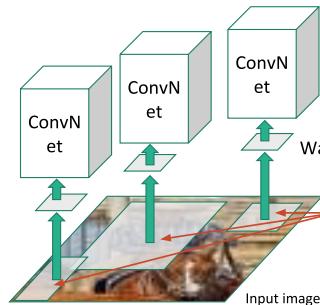
Regions of Interest (RoI) from a proposal method (~2k)

Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015. [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 55 May 10, 2018

R-CNN



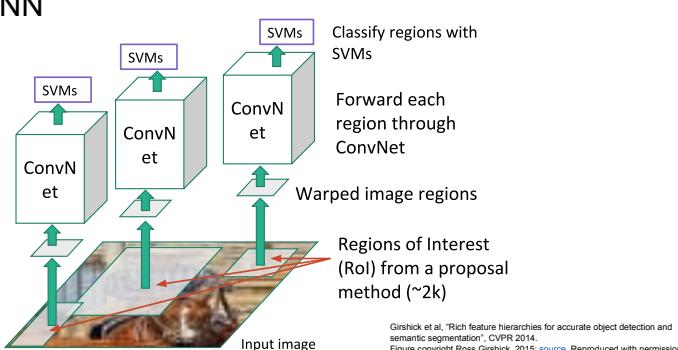
Forward each region through ConvNet

Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015. [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 56 May 10, 2018

R-CNN



Classify regions with SVMs

Forward each region through ConvNet

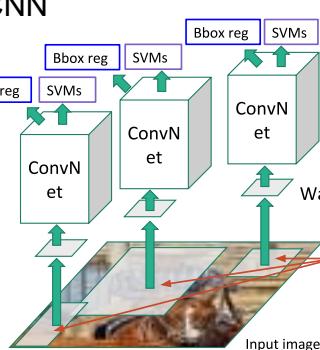
Regions of Interest (RoI) from a proposal method (~2k)

Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015. [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 57 May 10, 2018

R-CNN



Linear Regression for bounding box offsets

Classify regions with SVMs

Forward each region through ConvNet

Regions of Interest (RoI) from a proposal method (~2k)

Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
Figure copyright Ross Girshick, 2015. [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 58 May 10, 2018

R-CNN: Problems

- Ad hoc training objectives
 - Fine-tune network with softmax classifier (log loss)
 - Train post-hoc linear SVMs (hinge loss)
 - Train post-hoc bounding-box regressions (least squares)
- Training is slow (84h), takes a lot of disk space
- Inference (detection) is slow
 - 47s / image with VGG16 [Simonyan & Zisserman. ICLR15]
 - Fixed by SPP-net [He et al. ECCV14]



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 59 May 10, 2018

Fast R-CNN



Input image

Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 11 - 60 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

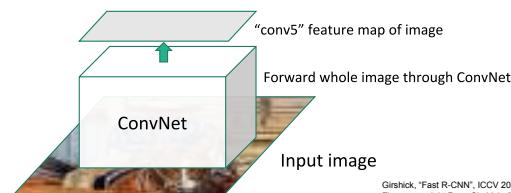
Lecture 11 - 59 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

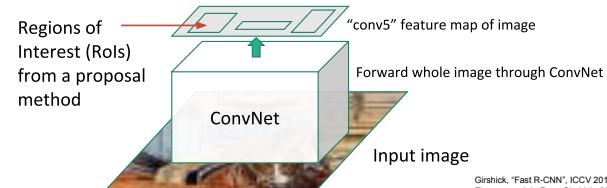
Lecture 11 - 60 May 10, 2018

Fast R-CNN

Fast R-CNN



Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.



Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

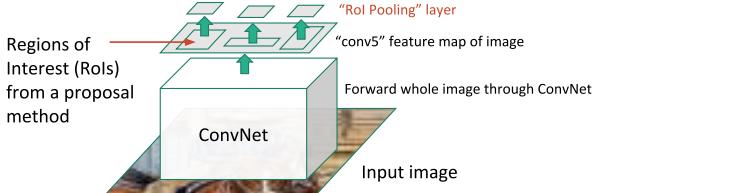
Lecture 11 - 61 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

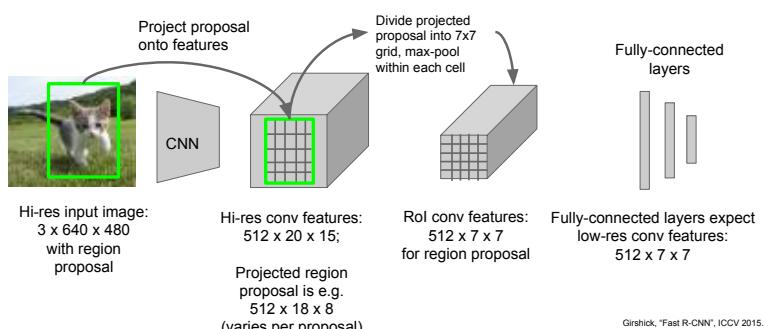
Lecture 11 - 62 May 10, 2018

Fast R-CNN

Fast R-CNN: RoI Pooling



Girshick, "Fast R-CNN", ICCV 2015.
Figure copyright Ross Girshick, 2015; [source](#). Reproduced with permission.



Girshick, "Fast R-CNN", ICCV 2015.

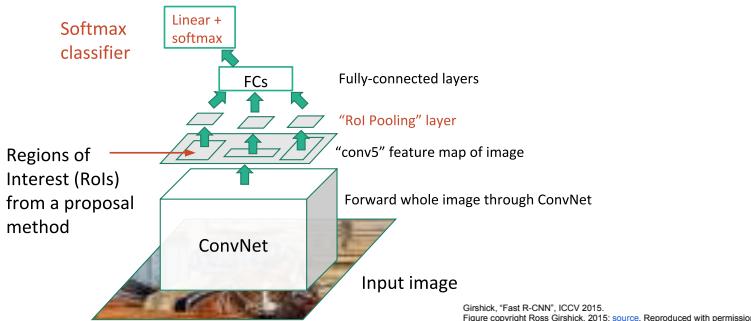
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 63 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 64 May 10, 2018

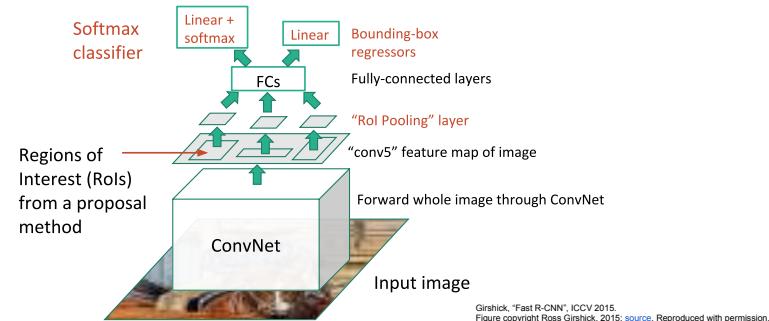
Fast R-CNN



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 65 May 10, 2018

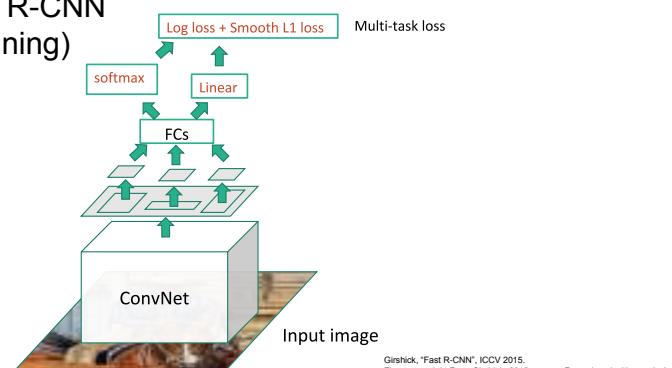
Fast R-CNN



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 66 May 10, 2018

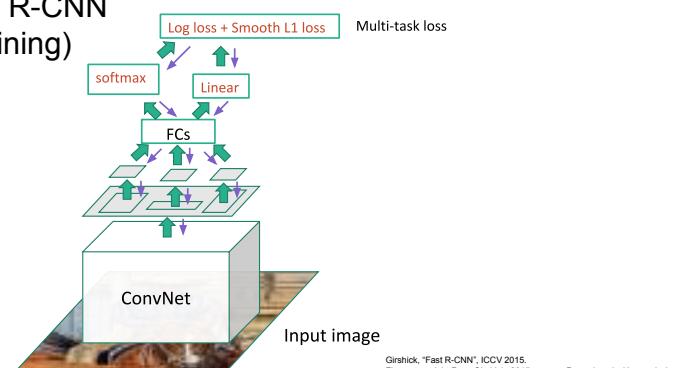
Fast R-CNN (Training)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 67 May 10, 2018

Fast R-CNN (Training)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 68 May 10, 2018

R-CNN vs SPP vs Fast R-CNN



R-CNN vs SPP vs Fast R-CNN



Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al., "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
Girshick, "Fast R-CNN", ICCV 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 69 May 10, 2018

Girshick et al., "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014.
He et al., "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014
Girshick, "Fast R-CNN", ICCV 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 70 May 10, 2018

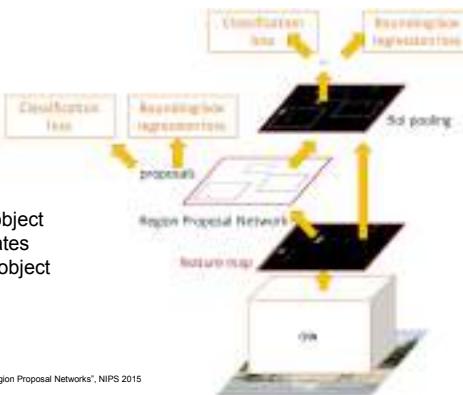
Faster R-CNN:

Make CNN do proposals!

Insert Region Proposal Network (RPN) to predict proposals from features

Jointly train with 4 losses:

1. RPN classify object / not object
2. RPN regress box coordinates
3. Final classification score (object classes)
4. Final box coordinates



Faster R-CNN:

Make CNN do proposals!

R-CNN Test-Time Speed



Ren et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015

Figure copyright 2015, Ross Girshick; reproduced with permission

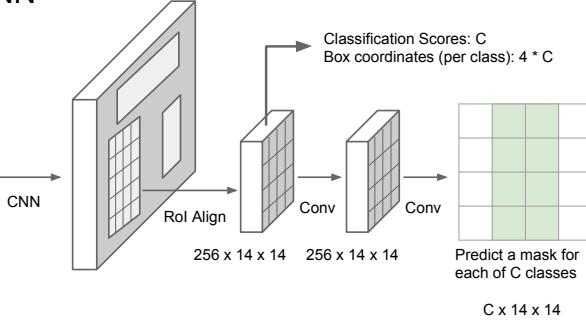
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 71 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 72 May 10, 2018

Mask R-CNN



Mask R-CNN: Very Good Results!



He et al., "Mask R-CNN", arXiv 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 73 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 74 May 10, 2018

Mask R-CNN

Also does pose



He et al., "Mask R-CNN", arXiv 2017

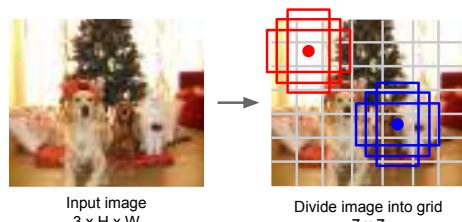
Figures copyright Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick, 2017.

Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 75 May 10, 2018

Detection without Proposals: YOLO / SSD



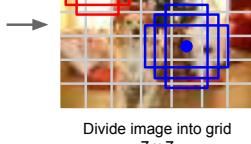
Within each grid cell:

- Regress from each of the B base boxes to a final box with 5 numbers: $(dx, dy, dh, dw, confidence)$
- Predict scores for each of C classes (including background as a class)

Output:
 $7 \times 7 \times (5 * B + C)$

Image a set of **base boxes** centered at each grid cell
Here $B = 3$

Input image
 $3 \times H \times W$



Redmon et al., "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016

Liu et al., "SSD: Single-Shot MultiBox Detector", ECCV 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

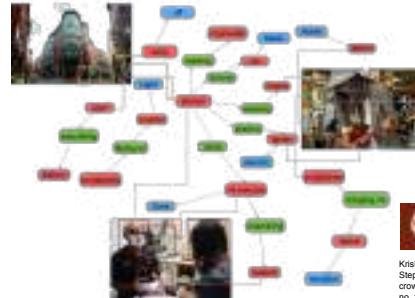
Lecture 11 - 76 May 10, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 83 May 10, 2018

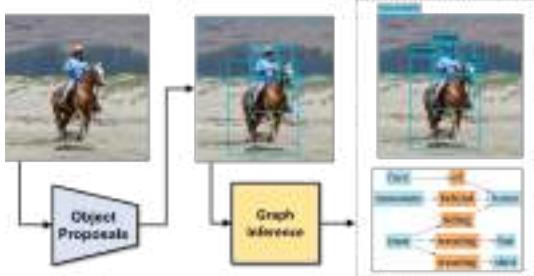
Aside: VisualGenome



108,077 Images
5.4 Million Region Descriptions
1.7 Million Visual Question Answers
3.8 Million Object Instances
2.8 Million Attributes
2.3 Million Relationships
Everything Mapped to Wordnet Synsets

Ranjay, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Alle Chen et al. "Visual genome: Connecting language and vision using uncrowded dense image annotations." International Journal of Computer Vision 123, 2017: 32-73.

Aside: Scene Graph Generation



Xu, Zhu, Choy, and Fei-Fei, "Scene Graph Generation by Iterative Message Passing", CVPR 2017
Figure copyright IEEE, 2018. Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 85 May 10, 2018

3D Object Detection



Object categories

This image is CC0 public domain

Simplified Camera Model

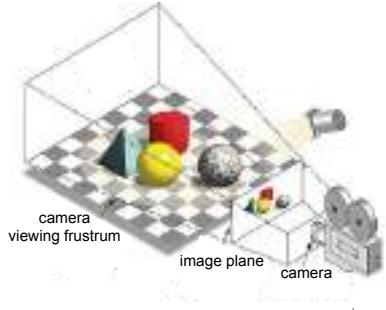
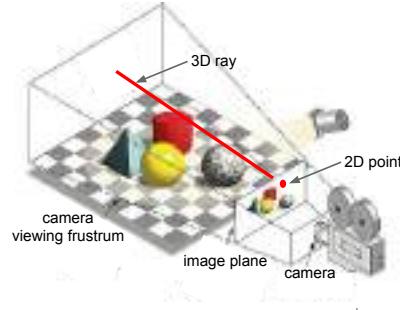


Image source: https://www.pcmag.com/encyclopedia_images/_FRUSTUM.GIF

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 87 May 10, 2018

Simplified Camera Model

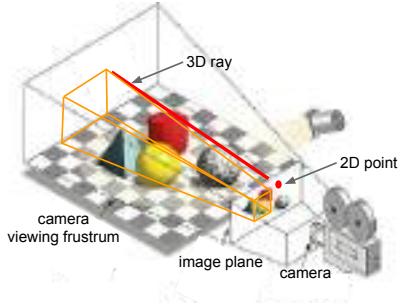


A point on the image plane corresponds to a **ray** in the 3D space

https://www.pcmag.com/encyclopedia_images/_FRUSTUM.GIF

Lecture 11 - 88 May 10, 2018

Simplified Camera Model



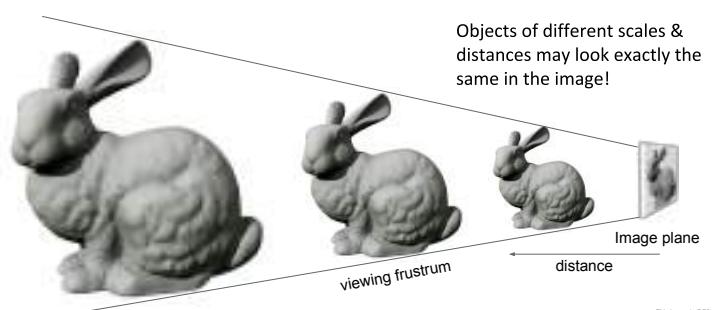
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 89 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

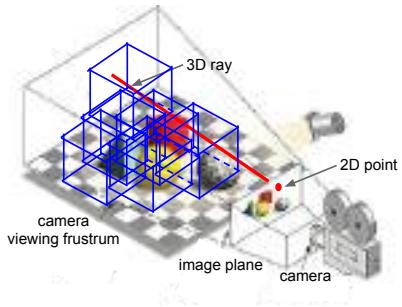
Lecture 11 - 90 May 10, 2018

Scale & Distance Ambiguity



This image is CC0 public domain

Simplified Camera Model



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 91 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 92 May 10, 2018

3D Object Detection



This image is CC0 public domain

3D Object Detection

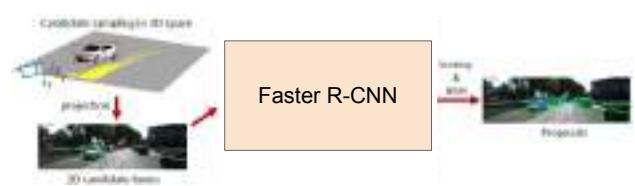


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 93 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

3D Object Detection: Monocular Camera

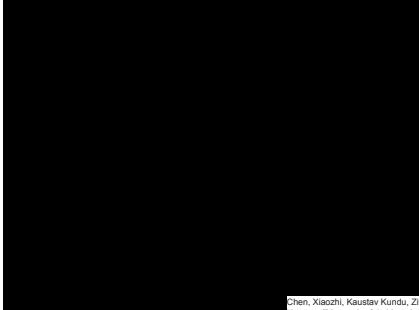


- Same idea as Faster RCNN, but proposals are in 3D
- 3D bounding box proposal, regress 3D box parameters + class score

Chen, Xiaozi, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. "Monocular 3d object detection for autonomous driving." CVPR 2016.

Lecture 11 - 94 May 10, 2018

3D Object Detection: Monocular View



Chen, Xiaozhi, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. "Monocular 3d object detection for autonomous driving." CVPR 2016.

3D Object Detection: Camera + LiDAR



Velodyne (HDL-64e)



3D Point Cloud

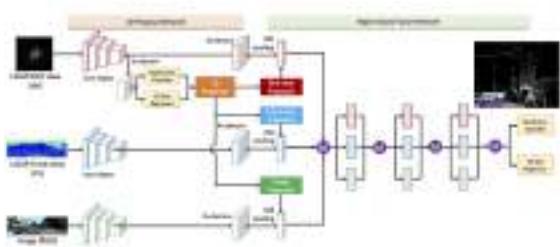
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 95 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 96 May 10, 2018

3D Object Detection: Camera + LiDAR



- Combine 3D proposals from multiple views & sensors
- regress 3D box parameters + class score

Chen, Xiaozhi, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. "Multi-view 3d object detection network for autonomous driving." CVPR 2017

3D Object Detection: Camera + LiDAR



Chen, Xiaozhi, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. "Multi-view 3d object detection network for autonomous driving." CVPR 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 97 May 10, 2018

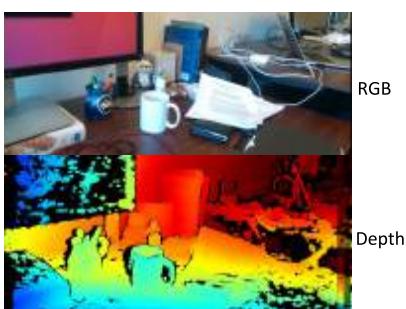
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 98 May 10, 2018

RGB-Depth Camera



Kinect (Xbox One)



RGB

Depth

RGB-Depth Camera



Registered RGB + depth point cloud

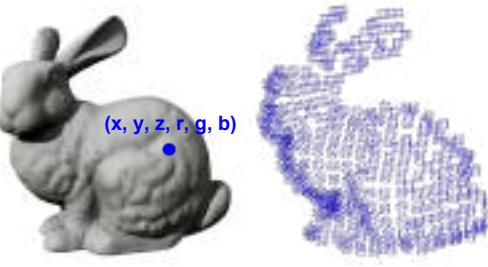
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 99 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 100 May 10, 2018

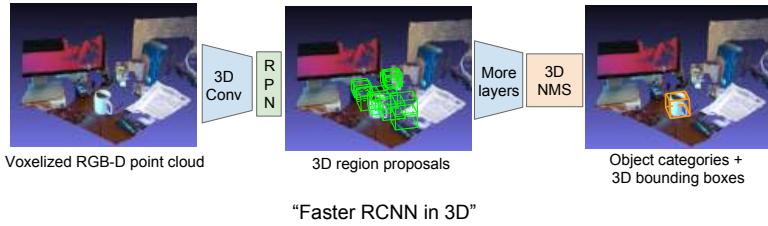
Point Cloud Voxelization



1. Capture RGB-D point cloud of a scene.
2. Partition the 3D space into a regular 3D grid.
3. For each grid cell that has a point fall into it, fill the cell with the RGB value of that point.

A bit like "3D image"

3D Object Detection: RGB-Depth Camera



S. Song, and J. Xiao. Deep Sliding Shapes for Amodal 3D Object Detection in RGB-D Images. CVPR 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 101 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 102 May 10, 2018

Recap



Next time:
Visualizing CNN features
DeepDream + Style Transfer

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 103 May 10, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 11 - 104 May 10, 2018

Administrative

Lecture 12: Generative Models

Project Milestone due tomorrow (Wed 5/16)

A3 due next Wed, 5/23

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 1 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 2 May 15, 2018

Overview

- Unsupervised Learning
- Generative Models
 - PixelRNN and PixelCNN
 - Variational Autoencoders (VAE)
 - Generative Adversarial Networks (GAN)

Supervised vs Unsupervised Learning

Supervised Learning

Data: (x, y)
 x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification,
regression, object detection,
semantic segmentation, image
captioning, etc.

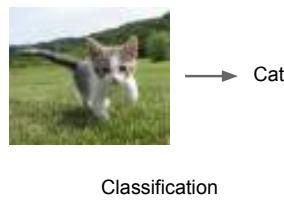
Supervised vs Unsupervised Learning

Supervised Learning

Data: (x, y)
 x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification,
regression, object detection,
semantic segmentation, image
captioning, etc.



Supervised vs Unsupervised Learning

Supervised Learning

Data: (x, y)
 x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification,
regression, object detection,
semantic segmentation, image
captioning, etc.



Supervised vs Unsupervised Learning

Supervised Learning

Data: (x, y)
 x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification,
regression, object detection,
semantic segmentation, image
captioning, etc.



Supervised vs Unsupervised Learning

Supervised Learning

Data: (x, y)
 x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification,
regression, object detection,
semantic segmentation, image
captioning, etc.



Supervised vs Unsupervised Learning

Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 9 May 15, 2018

Supervised vs Unsupervised Learning

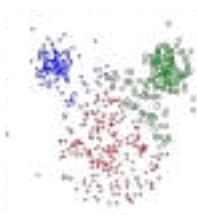
Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



K-means clustering

This image is CC0 public domain

Lecture 12 - 10 May 15, 2018

Supervised vs Unsupervised Learning

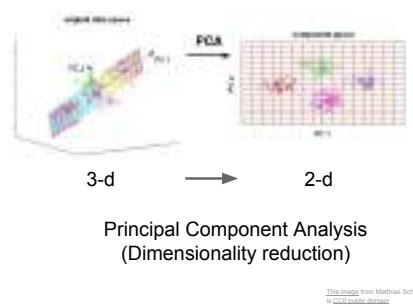
Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 11 May 15, 2018

Supervised vs Unsupervised Learning

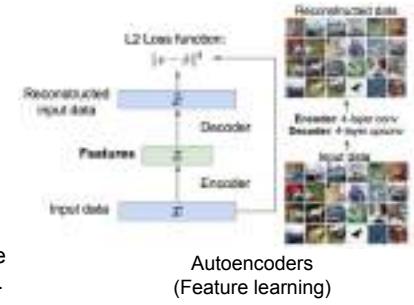
Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



Lecture 12 - 12 May 15, 2018

Supervised vs Unsupervised Learning

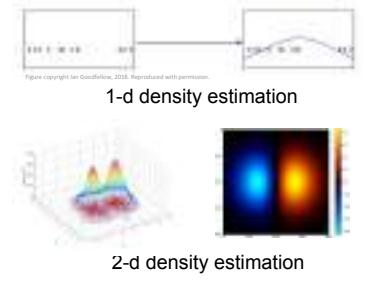
Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



This image is from Ian Goodfellow. 2010. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 13 May 15, 2018

Supervised vs Unsupervised Learning

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.

Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.

Lecture 12 - 14 May 15, 2018

Supervised vs Unsupervised Learning

Supervised Learning

Data: (x, y)
 x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$

Examples: Classification,
 regression, object detection,
 semantic segmentation, image
 captioning, etc.

Unsupervised Learning

Training data is cheap

Data: x

Just data, no labels!

Holy grail: Solve
 unsupervised learning
 \Rightarrow understand structure
 of visual world

Goal: Learn some underlying
 hidden *structure* of the data

Examples: Clustering,
 dimensionality reduction, feature
 learning, density estimation, etc.

Generative Models

Given training data, generate new samples from same distribution



Training data $\sim p_{\text{data}}(x)$



Generated samples $\sim p_{\text{model}}(x)$

Want to learn $p_{\text{model}}(x)$ similar to $p_{\text{data}}(x)$

Generative Models

Given training data, generate new samples from same distribution



Training data $\sim p_{\text{data}}(x)$



Generated samples $\sim p_{\text{model}}(x)$

Want to learn $p_{\text{model}}(x)$ similar to $p_{\text{data}}(x)$

Addresses density estimation, a core problem in unsupervised learning

Several flavors:

- Explicit density estimation: explicitly define and solve for $p_{\text{model}}(x)$
- Implicit density estimation: learn model that can sample from $p_{\text{model}}(x)$ w/o explicitly defining it

Why Generative Models?

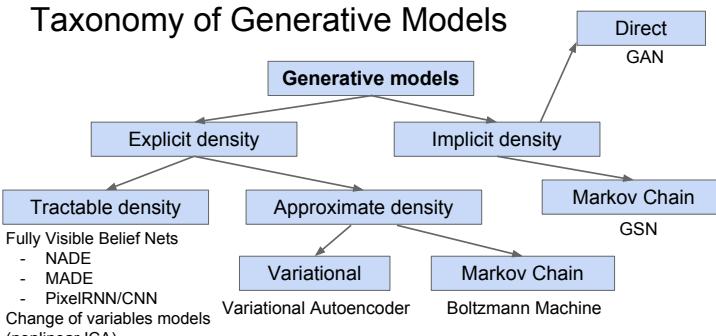
- Realistic samples for artwork, super-resolution, colorization, etc.



- Generative models of time-series data can be used for simulation and planning (reinforcement learning applications!)
- Training generative models can also enable inference of latent representations that can be useful as general features

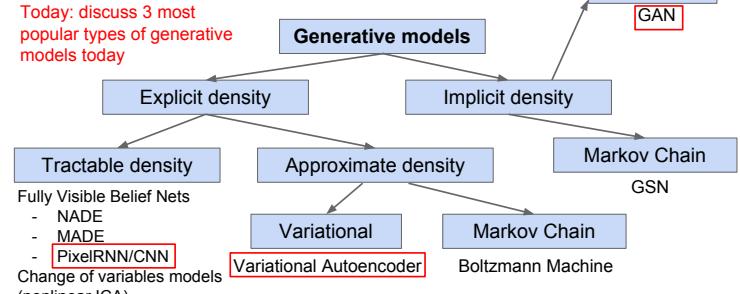
Figures from LR are copyright: (c) Alec Radford et al., 2015. (c) David Blei et al., 2011. (c) Phillip Isola et al., 2012. Reproduced with authors permission.

Taxonomy of Generative Models



Taxonomy of Generative Models

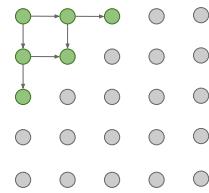
Today: discuss 3 most popular types of generative models today



PixelRNN [van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

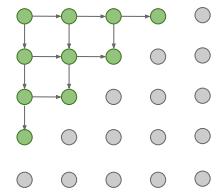


PixelRNN [van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

Drawback: sequential generation is slow!



PixelCNN [van der Oord et al. 2016]

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region

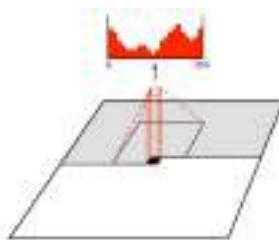


Figure copyright van der Oord et al., 2016. Reproduced with permission.

PixelCNN [van der Oord et al. 2016]

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region

Training: maximize likelihood of training images

$$p(x) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$

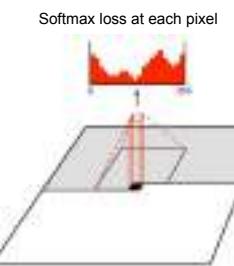


Figure copyright van der Oord et al., 2016. Reproduced with permission.

PixelCNN [van der Oord et al. 2016]

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region

Training is faster than PixelRNN
(can parallelize convolutions since context region values known from training images)

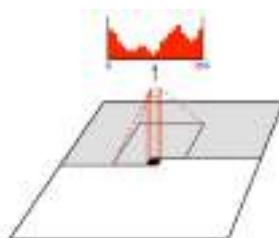


Figure copyright van der Oord et al., 2016. Reproduced with permission.

Generation Samples



32x32 CIFAR-10



32x32 ImageNet

Figures copyright Aaron van der Oord et al., 2016. Reproduced with permission.

PixelRNN and PixelCNN

Pros:

- Can explicitly compute likelihood $p(x)$
- Explicit likelihood of training data gives good evaluation metric
- Good samples

Con:

- Sequential generation => slow

Improving PixelCNN performance

- Gated convolutional layers
- Short-cut connections
- Discretized logistic loss
- Multi-scale
- Training tricks
- Etc...

See

- Van der Oord et al. NIPS 2016
- Salimans et al. 2017
(PixelCNN++)

Variational Autoencoders (VAE)

So far...

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i|x_1, \dots, x_{i-1})$$

So far...

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i|x_1, \dots, x_{i-1})$$

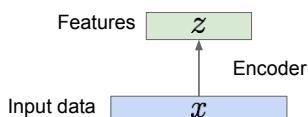
VAEs define intractable density function with latent z :

$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

Some background first: Autoencoders

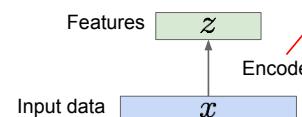
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data



Some background first: Autoencoders

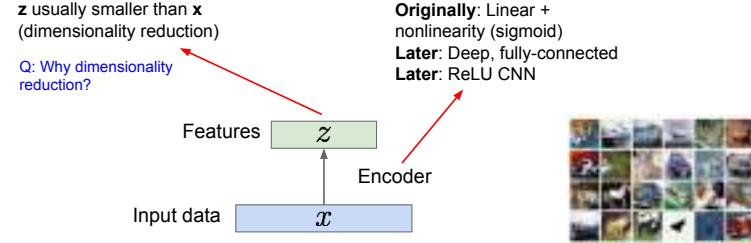
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

Originally: Linear + nonlinearity (sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN



Some background first: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

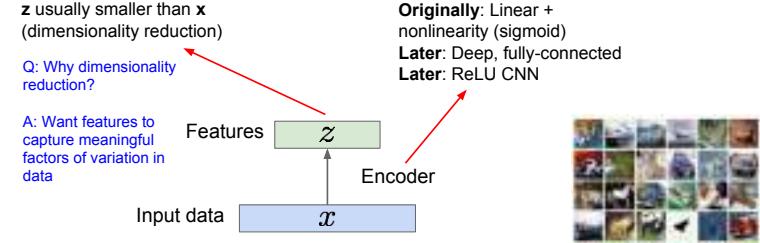


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 39 May 15, 2018

Some background first: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

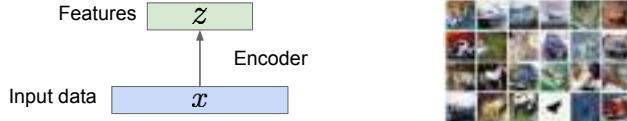


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 40 May 15, 2018

Some background first: Autoencoders

How to learn this feature representation?



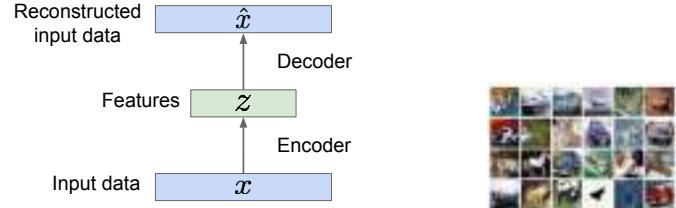
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 41 May 15, 2018

Some background first: Autoencoders

How to learn this feature representation?

Train such that features can be used to reconstruct original data
"Autoencoding" - encoding itself



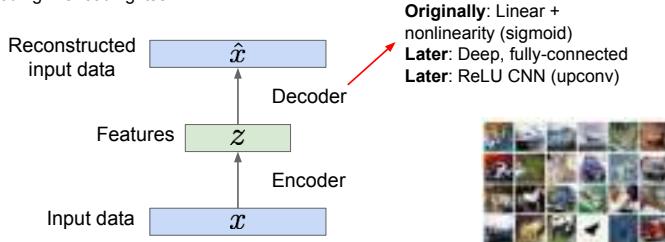
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 42 May 15, 2018

Some background first: Autoencoders

How to learn this feature representation?

Train such that features can be used to reconstruct original data
"Autoencoding" - encoding itself



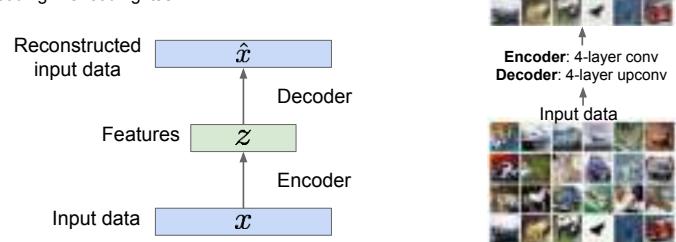
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 43 May 15, 2018

Some background first: Autoencoders

How to learn this feature representation?

Train such that features can be used to reconstruct original data
"Autoencoding" - encoding itself

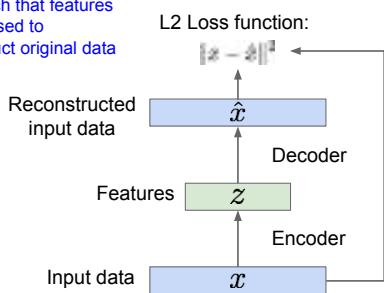


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 44 May 15, 2018

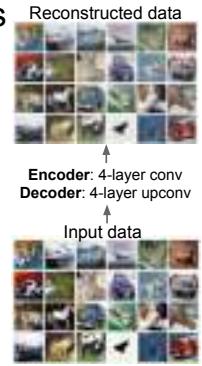
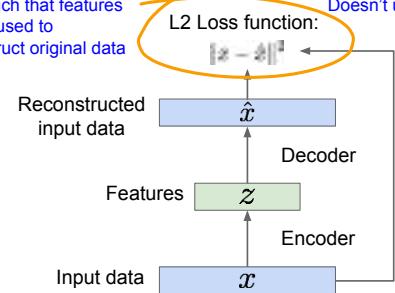
Some background first: Autoencoders

Train such that features can be used to reconstruct original data

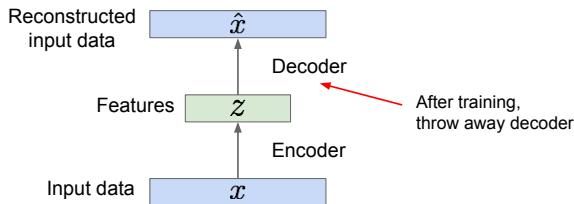


Some background first: Autoencoders

Train such that features can be used to reconstruct original data

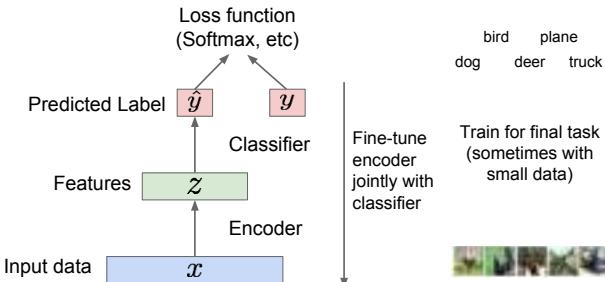


Some background first: Autoencoders

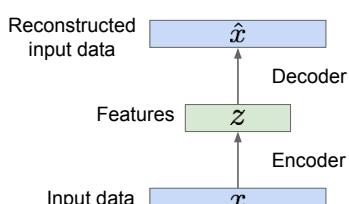


Some background first: Autoencoders

Encoder can be used to initialize a supervised model



Some background first: Autoencoders



Autoencoders can reconstruct data, and can learn features to initialize a supervised model

Features capture factors of variation in training data. Can we generate new images from an autoencoder?

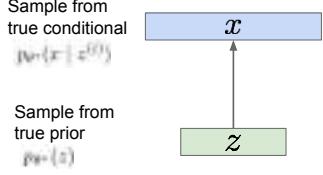
Variational Autoencoders

Probabilistic spin on autoencoders - will let us sample from the model to generate data!

Variational Autoencoders

Probabilistic spin on autoencoders - will let us sample from the model to generate data!

Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation z



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

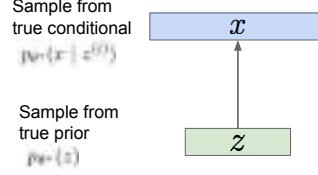
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 51 May 15, 2018

Variational Autoencoders

Probabilistic spin on autoencoders - will let us sample from the model to generate data!

Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from underlying unobserved (latent) representation z



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Intuition (remember from autoencoders!):
 x is an image, z is latent factors used to generate x : attributes, orientation, etc.

Fei-Fei Li & Justin Johnson & Serena Yeung

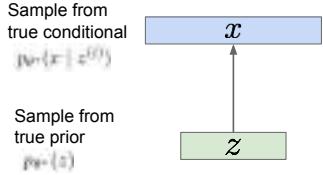
Lecture 12 - 51 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 52 May 15, 2018

Variational Autoencoders

We want to estimate the true parameters θ^* of this generative model.



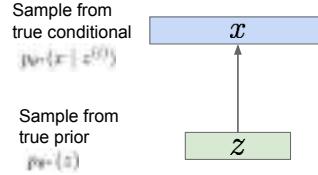
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 53 May 15, 2018

Variational Autoencoders

We want to estimate the true parameters θ^* of this generative model.



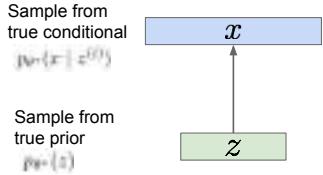
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 54 May 15, 2018

Variational Autoencoders

We want to estimate the true parameters θ^* of this generative model.



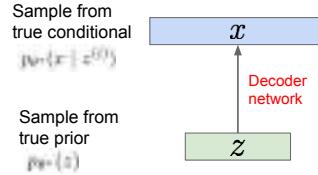
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 55 May 15, 2018

Variational Autoencoders

We want to estimate the true parameters θ^* of this generative model.

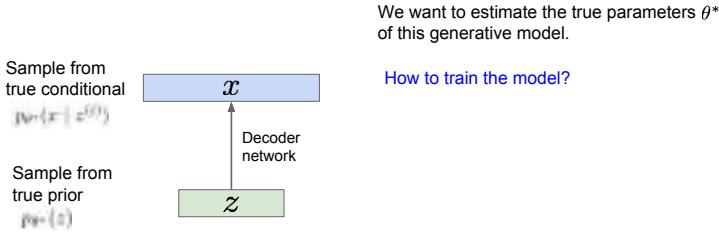


Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

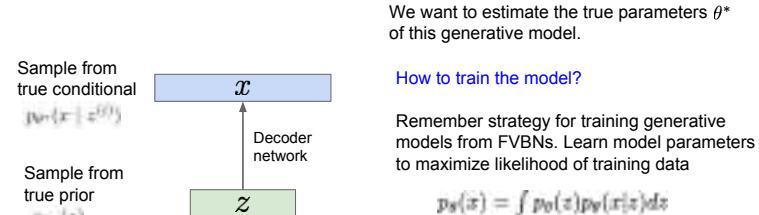
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 56 May 15, 2018

Variational Autoencoders



Variational Autoencoders



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

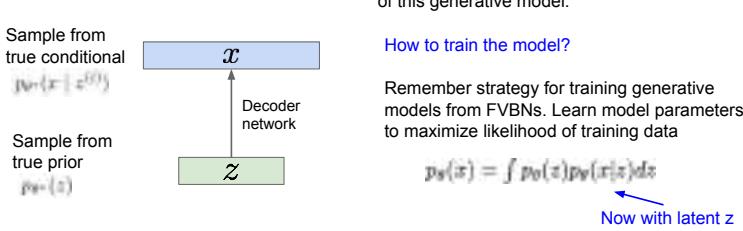
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 57 May 15, 2018

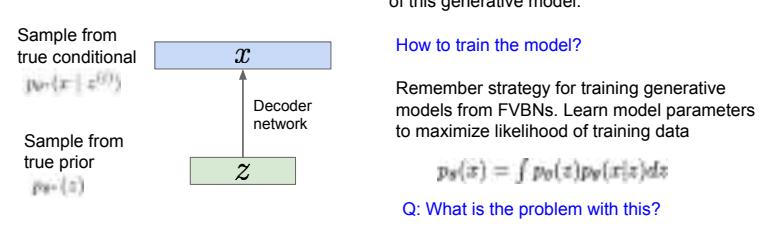
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 58 May 15, 2018

Variational Autoencoders



Variational Autoencoders



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

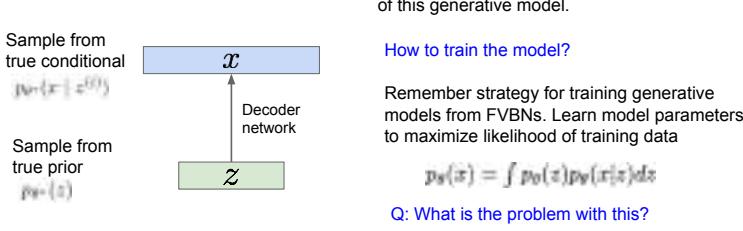
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 59 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 60 May 15, 2018

Variational Autoencoders



Variational Autoencoders: Intractability

Data likelihood: $p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 61 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 62 May 15, 2018

Variational Autoencoders: Intractability

Data likelihood: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

Simple Gaussian prior

Variational Autoencoders: Intractability

Data likelihood: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

Decoder neural network

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Autoencoders: Intractability

Data likelihood: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

Intractible to compute $p(x|z)$ for every z !

Variational Autoencoders: Intractability

Data likelihood: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

Posterior density also intractable: $p_\theta(z|x) = p_\theta(x|z)p_\theta(z)/p_\theta(x)$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Autoencoders: Intractability

Data likelihood: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

Posterior density also intractable: $p_\theta(z|x) = p_\theta(x|z)p_\theta(z)/p_\theta(x)$

Intractable data likelihood

Variational Autoencoders: Intractability

Data likelihood: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

Posterior density also intractable: $p_\theta(z|x) = p_\theta(x|z)p_\theta(z)/p_\theta(x)$

Solution: In addition to decoder network modeling $p_\theta(x|z)$, define additional encoder network $q_\phi(z|x)$ that approximates $p_\theta(z|x)$

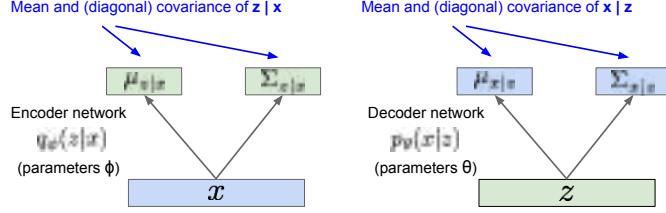
Will see that this allows us to derive a lower bound on the data likelihood that is tractable, which we can optimize

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Autoencoders

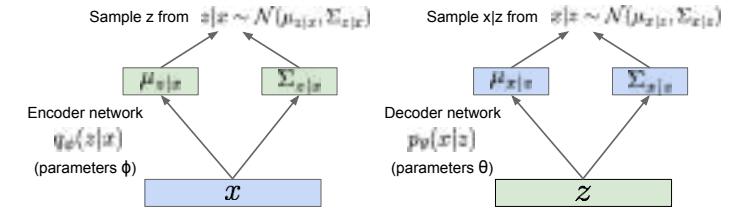
Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Autoencoders

Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

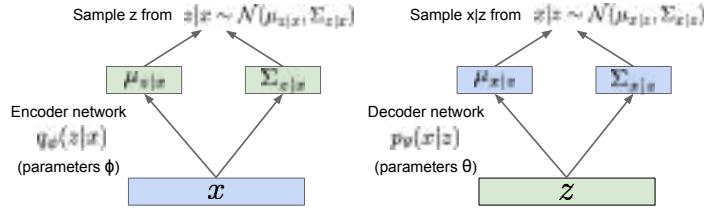
Lecture 12 - 69 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 70 May 15, 2018

Variational Autoencoders

Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic



Encoder and decoder networks also called
"recognition"/"inference" and "generation" networks

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 71 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 72 May 15, 2018

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

Taking expectation wrt. z
(using encoder network) will
come in handy later

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned} \log p_\theta(x^{(i)}) &= \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbb{E}_z \left[\log \frac{p_\theta(x^{(i)}|z) p_\phi(z)}{\int p_\theta(x^{(i)}|z') p_\phi(z') dz'} \right] \quad (\text{Bayes' Rule}) \end{aligned}$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 73 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 74 May 15, 2018

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned} \log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \end{aligned}$$

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned} \log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] + \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithm}) \end{aligned}$$

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned} \log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] + \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithm}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)})) \end{aligned}$$

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned} \log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] + \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithm}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)})) \end{aligned}$$

The expectation wrt. z (using encoder network) let us write nice KL terms

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned} \log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] + \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithm}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)})) \end{aligned}$$

↑
Decoder network gives $p_{\theta}(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)

↑
This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!

↑
 $p_{\theta}(z|x)$ intractable (saw earlier), can't compute this KL term. But we know KL divergence always ≥ 0 .

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned} \log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbb{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] + \mathbb{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithm}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)})) \end{aligned}$$

$\boxed{\mathcal{L}(p_{\theta}^{(i)}, \theta, \phi)}$
Tractable lower bound which we can take gradient of and optimize! ($p_{\theta}(x|z)$ differentiable, KL term differentiable)

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned}
 \log p_\theta(x^{(i)}) &= \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] - (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\
 &= \mathbb{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\
 &= \mathbb{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\
 &= \mathbb{E}_z \left[\log q_\phi(x^{(i)} | z) \right] - \mathbb{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] + \mathbb{E}_z \left[\log \frac{p_\theta(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\
 &= \mathbb{E}_z \left[\log q_\phi(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + D_{KL}(p_\theta(z | x^{(i)}) || q_\phi(z | x^{(i)}))
 \end{aligned}$$

$$\mathcal{L}(\phi^{(i)}, \theta, \phi)$$

Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned}
 \log p_\theta(x^{(i)}) &= \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] - (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\
 &= \mathbb{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\
 \text{Reconstruct} \quad \text{the input data} &= \mathbb{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\
 &= \mathbb{E}_z \left[\log q_\phi(x^{(i)} | z) \right] - \mathbb{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] + \mathbb{E}_z \left[\log \frac{p_\theta(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\
 &= \mathbb{E}_z \left[\log q_\phi(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + D_{KL}(p_\theta(z | x^{(i)}) || q_\phi(z | x^{(i)}))
 \end{aligned}$$

$$\mathcal{L}(\phi^{(i)}, \theta, \phi)$$

Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\mathbb{E}_z \left[\log q_\phi(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$$\mathcal{L}(\phi^{(i)}, \theta, \phi)$$

Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\mathbb{E}_z \left[\log q_\phi(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$$\mathcal{L}(\phi^{(i)}, \theta, \phi)$$

Let's look at computing the bound (forward pass) for a given minibatch of input data

Input Data x

Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\mathbb{E}_z \left[\log q_\phi(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$$\mathcal{L}(\phi^{(i)}, \theta, \phi)$$

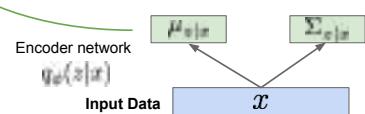
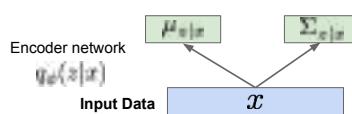
Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\mathbb{E}_z \left[\log q_\phi(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$$\mathcal{L}(\phi^{(i)}, \theta, \phi)$$

Make approximate posterior distribution close to prior

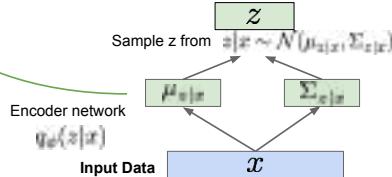


Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\frac{\mathbb{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x) || p_\theta(z))}{\mathcal{L}(\phi^{(i)}, \theta, \bar{n})}$$

Make approximate posterior distribution close to prior



Fei-Fei Li & Justin Johnson & Serena Yeung

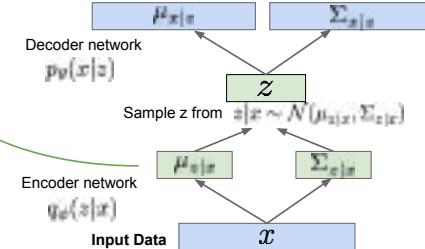
Lecture 12 - 87 May 15, 2018

Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\frac{\mathbb{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x) || p_\theta(z))}{\mathcal{L}(\phi^{(i)}, \theta, \bar{n})}$$

Make approximate posterior distribution close to prior



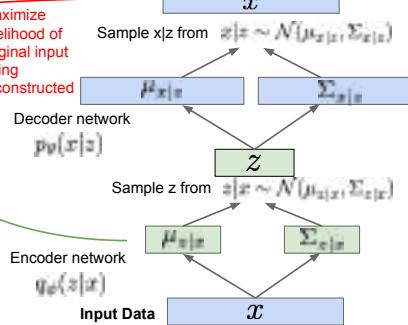
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 88 May 15, 2018

Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\frac{\mathbb{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x) || p_\theta(z))}{\mathcal{L}(\phi^{(i)}, \theta, \bar{n})}$$



Fei-Fei Li & Justin Johnson & Serena Yeung

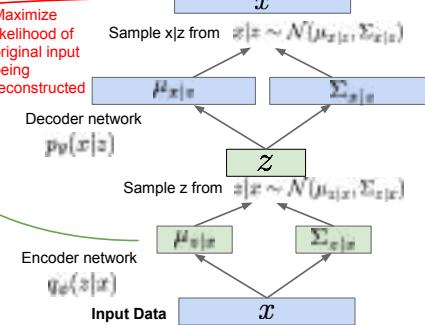
Lecture 12 - 89 May 15, 2018

Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\frac{\mathbb{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x) || p_\theta(z))}{\mathcal{L}(\phi^{(i)}, \theta, \bar{n})}$$

Make approximate posterior distribution close to prior

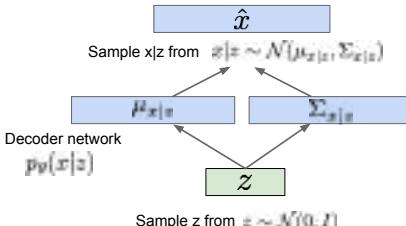


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 90 May 15, 2018

Variational Autoencoders: Generating Data!

Use decoder network. Now sample z from prior!



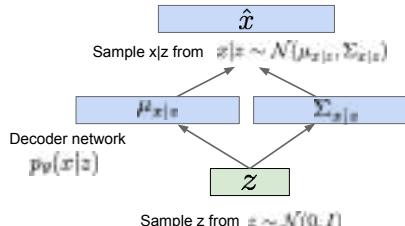
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

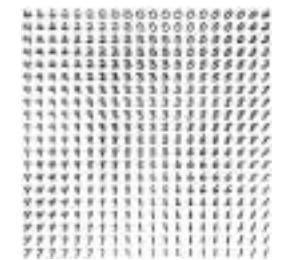
Lecture 12 - 91 May 15, 2018

Variational Autoencoders: Generating Data!

Use decoder network. Now sample z from prior!



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

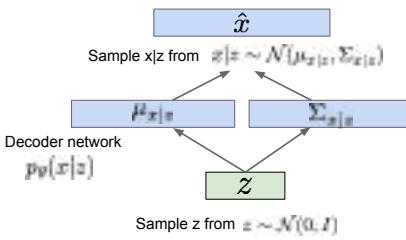


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 92 May 15, 2018

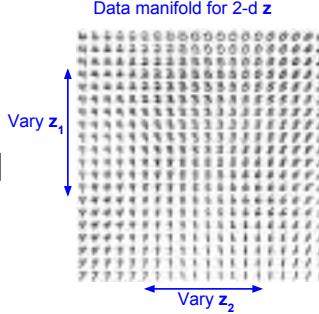
Variational Autoencoders: Generating Data!

Use decoder network. Now sample z from prior!



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung



Lecture 12 - 93 May 15, 2018

Variational Autoencoders: Generating Data!

Diagonal prior on z
=> independent latent variables

Different dimensions of z encode interpretable factors of variation

Degree of smile

Vary z_1



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 94 May 15, 2018

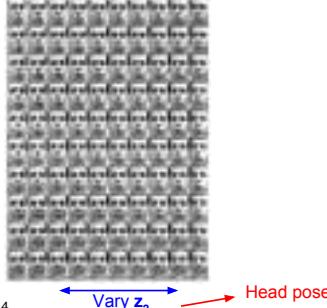
Variational Autoencoders: Generating Data!

Diagonal prior on z
=> independent latent variables

Different dimensions of z encode interpretable factors of variation

Also good feature representation that can be computed using $q_\phi(z|x)$!

Degree of smile
Vary z_1



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 95 May 15, 2018

Variational Autoencoders: Generating Data!



Figures copyright (L) Dirk Kingma et al. 2014; (R) Anders Larsen et al. 2017. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 96 May 15, 2018

Variational Autoencoders

Probabilistic spin to traditional autoencoders => allows generating data
Defines an intractable density => derive and optimize a (variational) lower bound

Pros:

- Principled approach to generative models
- Allows inference of $q(z|x)$, can be useful feature representation for other tasks

Cons:

- Maximizes lower bound of likelihood: okay, but not as good evaluation as PixelRNN/PixelCNN
- Samples blurrier and lower quality compared to state-of-the-art (GANs)

Active areas of research:

- More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian
- Incorporating structure in latent variables

Generative Adversarial Networks (GAN)

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 97 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 98 May 15, 2018

So far...

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(\mathbf{z}) p_{\theta}(x|\mathbf{z}) d\mathbf{z}$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

So far...

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(\mathbf{z}) p_{\theta}(x|\mathbf{z}) d\mathbf{z}$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

What if we give up on explicitly modeling density, and just want ability to sample?

So far...

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent \mathbf{z} :

$$p_{\theta}(x) = \int p_{\theta}(\mathbf{z}) p_{\theta}(x|\mathbf{z}) d\mathbf{z}$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

What if we give up on explicitly modeling density, and just want ability to sample?

GANs: don't work with any explicit density function!

Instead, take game-theoretic approach: learn to generate from training distribution through 2-player game

Generative Adversarial Networks

Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution, e.g. random noise. Learn transformation to training distribution.

Q: What can we use to represent this complex transformation?

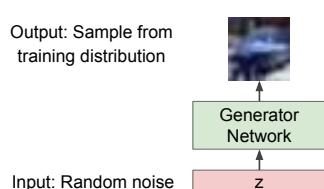
Generative Adversarial Networks

Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution, e.g. random noise. Learn transformation to training distribution.

Q: What can we use to represent this complex transformation?

A: A neural network!



Training GANs: Two-player game

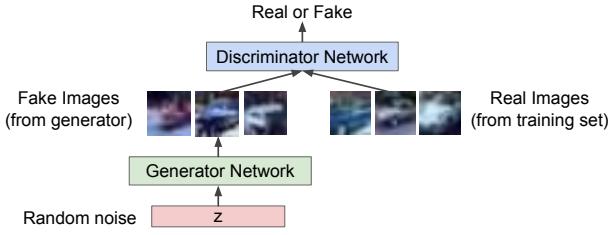
Generator network: try to fool the discriminator by generating real-looking images

Discriminator network: try to distinguish between real and fake images

Training GANs: Two-player game

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Generator network: try to fool the discriminator by generating real-looking images
Discriminator network: try to distinguish between real and fake images



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 10 May 15, 2018

Training GANs: Two-player game

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Generator network: try to fool the discriminator by generating real-looking images
Discriminator network: try to distinguish between real and fake images

Train jointly in **minimax game**

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 10 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 10 May 15, 2018

Training GANs: Two-player game

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Generator network: try to fool the discriminator by generating real-looking images
Discriminator network: try to distinguish between real and fake images

Train jointly in **minimax game**

Discriminator outputs likelihood in (0,1) of real image

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\underbrace{\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x)}_{\text{Discriminator output for real data } x} + \underbrace{\mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))}_{\text{Discriminator output for generated fake data } G(z)} \right]$$

Training GANs: Two-player game

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Generator network: try to fool the discriminator by generating real-looking images
Discriminator network: try to distinguish between real and fake images

Train jointly in **minimax game**

Discriminator outputs likelihood in (0,1) of real image

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\underbrace{\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x)}_{\text{Discriminator output for real data } x} + \underbrace{\mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))}_{\text{Discriminator output for generated fake data } G(z)} \right]$$

- Discriminator (θ_d) wants to **maximize objective** such that $D(x)$ is close to 1 (real) and $D(G(z))$ is close to 0 (fake)
- Generator (θ_g) wants to **minimize objective** such that $D(G(z))$ is close to 1 (discriminator is fooled into thinking generated $G(z)$ is real)

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 10 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 10 May 15, 2018

Training GANs: Two-player game

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. **Gradient ascent** on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. **Gradient descent** on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

Training GANs: Two-player game

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. **Gradient ascent** on discriminator

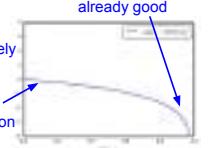
$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. **Gradient descent** on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

In practice, optimizing this generator objective does not work well!

When sample is likely fake, want to learn from it to improve generator. But gradient in this region is relatively flat!



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 10 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 11 May 15, 2018

Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Instead: Gradient ascent on generator, different objective

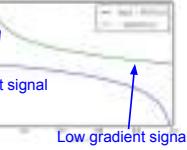
$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log [D_{\theta_d}(G_{\theta_g}(z))]$$

Instead of minimizing likelihood of discriminator being correct, now maximize likelihood of discriminator being wrong.

Same objective of fooling discriminator, but now higher gradient signal for bad samples => works much better! Standard in practice.

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014



Aside: Jointly training two networks is challenging, can be unstable. Choosing objectives with better loss landscapes helps training, is an active area of research.

Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

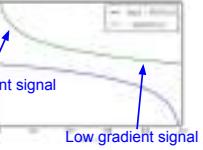
Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Instead: Gradient ascent on generator, different objective

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log [D_{\theta_d}(G_{\theta_g}(z))]$$



Instead of minimizing likelihood of discriminator being correct, now maximize likelihood of discriminator being wrong.

Same objective of fooling discriminator, but now higher gradient signal for bad samples => works much better! Standard in practice.

Training GANs: Two-player game

Putting it together: GAN training algorithm

```
for number of training iterations do
    for i steps do
        • Sample minibatch of m noise samples {z(1), ..., z(m)} from noise prior pz(z).
        • Sample minibatch of n examples {x(1), ..., x(n)} from data generating distribution pgen(x).
        • Update the discriminator by ascending its stochastic gradient:
             $\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D_{\theta_d}(x^{(i)}) + \log (1 - D_{\theta_d}(G_{\theta_g}(z^{(i)})))]$ 
    and for
        • Sample minibatch of m noise samples {z(1), ..., z(m)} from noise prior pz(z).
        • Update the generator by ascending its stochastic gradient. (improved objective)
             $\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m [\log (1 - D_{\theta_d}(G_{\theta_g}(z^{(i)})))]$ 
    end for
```

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Training GANs: Two-player game

Putting it together: GAN training algorithm

```
for number of training iterations do
    for i steps do
        • Sample minibatch of m noise samples {z(1), ..., z(m)} from noise prior pz(z).
        • Sample minibatch of n examples {x(1), ..., x(n)} from data generating distribution pgen(x).
        • Update the discriminator by ascending its stochastic gradient:
             $\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D_{\theta_d}(x^{(i)}) + \log (1 - D_{\theta_d}(G_{\theta_g}(z^{(i)})))]$ 
Some find k=1 more stable, others use k > 1, no best rule.
Recent work (e.g. Wasserstein GAN) alleviates this problem, better stability!
    and for
        • Sample minibatch of m noise samples {z(1), ..., z(m)} from noise prior pz(z).
        • Update the generator by ascending its stochastic gradient. (improved objective)
             $\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m [\log (1 - D_{\theta_d}(G_{\theta_g}(z^{(i)})))]$ 
    end for
```

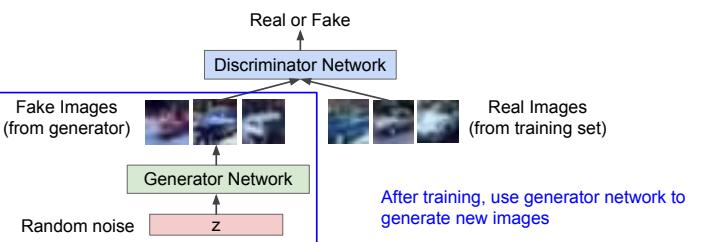
Training GANs: Two-player game

Generator network: try to fool the discriminator by generating real-looking images

Discriminator network: try to distinguish between real and fake images

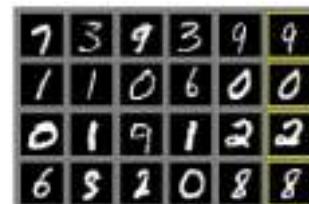
Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014



Generative Adversarial Nets

Generated samples



Nearest neighbor from training set

Figures copyright Ian Goodfellow et al., 2014. Reproduced with permission.

Generative Adversarial Nets

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Generated samples (CIFAR-10)



Figures copyright Ian Goodfellow et al., 2014. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 11
May 15, 2018

Generative Adversarial Nets: Convolutional Architectures

Generator is an upsampling network with fractionally-strided convolutions
Discriminator is a convolutional network

Architecture guidelines for stable Deep Convolutional GANs

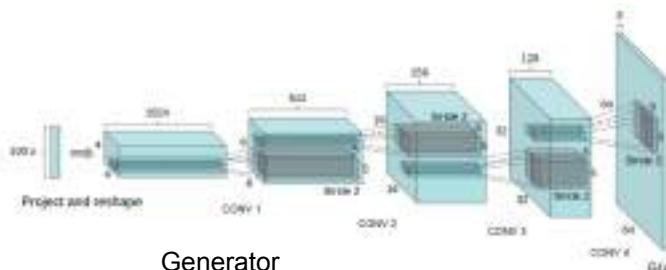
- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 11
May 15, 2018

Generative Adversarial Nets: Convolutional Architectures



Generator

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 11
May 15, 2018

Generative Adversarial Nets: Convolutional Architectures



Samples from the model look much better!

Radford et al, ICLR 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 12
May 15, 2018

Generative Adversarial Nets: Convolutional Architectures



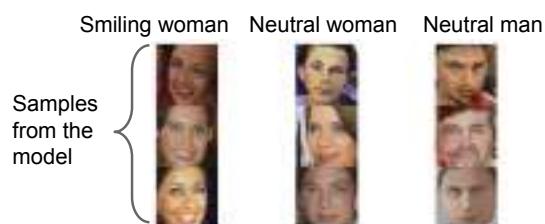
Interpolating between random points in latent space

Radford et al, ICLR 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 12
May 15, 2018

Generative Adversarial Nets: Interpretable Vector Math



Radford et al, ICLR 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 12
May 15, 2018

2017: Explosion of GANs

Better training and generation



LSGAN, Zhu 2017.



Wasserstein GAN,
Arjovsky 2017.
Improved Wasserstein
GAN, Gulrajani 2017.



Progressive GAN, Karras 2018.

2017: Explosion of GANs

Source->Target domain transfer



CycleGAN, Zhu et al. 2017.

Text -> Image Synthesis



Reed et al. 2017.

Many GAN applications



Pix2pix, Isola 2017. Many examples at
<https://phillipi.github.io/pix2pix/>

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 12 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 13 May 15, 2018

GANs

Don't work with an explicit density function

Take game-theoretic approach: learn to generate from training distribution through 2-player game

Pros:

- Beautiful, state-of-the-art samples!

Cons:

- Trickier / more unstable to train
- Can't solve inference queries such as $p(x)$, $p(z|x)$

Active areas of research:

- Better loss functions, more stable training (Wasserstein GAN, LSGAN, many others)
- Conditional GANs, GANs for all kinds of applications

Recap

Generative Models

- PixelRNN and PixelCNN Explicit density model, optimizes exact likelihood, good samples. But inefficient sequential generation.
- Variational Autoencoders (VAE) Optimize variational lower bound on likelihood. Useful latent representation, inference queries. But current sample quality not the best.
- Generative Adversarial Networks (GANs) Game-theoretic approach, best samples! But can be tricky and unstable to train, no inference queries.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 13 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 13 May 15, 2018

Recap

Generative Models

- PixelRNN and PixelCNN Explicit density model, optimizes exact likelihood, good samples. But inefficient sequential generation.
- Variational Autoencoders (VAE) Optimize variational lower bound on likelihood. Useful latent representation, inference queries. But current sample quality not the best.
- Generative Adversarial Networks (GANs) Game-theoretic approach, best samples! But can be tricky and unstable to train, no inference queries.

Also recent work in combinations of
these types of models! E.g., Adversarial
Autoencoders (Makhzani 2015) and
PixelVAE (Gulrajani 2016)

Recap

Generative Models

- PixelRNN and PixelCNN Explicit density model, optimizes exact likelihood, good samples. But inefficient sequential generation.
- Variational Autoencoders (VAE) Optimize variational lower bound on likelihood. Useful latent representation, inference queries. But current sample quality not the best.
- Generative Adversarial Networks (GANs) Game-theoretic approach, best samples! But can be tricky and unstable to train, no inference queries.

Next time: Visualizing and understanding networks

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 13 May 15, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 12 - 13 May 15, 2018

Administrative

- Project milestone was due yesterday 5/16
- Assignment 3 due a week from yesterday 5/23
- Google Cloud:
 - See pinned Piazza post about GPUs
 - Make a private post on Piazza if you need more credits

Lecture 13: Visualizing and Understanding

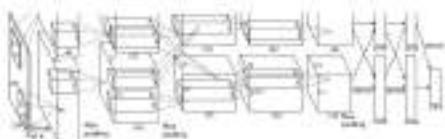
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 1 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 2 May 17, 2018

What's going on inside ConvNets?

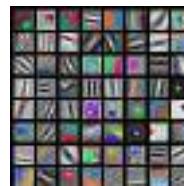


Input Image:
 $3 \times 224 \times 224$

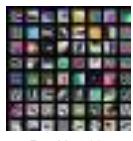
What are the intermediate features looking for?

Class Scores:
1000 numbers

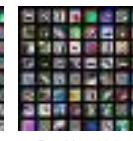
First Layer: Visualize Filters



AlexNet:
 $64 \times 3 \times 11 \times 11$



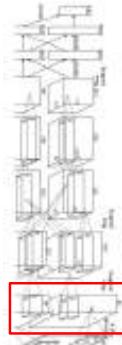
ResNet-18:
 $64 \times 3 \times 7 \times 7$



ResNet-101:
 $64 \times 3 \times 7 \times 7$



DenseNet-121:
 $64 \times 3 \times 7 \times 7$



Krizhevsky et al., "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.
Figure reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 3 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 4 May 17, 2018

Visualize the filters/kernels (raw weights)



layer 1 weights
 $16 \times 3 \times 7 \times 7$

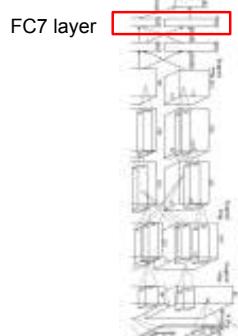
layer 2 weights
 $20 \times 16 \times 7 \times 7$

layer 3 weights
 $20 \times 20 \times 7 \times 7$

Last Layer

4096-dimensional feature vector for an image
(layer immediately before the classifier)

Run the network on many images, collect the
feature vectors



We can visualize
filters at higher
layers, but not
that interesting

(these are taken
from ConvNetJS
CIFAR-10
demo)

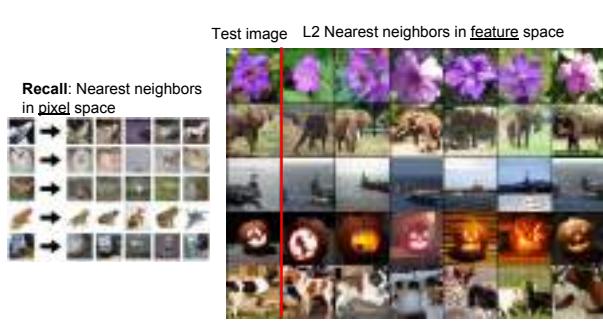
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 5 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

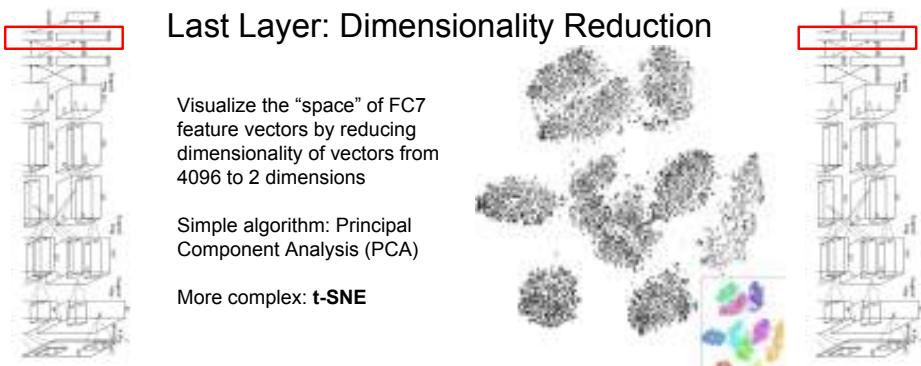
Lecture 13 - 6 May 17, 2018

Last Layer: Nearest Neighbors



Fei-Fei Li & Justin Johnson & Serena Yeung

Last Layer: Dimensionality Reduction



Lecture 13 - 7 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 8 May 17, 2018

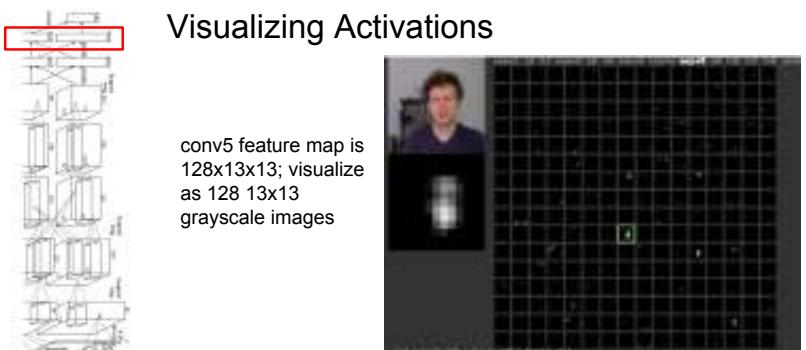
Last Layer: Dimensionality Reduction



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 9 May 17, 2018

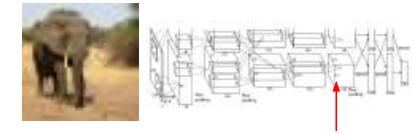
Visualizing Activations



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 10 May 17, 2018

Maximally Activating Patches



Pick a layer and a channel; e.g. conv5 is 128 x 13 x 13, pick channel 17/128

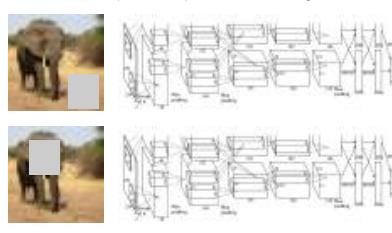
Run many images through the network, record values of chosen channel

Visualize image patches that correspond to maximal activations



Which pixels matter: Saliency vs Occlusion

Mask part of the image before feeding to CNN, check how much predicted probabilities change



Fei-Fei Li & Justin Johnson & Serena Yeung

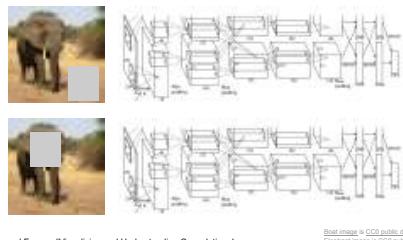
Lecture 13 - 11 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

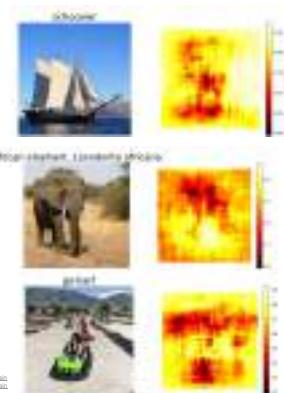
Lecture 13 - 12 May 17, 2018

Which pixels matter: Saliency vs Occlusion

Mask part of the image before feeding to CNN,
check how much predicted probabilities change



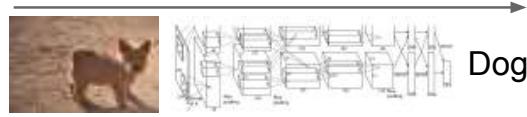
Boat image in CC-BY public domain
Background image in CC-BY license from
pckaren image in CC-BY public domain



Boat image in CC-BY public domain
Background image in CC-BY license from
pckaren image in CC-BY public domain

Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities



Dog

Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

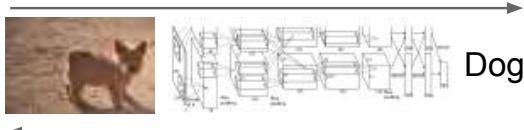
Lecture 13 - 13 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 14 May 17, 2018

Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities



Dog

Compute gradient of (unnormalized) class score with respect to image pixels, take absolute value and max over RGB channels



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Saliency Maps



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

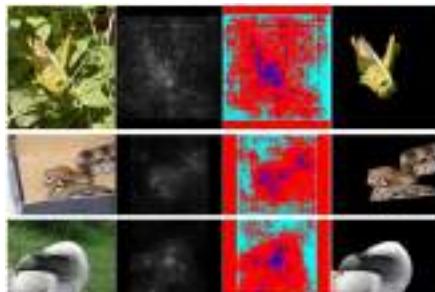
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 15 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 16 May 17, 2018

Saliency Maps: Segmentation without supervision



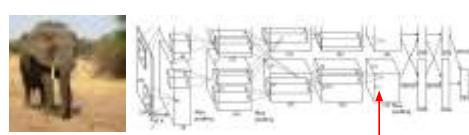
Use GrabCut on
saliency map

Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 17 May 17, 2018

Intermediate Features via (guided) backprop



Pick a single intermediate neuron, e.g. one value in $128 \times 13 \times 13$ conv5 feature map

Compute gradient of neuron value with respect to image pixels

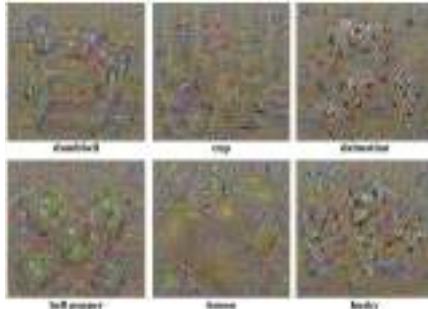
Zeller and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014.
Springenberg et al., "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015.

Lecture 13 - 18 May 17, 2018

Visualizing CNN features: Gradient Ascent

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

Simple regularizer: Penalize L2 norm of generated image



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.
Figure copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

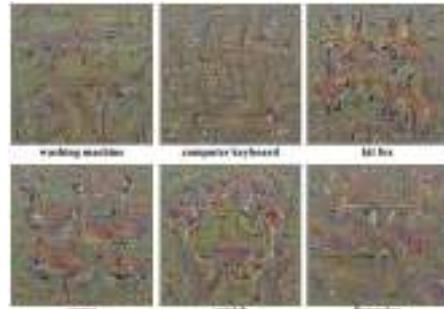
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 25 May 17, 2018

Visualizing CNN features: Gradient Ascent

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

Simple regularizer: Penalize L2 norm of generated image



Yosinski et al., "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.
Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014; Reproduced with permission.

Lecture 13 - 26 May 17, 2018

Visualizing CNN features: Gradient Ascent

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

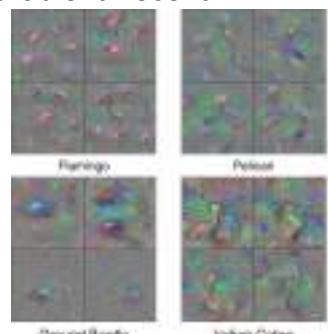
- (1) Gaussian blur image
- (2) Clip pixels with small values to 0
- (3) Clip pixels with small gradients to 0

Visualizing CNN features: Gradient Ascent

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

- (1) Gaussian blur image
- (2) Clip pixels with small values to 0
- (3) Clip pixels with small gradients to 0



Yosinski et al., "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.
Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014; Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 27 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

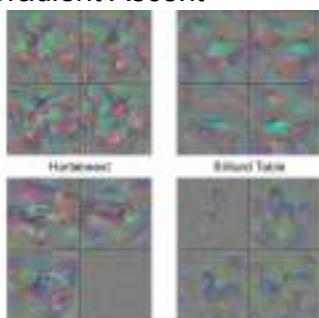
Lecture 13 - 28 May 17, 2018

Visualizing CNN features: Gradient Ascent

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

- (1) Gaussian blur image
- (2) Clip pixels with small values to 0
- (3) Clip pixels with small gradients to 0



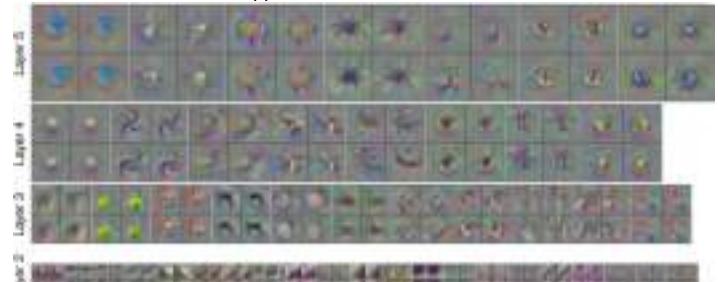
Yosinski et al., "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.
Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014; Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 29 May 17, 2018

Visualizing CNN features: Gradient Ascent

Use the same approach to visualize intermediate features



Yosinski et al., "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.
Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014; Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 30 May 17, 2018

Visualizing CNN features: Gradient Ascent

Adding "multi-faceted" visualization gives even nicer results:
(Plus more careful regularization, center-bias)



Nguyen et al., "Multi-faceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks", ICML Visualization for Deep Learning Workshop 2016.

Figure copyright Anh Nguyen, Jason Yosinski, and Jeff Clune, 2016; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 31 May 17, 2018

Visualizing CNN features: Gradient Ascent



Nguyen et al., "Multi-faceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks", ICML Visualization for Deep Learning Workshop 2016.

Figure copyright Anh Nguyen, Jason Yosinski, and Jeff Clune, 2016; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 32 May 17, 2018

Visualizing CNN features: Gradient Ascent

Optimize in FC6 latent space instead of pixel space:



Nguyen et al., "Synthesizing the preferred inputs for neurons in neural networks via deep generator networks", NIPS 2016

Figure copyright Nguyen et al., 2016; reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 33 May 17, 2018

Fooling Images / Adversarial Examples

- (1) Start from an arbitrary image
- (2) Pick an arbitrary class
- (3) Modify the image to maximize the class
- (4) Repeat until network is fooled

Fooling Images / Adversarial Examples



Boat image is CC0 public domain

Elephant image is CC0 public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 35 May 17, 2018

Fooling Images / Adversarial Examples



Check out [Ian Goodfellow's lecture](#) from last year

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 36 May 17, 2018



Image is licensed under CC BY 3.0

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 43 May 17, 2018



Image is licensed under CC BY 4.0

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 44 May 17, 2018



Image is licensed under CC BY 4.0

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 45 May 17, 2018



Image is licensed under CC BY 3.0

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 46 May 17, 2018



Image is licensed under CC BY 3.0

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 47 May 17, 2018



Image is licensed under CC BY 4.0

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 48 May 17, 2018

Feature Inversion

Given a CNN feature vector for an image, find a new image that:

- Matches the given feature vector
- "looks natural" (image prior regularization)

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Given feature vector

$$\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2$$

$$\mathcal{R}_V(\mathbf{x}) = \sum_{i,j} \left((x_{i,j+1} - x_{i,j})^2 + (x_{i+1,j} - x_{i,j})^2 \right)^{\frac{1}{2}}$$

Features of new image

Total Variation regularizer
(encourages spatial smoothness)

Mahendran and Vedaldi, "Understanding Deep Image Representations by Inverting Them", CVPR 2015

Feature Inversion

Reconstructing from different layers of VGG-16



Mahendran and Vedaldi, "Understanding Deep Image Representations by Inverting Them", CVPR 2015
Figure from Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016. Copyright Springer, 2016.
Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

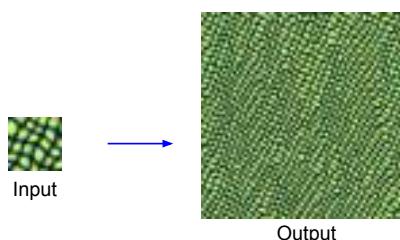
Lecture 13 - 49 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 50 May 17, 2018

Texture Synthesis

Given a sample patch of some texture, can we generate a bigger image of the same texture?



[Output image](#) is licensed under the [MIT license](#)

Texture Synthesis: Nearest Neighbor

Generate pixels one at a time in scanline order; form neighborhood of already generated pixels and copy nearest neighbor from input



Fei-Fei Li & Justin Johnson & Serena Yeung

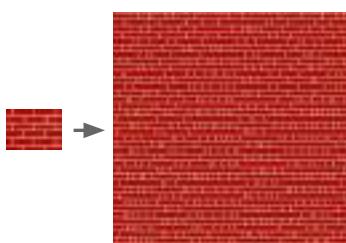
Lecture 13 - 51 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

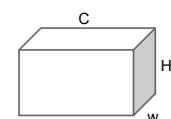
Lecture 13 - 52 May 17, 2018

Texture Synthesis: Nearest Neighbor

Neural Texture Synthesis: Gram Matrix



[Image](#) licensed under the [MIT license](#)



Each layer of CNN gives $C \times H \times W$ tensor of features; $H \times W$ grid of C -dimensional vectors

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 53 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

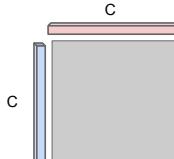
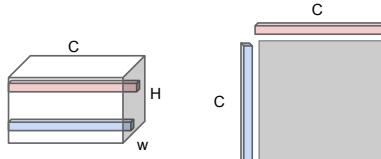
Lecture 13 - 54 May 17, 2018

Neural Texture Synthesis: Gram Matrix

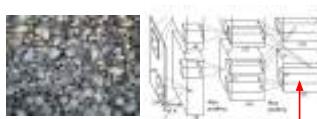


Each layer of CNN gives $C \times H \times W$ tensor of features; $H \times W$ grid of C -dimensional vectors

Outer product of two C -dimensional vectors gives $C \times C$ matrix measuring co-occurrence



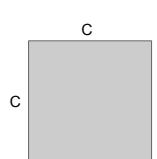
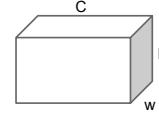
Neural Texture Synthesis: Gram Matrix



Each layer of CNN gives $C \times H \times W$ tensor of features; $H \times W$ grid of C -dimensional vectors

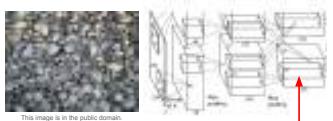
Outer product of two C -dimensional vectors gives $C \times C$ matrix measuring co-occurrence

Average over all HW pairs of vectors, giving **Gram matrix** of shape $C \times C$



Gram Matrix

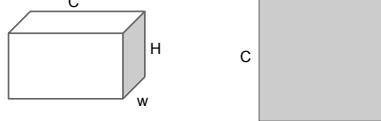
Neural Texture Synthesis: Gram Matrix



Each layer of CNN gives $C \times H \times W$ tensor of features; $H \times W$ grid of C -dimensional vectors

Outer product of two C -dimensional vectors gives $C \times C$ matrix measuring co-occurrence

Average over all HW pairs of vectors, giving **Gram matrix** of shape $C \times C$

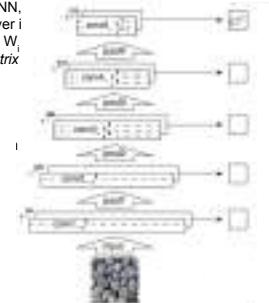


Efficient to compute; reshape features from $C \times H \times W$ to $=C \times HW$
then compute $G = FFT$

Neural Texture Synthesis

1. Pretrain a CNN on ImageNet (VGG-19)
2. Run input texture forward through CNN, record activations on every layer; layer i gives feature map of shape $C_i \times H_i \times W_i$
3. At each layer compute the **Gram matrix** giving outer product of features:

$$G_{ij} = \sum_k F_{ik}^T F_{jk} \quad (\text{shape } C_i \times C_j)$$

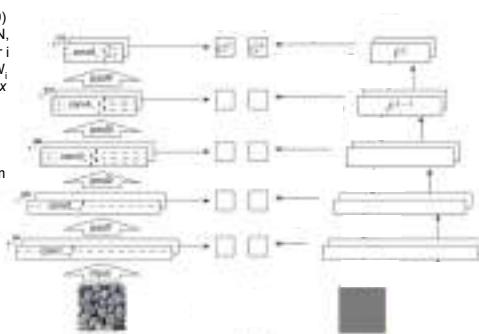


Gatys, Ecker, and Bethge, "Texture Synthesis Using Convolutional Neural Networks", NIPS 2015
Figure copyright Leon Gatys, Alexander S. Ecker, and Matthias Bethge, 2015. Reproduced with permission.

Neural Texture Synthesis

1. Pretrain a CNN on ImageNet (VGG-19)
2. Run input texture forward through CNN, record activations on every layer; layer i gives feature map of shape $C_i \times H_i \times W_i$
3. At each layer compute the **Gram matrix** giving outer product of features:

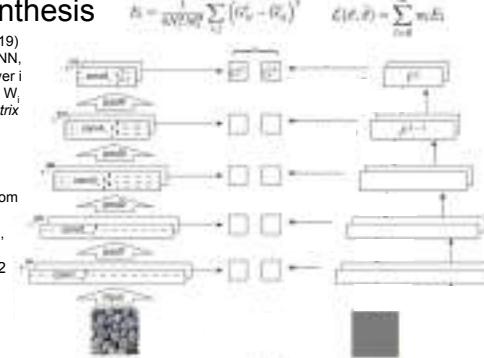
$$G_{ij} = \sum_k F_{ik}^T F_{jk} \quad (\text{shape } C_i \times C_j)$$



Neural Texture Synthesis

1. Pretrain a CNN on ImageNet (VGG-19)
2. Run input texture forward through CNN, record activations on every layer; layer i gives feature map of shape $C_i \times H_i \times W_i$
3. At each layer compute the **Gram matrix** giving outer product of features:

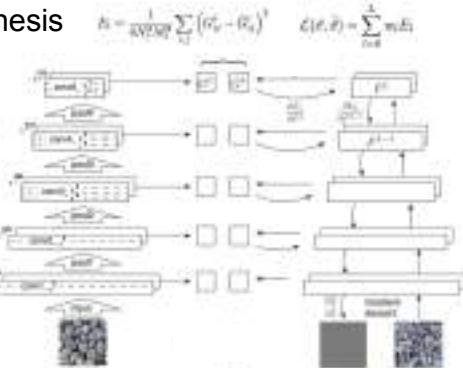
$$G_{ij} = \sum_k F_{ik}^T F_{jk} \quad (\text{shape } C_i \times C_j)$$



Gatys, Ecker, and Bethge, "Texture Synthesis Using Convolutional Neural Networks", NIPS 2015
Figure copyright Leon Gatys, Alexander S. Ecker, and Matthias Bethge, 2015. Reproduced with permission.

Neural Texture Synthesis

1. Pretrain a CNN on ImageNet (VGG-19)
2. Run input texture forward through CNN, record activations on every layer; layer i gives feature map of shape $C_i \times H_i \times W_i$
3. At each layer compute the *Gram matrix* giving outer product of features:
$$G_i = \sum_{k=1}^{C_i} F_k F_k^T \quad (\text{shape } C_i \times C_i)$$
4. Initialize generated image from random noise
5. Pass generated image through CNN, compute Gram matrix on each layer
6. Compute loss: weighted sum of L2 distance between Gram matrices
7. Backprop to get gradient on image
8. Make gradient step on image
9. GOTO 5



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 61 May 17, 2018

Neural Texture Synthesis

Reconstructing texture from higher layers recovers larger features from the input texture



Gatys, Ecker, and Bethge, "Texture Synthesis Using Convolutional Neural Networks", NIPS 2015. Copyright Leon Gatys, Alexander S. Ecker, and Matthias Bethge, 2015. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 62 May 17, 2018

Neural Texture Synthesis: Texture = Artwork

Texture synthesis
(Gram reconstruction)

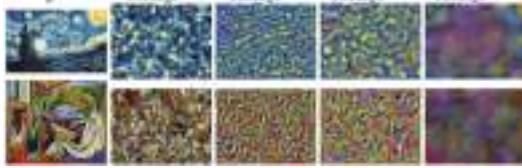


Figure from Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016. Copyright Springer, 2016. Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 63 May 17, 2018

Neural Style Transfer: Feature + Gram Reconstruction

Texture synthesis
(Gram reconstruction)

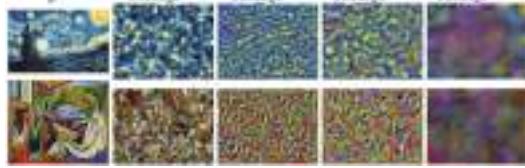


Figure from Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016. Copyright Springer, 2016. Reproduced for educational purposes.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 64 May 17, 2018

Neural Style Transfer

Content Image



Style Image



+

Starry Night by Van Gogh is in the public domain

This image is licensed under CC-BY 3.0

Neural Style Transfer

Content Image



+

Style Image



=

Style Transfer!



This image, copyright Justin Johnson, 2015. Reproduced with permission.

Gatys, Ecker, and Bethge, "Texture Synthesis Using Convolutional Neural Networks", NIPS 2015.

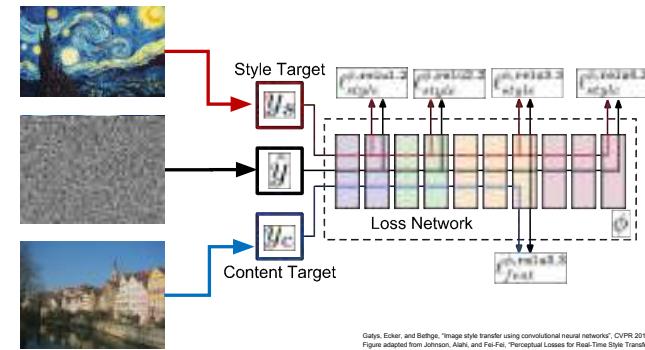
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 65 May 17, 2018

Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 66 May 17, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

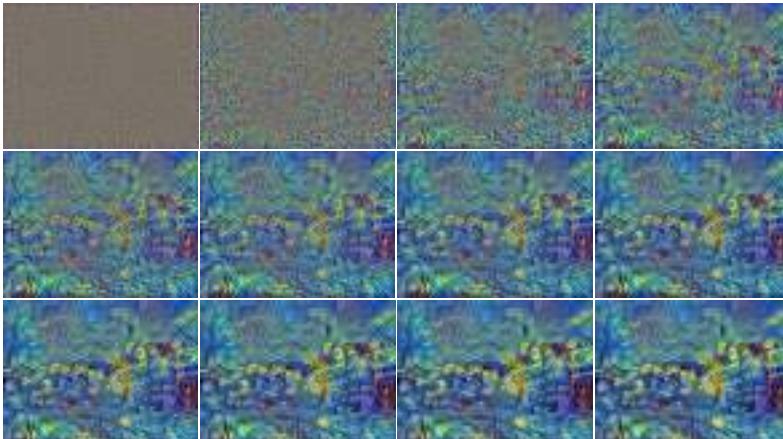
Lecture 13 - 67 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 68 May 17, 2018

Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016
Figure adapted from Johnson, Ashish, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016. Copyright Springer, 2016. Reproduced for educational purposes.

Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016
Figure adapted from Johnson, Ashish, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016. Copyright Springer, 2016. Reproduced for educational purposes.



Neural Style Transfer

Example outputs from [my implementation](#) (in Torch)



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 70 May 17, 2018

Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016
Figure copyright Justin Johnson, 2015.

Neural Style Transfer



More weight to content loss More weight to style loss

Neural Style Transfer

Resizing style image before running style transfer algorithm can transfer different types of features



Larger style image Smaller style image

Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016
Figure copyright Justin Johnson, 2015.

Fei-Fei Li & Justin Johnson & Serena Yeung

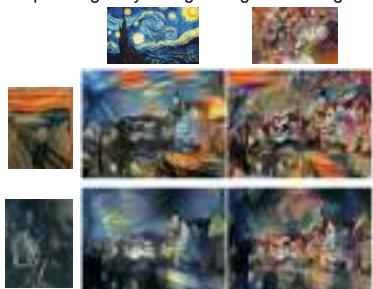
Lecture 13 - 71 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 72 May 17, 2018

Neural Style Transfer: Multiple Style Images

Mix style from multiple images by taking a weighted average of Gram matrices



Gatys, Ecker, and Bethge, "Image style transfer using convolutional neural networks", CVPR 2016.

Figure copyright Justin Johnson.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 73 May 17, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

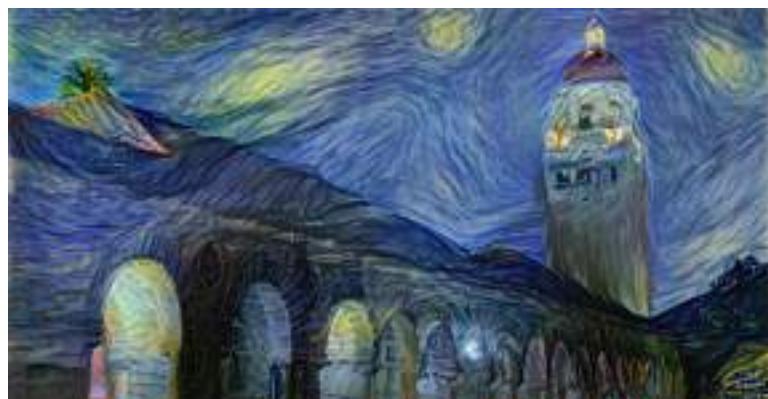
Lecture 13 - 74 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 73 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 74 May 17, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 75 May 17, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 76 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 75 May 17, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 76 May 17, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 77 May 17, 2018

Neural Style Transfer

Problem: Style transfer requires many forward / backward passes through VGG; very slow!

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 13 - 78 May 17, 2018

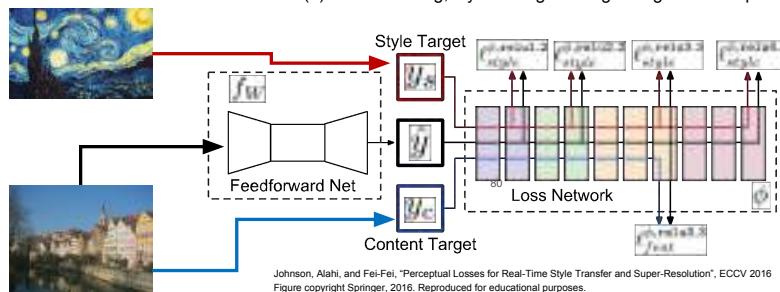
Neural Style Transfer

Problem: Style transfer requires many forward / backward passes through VGG; very slow!

Solution: Train another neural network to perform style transfer for us!

Fast Style Transfer

- (1) Train a feedforward network for each style
- (2) Use pretrained CNN to compute same losses as before
- (3) After training, stylize images using a single forward pass

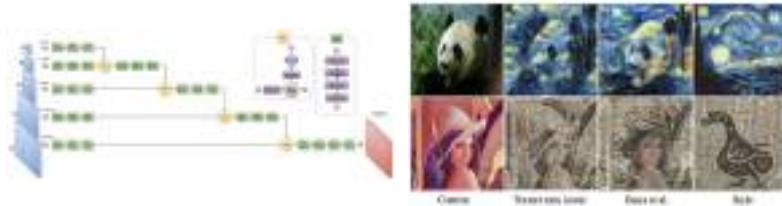


Fast Style Transfer



<https://github.com/jcjohnson/fast-neural-style>

Fast Style Transfer



Fast Style Transfer



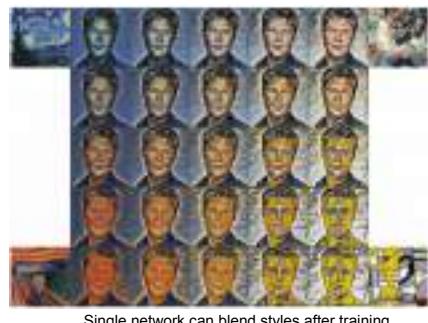
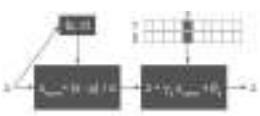
Replacing batch normalization with Instance Normalization improves results

One Network, Many Styles



One Network, Many Styles

Use the same network for multiple styles using *conditional instance normalization*: learn separate scale and shift parameters per style



Single network can blend styles after training

Dumoulin, Shlens, and Kudlur, "A Learned Representation for Artistic Style", ICLR 2017.

Figure copyright Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur, 2016; reproduced with permission.

Summary

Many methods for understanding CNN representations

Activations: Nearest neighbors, Dimensionality reduction, maximal patches, occlusion

Gradients: Saliency maps, class visualization, fooling images, feature inversion

Fun: DeepDream, Style Transfer.

Next time: **(Deep) Reinforcement Learning**

Lecture 14: Reinforcement Learning

Administrative

- All projects must be registered online (see Piazza for instructions), due Fri 11:59pm
- Midterm and A2 regrade requests also due Fri 11:59pm

So far... Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn a *function* to map $x \rightarrow y$



→ Cat

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.

Classification

So far... Unsupervised Learning

Data: x

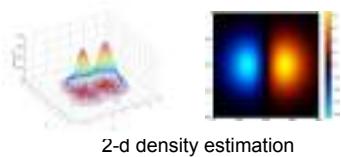
Just data, no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



1-d density estimation



2-d density images ([1](#) and [2](#)) are CC-BY public domain.

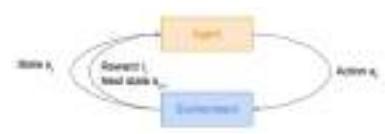
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 4 May 22, 2018

Today: Reinforcement Learning

Problems involving an **agent** interacting with an **environment**, which provides numeric **reward** signals

Goal: Learn how to take actions in order to maximize reward



Atari games figure copyright Volodymyr Mnih et al., 2013. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 5 May 22, 2018

Overview

Reinforcement Learning

- What is Reinforcement Learning?
- Markov Decision Processes
- Q-Learning
- Policy Gradients

Agent

Environment

Fei-Fei Li & Justin Johnson & Serena Yeung

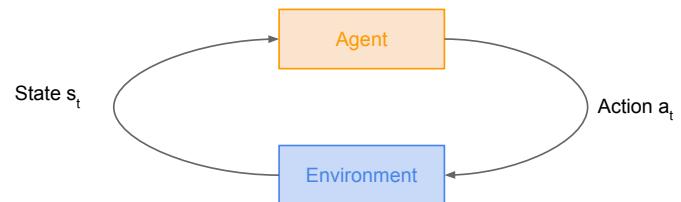
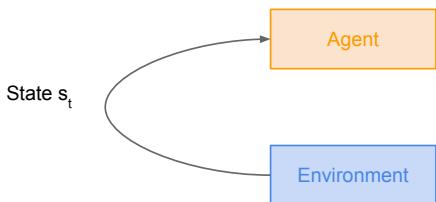
Lecture 14 - 6 May 22, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 7 May 22, 2018

Reinforcement Learning

Reinforcement Learning



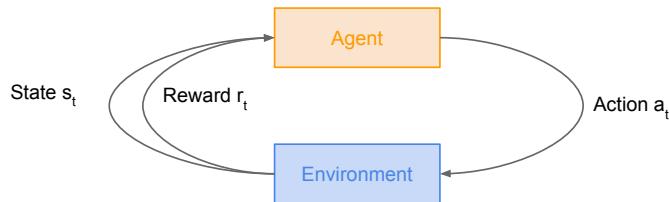
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 8 May 22, 2018

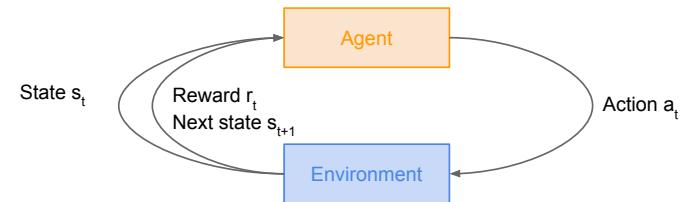
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 9 May 22, 2018

Reinforcement Learning



Reinforcement Learning



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 10

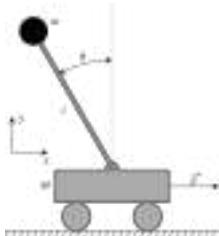
May 22, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 11

May 22, 2018

Cart-Pole Problem

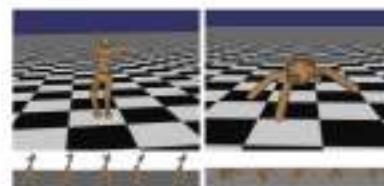


Objective: Balance a pole on top of a movable cart

State: angle, angular speed, position, horizontal velocity
Action: horizontal force applied on the cart

Reward: 1 at each time step if the pole is upright

Robot Locomotion



Objective: Make the robot move forward

State: Angle and position of the joints
Action: Torques applied on joints
Reward: 1 at each time step upright + forward movement

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 12 May 22, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 13 May 22, 2018

This image is CC0 public domain.

Figures copyright John Schulman et al., 2016. Reproduced with permission.

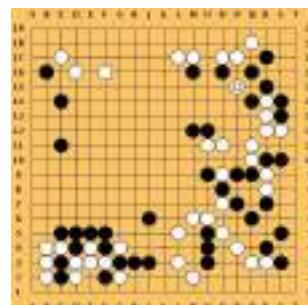
Atari Games



Objective: Complete the game with the highest score

State: Raw pixel inputs of the game state
Action: Game controls e.g. Left, Right, Up, Down
Reward: Score increase/decrease at each time step

Go



Objective: Win the game!

State: Position of all pieces
Action: Where to put the next piece down
Reward: 1 if win at the end of the game, 0 otherwise

Fei-Fei Li & Justin Johnson & Serena Yeung

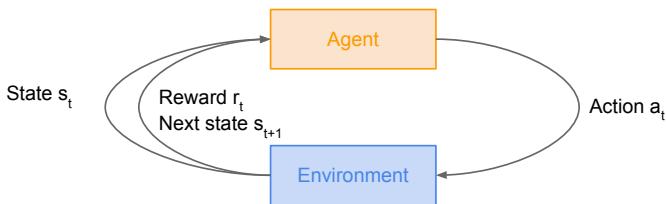
Lecture 14 - 14 May 22, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 15 May 22, 2018

This image is CC0 public domain.

How can we mathematically formalize the RL problem?



Markov Decision Process

- Mathematical formulation of the RL problem
- **Markov property:** Current state completely characterises the state of the world

Defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

\mathcal{S} : set of possible states
 \mathcal{A} : set of possible actions
 \mathcal{R} : distribution of reward given (state, action) pair
 \mathbb{P} : transition probability i.e. distribution over next state given (state, action) pair
 γ : discount factor

Markov Decision Process

- At time step $t=0$, environment samples initial state $s_0 \sim p(s_0)$

- Then, for $t=0$ until done:

- Agent selects action a_t
- Environment samples reward $r_t \sim R(\cdot | s_t, a_t)$
- Environment samples next state $s_{t+1} \sim P(\cdot | s_t, a_t)$
- Agent receives reward r_t and next state s_{t+1}

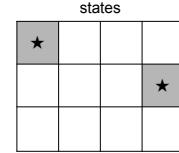
- A policy π is a function from S to A that specifies what action to take in each state

- **Objective:** find policy π^* that maximizes cumulative discounted reward:

$$\sum_{t \geq 0} \gamma^t r_t$$

A simple MDP: Grid World

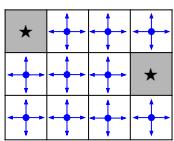
actions = {
 1. right →
 2. left ←
 3. up ↑
 4. down ↓
 }



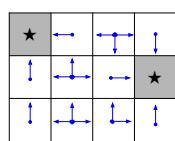
Set a negative "reward"
for each transition
(e.g. $r = -1$)

Objective: reach one of terminal states (greyed out) in least number of actions

A simple MDP: Grid World



Random Policy



Optimal Policy

The optimal policy π^*

We want to find optimal policy π^* that maximizes the sum of rewards.

How do we handle the randomness (initial state, transition probability...)?

The optimal policy π^*

We want to find optimal policy π^* that maximizes the sum of rewards.

How do we handle the randomness (initial state, transition probability...)?
Maximize the **expected sum of rewards!**

Formally: $\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \right] \text{ with } s_0 \sim p(s_0), a_0 \sim \pi(\cdot | s_0), s_{t+1} \sim p(\cdot | s_t, a_t)$

Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) $s_0, a_0, r_0, s_1, a_1, r_1, \dots$

Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) $s_0, a_0, r_0, s_1, a_1, r_1, \dots$

How good is a state?

The **value function** at state s , is the expected cumulative reward from following the policy from state s :

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, \pi \right]$$

Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) $s_0, a_0, r_0, s_1, a_1, r_1, \dots$

How good is a state?

The **value function** at state s , is the expected cumulative reward from following the policy from state s :

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, \pi \right]$$

How good is a state-action pair?

The **Q-value function** at state s and action a , is the expected cumulative reward from taking action a in state s and then following the policy:

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]$$

Bellman equation

The optimal Q-value function Q^* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]$$

Bellman equation

The optimal Q-value function Q^* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]$$

Q^* satisfies the following **Bellman equation**:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$$

Intuition: if the optimal state-action values for the next time-step $Q^*(s', a')$ are known, then the optimal strategy is to take the action that maximizes the expected value of $r + \gamma Q^*(s', a')$

Bellman equation

The optimal Q-value function Q^* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

Q^* satisfies the following **Bellman equation**:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q^*(s', a') | s, a]$$

Intuition: if the optimal state-action values for the next time-step $Q^*(s', a')$ are known, then the optimal strategy is to take the action that maximizes the expected value of $r + \gamma Q^*(s', a')$

The optimal policy π^* corresponds to taking the best action in any state as specified by Q^*

Solving for the optimal policy

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s, a) = \mathbb{E} [r + \gamma \max_{a'} Q_i(s', a') | s, a]$$

Q_i will converge to Q^* as $i \rightarrow \infty$

[What's the problem with this?](#)

Solving for the optimal policy

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s, a) = \mathbb{E} [r + \gamma \max_{a'} Q_i(s', a') | s, a]$$

Q_i will converge to Q^* as $i \rightarrow \infty$

[What's the problem with this?](#)

Not scalable. Must compute $Q(s, a)$ for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!

Solving for the optimal policy

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s, a) = \mathbb{E} [r + \gamma \max_{a'} Q_i(s', a') | s, a]$$

Q_i will converge to Q^* as $i \rightarrow \infty$

[What's the problem with this?](#)

Not scalable. Must compute $Q(s, a)$ for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!

Solution: use a function approximator to estimate $Q(s, a)$. E.g. a neural network!

Solving for the optimal policy: Q-learning

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s, a; \theta) \approx Q^*(s, a)$$

Solving for the optimal policy: Q-learning

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s, a; \theta) \approx Q^*(s, a)$$

If the function approximator is a deep neural network => **deep q-learning!**

Solving for the optimal policy: Q-learning

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s, a; \theta) \approx Q^*(s, a)$$

function parameters (weights)

If the function approximator is a deep neural network => **deep q-learning!**

Solving for the optimal policy: Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q^*(s', a') | s, a]$$

Solving for the optimal policy: Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q^*(s', a') | s, a]$$

Forward Pass

Loss function: $L_i(\theta_i) = \mathbb{E}_{s, a \sim g(\cdot)} [(y_i - Q(s, a; \theta_i))^2]$

where $y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$

Solving for the optimal policy: Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q^*(s', a') | s, a]$$

Forward Pass

Loss function: $L_i(\theta_i) = \mathbb{E}_{s, a \sim g(\cdot)} [(y_i - Q(s, a; \theta_i))^2]$

where $y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$

Backward Pass

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim g(\cdot), s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)] \nabla_{\theta_i} Q(s, a; \theta_i)$$

Solving for the optimal policy: Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q^*(s', a') | s, a]$$

Forward Pass

Loss function: $L_i(\theta_i) = \mathbb{E}_{s, a \sim g(\cdot)} [(y_i - Q(s, a; \theta_i))^2]$

where $y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$

Iteratively try to make the Q-value close to the target value (y) it should have, if Q-function corresponds to optimal Q^* (and optimal policy π^*)

Backward Pass

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim g(\cdot), s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)] \nabla_{\theta_i} Q(s, a; \theta_i)$$

Case Study: Playing Atari Games



Objective: Complete the game with the highest score

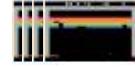
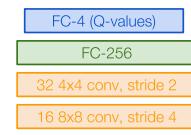
State: Raw pixel inputs of the game state

Action: Game controls e.g. Left, Right, Up, Down

Reward: Score increase/decrease at each time step

Q-network Architecture

$Q(s, a; \theta)$:
neural network
with weights θ

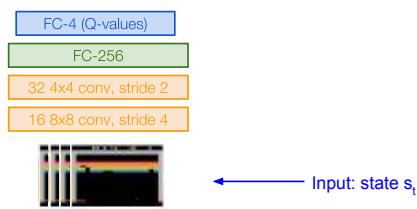


Current state s_t : 84x84x4 stack of last 4 frames
(after RGB->grayscale conversion, downsampling, and cropping)

Figures copyright Volodymyr Mnih et al., 2013. Reproduced with permission.

Q-network Architecture

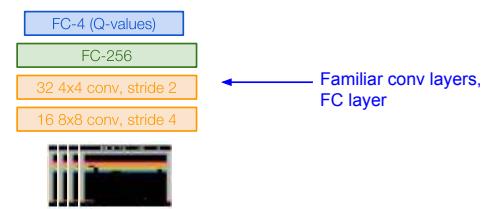
$Q(s, a; \theta)$:
neural network
with weights θ



Current state s_t : 84x84x4 stack of last 4 frames
(after RGB->grayscale conversion, downsampling, and cropping)

Q-network Architecture

$Q(s, a; \theta)$:
neural network
with weights θ



Current state s_t : 84x84x4 stack of last 4 frames
(after RGB->grayscale conversion, downsampling, and cropping)

[Mnih et al. NIPS Workshop 2013; Nature 2015]

[Mnih et al. NIPS Workshop 2013; Nature 2015]

Q-network Architecture

$Q(s, a; \theta)$:
neural network
with weights θ

Last FC layer has 4-d
output (if 4 actions),
corresponding to $Q(s_t, a_1), Q(s_t, a_2), Q(s_t, a_3), Q(s_t, a_4)$



Current state s_t : 84x84x4 stack of last 4 frames
(after RGB->grayscale conversion, downsampling, and cropping)

Q-network Architecture

$Q(s, a; \theta)$:
neural network
with weights θ

Last FC layer has 4-d
output (if 4 actions),
corresponding to $Q(s_t, a_1), Q(s_t, a_2), Q(s_t, a_3), Q(s_t, a_4)$

Number of actions between 4-18
depending on Atari game

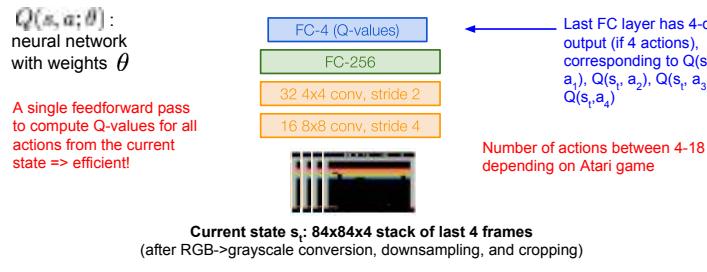


Current state s_t : 84x84x4 stack of last 4 frames
(after RGB->grayscale conversion, downsampling, and cropping)

[Mnih et al. NIPS Workshop 2013; Nature 2015]

[Mnih et al. NIPS Workshop 2013; Nature 2015]

Q-network Architecture



Training the Q-network: Loss function (from before)

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q^*(s', a') | s, a]$$

Forward Pass

$$\text{Loss function: } L_t(\theta_t) = \mathbb{E}_{s, a \sim \rho(\cdot)} [(y_t - Q(s, a; \theta_t))^2]$$

$$\text{where } y_t = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{t-1}) | s, a]$$

Iteratively try to make the Q-value
close to the target value (y_t) it
should have, if Q-function
corresponds to optimal Q^* (and
optimal policy π^*)

Backward Pass

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_t} L_t(\theta_t) = \mathbb{E}_{s, a \sim \rho(\cdot)} \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{t-1}) - Q(s, a; \theta_t)] \nabla_{\theta_t} Q(s, a; \theta_t)$$

Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:

- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:

- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand side) => can lead to bad feedback loops

Address these problems using **experience replay**

- Continually update a **replay memory** table of transitions (s_t, a_t, r_t, s_{t+1}) as game (experience) episodes are played
- Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples

Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:

- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 2 Deep Q-learning with Experience Replay

```

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights
for episode = 1, M do
    Initialize sequence  $s_0 = \{s_0\}$  and preprocessed sequence  $\phi_0 = \phi(s_0)$ 
    for t = 1, T do
        With probability  $\epsilon$  select  $a_t = \arg\max_a Q^*(\phi(s'_t), a; \theta)$ 
        otherwise select  $a_t = \arg\max_a Q^*(\phi(s'_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $s_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, r_t, s_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_i, a_i, r_i, \phi_{i+1})$  from  $\mathcal{D}$ 
        Sci  $\beta_i = \begin{cases} r_i & \text{for terminal } \phi_{i+1} \\ r_i + \gamma \max_{a'} Q(\phi_{i+1}, a'; \theta) & \text{for non-terminal } \phi_{i+1} \end{cases}$ 
        Perform a gradient descent step on  $|y_i - Q(\phi_i, a_i; \theta)|^2$  according to equation 3
    end for
end for
  
```

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 1 Deep Q-learning with Experience Replay

```

Initializes replay memory  $\mathcal{D}$  to capacity  $N$ 
Initializes action-value function  $Q$  with random weights
for episode = 1, M do
    Initialize sequence  $s_0 = [x_0]$  and preprocessed request  $\phi_0 = \phi(s_0)$ 
    for t = 1, T do
        With probability  $\epsilon$  select a random action  $a_t$ ;
        otherwise select  $a_t = \max_{a \in A} Q^*(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $s_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, r_t, s_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_i, a_i, r_i, \phi_{i+1})$  from  $\mathcal{D}$ 
        Set  $\hat{y}_i = \begin{cases} r_i & \text{for terminal } \theta_{i+1} \\ r_i + \gamma \max_a Q(\phi_{i+1}, a'; \theta) & \text{for non-terminal } \theta_{i+1} \end{cases}$ 
        Perform a gradient descent step on  $[\hat{y}_i - Q(\phi_i, a_i; \theta)]^2$  according to equation 3
    end for
end for

```

Initialize replay memory, Q-network

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 2 Deep Q-learning with Experience Replay

```

Initializes replay memory  $\mathcal{D}$  to capacity  $N$ 
Initializes action-value function  $Q$  with random weights
for episode = 1, M do
    Initialize sequence  $s_0 = [x_0]$  and preprocessed request  $\phi_0 = \phi(s_0)$ 
    for t = 1, T do
        With probability  $\epsilon$  select a random action  $a_t$ ;
        otherwise select  $a_t = \max_{a \in A} Q^*(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $s_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, r_t, s_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_i, a_i, r_i, \phi_{i+1})$  from  $\mathcal{D}$ 
        Set  $\hat{y}_i = \begin{cases} r_i & \text{for terminal } \theta_{i+1} \\ r_i + \gamma \max_a Q(\phi_{i+1}, a'; \theta) & \text{for non-terminal } \theta_{i+1} \end{cases}$ 
        Perform a gradient descent step on  $[\hat{y}_i - Q(\phi_i, a_i; \theta)]^2$  according to equation 3
    end for
end for

```

Play M episodes (full games)

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 1 Deep Q-learning with Experience Replay

```

Initializes replay memory  $\mathcal{D}$  to capacity  $N$ 
Initializes action-value function  $Q$  with random weights
for episode = 1, M do
    Initialize sequence  $s_0 = [x_0]$  and preprocessed request  $\phi_0 = \phi(s_0)$ 
    for t = 1, T do
        With probability  $\epsilon$  select a random action  $a_t$ ;
        otherwise select  $a_t = \max_{a \in A} Q^*(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $s_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, r_t, s_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_i, a_i, r_i, \phi_{i+1})$  from  $\mathcal{D}$ 
        Set  $\hat{y}_i = \begin{cases} r_i & \text{for terminal } \theta_{i+1} \\ r_i + \gamma \max_a Q(\phi_{i+1}, a'; \theta) & \text{for non-terminal } \theta_{i+1} \end{cases}$ 
        Perform a gradient descent step on  $[\hat{y}_i - Q(\phi_i, a_i; \theta)]^2$  according to equation 3
    end for
end for

```

Initialize state
(starting game
screen pixels) at the
beginning of each
episode

For each timestep t
of the game

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 2 Deep Q-learning with Experience Replay

```

Initializes replay memory  $\mathcal{D}$  to capacity  $N$ 
Initializes action-value function  $Q$  with random weights
for episode = 1, M do
    Initialize sequence  $s_0 = [x_0]$  and preprocessed request  $\phi_0 = \phi(s_0)$ 
    for t = 1, T do
        With probability  $\epsilon$  select a random action  $a_t$ ;
        otherwise select  $a_t = \max_{a \in A} Q^*(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $s_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, r_t, s_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_i, a_i, r_i, \phi_{i+1})$  from  $\mathcal{D}$ 
        Set  $\hat{y}_i = \begin{cases} r_i & \text{for terminal } \theta_{i+1} \\ r_i + \gamma \max_a Q(\phi_{i+1}, a'; \theta) & \text{for non-terminal } \theta_{i+1} \end{cases}$ 
        Perform a gradient descent step on  $[\hat{y}_i - Q(\phi_i, a_i; \theta)]^2$  according to equation 3
    end for
end for

```

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 1 Deep Q-learning with Experience Replay

```

Initializes replay memory  $\mathcal{D}$  to capacity  $N$ 
Initializes action-value function  $Q$  with random weights
for episode = 1, M do
    Initialize sequence  $s_0 = [x_0]$  and preprocessed request  $\phi_0 = \phi(s_0)$ 
    for t = 1, T do
        With small probability,   
With probability  $\epsilon$ , select a random action  $a_t$ ;   
otherwise select a random action (explore),   
otherwise select greedy action from current policy
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $s_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, r_t, s_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_i, a_i, r_i, \phi_{i+1})$  from  $\mathcal{D}$ 
        Set  $\hat{y}_i = \begin{cases} r_i & \text{for terminal } \theta_{i+1} \\ r_i + \gamma \max_a Q(\phi_{i+1}, a'; \theta) & \text{for non-terminal } \theta_{i+1} \end{cases}$ 
        Perform a gradient descent step on  $[\hat{y}_i - Q(\phi_i, a_i; \theta)]^2$  according to equation 3
    end for
end for

```

With small probability,
select a random
action (explore),
otherwise select
greedy action from
current policy

Take the action (a_t) ,
and observe the
reward r_t and next
state s_{t+1}

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 2 Deep Q-learning with Experience Replay

```

Initializes replay memory  $\mathcal{D}$  to capacity  $N$ 
Initializes action-value function  $Q$  with random weights
for episode = 1, M do
    Initialize sequence  $s_0 = [x_0]$  and preprocessed request  $\phi_0 = \phi(s_0)$ 
    for t = 1, T do
        With probability  $\epsilon$  select a random action  $a_t$ ;
        otherwise select  $a_t = \max_{a \in A} Q^*(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $s_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, r_t, s_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_i, a_i, r_i, \phi_{i+1})$  from  $\mathcal{D}$ 
        Set  $\hat{y}_i = \begin{cases} r_i & \text{for terminal } \theta_{i+1} \\ r_i + \gamma \max_a Q(\phi_{i+1}, a'; \theta) & \text{for non-terminal } \theta_{i+1} \end{cases}$ 
        Perform a gradient descent step on  $[\hat{y}_i - Q(\phi_i, a_i; \theta)]^2$  according to equation 3
    end for
end for

```

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 1 Deep Q-learning with Experience Replay

```

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights
for episode = 1.. $M$  do
    Initialize sequence  $s_1 = [x_1]$  and preprocessed sequence  $\phi_1 = \phi(x_1)$ 
    for t = 1.. $T$  do
        With probability  $c$  select a random action  $a_t$ 
        otherwise select  $a_t = \max_{a \in A} Q^{\pi}(s_t, a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, r_t$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(s_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ 
         $\text{Sci}[j] = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a' \in A} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ 
        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3
    end for

```

Store transition in replay memory

Putting it together: Deep Q-Learning with Experience Replay

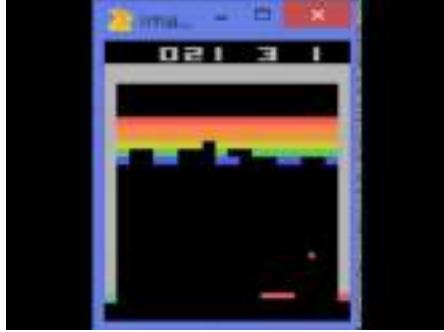
Algorithm 2 Deep Q-learning with Experience Replay

```

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights
for episode = 1.. $M$  do
    Initialize sequence  $s_1 = [x_1]$  and preprocessed sequence  $\phi_1 = \phi(x_1)$ 
    for t = 1.. $T$  do
        With probability  $c$  select a random action  $a_t$ 
        otherwise select  $a_t = \max_{a \in A} Q^{\pi}(s_t, a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, r_t$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(s_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ 
         $\text{Sci}[j] = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a' \in A} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ 
        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3
    end for

```

Experience Replay:
Sample a random minibatch of transitions from replay memory and perform a gradient descent step



<https://www.youtube.com/watch?v=V1eYnJ0Rok>

Video by Károly Zsolnai-Fehér. Reproduced with permission.

Policy Gradients

What is a problem with Q-learning?
The Q-function can be very complicated!

Example: a robot grasping an object has a very high-dimensional state => hard to learn exact value of every (state, action) pair

Policy Gradients

What is a problem with Q-learning?
The Q-function can be very complicated!

Example: a robot grasping an object has a very high-dimensional state => hard to learn exact value of every (state, action) pair

But the policy can be much simpler: just close your hand
Can we learn a policy directly, e.g. finding the best policy from a collection of policies?

Policy Gradients

Formally, let's define a class of parameterized policies: $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi_\theta \right]$$

Policy Gradients

Formally, let's define a class of parameterized policies: $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi_\theta \right]$$

We want to find the optimal policy $\hat{\theta}^* = \arg \max_{\theta} J(\theta)$

How can we do this?

Policy Gradients

Formally, let's define a class of parameterized policies: $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi_\theta \right]$$

We want to find the optimal policy $\hat{\theta}^* = \arg \max_{\theta} J(\theta)$

How can we do this?

[Gradient ascent on policy parameters!](#)

REINFORCE algorithm

Mathematically, we can write:

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)] \\ &= \int_{\tau} r(\tau)p(\tau; \theta)d\tau \end{aligned}$$

Where $r(\tau)$ is the reward of a trajectory $\tau = (s_0, a_0, r_0, s_1, \dots)$

REINFORCE algorithm

$$\begin{aligned} \text{Expected reward: } J(\theta) &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)] \\ &= \int_{\tau} r(\tau)p(\tau; \theta)d\tau \end{aligned}$$

REINFORCE algorithm

$$\begin{aligned} \text{Expected reward: } J(\theta) &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)] \\ &= \int_{\tau} r(\tau)p(\tau; \theta)d\tau \end{aligned}$$

Now let's differentiate this: $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta)d\tau$

REINFORCE algorithm

$$\begin{aligned} \text{Expected reward: } J(\theta) &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)] \\ &= \int_{\tau} r(\tau)p(\tau; \theta)d\tau \end{aligned}$$

Now let's differentiate this: $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta)d\tau$ Intractable! Gradient of an expectation is problematic when p depends on θ

REINFORCE algorithm

Expected reward:

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)] = \int_{\tau} r(\tau)p(\tau; \theta)d\tau$$

Now let's differentiate this: $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$ Intractable! Gradient of an expectation is problematic when p depends on θ

However, we can use a nice trick: $\nabla_{\theta} p(\tau; \theta) = p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta)$

REINFORCE algorithm

Expected reward:

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)] = \int_{\tau} r(\tau)p(\tau; \theta)d\tau$$

Now let's differentiate this: $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$ Intractable! Gradient of an expectation is problematic when p depends on θ

However, we can use a nice trick: $\nabla_{\theta} p(\tau; \theta) = p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta)$ If we inject this back:

$$\nabla_{\theta} J(\theta) = \int_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau) \nabla_{\theta} \log p(\tau; \theta)]$$

Can estimate with Monte Carlo sampling

REINFORCE algorithm

Can we compute those quantities without knowing the transition probabilities?

We have: $p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t)$

REINFORCE algorithm

Can we compute those quantities without knowing the transition probabilities?

We have: $p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t)$

Thus: $\log p(\tau; \theta) = \sum_{t \geq 0} \log p(s_{t+1}|s_t, a_t) + \log \pi_{\theta}(a_t|s_t)$

REINFORCE algorithm

Can we compute those quantities without knowing the transition probabilities?

We have: $p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t)$

Thus: $\log p(\tau; \theta) = \sum_{t \geq 0} \log p(s_{t+1}|s_t, a_t) + \log \pi_{\theta}(a_t|s_t)$

And when differentiating: $\nabla_{\theta} \log p(\tau; \theta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$ Doesn't depend on transition probabilities!

REINFORCE algorithm

Can we compute those quantities without knowing the transition probabilities?

We have: $p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1}|s_t, a_t) \pi_{\theta}(a_t|s_t)$

Thus: $\log p(\tau; \theta) = \sum_{t \geq 0} \log p(s_{t+1}|s_t, a_t) + \log \pi_{\theta}(a_t|s_t)$

And when differentiating: $\nabla_{\theta} \log p(\tau; \theta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$ Doesn't depend on transition probabilities!

Therefore when sampling a trajectory τ , we can estimate $J(\theta)$ with

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$$

Intuition

$$\text{Gradient estimator: } \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen
- If $r(\tau)$ is low, push down the probabilities of the actions seen

Intuition

$$\text{Gradient estimator: } \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen
- If $r(\tau)$ is low, push down the probabilities of the actions seen

Might seem simplistic to say that if a trajectory is good then all its actions were good. **But in expectation, it averages out!**

Intuition

$$\text{Gradient estimator: } \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen
- If $r(\tau)$ is low, push down the probabilities of the actions seen

Might seem simplistic to say that if a trajectory is good then all its actions were good. **But in expectation, it averages out!**

However, this also suffers from high variance because credit assignment is really hard. Can we help the estimator?

Variance reduction

$$\text{Gradient estimator: } \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Variance reduction

$$\text{Gradient estimator: } \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

First idea: Push up probabilities of an action seen, only by the cumulative future reward from that state

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Variance reduction

$$\text{Gradient estimator: } \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

First idea: Push up probabilities of an action seen, only by the cumulative future reward from that state

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Second idea: Use discount factor γ to ignore delayed effects

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Variance reduction: Baseline

Problem: The raw value of a trajectory isn't necessarily meaningful. For example, if rewards are all positive, you keep pushing up probabilities of actions.

What is important then? Whether a reward is better or worse than what you expect to get

Idea: Introduce a baseline function dependent on the state.

Concretely, estimator is now:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_{t'} - b(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

How to choose the baseline?

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_{t'} - b(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

A simple baseline: constant moving average of rewards experienced so far from all trajectories

How to choose the baseline?

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_{t'} - b(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

A simple baseline: constant moving average of rewards experienced so far from all trajectories

Variance reduction techniques seen so far are typically used in "Vanilla REINFORCE"

How to choose the baseline?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

How to choose the baseline?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

How to choose the baseline?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

Intuitively, we are happy with an action a_t in a state s_t if $Q^{\pi}(s_t, a_t) - V^{\pi}(s_t)$ is large. On the contrary, we are unhappy with an action if it's small.

How to choose the baseline?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

Intuitively, we are happy with an action a_t in a state s_t if $Q^\pi(s_t, a_t) - V^\pi(s_t)$ is large. On the contrary, we are unhappy with an action if it's small.

Using this, we get the estimator: $\nabla_\theta J(\theta) \approx \sum_{t \geq 0} (Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t)) \nabla_\theta \log \pi_\theta(a_t | s_t)$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 88 May 22, 2018

Actor-Critic Algorithm

Problem: we don't know Q and V. Can we learn them?

Yes, using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q-function).

- The actor decides which action to take, and the critic tells the actor how good its action was and how it should adjust
- Also alleviates the task of the critic as it only has to learn the values of (state, action) pairs generated by the policy
- Can also incorporate Q-learning tricks e.g. experience replay
- **Remark:** we can define by the **advantage function** how much an action was better than expected

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$

Lecture 14 - 89 May 22, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 88 May 22, 2018

Actor-Critic Algorithm

```
Initialize policy parameters  $\theta$ , critic parameters  $\phi$ 
For iteration=1, 2 ... do
    Sample m trajectories under the current policy
     $\Delta\theta \leftarrow 0$ 
    For i=1, ..., m do
        For t=1, ..., T do
             $A_t^i = \sum_{t' \geq t} \gamma^{t'-t} r_t^i - V_\theta(s_t^i)$ 
             $\Delta\theta \leftarrow \Delta\theta + A_t^i \nabla_\theta \log \pi_\theta(a_t^i | s_t^i)$ 
             $\Delta\phi \leftarrow \sum_t \sum_i \nabla_\phi [A_t^i]^2$ 
             $\theta \leftarrow \alpha \Delta\theta$ 
             $\phi \leftarrow \beta \Delta\phi$ 
    End for
End for
```

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 90 May 22, 2018

REINFORCE in action: Recurrent Attention Model (RAM)

Objective: Image Classification

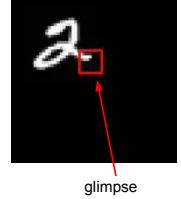
Take a sequence of "glimpses" selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image

State: Glimpses seen so far

Action: (x,y) coordinates (center of glimpse) of where to look next in image

Reward: 1 at the final timestep if image correctly classified, 0 otherwise



[Mnih et al. 2014]

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 91 May 22, 2018

REINFORCE in action: Recurrent Attention Model (RAM)

Objective: Image Classification

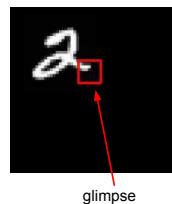
Take a sequence of "glimpses" selectively focusing on regions of the image, to predict class

- Inspiration from human perception and eye movements
- Saves computational resources => scalability
- Able to ignore clutter / irrelevant parts of image

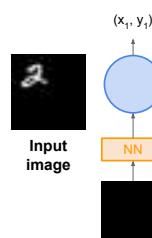
State: Glimpses seen so far

Action: (x,y) coordinates (center of glimpse) of where to look next in image

Reward: 1 at the final timestep if image correctly classified, 0 otherwise



REINFORCE in action: Recurrent Attention Model (RAM)



[Mnih et al. 2014]

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 92 May 22, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

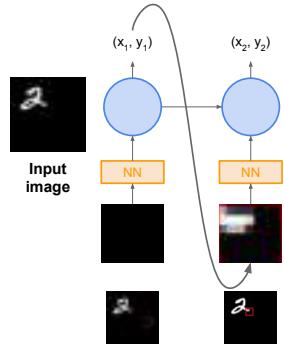
Lecture 14 - 93 May 22, 2018

Glimpsing is a non-differentiable operation => learn policy for how to take glimpse actions using REINFORCE

Given state of glimpses seen so far, use RNN to model the state and output next action

[Mnih et al. 2014]

REINFORCE in action: Recurrent Attention Model (RAM)

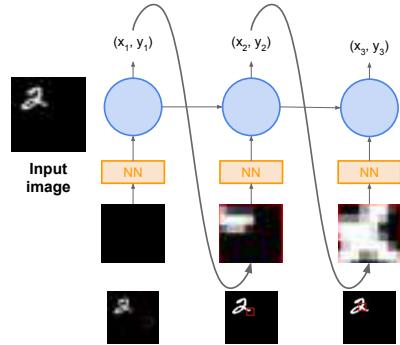


[Mnih et al. 2014]

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 94 May 22, 2018

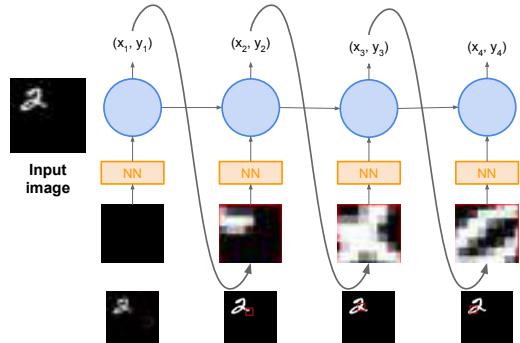
REINFORCE in action: Recurrent Attention Model (RAM)



[Mnih et al. 2014]

Lecture 14 - 95 May 22, 2018

REINFORCE in action: Recurrent Attention Model (RAM)

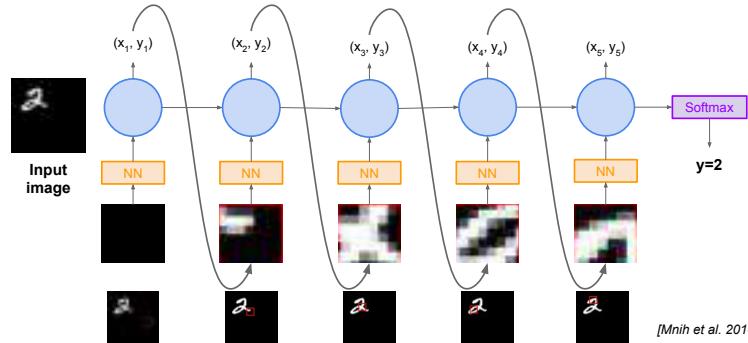


[Mnih et al. 2014]

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 96 May 22, 2018

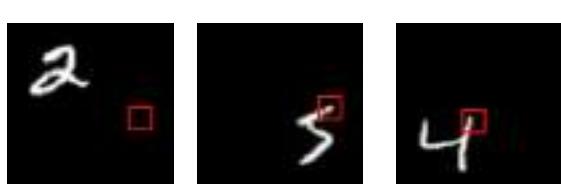
REINFORCE in action: Recurrent Attention Model (RAM)



[Mnih et al. 2014]

Lecture 14 - 97 May 22, 2018

REINFORCE in action: Recurrent Attention Model (RAM)



Has also been used in many other tasks including fine-grained image recognition, image captioning, and visual question-answering!

More policy gradients: AlphaGo

AlphaGo [Nature 2016]:

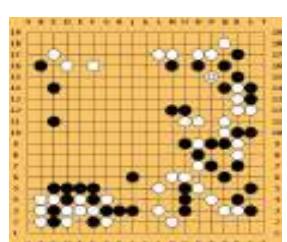
- Required many engineering tricks
- Bootstrapped from human play
- Beat 18-time world champion Lee Sedol

AlphaGo Zero [Nature 2017]:

- Simplified and elegant version of AlphaGo
- No longer bootstrapped from human play
- Beat (at the time) #1 world ranked Ke Jie

AlphaZero: Dec. 2017

- Generalized to beat world champion programs on chess and shogi as well



This image is CC-BY public domain

Figures copyright Daniel Levy, 2017. Reproduced with permission.

[Mnih et al. 2014]

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 98 May 22, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 99 May 22, 2018

Summary

- **Policy gradients:** very general but suffer from high variance so requires a lot of samples. **Challenge:** sample-efficiency
- **Q-learning:** does not always work but when it works, usually more sample-efficient. **Challenge:** exploration
- Guarantees:
 - **Policy Gradients:** Converges to a local minima of $J(\theta)$, often good enough!
 - **Q-learning:** Zero guarantees since you are approximating Bellman equation with a complicated function approximator

Next Time

Guest Lecture: Andrej Karpathy