



# Deep Learning & Neural Networks

Machine Learning – CSE4546  
Sham Kakade  
University of Washington  
December 1, 2016

©Sham Kakade

1

## Announcements:



- HW4 posted
- Poster Session Thurs, Dec 8
  - TAs (or your CSE friends) can help with printing
- Today:
  - Review: Deep Learning
  - Convolutional Neural Nets (+ RNNs?)
  - Start: RL
  - Also: MusicNet is out!

©Sham Kakade

2

# Poster Session

- Thursday Dec 8, 9-11:30am
  - Please arrive 20 mins early to set up
- Everyone is expected to attend
- Prepare a poster
  - We provide poster board and pins
  - Both one large poster (recommended) and several pinned pages are OK.
- Capture
  - Problem you are solving
  - Data you used
  - ML methodology
  - Results
- **Prepare a ~~1 minute speech~~ *2-min summary* about your project**
- ~~Two instructors~~ *people* will visit your poster separately
- Project Grading: scope, depth, data

©Sham Kakade

3

## Review

Machine Learning – CSE4546  
Sham Kakade  
University of Washington  
December 1, 2016

©Sham Kakade

4

# Hidden layer

- 1-hidden layer:

$$out(\mathbf{x}) = g \left( w_0 + \sum_k w_k g \left( w_0^k + \sum_i w_i^k x_i \right) \right)$$



©Sham Kakade

5

# Forward propagation for 1-hidden layer - Prediction

- 1-hidden layer:

$$out(\mathbf{x}) = g \left( w_0 + \sum_k w_k g \left( w_0^k + \sum_i w_i^k x_i \right) \right)$$

compute/cache  
output/activation at all  
layers.

©Sham Kakade

6

## Gradient descent for 1-hidden layer – Back-propagation: Computing $\frac{\partial \ell(W)}{\partial w_k}$

$$\ell(W) = \frac{1}{2} \sum_j [y^j - \text{out}(\mathbf{x}^j)]^2$$

Dropped  $w_0$  to make derivation simpler

$$\text{out}(\mathbf{x}) = g \left( \sum_{k'} w_{k'} g \left( \sum_{i'} w_{i'}^{k'} x_{i'} \right) \right)$$

$$\frac{\partial \ell(W)}{\partial w_k} = \sum_{j=1}^m -[y^j - \text{out}(\mathbf{x}^j)] \frac{\partial \text{out}(\mathbf{x}^j)}{\partial w_k}$$

*1st-layer*  
 $g'(\text{2nd layer activation}) \times \text{output}$

©Sham Kakade

7

## Gradient descent for 1-hidden layer – Back-propagation: Computing $\frac{\partial \ell(W)}{\partial w_i^k}$

$$\ell(W) = \frac{1}{2} \sum_j [y^j - \text{out}(\mathbf{x}^j)]^2$$

Dropped  $w_0$  to make derivation simpler

$$\text{out}(\mathbf{x}) = g \left( \sum_{k'} w_{k'} g \left( \sum_{i'} w_{i'}^{k'} x_{i'} \right) \right)$$

$$\frac{\partial \ell(W)}{\partial w_i^k} = \sum_{j=1}^m -[y^j - \text{out}(\mathbf{x}^j)] \frac{\partial \text{out}(\mathbf{x}^j)}{\partial w_i^k}$$

©Sham Kakade

8

# Theory, Optimization, and Debugging

Machine Learning – CSE4546

Sham Kakade

University of Washington

December 1, 2016

©Sham Kakade

9

## Architecture Selection

### ■ Feed-forward nets

- These are fully interconnected nets
- Try wider nets
- (empirical question) When do deeper nets help?
- (empirical question) Do feed-forward nets perform better than random features?

### ■ Structured Nets

- ConvNets are a great idea
  - Some idea of how to choose architecture
- Recurrent nets
  - Architecture chosen optimization

©Sham Kakade

10

# Optimization Issues

- Initialization
  - Want non-zero gradients
  - Init with a 'sensitivity analysis'
  - Want to start with a point not too far from some local opt
- Needs lots of Training data?
- Learning rates
  - Set by hand
  - Turn down when learning slows down
- Tensor Flow Defaults?

©Sham Kakade

11

# Regularization

- Needs lots of Training data?
  - Sometimes
  - (briefly) Share MusicNet case study
- Regularization (sometimes important?)
  - L2?
  - Dropout?

↳ very similar  
to weight decay /  
L2

©Sham Kakade

12

## “Theory”

$$w \leftarrow w - \eta \nabla L(w)$$

- $L(w)$  is our total loss on  $N$  data points
- Suppose  $L(w)$  is  $R$ -smooth (and  $L$  is bounded)
- Let's do batch gradient descent.
- What can we say?

find  $w$  s.t.  $\|\nabla L(w)\|^2 \leq \epsilon$   
in  $\mathcal{O}(\frac{1}{\epsilon})$  steps.

©Sham Kakade

13

## Convolutional Neural Networks & Application to Computer Vision

Machine Learning – CSE4546

Sham Kakade

University of Washington

December 1, 2016

©Sham Kakade

14

## Contains slides from...

- LeCun & Ranzato
- Russ Salakhutdinov
- Honglak Lee

©Sham Kakade

15

## Neural Networks in Computer Vision

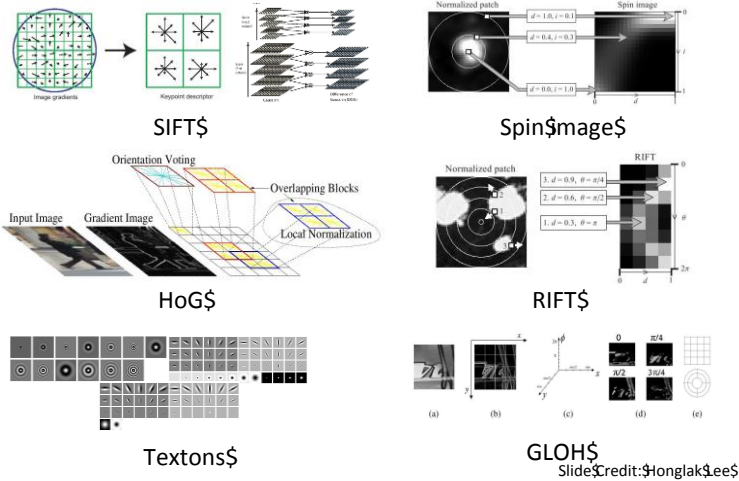
- Neural nets have made an amazing come back
  - Used to engineer high-level features of images
- Image features:

©Sham Kakade

16



# Some hand-created image features



©Sham Kakade

17

# Scanning an image with a detector

- Detector = Classifier from image patches:

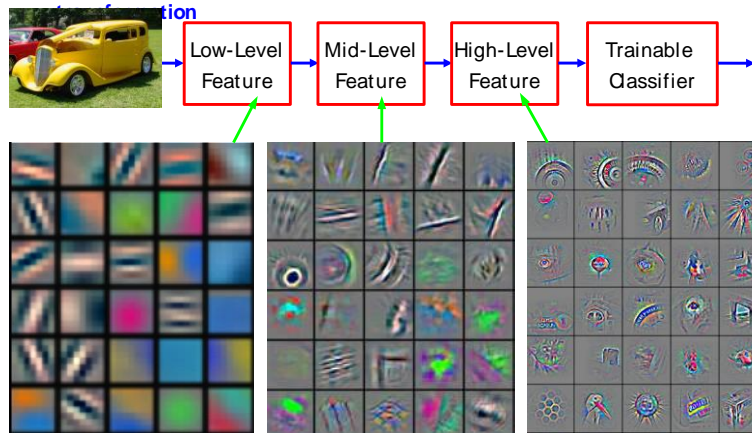
- Typically scan image with detector:



©Sham Kakade

18

# Using neural nets to learn non-linear features

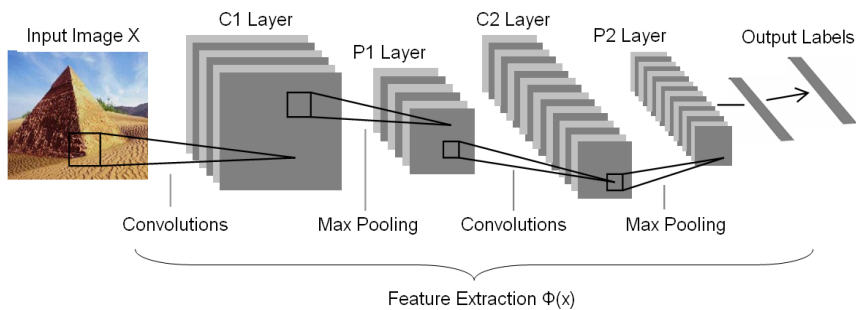


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

©Sham Kakade

19

But, many tricks needed to work well...



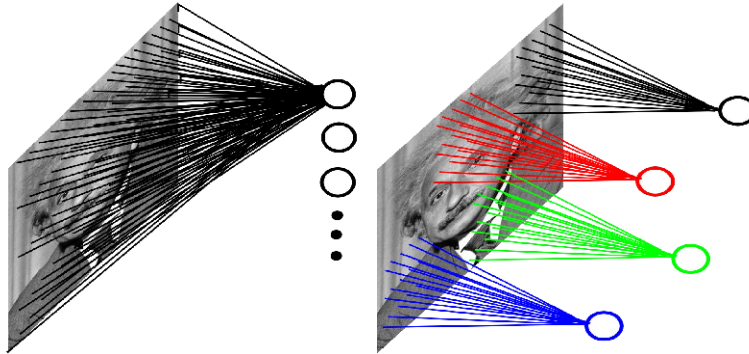
©Sham Kakade

20

# Convolution Layer

## Example: 200x200 image

- ▶ Fully-connected, 400,000 hidden units = 16 billion parameters
- ▶ Locally-connected, 400,000 hidden units 10x10 fields = 40 million params
- ▶ Local connections capture local dependencies



©Sham Kakade

21

# Parameter sharing

- Fundamental technique used throughout ML
- Neural net without parameter sharing:

- Sharing parameters:

(implementing convolutions).

↳ nodes have the same weights.

(

©Sham Kakade

22

# Pooling/Subsampling

- Convolutions act like detectors:

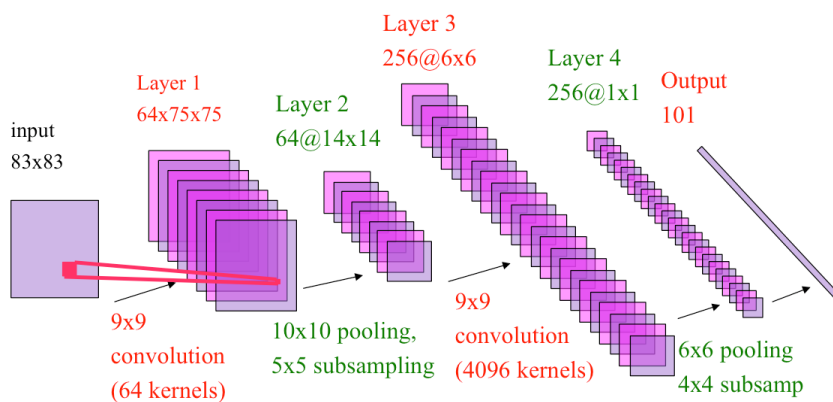


- But we don't expect true detections in every patch
- Pooling/subsampling nodes:

©Sham Kakade

23

## Example neural net architecture



©Sham Kakade

24

# Sample results

## Traffic Sign Recognition (GTSRB)

- ▶ German Traffic Sign Reco Bench
- ▶ 99.2% accuracy



## House Number Recognition (Google)

- ▶ Street View House Numbers
- ▶ 94.3 % accuracy



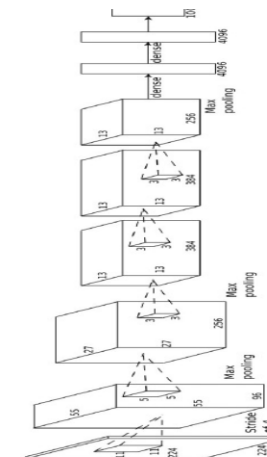
©Sham Kakade

25

# Example from Krizhevsky, Sutskever, Hinton 2012

## Won the 2012 ImageNet LSVRC. 60 Million parameters, 832M MAC ops

4M	<b>FULL CONNECT</b>	4Mflop
16M	<b>FULL 4096/ReLU</b>	16M
37M	<b>FULL 4096/ReLU</b>	37M
	<b>MAX POOLING</b>	
442K	<b>CONV 3x3/ReLU 256fm</b>	74M
1.3M	<b>CONV 3x3/ReLU 384fm</b>	224M
884K	<b>CONV 3x3/ReLU 384fm</b>	149M
	<b>MAX POOLING 2x2sub</b>	
	<b>LOCAL CONTRAST NORM</b>	
307K	<b>CONV 11x11/ReLU 256fm</b>	223M
	<b>MAX POOL 2x2sub</b>	
	<b>LOCAL CONTRAST NORM</b>	
35K	<b>CONV 11x11/ReLU 96fm</b>	105M

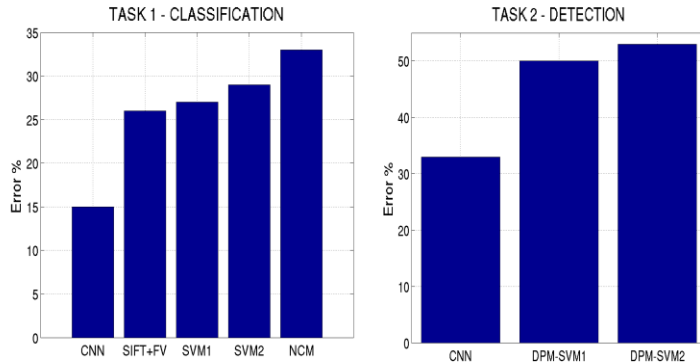


©Sham Kakade

26

# Results by Krizhevsky, Sutskever, Hinton 2012

- ImageNet Large Scale Visual Recognition Challenge
- 1000 categories, 1.5 Million labeled training samples



©Sham Kakade

27



©Sham Kakade

28

## TEST IMAGE



## RETRIEVED IMAGES



©Sham Kakade

29

## RNNs and LSTMs



©Sham Kakade

30