

DEEP LEARNING


David Kauchak  
CS158 – Fall 2019

Admin

Assignment 7

Assignment 8

Deep learning



WIKIPEDIA

Deep learning is a branch of machine learning based on a set of algorithms that attempt to model high level abstractions in data by using a deep graph with multiple processing layers, composed of multiple linear and non-linear transformations.

Deep learning is part of a broader family of machine learning methods based on learning representations of data.

Deep learning

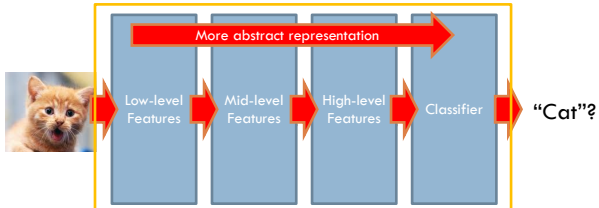
Key: learning better features that abstract from the “raw” data

Using **learned** feature representations based on large amounts of data, generally unsupervised

Using classifiers with multiple layers of learning

## Deep learning

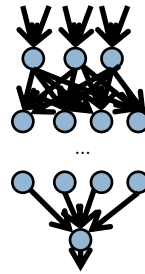
- Train *multiple layers* of features/abstractions from data.
- Try to discover *representation* that makes decisions easy.



Deep Learning: train layers of features so that classifier works well.

Slide adapted from: Adam Coates

## Deep learning for neural networks



Traditional NN models: 1-2 hidden layers

Deep learning NN models: 3+ hidden layers

## Importance of features

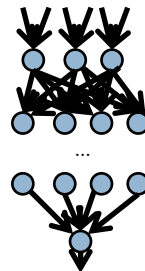
Once you have the right features, the algorithm you pick is relatively unimportant

Normal process = hand-crafted features

Deep learning: find algorithms to automatically discover features from the data

## Challenges

What makes “deep learning” hard for NNs?



$$w = w + input * \Delta_{current}$$

$$\Delta_{current} = f'(current\_input) \sum w_{output} \Delta_{output}$$

### Challenges

What makes "deep learning" hard for NNs?

$$w = w + input * \Delta_{current}$$

$$\Delta_{current} = f'(current\_input) \sum w_{output} \Delta_{output}$$

$$w = w + input * \Delta_{current}$$

$$\Delta_{current} = f'(current\_input) \sum w_{output} \Delta_{output}$$

...

...

What will happen to the modified errors further up?

### Challenges

What makes "deep learning" hard for NNs?

$$w = w + input * \Delta_{current}$$

$$\Delta_{current} = f'(current\_input) \sum w_{output} \Delta_{output}$$

$$w = w + input * \Delta_{current}$$

$$\Delta_{current} = f'(current\_input) \sum w_{output} \Delta_{output}$$

...

...

Modified errors tend to get diluted as they get combined with many layers of weight corrections

### Deep learning

Growing field

Driven by:

- ▣ Increase in data availability
- ▣ Increase in computational power
- ▣ Parallelizability of many of the algorithms

Involves more than just neural networks (though, they're a very popular model)

### word2vec

How many people have heard of it?

What is it?

## Word representations

Wine data uses word occurrences as a feature

What does this miss?

## Word representations

Wine data uses word occurrences as a feature

What does this miss?

"The wine had a dark red color" Zinfandel

"The wine was a deep crimson color" label?

"The wine was a deep yellow color" label?

Would like to recognize that words have similar meaning even though they aren't lexically the same

## Word representations

Key idea: project words into a multi-dimensional "meaning" space

word  $\rightarrow [x_1, x_2, \dots, x_d]$

## Word representations

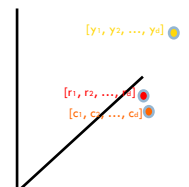
Key idea: project words into a multi-dimensional "meaning" space

word  $\rightarrow [x_1, x_2, \dots, x_d]$

red  $\rightarrow [r_1, r_2, \dots, r_d]$

crimson  $\rightarrow [c_1, c_2, \dots, c_d]$

yellow  $\rightarrow [y_1, y_2, \dots, y_d]$



## Word representations

Key idea: project words into a multi-dimensional “meaning” space

word  $\rightarrow [x_1, x_2, \dots, x_d]$

The idea of word representations is not new:

- Co-occurrence matrices
- Latent Semantic Analysis (LSA)

New idea: learn word representation using a task-driven approach

## A prediction problem

I like to eat bananas with cream cheese

Given a context of words

Predict what words are likely to occur in that context

## A prediction problem

Given text, can generate lots of **positive** examples:

I like to eat bananas with cream cheese

**input**

\_\_\_ like to eat

I \_\_\_ to eat bananas

I like \_\_\_ eat bananas with

I like to \_\_\_ bananas with cream

...

**prediction**

I

like

to

eat

...

## A prediction problem

Use data like this to learn a distribution:

$$p(\text{word} \mid \text{context})$$

$$p(w_i \mid w_{i-2} w_{i-1} w_{i+1} w_{i+2})$$

words before      words after

### A prediction problem

Any problems with using only positive examples?

$$p(w_i | w_{i-2} w_{i-1} w_{i+1} w_{i+2})$$

input

\_\_\_ like to eat  
I \_\_\_ to eat bananas  
I like \_\_\_ eat bananas with  
I like to \_\_\_ bananas with cream  
...

prediction

I  
like  
to  
eat  
...

### A prediction problem

Want to learn a distribution over **all** words

$$p(w_i | w_{i-2} w_{i-1} w_{i+1} w_{i+2})$$

input

\_\_\_ like to eat  
I \_\_\_ to eat bananas  
I like \_\_\_ eat bananas with  
I like to \_\_\_ bananas with cream  
...

prediction

I  
like  
to  
eat  
...

### A prediction problem

Negative examples?

I like to eat bananas with cream cheese

input

\_\_\_ like to eat  
I \_\_\_ to eat bananas  
I like \_\_\_ eat bananas with  
I like to \_\_\_ bananas with cream  
...

prediction

I  
like  
to  
eat  
...

### A prediction problem

Use random words to generate negative examples

I like to eat bananas with cream cheese

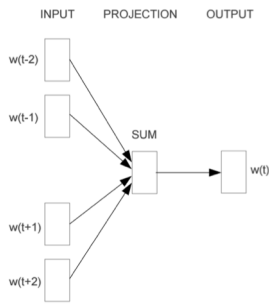
input

\_\_\_ like to eat  
I \_\_\_ to eat bananas  
I like \_\_\_ eat bananas with  
I like to \_\_\_ bananas with cream  
...

prediction (negative)

car  
snoopy  
run  
sloth  
...

### Train a neural network on this problem



<https://arxiv.org/pdf/1301.3781v3.pdf>

### Encoding words

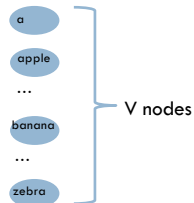
How can we input a "word" into a network?



### "One-hot" encoding

For a vocabulary of  $V$  words, have  $V$  input nodes

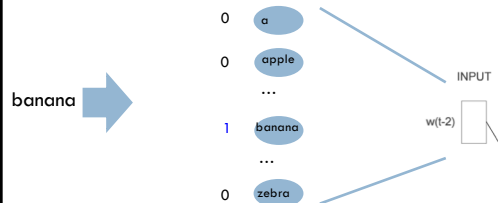
All inputs are 0 except the for the one corresponding to the word



### "One-hot" encoding

For a vocabulary of  $V$  words, have  $V$  input nodes

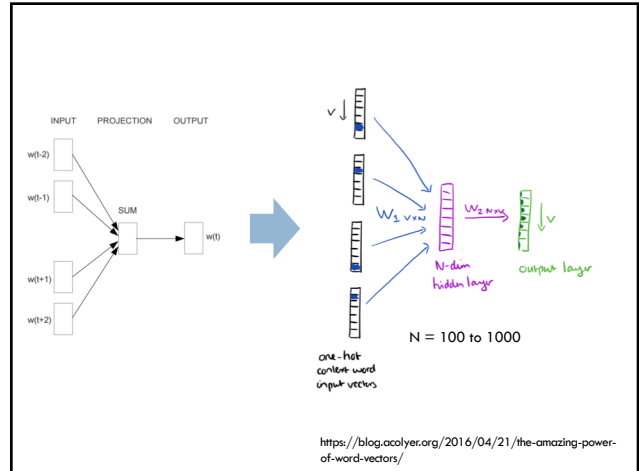
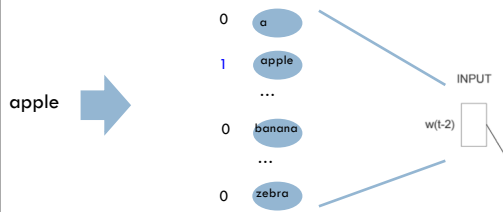
All inputs are 0 except the for the one corresponding to the word



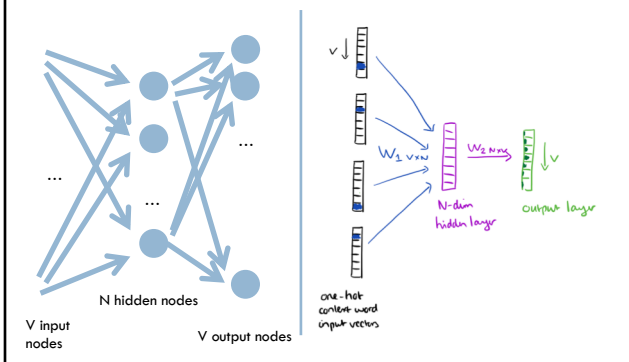
## “One-hot” encoding

For a vocabulary of  $V$  words, have  $V$  input nodes

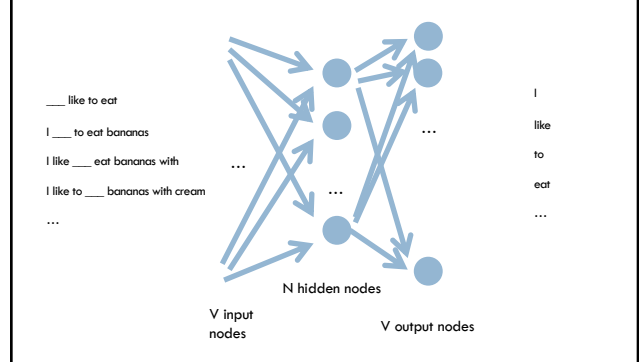
All inputs are 0 except the one corresponding to the word



## Another view

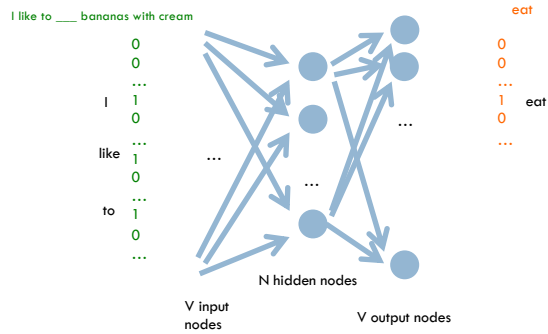


## Training: backpropagation

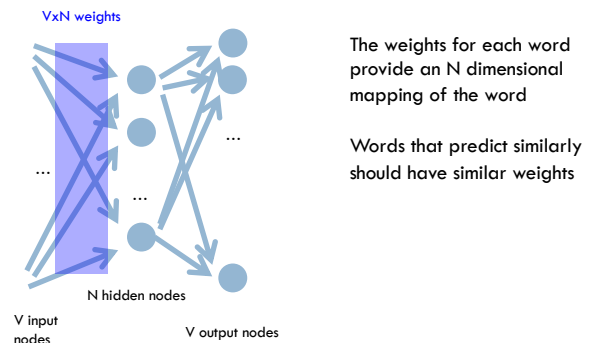




### Training: backpropagation



### Word representation



### Results

$$\text{vector}(\text{word1}) - \text{vector}(\text{word2}) = \text{vector}(\text{word3}) - X$$

word1 is to word2 as word3 is to X

Type of relationship	Word Pair 1		Word Pair 2	
Common capital city	Athens	Greece	Oslo	Norway
All capital cities	Astana	Kazakhstan	Harare	Zimbabwe
Currency	Angola	kwanza	Iran	rial
City-in-state	Chicago	Illinois	Stockton	California
Man-Woman	brother	sister	grandson	granddaughter

### Results

$$\text{vector}(\text{word1}) - \text{vector}(\text{word2}) = \text{vector}(\text{word3}) - X$$

word1 is to word2 as word3 is to X

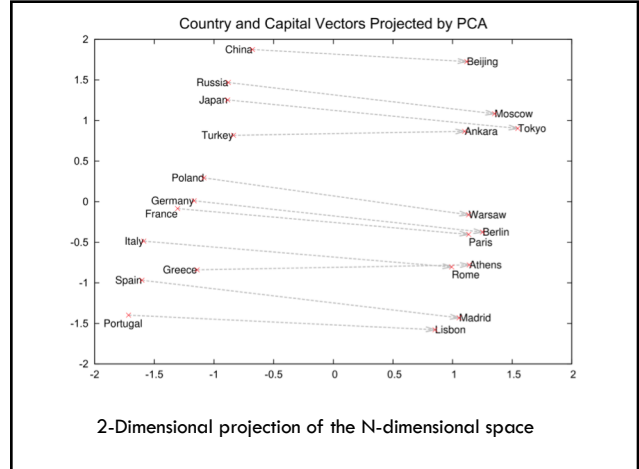
Type of relationship	Word Pair 1		Word Pair 2	
Adjective to adverb	apparent	apparently	rapid	rapidly
Opposite	possibly	impossibly	ethical	unethical
Comparative	great	greater	tough	tougher
Superlative	easy	easiest	lucky	luckiest
Present Participle	think	thinking	read	reading
Nationality adjective	Switzerland	Swiss	Cambodia	Cambodian
Past tense	walking	walked	swimming	swam
Plural nouns	mouse	mice	dollar	dollars
Plural verbs	work	works	speak	speaks

## Results

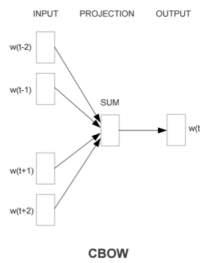
$$\text{vector}(\text{word1}) - \text{vector}(\text{word2}) = \text{vector}(\text{word3}) - X$$

word1 is to word2 as word3 is to X

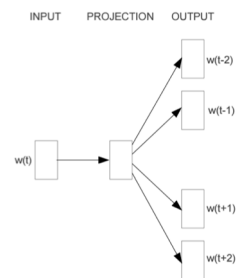
Newspapers			
New York San Jose	New York Times San Jose Mercury News	Baltimore Cincinnati	Baltimore Sun Cincinnati Enquirer
NHL Teams			
Boston Phoenix	Boston Bruins Phoenix Coyotes	Montreal Nashville	Montreal Canadiens Nashville Predators
NBA Teams			
Detroit Oakland	Detroit Pistons Golden State Warriors	Toronto Memphis	Toronto Raptors Memphis Grizzlies
Airlines			
Austria Belgium	Austrian Airlines Brussels Airlines	Spain Greece	Spainair Aegean Airlines
Company executives			
Steve Ballmer Samuel J. Palmisano	Microsoft IBM	Larry Page Werner Vogels	Google Amazon



## Continuous Bag Of Words



## Other models: skip-gram



## word2vec

A model for learning word representations from large amounts of data

Has become a popular pre-processing step for learning a more robust feature representation

Models like word2vec have also been incorporated into other learning approaches (e.g. translation tasks)

## word2vec resources

- <https://blog.acolyer.org/2016/04/21/the-amazing-power-of-word-vectors/>
- <https://code.google.com/archive/p/word2vec/>
- <https://deeplearning4j.org/word2vec>
- <https://arxiv.org/pdf/1301.3781v3.pdf>

## Image classification

Input: pixels of the image

Output: what's in the image

Ideally: have some features that identify "parts"



## Challenge: many different features



### Challenge: many different features

Round, elongated head with orange or black beak, can be turned backwards

Long white neck, can bend around, not necessarily straight

Small black circles, can be facing the camera, sometimes can see both

Black triangular shaped form, on the head, can have different sizes

White tail, generally far from the head, looks feathery

White, oval shaped body, with or without wings visible

Black feet, under body, can have different shapes

White elongated piece, can be squared or more triangular, can be obstructed sometimes

Luckily, the color is consistent...

<https://towardsdatascience.com/simple-introduction-to-convolutional-neural-networks-part-3-3077b0e>

### Image kernels

0	25	75	80	80
0	75	80	80	80
0	75	80	80	80
0	70	75	80	80
0	0	0	0	0

 $\times$ 

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

 $=$ 

0	0	75
0	0	80
0	0	80

 $\Sigma$ 

		75		
		80		
		80		

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

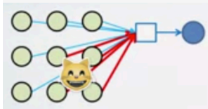
-1	0	1
-2	0	2
-1	0	1

Idea: learn kernels

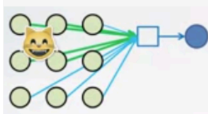
-1	-2	-1
0	0	0
1	2	1

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

## Traditional NN approach doesn't work



In this case, the **red weights** will be modified to better recognize cats



In this case, the **green weights** will be modified.

<https://towardsdatascience.com/simple-introduction-to-convolutional-neural-networks-eef843072bce>

## Traditional NN approach doesn't work

The information of image is the pixel

If we're dealing with a  $512 \times 512$  RGB image, we have  $512 \times 512 \times 3 = 786,432$  inputs

How many weights will we have with 5 hidden nodes?

For example, a  $512 \times 512$  RGB image has  $512 \times 512 \times 3 = 786,432$  and therefore **786,432 weights** in the next layer **per neuron**

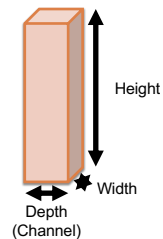
## Traditional NN approach doesn't work

The information of image is the pixel

If we're dealing with a  $512 \times 512$  RGB image, we have  $512 \times 512 \times 3 = 786,432$  inputs

**786,432 weights per neuron = ~4M weights!**

## Convolutional Neural Networks (CNNs)



We'll draw layers as blocks of nodes/inputs, e.g.,  $512 \times 512 \times 3$

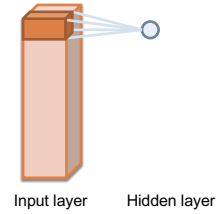
## Locally connected

image features are usually local

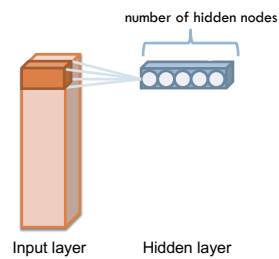
reduce the fully-connected network to locally-connected network.

For example, if we set window size 5, we only need  $5 \times 5 \times 3 = 75$

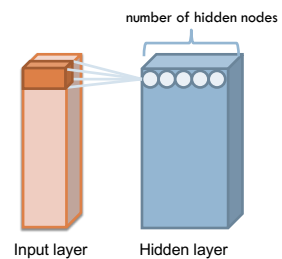
## fully connected -> locally connected



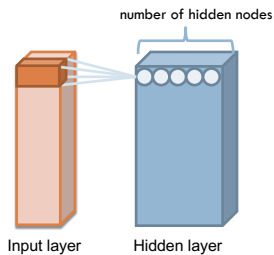
## hidden nodes



## apply across entire image



### apply across entire image



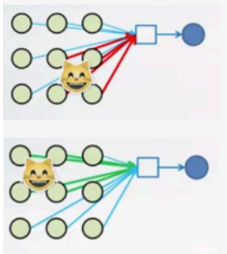
How many weights assuming: 512x512x3 images, 5 x 5 x 3 locally connected, and 5 hidden nodes?

### Too many weights!

Despite we locally-connected the layer, there are still too many weights

512x512x5 neurons in the next layer, we have 5x5x3 local connections = 98 million weights

### Share weights:



All weights to a given hidden node are the same for the locally-connected edges

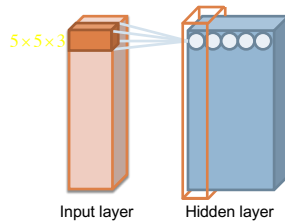
During classifying, we treat it like we have different edges, just with the same weight

During training, we update the weights as normal except we update *the same weights* for a given hidden node

Solves the positional issue!

### Share weights

We share parameter in the **same depth**.



## Parameter sharing

We share parameter in the **same depth**

Now we only have  $75 \times 5 = 375$  weights

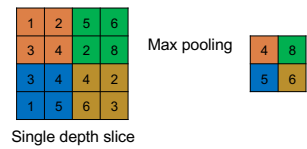
We call these layers "**convolution layers**".

What is learned can be considered as the convolution filters (like a kernel)

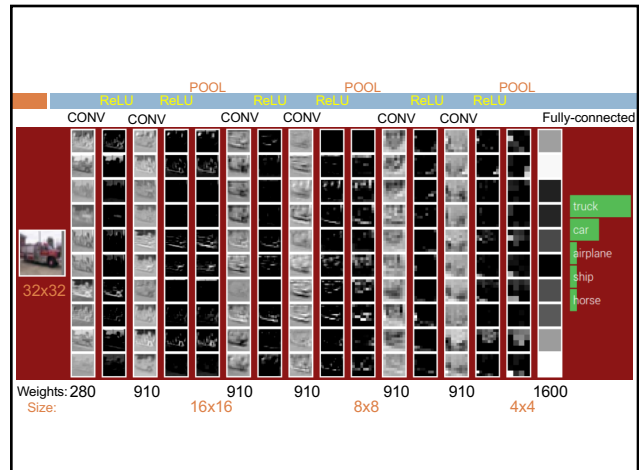
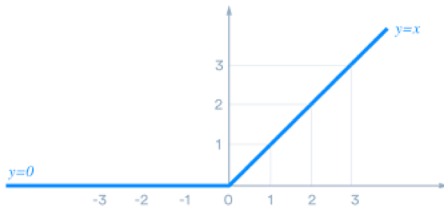
## Pool layers

Convolution layers are often followed by pool layers

Reduce the weights without losing too much information



## Rectified Linear Unit (ReLU)





## Another example

<http://scs.ryerson.ca/~aharley/vis/conv/>