

Deep Learning for NLP

(without Magic)



Université 
de Montréal

Richard Socher, Yoshua Bengio and Chris Manning

ACL 2012



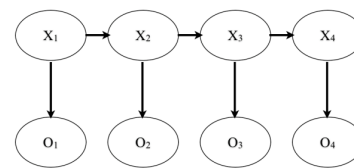
Deep Learning

Most current machine learning works well because of human-designed representations and input features

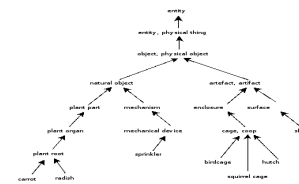
Machine learning becomes just optimizing weights to best make a final prediction

Representation learning attempts to automatically learn good features or representations

Deep learning algorithms attempt to learn multiple levels of representation of increasing complexity/abstraction



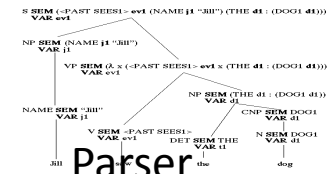
NER



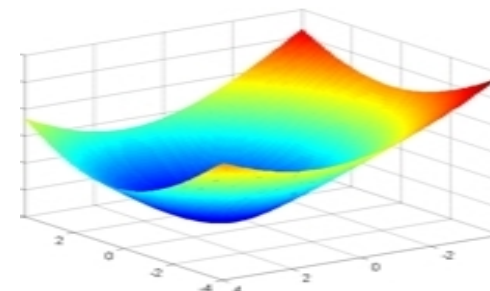
WordNet

```
<DOCID= wa94.008.0212 </DOCID>
<DOCNO= 940413-0062 </DOCNO>
<Q= Who's News?
@ 11:59 PM EDT </HL>
<ED= 11/19/94 </ED>
<EO= WALL STREET JOURNAL (J), PAGE B10 </SO>
<CO= MER </CO>
<IC= SECURITIES (SCR) </ID>
<TXT>
92
<BNSP FRY 11c1 <corcora> --- Donald Wright, 46 years old, was named executive vice president and director of fixed income at this brokerage firm. Mr. Wright resigned as president of Corcoran Group, a unit of Merrill Lynch & Co., to succeed Mark ...
Wright, 49, who left on Oct. 29 last month. A Corcoran Group spokesman said it hasn't named a successor to Mr. Wright, who is expected to begin his new job in the end of the month.
</P>
</TEXT>
</DOC>
```

SRL



Parser



A Deep Architecture

Mainly, work has explored [deep belief networks \(DBNs\)](#), Markov Random Fields with multiple layers, and various types of multiple-layer neural networks

Output layer

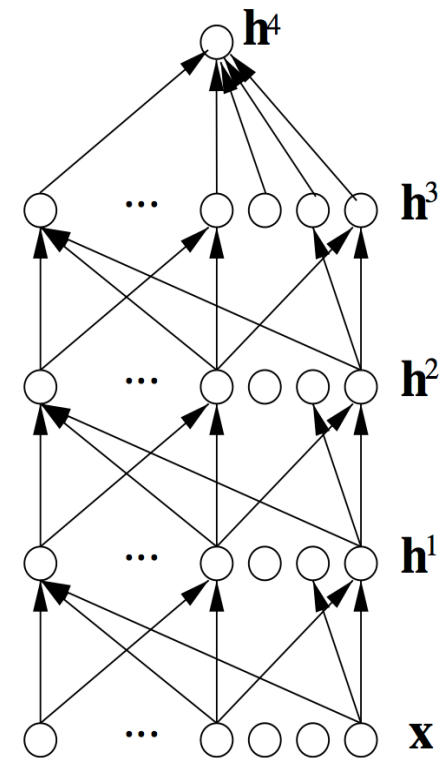
Here predicting a supervised target

Hidden layers

These learn more abstract representations as you head up

Input layer

3 Raw sensory inputs (roughly)



Part 1.1: The Basics

Five Reasons to Explore Deep Learning

#1 Learning representations

Handcrafting features is time-consuming

The features are often both over-specified and incomplete

The work has to be done again for each task/domain/...

We must move beyond handcrafted features and simple ML

Humans develop representations for learning and reasoning

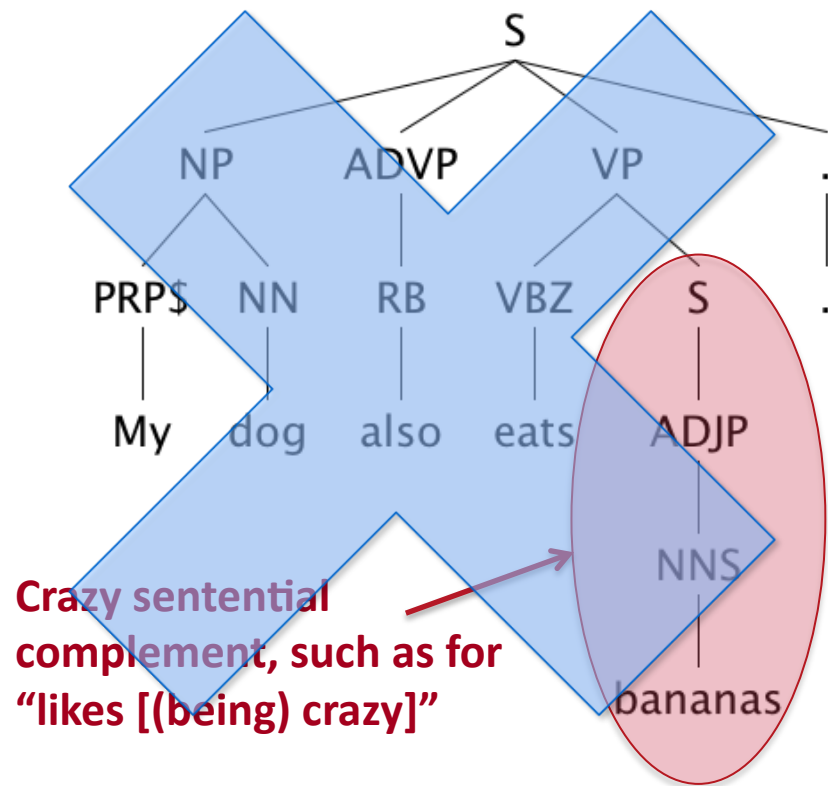
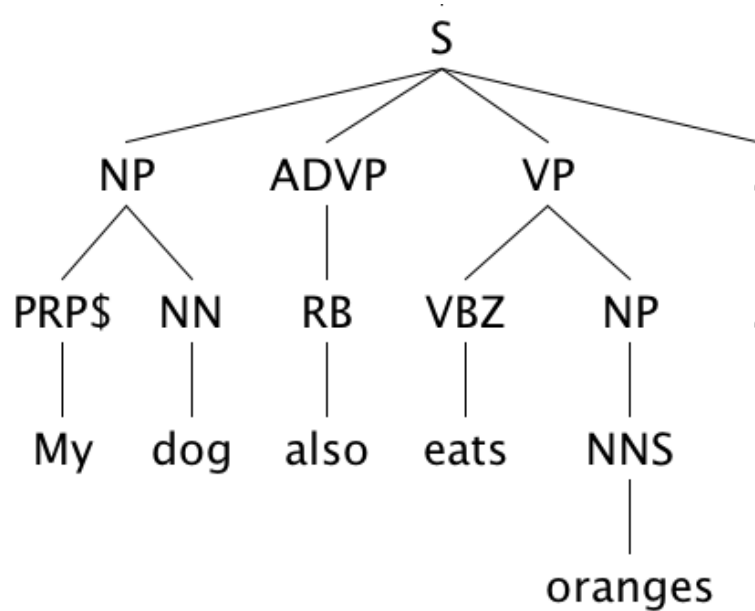
Our computers should do the same

Deep learning provides a way of doing this



#2 The need for distributed representations

Current NLP systems are incredibly fragile because of their atomic symbol representations



#2 The need for distributed representations

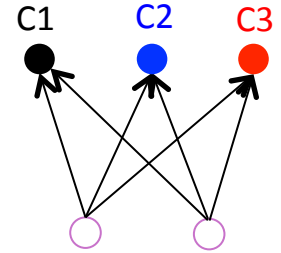
Learned word representations that model similarities help enormously in NLP

Distributional similarity based word clusters greatly help most applications

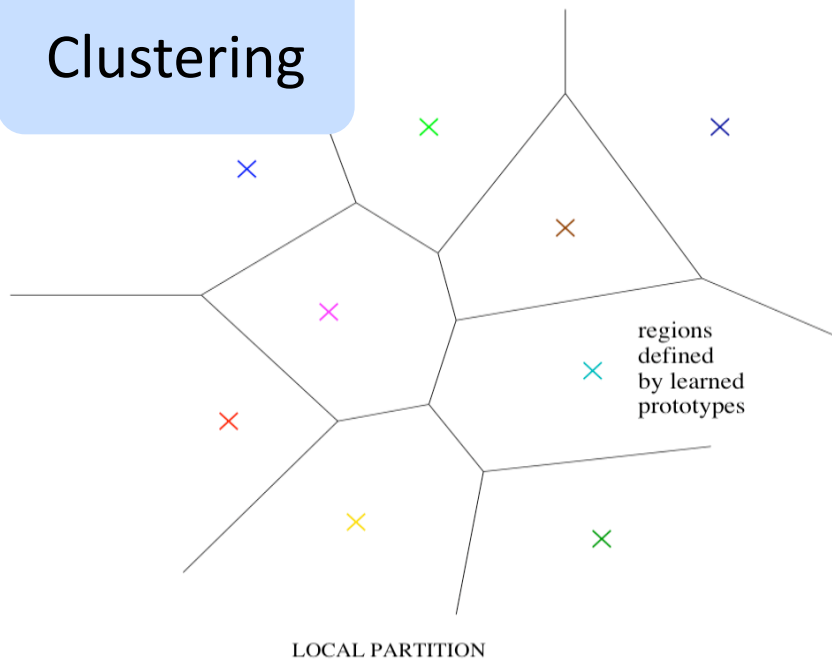
+1.4% F1 Dependency Parsing **15.2% error reduction** (Koo & Collins 2008, Brown clustering)

+3.4% F1 Named Entity Recognition **23.7% error reduction** (Stanford NER, exchange clustering)

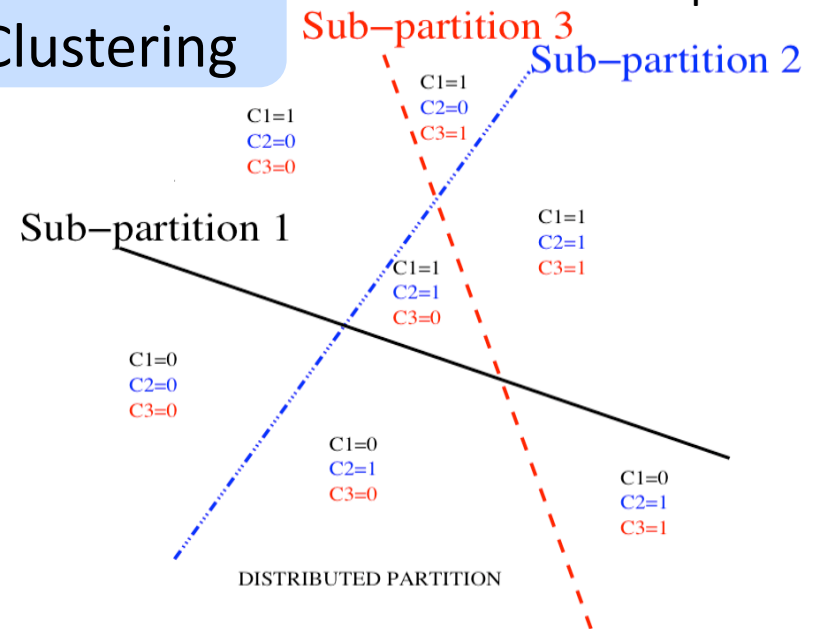
#2 The need for distributed representations



Clustering



Multi-Clustering



Learning features that are not mutually exclusive can be **exponentially more efficient** than nearest-neighbor-like or clustering-like models

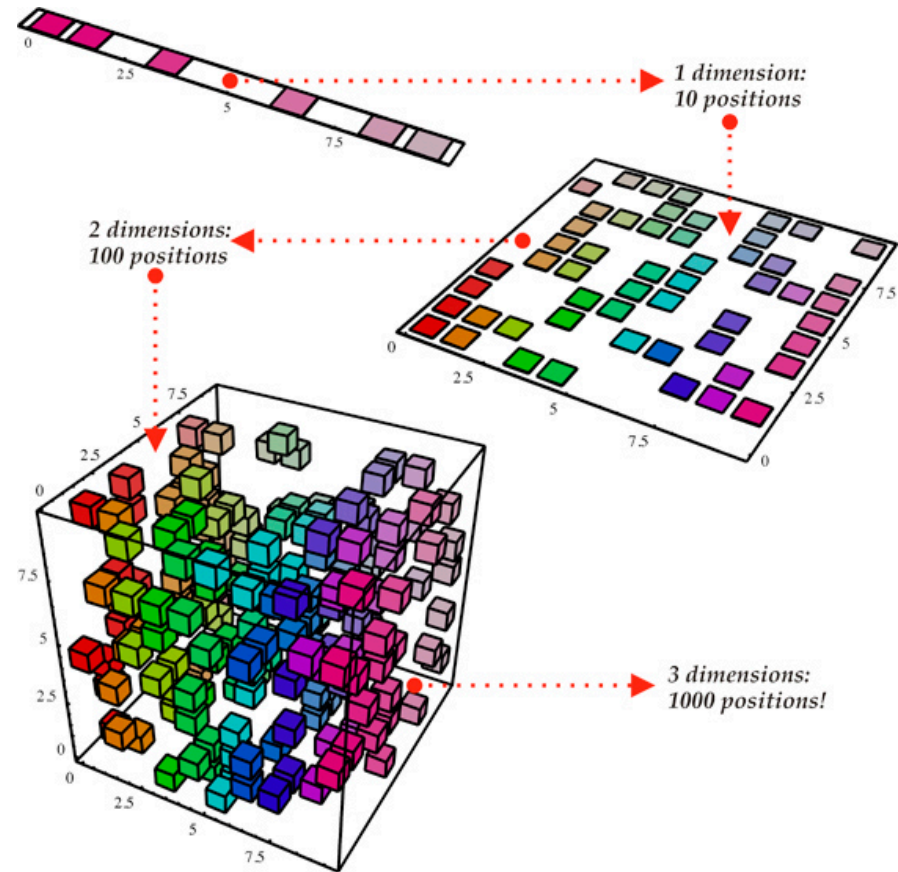
Distributed representations deal with the curse of dimensionality

Generalizing locally (e.g., nearest neighbors) requires representative examples for all relevant variations!

Classic solutions:

- Manual feature design
- Assuming a smooth target function (e.g., linear models)
- Kernel methods (linear in terms of kernel based on data points)

Neural networks parameterize and learn a “similarity” kernel



#3 Unsupervised feature and weight learning

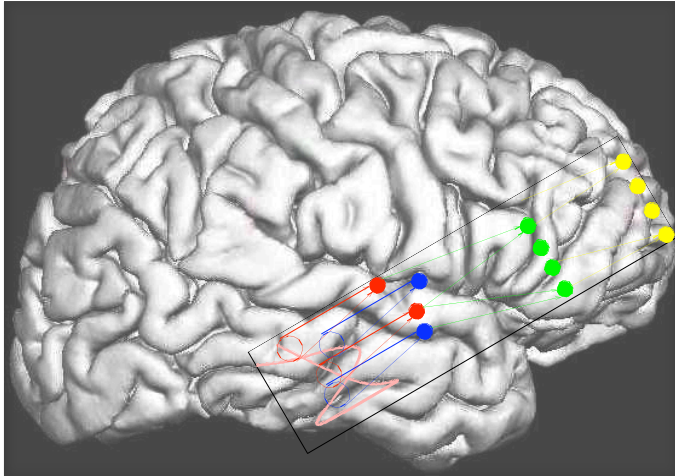
Today, most practical, good NLP& ML methods require labeled training data (i.e., supervised learning)

But almost all data is unlabeled

Most information must be acquired **unsupervised**

Fortunately, a good model of observed data can really help you learn classification decisions

#4 Learning multiple levels of representation



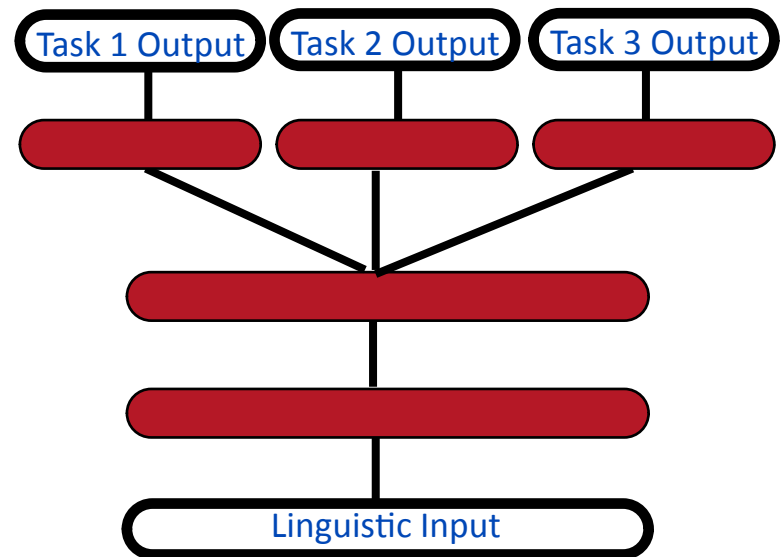
Biologically inspired learning

- The cortex seems to have a generic learning algorithm
- The brain has a deep architecture

We need good intermediate representations that can be shared across tasks

Multiple levels of latent variables allow combinatorial sharing of statistical strength

- Insufficient model depth can be exponentially inefficient

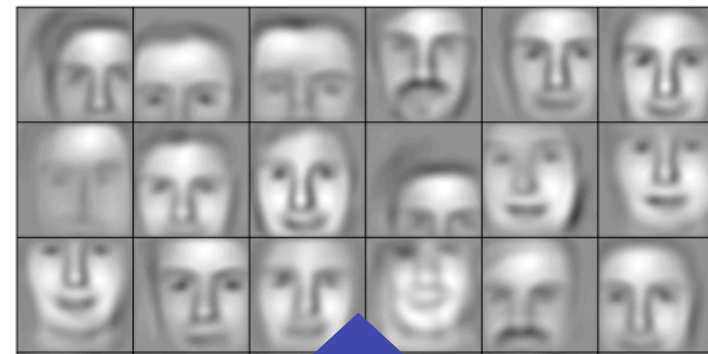


#4 Learning multiple levels of representation

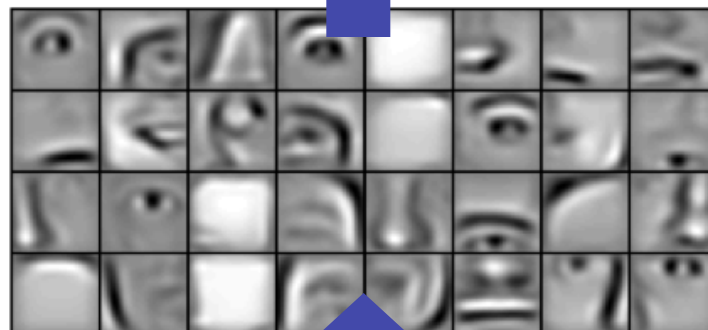


[Lee et al. ICML 2009; Lee et al. NIPS 2009]

Successive model layers learn deeper intermediate representations



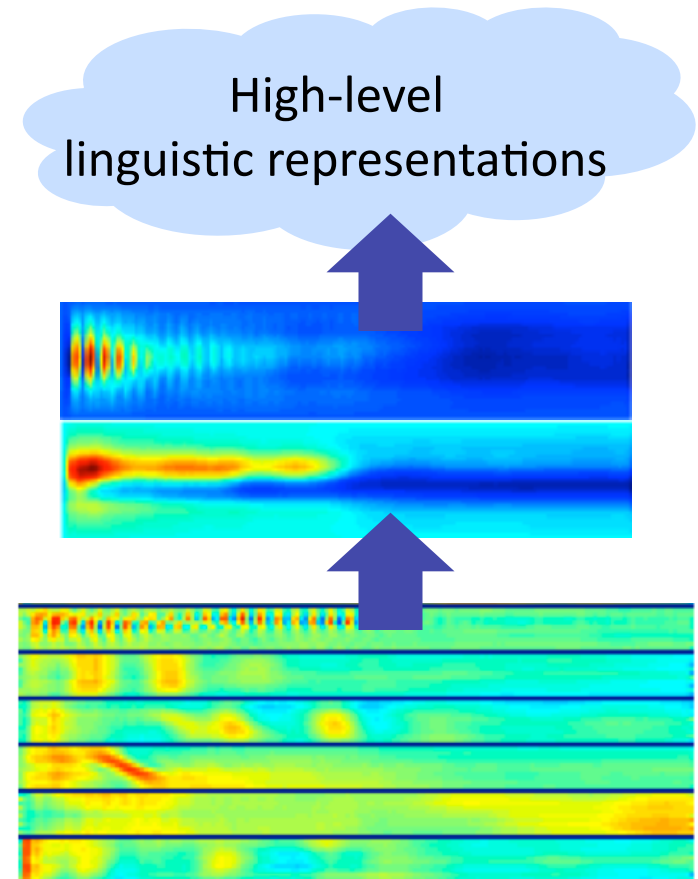
Layer 3



Layer 2



Layer 1

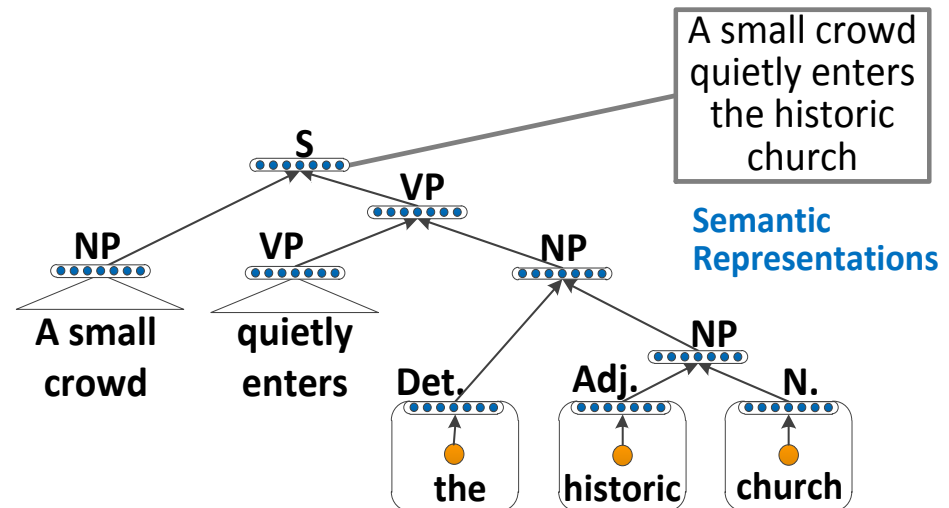
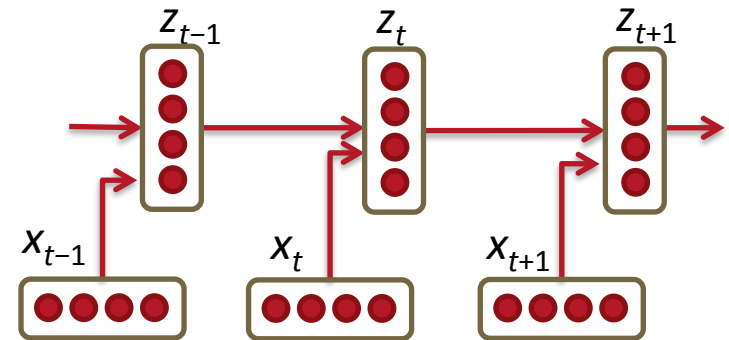


Handling the recursivity of human language

Human sentences are composed from words and phrases

We need **compositionality** in our ML models

Recursion: the same operator (same parameters) is applied repeatedly on different components



#5 Why now?

Despite prior investigation and understanding of many of the algorithmic techniques ...

Before 2006 training deep architectures was **unsuccessful** 😞

What has changed?

- New methods for unsupervised pre-training have been developed (Restricted Boltzmann Machines = RBMs, autoencoders, contrastive estimation, etc.)
- More efficient parameter estimation methods
- Better understanding of model regularization

Deep Learning models have already achieved impressive results for HLT

Neural Language Model [Mikolov et al. Interspeech 2011]



Model \ WSJ task	Eval WER
KN5 Baseline	17.2
Discriminative LM	16.9
Recurrent NN combination	14.4

MSR MAVIS Speech System [Dahl et al. 2012; Seide et al. 2011; following Mohamed et al. 2011]



“The algorithms represent the first time a company has released a deep-neural-networks (DNN)-based speech-recognition algorithm in a commercial product.”

Acoustic model & training	Recog \ WER	RT03S FSH	Hub5 SWB
GMM 40-mix, BMMI, SWB 309h	1-pass -adapt	27.4	23.6
CD-DNN 7 layer x 2048, SWB 309h	1-pass -adapt	18.5 (-33%)	16.1 (-32%)
GMM 72-mix, BMMI, FSH 2000h	k-pass +adapt	18.6	17.1

Deep Learn Models Have Interesting Performance Characteristics

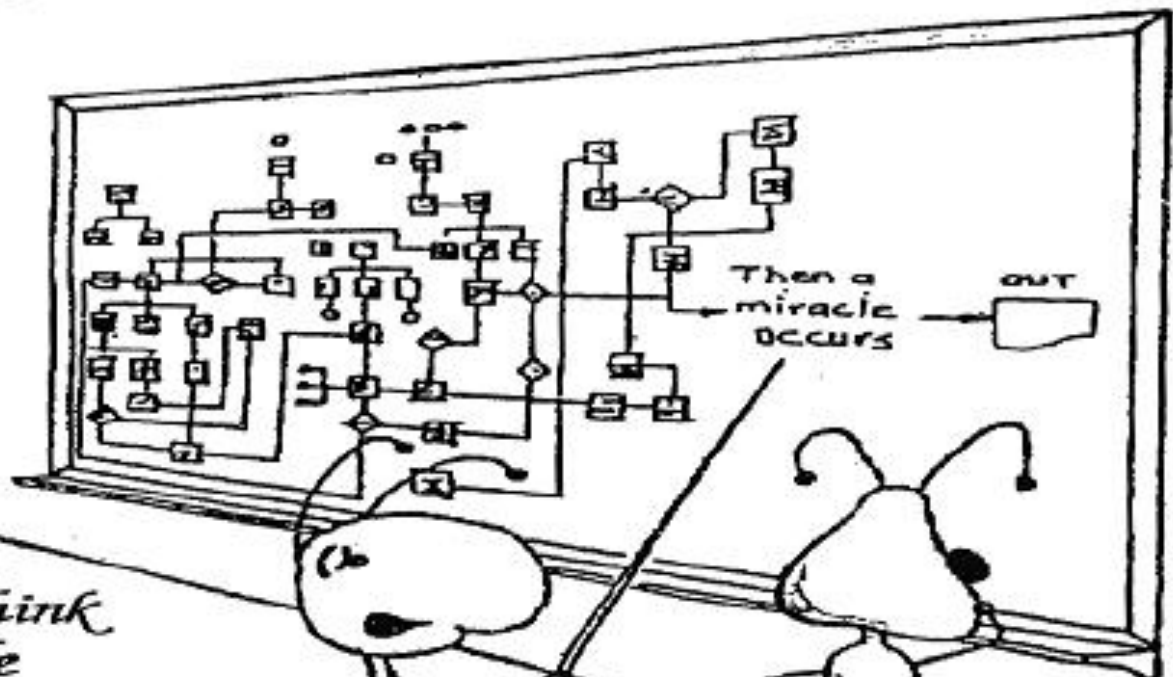
Deep learning models can now be very fast *in some circumstances*

- SENNA [Collobert et al. 2011] can do POS or NER faster than other SOTA taggers (16x to 122x), using 25x less memory
- WSJ POS 97.29% acc; CoNLL NER 89.59% F1; CoNLL Chunking 94.32% F1

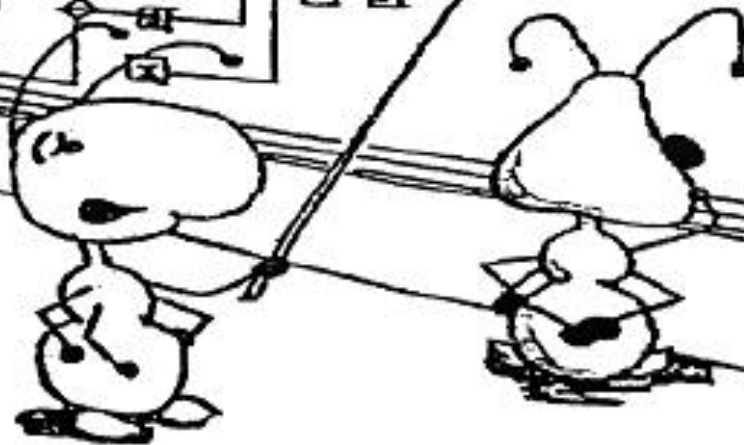
Changes in computing technology favor deep learning

- In NLP, speed has traditionally come from exploiting sparsity
- But with modern machines, branches and widely spaced memory accesses are costly
- Uniform parallel operations on dense vectors are faster

These trends are even stronger with multi-core CPUs and GPUs



Good work -- but I think we might need a little more detail right here.



Outline of the Tutorial

1. The Basics
 1. Motivations
 2. From logistic regression to neural networks
 3. Word representations
 4. Unsupervised word vector learning
 5. Backpropagation Training
 6. Learning word-level classifiers: POS and NER
 7. Sharing statistical strength
2. Recursive Neural Networks
3. Applications, Discussion, and Resources

Outline of the Tutorial

1. The Basics
2. Recursive Neural Networks
 1. Motivation
 2. Recursive Neural Networks for Parsing
 3. Theory: Backpropagation Through Structure
 4. Recursive Autoencoders
 5. Application to Sentiment Analysis and Paraphrase Detection
 6. Compositionality Through Recursive Matrix-Vector Spaces
 7. Relation classification
3. Applications, Discussion, and Resources

Outline of the Tutorial

1. The Basics
2. Recursive Neural Networks
3. Applications, Discussion, and Resources
 1. Applications
 1. Neural language models
 2. Structured embedding of knowledge bases
 3. Assorted other speech and NLP applications
 2. Resources (readings, code, ...)
 3. Tricks of the trade
 4. Discussion: Limitations, advantages, future directions

Part 1.2: The Basics

From logistic regression to neural nets

Demystifying neural networks

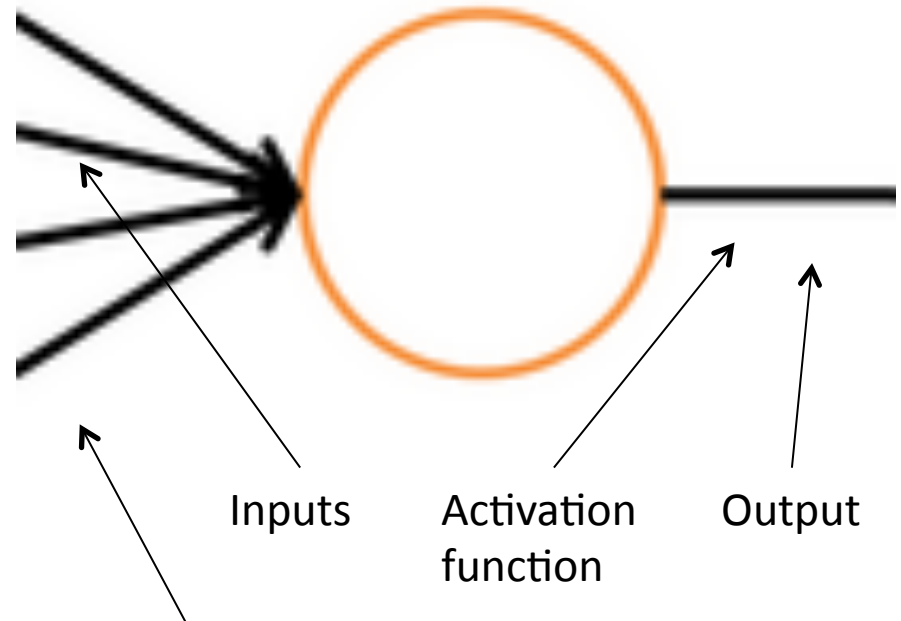
Neural networks come with their own terminological baggage

... just like SVMs

But if you understand how maxent/logistic regression models work

Then **you already understand** the operation of a basic neural network neuron!

A single neuron
A computational unit with n (3) inputs and 1 output and parameters W, b



Bias unit corresponds to intercept term

From Maxent Classifiers to Neural Networks

In NLP, a maxent classifier is normally written as:

$$P(c | d, \lambda) = \frac{\exp \sum_i \lambda_i f_i(c, d)}{\sum_{c' \in C} \exp \sum_i \lambda_i f_i(c', d)}$$

Supervised learning gives us a distribution for datum d over classes in C

Vector form:
$$P(c | d, \lambda) = \frac{e^{\lambda^\top f(c, d)}}{\sum_{c'} e^{\lambda^\top f(c', d)}}$$

Such a classifier is used as-is in a neural network (“a softmax layer”)

- Often as the top layer

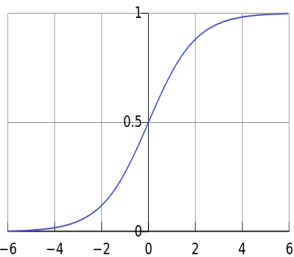
But for now we’ll derive a two-class logistic model for one neuron

From Maxent Classifiers to Neural Networks

Vector form:
$$P(c | d, \lambda) = \frac{e^{\lambda^T f(c, d)}}{\sum_{c'} e^{\lambda^T f(c', d)}}$$

Make two class:

$$\begin{aligned} P(c_1 | d, \lambda) &= \frac{e^{\lambda^T f(c_1, d)}}{e^{\lambda^T f(c_1, d)} + e^{\lambda^T f(c_2, d)}} = \frac{e^{\lambda^T f(c_1, d)}}{e^{\lambda^T f(c_1, d)} + e^{\lambda^T f(c_2, d)}} \cdot \frac{e^{-\lambda^T f(c_1, d)}}{e^{-\lambda^T f(c_1, d)}} \\ &= \frac{1}{1 + e^{\lambda^T [f(c_2, d) - f(c_1, d)]}} = \frac{1}{1 + e^{-\lambda^T x}} \quad \text{for } x = f(c_1, d) - f(c_2, d) \\ &= f(\lambda^T x) \end{aligned}$$



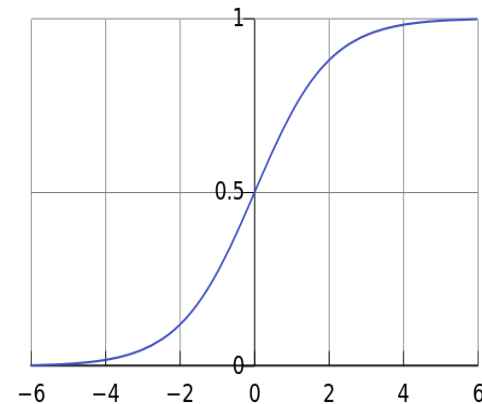
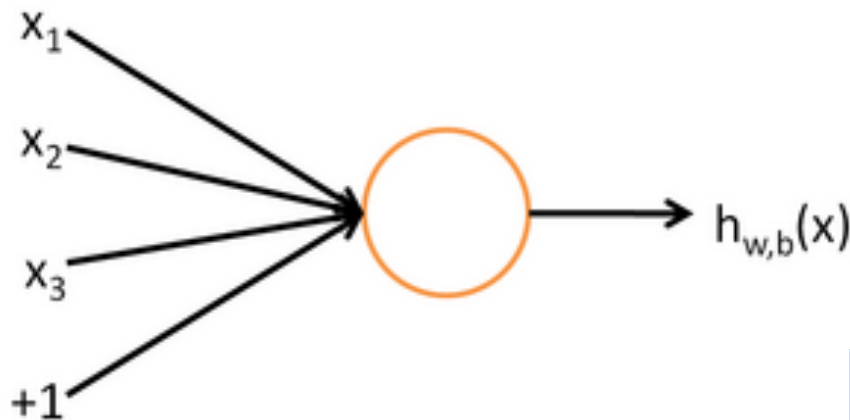
for $f(z) = 1/(1 + \exp(-z))$, the logistic function – a sigmoid **non-linearity**.

This is exactly what a neuron computes

$$h_{w,b}(x) = f(w^T x + b)$$

b : We can have an “always on” feature, which gives a class prior, or separate it out, as a bias term

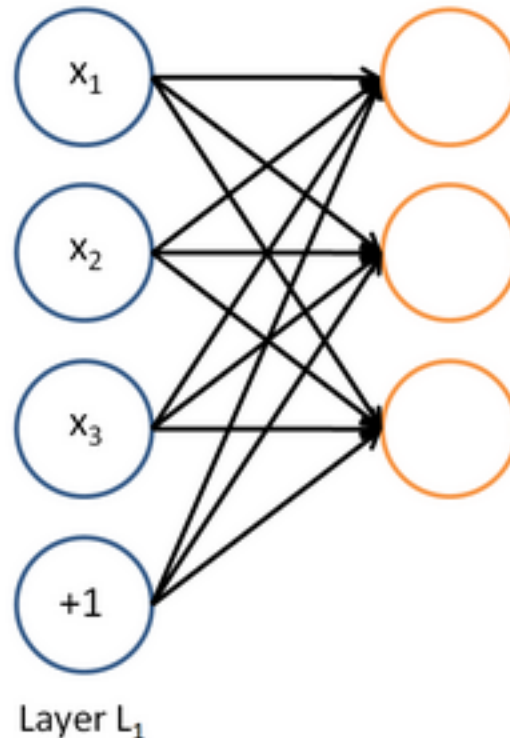
$$f(z) = \frac{1}{1 + e^{-z}}$$



w , b are the parameters of this neuron
i.e., this logistic regression model

A neural network = running several logistic regressions at the same time

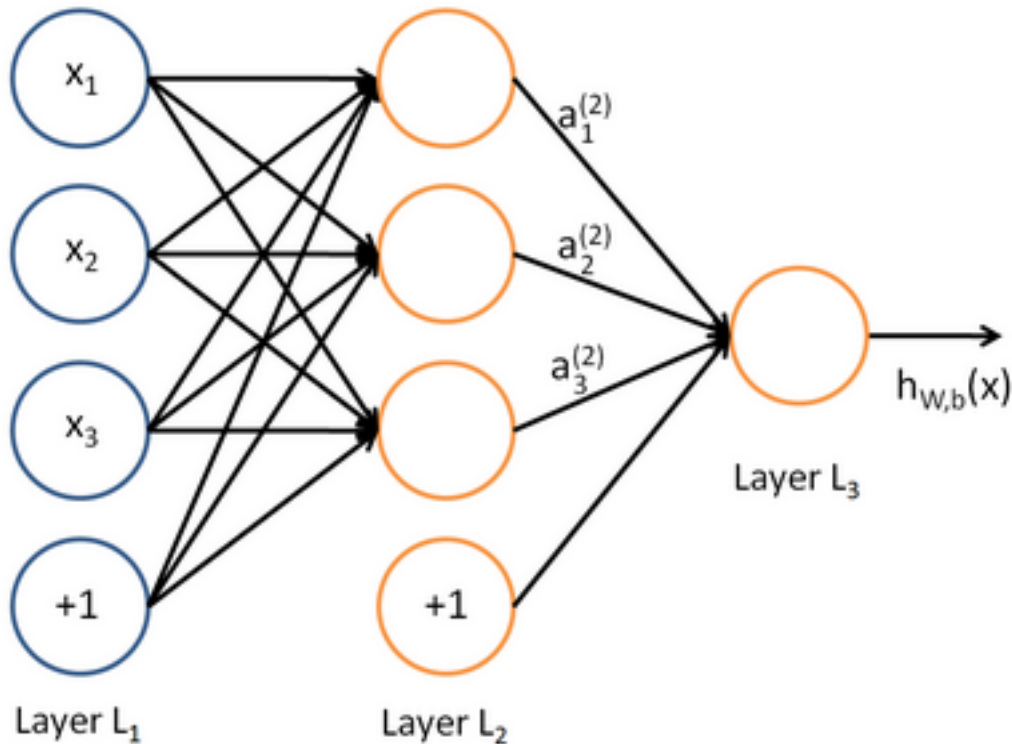
If we feed a vector of inputs through a bunch of logistic regression functions, then we get a vector of outputs



But we don't have to decide ahead of time what variables these logistic regressions are trying to predict!

A neural network = running several logistic regressions at the same time

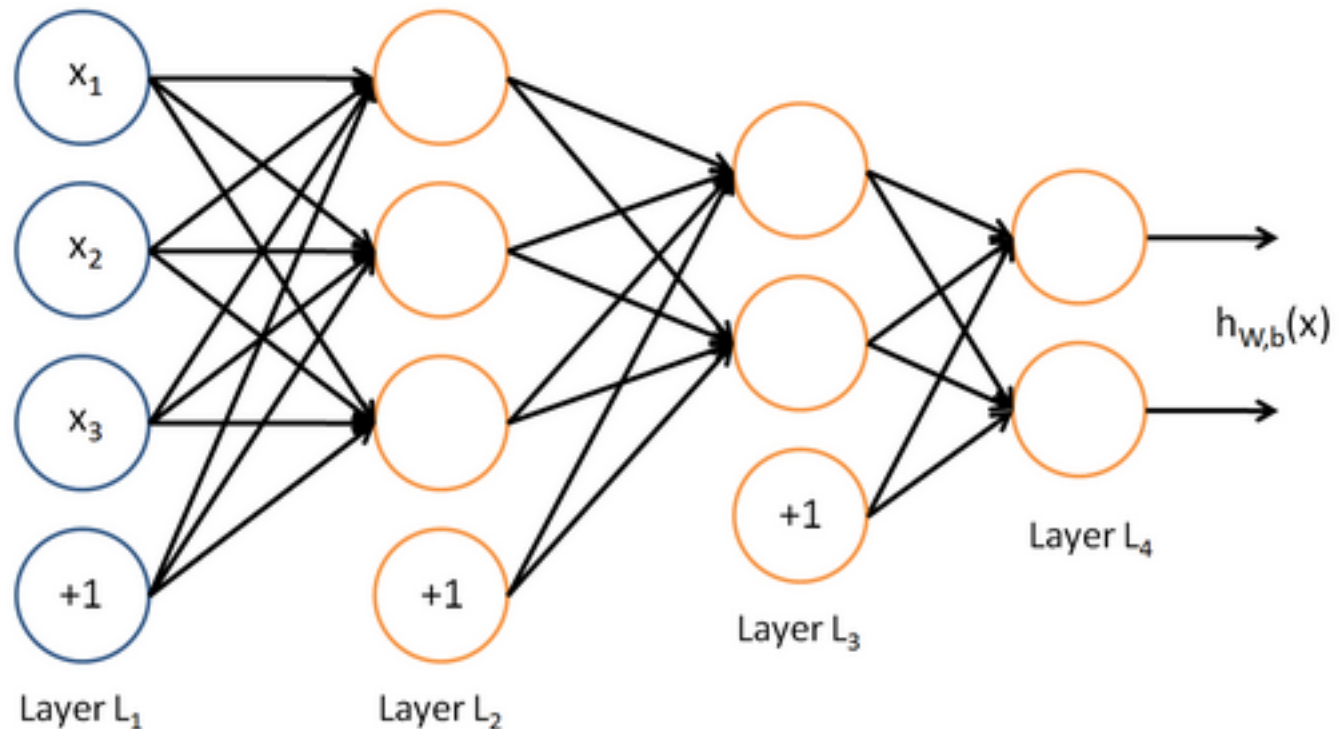
... which we can feed into another logistic regression function



and it is the training criterion that will decide what those intermediate binary target variables should be, so as to do a good job of predicting the targets for the next layer, etc.

A neural network = running several logistic regressions at the same time

Before we know it, we have a multilayer neural network....



Matrix notation for a layer

We have

$$a_1 = f(W_{11}x_1 + W_{12}x_2 + W_{13}x_3 + b_1)$$

$$a_2 = f(W_{21}x_1 + W_{22}x_2 + W_{23}x_3 + b_2)$$

etc.

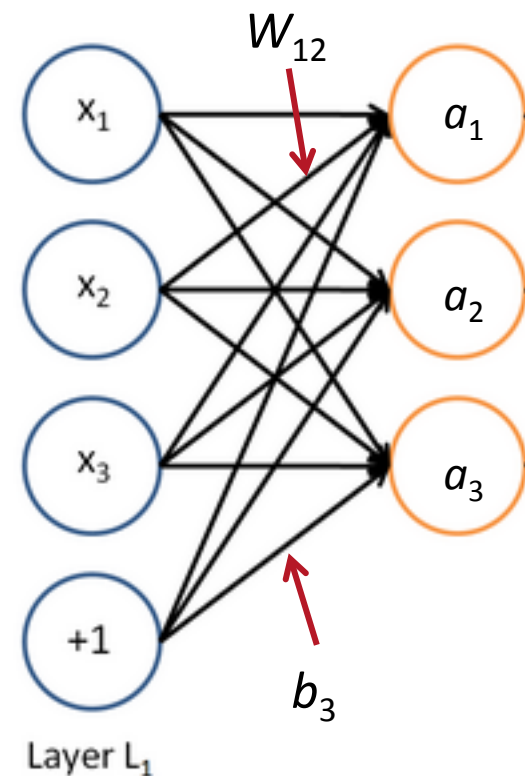
In matrix notation

$$z = Wx + b$$

$$a = f(z)$$

where f is applied element-wise:

$$f([z_1, z_2, z_3]) = [f(z_1), f(z_2), f(z_3)]$$

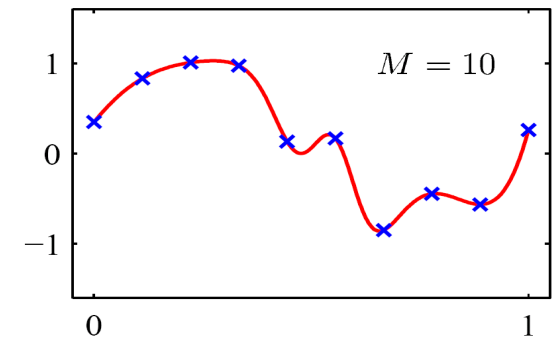
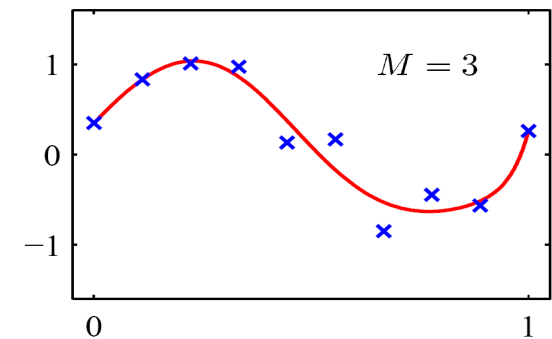
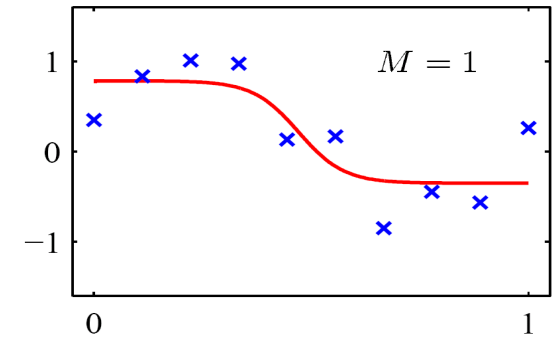


How do we train the weights W ?

- For a supervised single layer neural net, we can train the model just like a maxent model – we calculate and use gradients
 - Stochastic gradient descent (SGD)
 - Conjugate gradient or L-BFGS
- A multilayer net could be more complex because the internal (“hidden”) logistic units make the function non-convex ... just as for hidden CRFs [Quattoni et al. 2005, Gunawardana et al. 2005]
 - But we can use the same ideas and techniques
 - Just without guarantees ...
 - This leads into “**backpropagation**”, which we cover later

Non-linearities: Why they're needed

- For logistic regression: map to probabilities
- Here: function approximation, e.g., regression or classification
 - Without non-linearities, deep neural networks can't do anything more than a linear transform
 - Extra layers could just be compiled down into a single linear transform
 - Probabilistic interpretation unnecessary except in the Boltzmann machine/graphical models
 - People often use other non-linearities, such as **tanh**, as we'll discuss in part 3



Summary

Knowing the meaning of words!

You now understand the basics and the relation to other models

- Neuron = logistic regression or similar function
- Input layer = input training/test vector
- Bias unit = intercept term/always on feature
- Activation = response
- Activation function is a logistic (or similar “sigmoid” nonlinearity)
- Backpropagation = running stochastic gradient descent across a multilayer network
- Weight decay = regularization / Bayesian prior

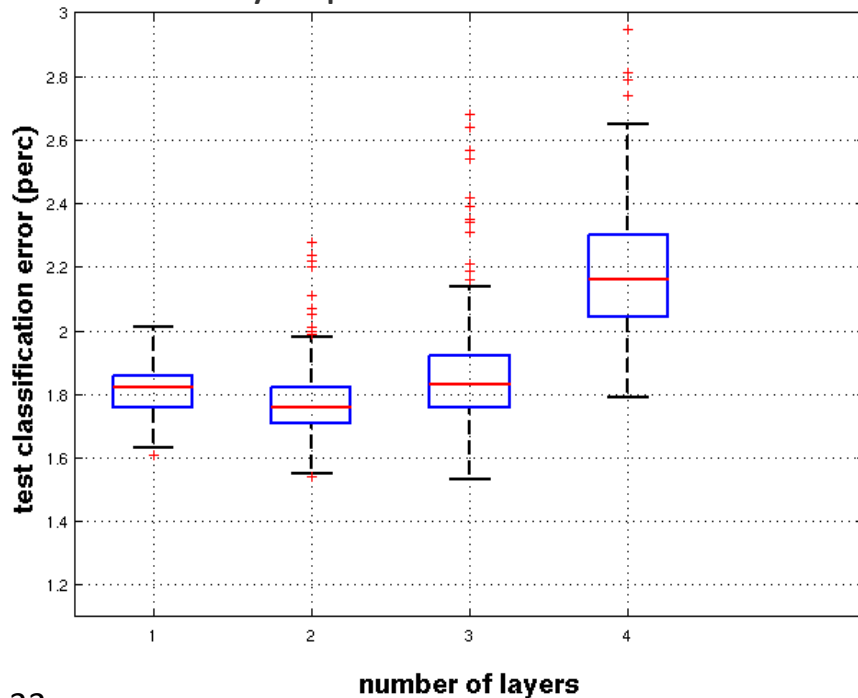
Effective deep learning became possible through unsupervised pre-training

[Erhan et al., JMLR 2010]

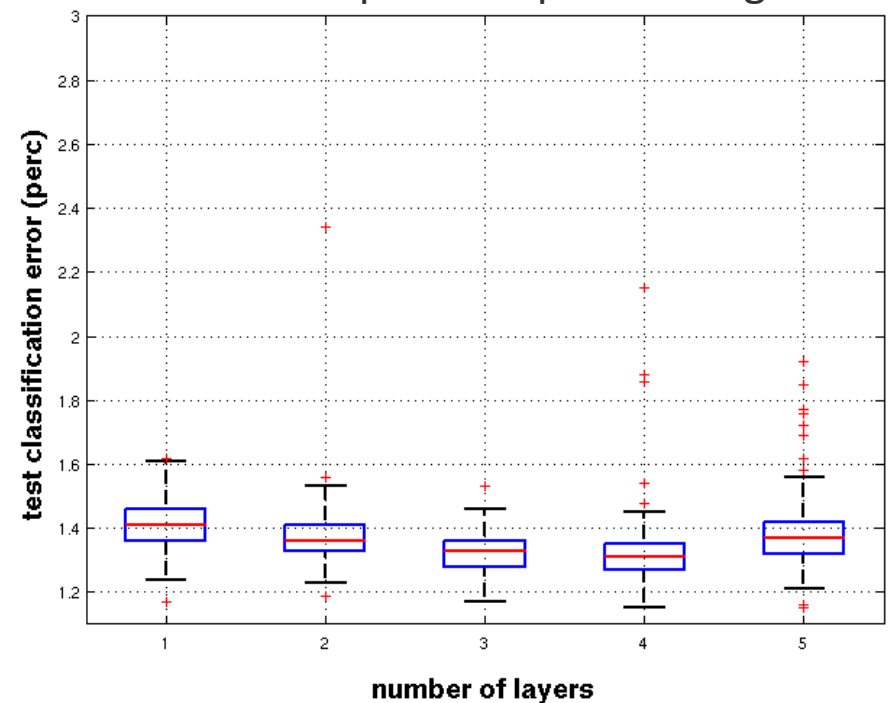


(with RBMs and Denoising Auto-Encoders)

Purely supervised neural net



With unsupervised pre-training



Part 1.3: The Basics

Word Representations

The standard word representation

The vast majority of rule-based **and** statistical NLP work regards words as atomic symbols: *hotel, conference, walk*

In vector space terms, this is a vector with one 1 and a lot of zeroes

$[0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$

Dimensionality: 20K (speech) – 50K (PTB) – 500K (big vocab) – 13M (Google 1T)

We call this a “*one-hot*” representation. Its problem:

motel $[0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$ AND
hotel $[0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$ = 0

Distributional similarity based representations

You can get a lot of value by representing a word by means of its neighbors

“You shall know a word by the company it keeps”

(J. R. Firth 1957: 11)

One of the most successful ideas of modern statistical NLP

government debt problems turning into banking crises as has happened in

saying that Europe needs unified banking regulation to replace the hodgepodge

↖ These words will represent *banking* ↗

You can vary whether you use local or large context to get a more syntactic or semantic clustering

Class-based (hard) and soft clustering word representations

Class based models learn word classes of similar words based on distributional information (\sim class HMM)

- Brown clustering (Brown et al. 1992)
- Exchange clustering (Martin et al. 1998, Clark 2003)
- Desparsification and great example of unsupervised pre-training

Soft clustering models learn for each cluster/topic a distribution over words of how likely that word is in each cluster

- Latent Semantic Analysis (LSA/LSI), Random projections
- Latent Dirichlet Analysis (LDA), HMM clustering

Neural word embeddings as a distributed representation

Similar idea

Combine vector space semantics with the prediction of probabilistic models (Bengio et al. 2003, Collobert & Weston 2008, Turian et al. 2010)

In all of these approaches, including deep learning models, a word is represented as a dense vector

linguistics =

$$\begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}$$

Neural word embeddings - visualization



Advantages of the neural word embedding approach

Compared to a method like LSA, neural word embeddings can become **more meaningful** through adding supervision from one or multiple tasks

For instance, sentiment is usually not captured in unsupervised word embeddings but can be in neural word vectors

We can build representations for large linguistic units

See part 2

Part 1.4: The Basics

Unsupervised word vector Learning

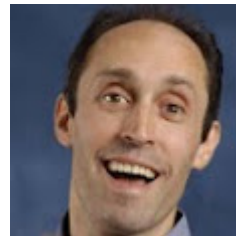
A neural network for learning word vectors

(Collobert et al. JMLR 2011)

Idea: A word and its context is a positive training sample; a random word in that same context gives a negative training sample:

+ cat chills on a mat - cat chills Jeju a mat

Similar: Implicit negative evidence in Contrastive Estimation, (Smith and Eisner 2005)



A neural network for learning word vectors

How do we formalize this idea? Ask that

score(cat chills on a mat) > score(cat chills Jeju a mat)

How do we compute the score?

- With a neural network
- Each word is associated with an n -dimensional vector



Word embedding matrix

- Initialize all word vectors randomly to form a word embedding matrix $L \in \mathbb{R}^{n \times |V|}$

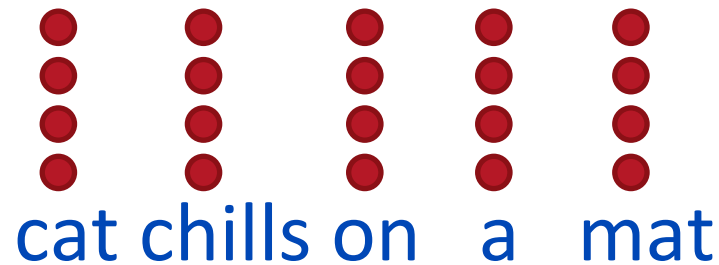
$$L = \begin{bmatrix} \bullet & \bullet & \bullet & \dots & \bullet & \bullet \\ \bullet & \bullet & \bullet & \dots & \bullet & \bullet \\ \bullet & \bullet & \bullet & \dots & \bullet & \bullet \\ \bullet & \bullet & \bullet & \dots & \bullet & \bullet \end{bmatrix}_n$$

the cat mat ...

- These are the word features we want to learn
- Also called a look-up table
 - Conceptually you get a word's vector by left multiplying a one-hot vector e by L : $x = Le$

Word vectors as input to a neural network

- $\text{score}(\text{cat chills on a mat})$
- To describe a phrase, retrieve (via index) the corresponding vectors from L



- Then concatenate them to $5n$ vector:
- $x = [\text{●●●●} \text{●●●●} \text{●●●●} \text{●●●●} \text{●●●●}]$
- How do we then compute $\text{score}(x)$?

A Single Layer Neural Network

- A single layer is a combination of a linear layer and a nonlinearity:

$$z = Wx + b$$

$$a = f(z)$$

- The neural activations can then be used to compute some function.
- For instance, the score we care about:

$$\text{score}(x) = U^T a \in \mathbb{R}$$

Summary: Feed-forward Computation

Computing a window's score with a 3-layer Neural Net: $s = \text{score}(\text{cat chills on a mat})$

$$s = U^T f(Wx + b) \quad x \in \mathbb{R}^{20 \times 1}, W \in \mathbb{R}^{8 \times 20}, U \in \mathbb{R}^{8 \times 1}$$

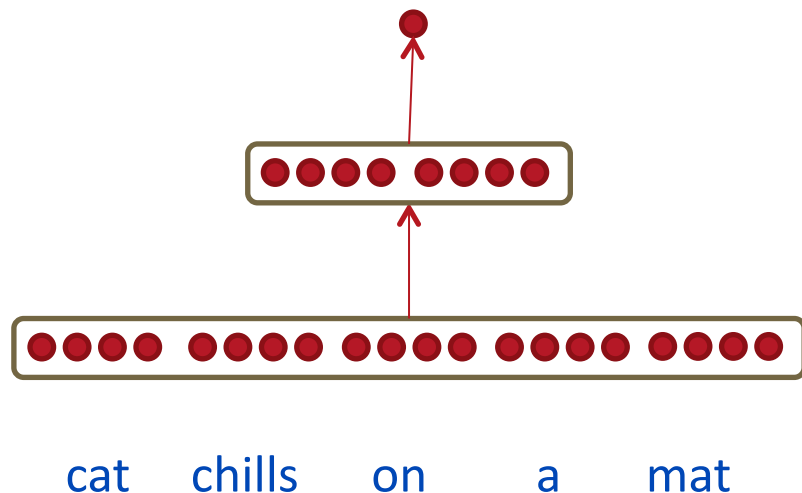
$$s = U^T a$$

$$a = f(z)$$

$$z = Wx + b$$

$$x = [x_{\text{cat}} \quad x_{\text{chills}} \quad x_{\text{on}} \quad x_{\text{a}} \quad x_{\text{mat}}]$$

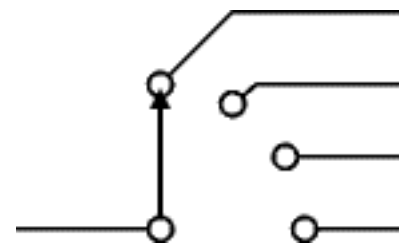
$$L \in \mathbb{R}^{n \times |V|}$$



Summary: Feed-forward Computation

- $s = \text{score}(\text{cat chills on a mat})$
- $s_c = \text{score}(\text{cat chills Jeju a mat})$
- Idea for training objective: make score of true window larger and corrupt window's score lower (until they're good enough): minimize

$$J = \max(0, 1 - s + s_c)$$



- This is continuous, can perform SGD

Training with Backpropagation

$$J = \max(0, 1 - s + s_c)$$

$$s = U^T f(Wx + b)$$

$$s_c = U^T f(Wx_c + b)$$

Assuming cost J is > 0 , it is simple to see that we can compute the derivatives of s and s_c wrt all the involved variables: U, W, b, x

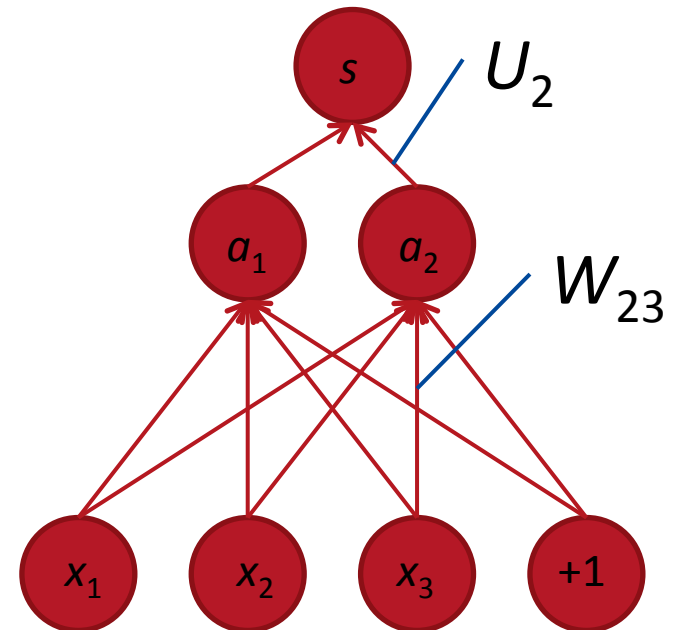
$$\frac{\partial s}{\partial U} = \frac{\partial}{\partial U} U^T a \qquad \frac{\partial s}{\partial U} = a$$

Training with Backpropagation

- Let's consider the derivative of a single weight W_{ij}

$$\frac{\partial s}{\partial W} = \frac{\partial}{\partial W} U^T a = \frac{\partial}{\partial W} U^T f(z) = \frac{\partial}{\partial W} U^T f(Wx + b)$$

- This only appears inside a_i
- For example: W_{23} is only used to compute a_2



Training with Backpropagation

$$\frac{\partial s}{\partial W} = \frac{\partial}{\partial W} U^T a = \frac{\partial}{\partial W} U^T f(z) = \frac{\partial}{\partial W} U^T f(Wx + b)$$

Derivative of weight W_{ij} :

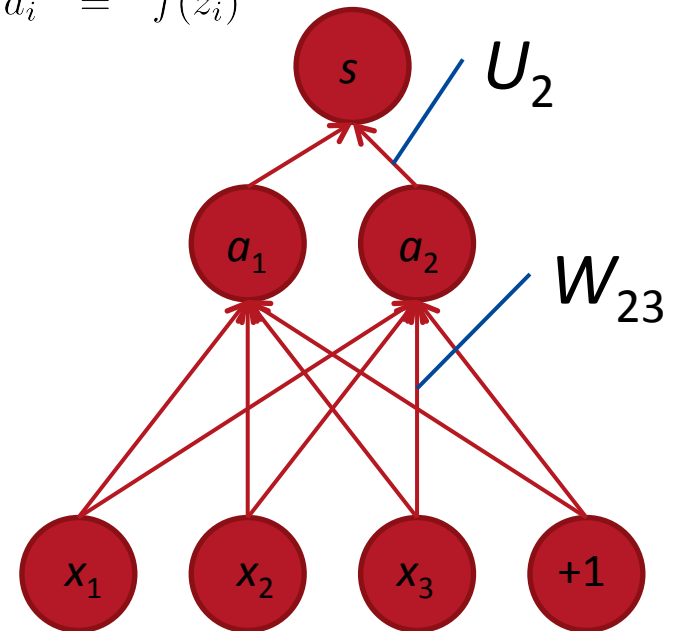
$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \frac{\partial u}{\partial x}$$

$$\frac{\partial}{\partial W_{ij}} U^T a \rightarrow \frac{\partial}{\partial W_{ij}} U_i a_i$$

$$z_i = W_{i \cdot} x + b_i = \sum_{j=1}^3 W_{ij} x_j + b_i$$

$$a_i = f(z_i)$$

$$\begin{aligned} U_i \frac{\partial}{\partial W_{ij}} a_i &= U_i \frac{\partial a_i}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} \\ &= U_i \frac{\partial f(z_i)}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} \\ &= U_i f'(z_i) \frac{\partial z_i}{\partial W_{ij}} \\ &= U_i f'(z_i) \frac{\partial W_{i \cdot} x + b_i}{\partial W_{ij}} \end{aligned}$$



Training with Backpropagation

Derivative of single weight W_{ij} :

$$= U_i f'(z_i) \frac{\partial W_i \cdot x + b_i}{\partial W_{ij}}$$

$$= U_i f'(z_i) \frac{\partial}{\partial W_{ij}} \sum_k W_{ik} x_k$$

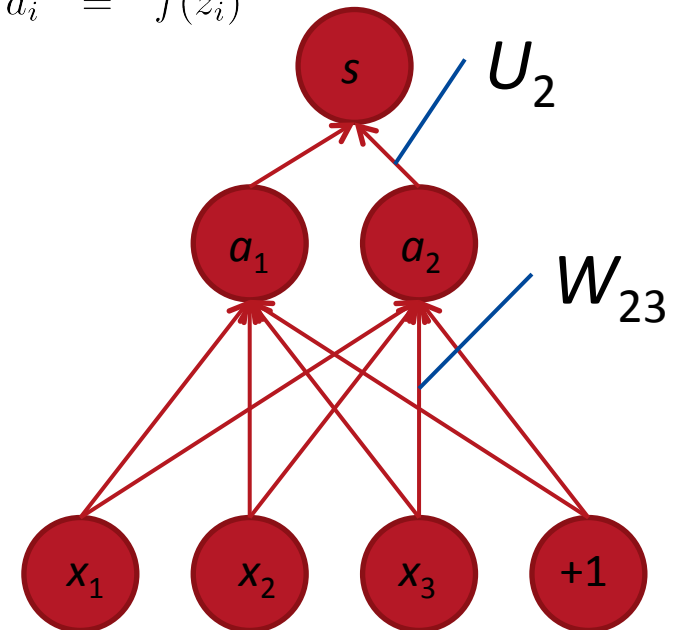
$$= \underbrace{U_i f'(z_i)}_{\delta_i} x_j$$

$$= \underbrace{\delta_i}_{\text{Local error signal}} \underbrace{x_j}_{\text{Local input signal}}$$

$$U_i \frac{\partial}{\partial W_{ij}} a_i$$

$$z_i = W_i \cdot x + b_i = \sum_{j=1}^3 W_{ij} x_j + b_i$$

$$a_i = f(z_i)$$



Training with Backpropagation

- From single weight W_{ij} to full W :

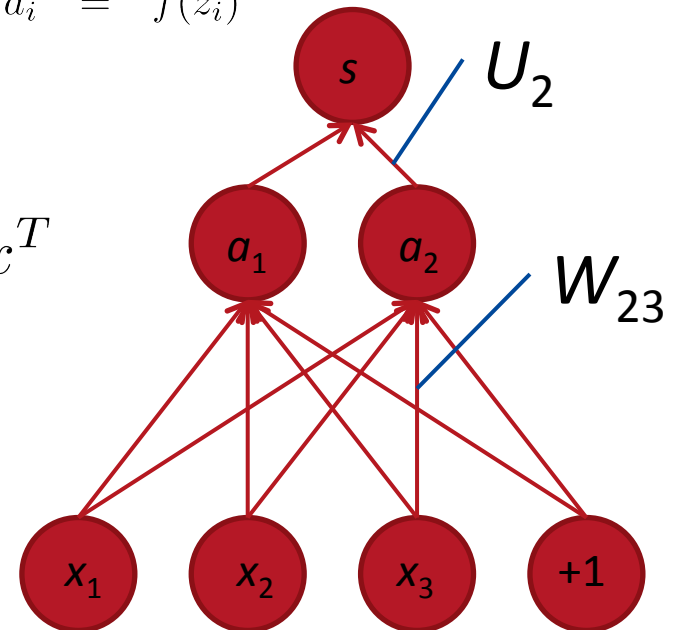
$$\begin{aligned}\frac{\partial J}{\partial W_{ij}} &= \underbrace{U_i f'(z_i)}_{\delta_i} x_j \\ &= \delta_i x_j\end{aligned}$$

$$z_i = W_i \cdot x + b_i = \sum_{j=1}^3 W_{ij} x_j + b_i$$

$$a_i = f(z_i)$$

- We want all combinations of $i = 1, 2$ and $j = 1, 2, 3$

- Solution: Outer product: $\frac{\partial J}{\partial W} = \delta x^T$
where $\delta \in \mathbb{R}^{2 \times 1}$ is the “responsibility” coming from each activation a



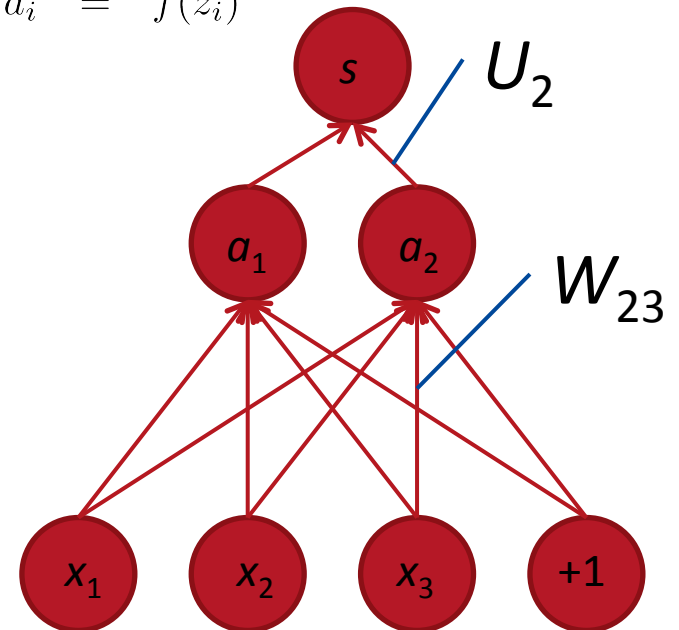
Training with Backpropagation

- For biases b , we get:

$$\begin{aligned} & U_i \frac{\partial}{\partial b_i} a_i \\ = & U_i f'(z_i) \frac{\partial W_i \cdot x + b_i}{\partial b_i} \\ = & \delta_i \end{aligned}$$

$$z_i = W_i \cdot x + b_i = \sum_{j=1}^3 W_{ij} x_j + b_i$$

$$a_i = f(z_i)$$



Training with Backpropagation

That's almost backpropagation

It's simply taking derivatives and using the chain rule!

Remaining trick: we can re-use derivatives computed for higher layers in computing derivatives for lower layers

Example: last derivatives of model, the word vectors in x

Training with Backpropagation

- Take derivative of score with respect to single word vector (for simplicity a 1d vector, but same if it was longer)
- Now, we cannot just take into consideration one a_i because each x_j is connected to all the neurons above and hence x_j influences the overall score through all of these, hence:

$$\begin{aligned}\frac{\partial s}{\partial x_j} &= \sum_{i=1}^2 \frac{\partial s}{\partial a_i} \frac{\partial a_i}{\partial x_j} \\ &= \sum_{i=1}^2 \frac{\partial U^T a}{\partial a_i} \frac{\partial a_i}{\partial x_j} \\ &= \sum_{i=1}^2 U_i \frac{\partial f(W_i \cdot x + b)}{\partial x_j} \\ &= \sum_{i=1}^2 \underbrace{U_i f'(W_i \cdot x + b)}_{\delta_i} \frac{\partial W_i \cdot x}{\partial x_j} \\ &= \sum_{i=1}^2 \delta_i W_{ij} \\ &= \underbrace{\delta^T}_{\text{Re-used part of previous derivative}} W_{\cdot j}\end{aligned}$$

Training with Backpropagation: softmax

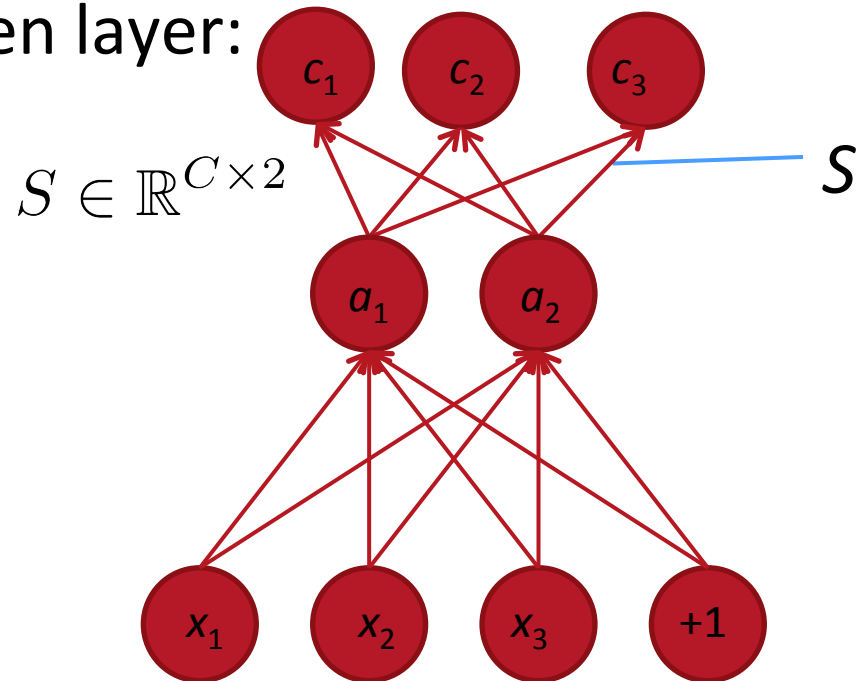
What is the major benefit of learned word vectors?

Ability to also propagate labeled information into them, via softmax/maxent and hidden layer:

$$P(c|d, \lambda) = \frac{e^{\lambda^T f(c,d)}}{\sum_{c'} e^{\lambda^T f(c',d)}}$$



$$p(c|x) = \frac{\exp(S_c \cdot a)}{\sum_{c'} \exp(S_{c'} \cdot a)}$$



Part 1.5: The Basics

Backpropagation Training

Back-Prop

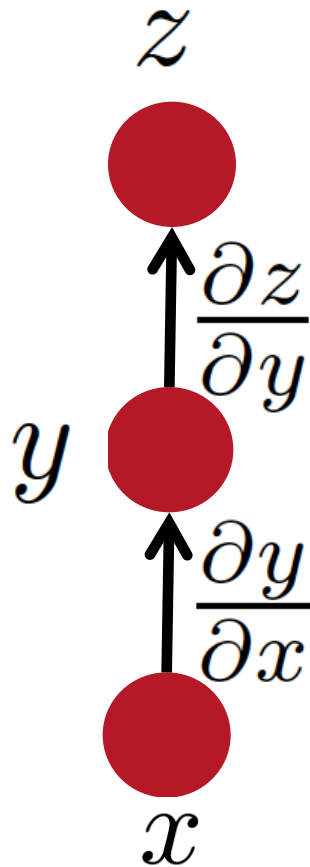
- Compute gradient of example-wise loss wrt parameters

- Simply applying the derivative chain rule wisely

$$z = f(y) \quad y = g(x) \quad \frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

- If computing the loss(example, parameters) is $O(n)$ computation, then so is computing the gradient

Simple Chain Rule



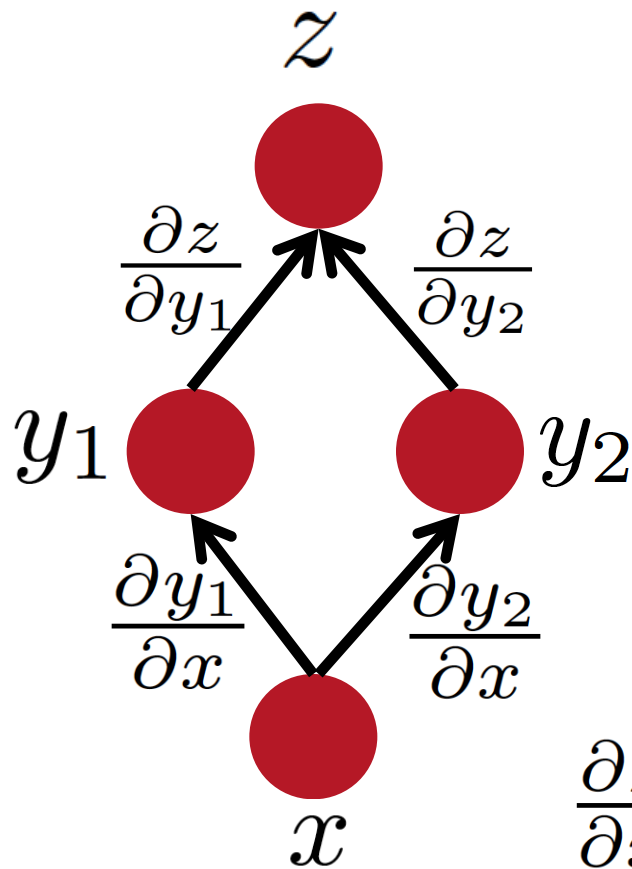
$$\Delta z = \frac{\partial z}{\partial y} \Delta y$$

$$\Delta y = \frac{\partial y}{\partial x} \Delta x$$

$$\Delta z = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \Delta x$$

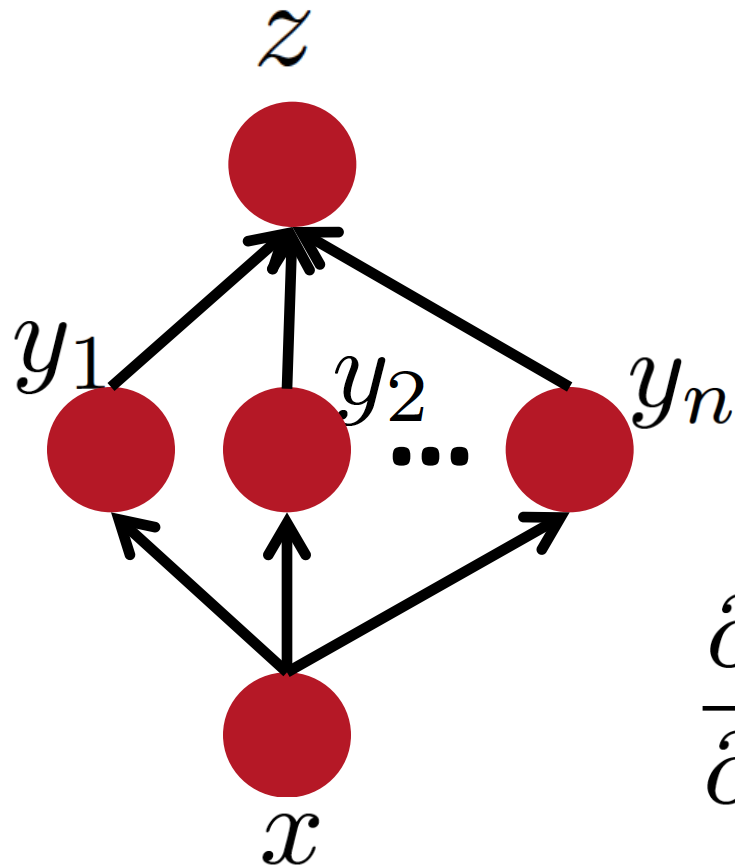
$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

Multiple Paths Chain Rule



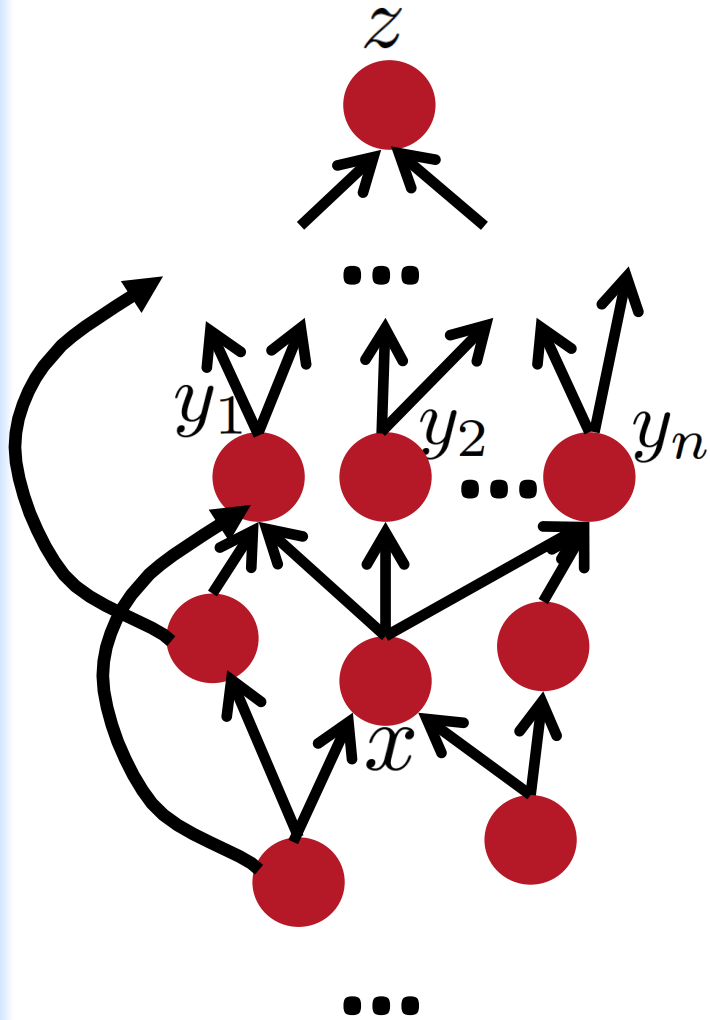
$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x} + \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial x}$$

Multiple Paths Chain Rule - General



$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$$

Chain Rule in Flow Graph

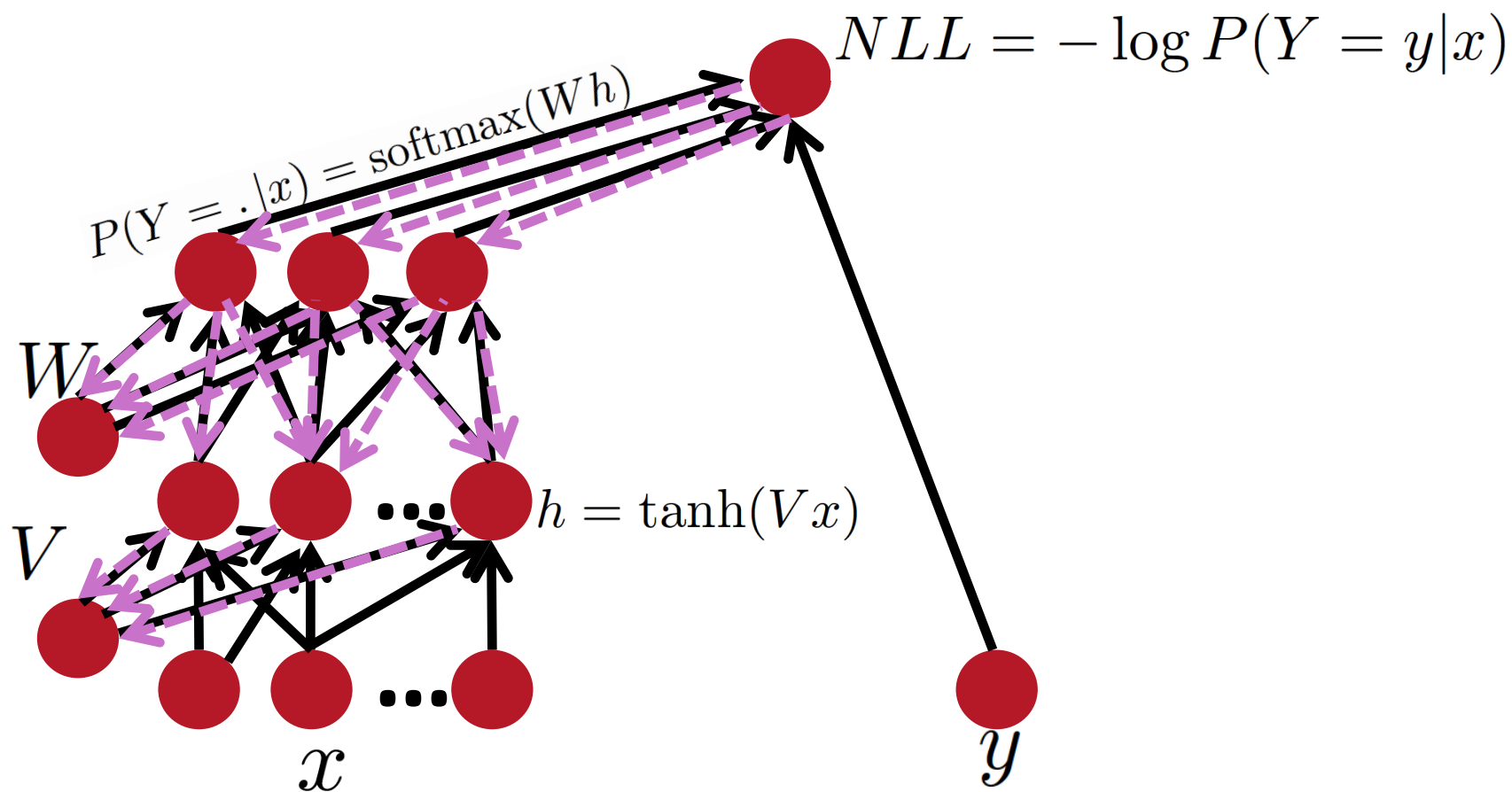


Flow graph: any directed acyclic graph
node = computation result
arc = computation dependency

$\{y_1, y_2, \dots, y_n\}$ = successors of x

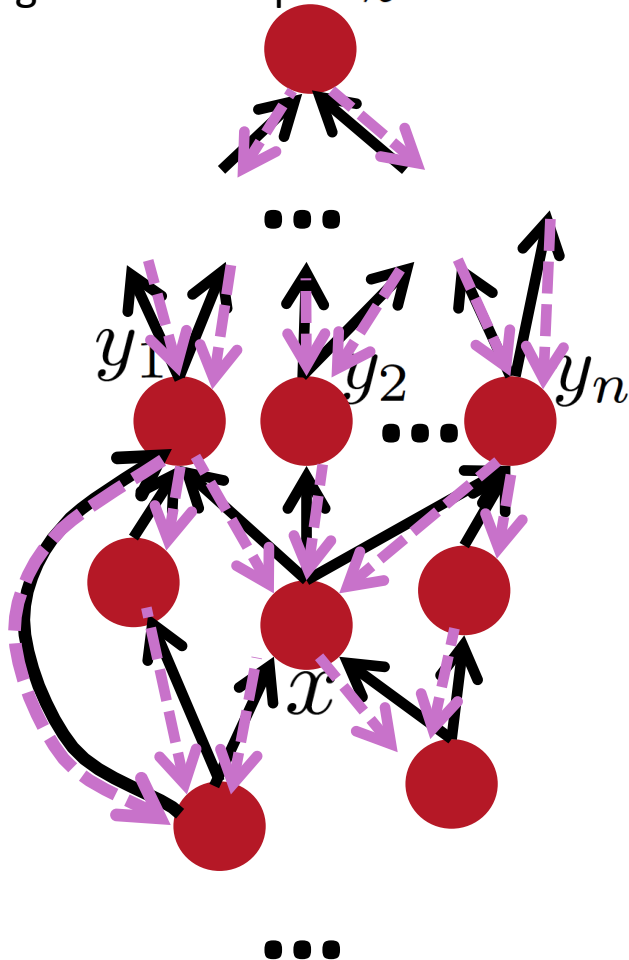
$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$$

Back-Prop in Multi-Layer Net



Back-Prop in General Flow Graph

Single scalar output z

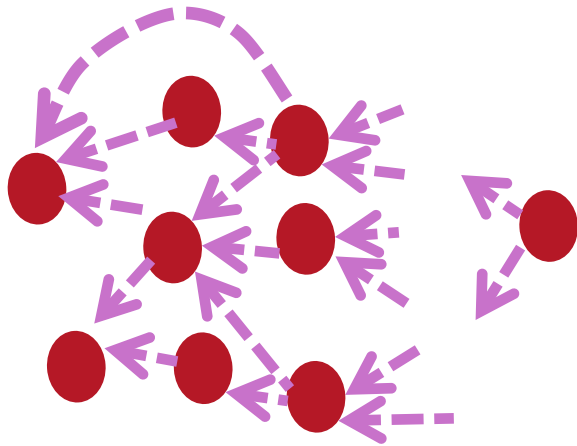
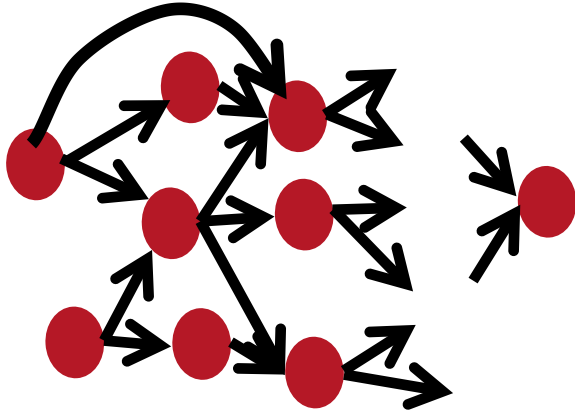


1. Fprop: visit nodes in topo-sort order
 - Compute value of node given predecessors
2. Bprop:
 - initialize output gradient = 1
 - visit nodes in reverse order:
 - Compute gradient wrt each node using gradient wrt successors

$\{y_1, y_2, \dots, y_n\}$ = successors of x

$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$$

Automatic Differentiation



- The gradient computation can be automatically inferred from the symbolic expression of the fprop.
- Each node type needs to know how to compute its output and how to compute the gradient wrt its inputs given the gradient wrt its output.
- Easy and fast prototyping

Part 1.6: The Basics

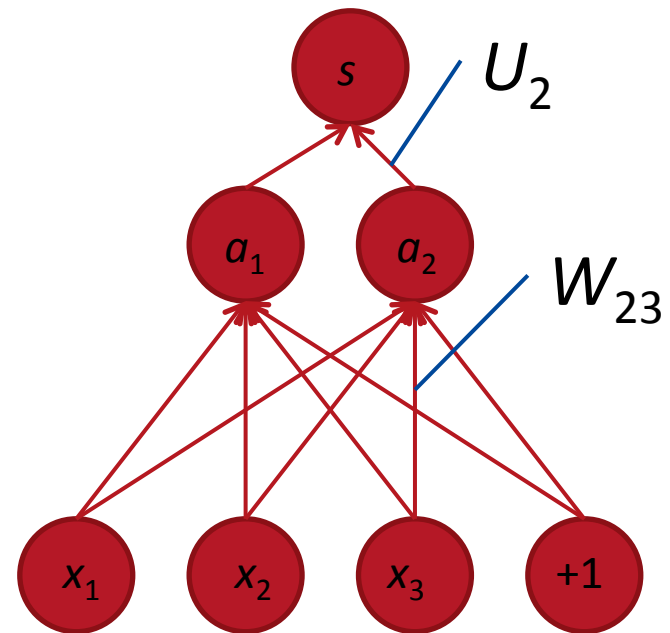
Learning word-level classifiers: POS and NER

The Model

(Collobert & Weston 2008;
Collobert et al. 2011)

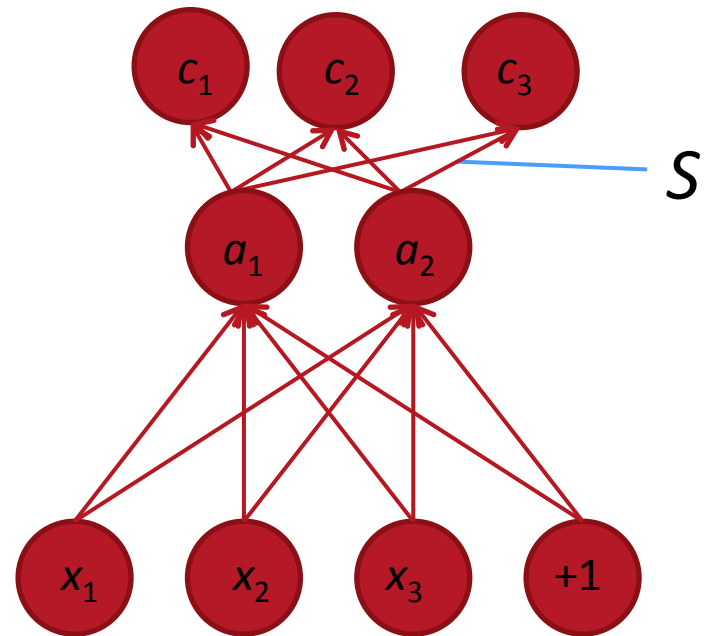
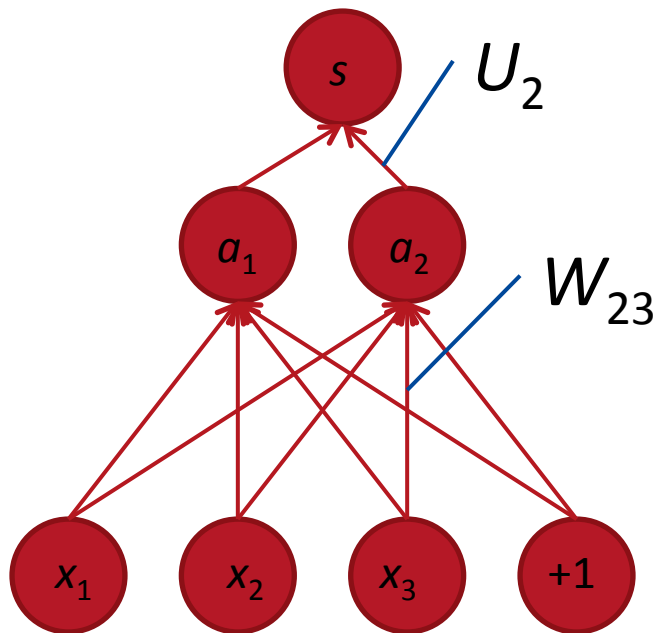


- Similar to word vector learning but replaces the single scalar score with a *Softmax/Maxent* classifier
- Training is again done via backpropagation which gives an error similar to the score in the unsupervised word vector learning model



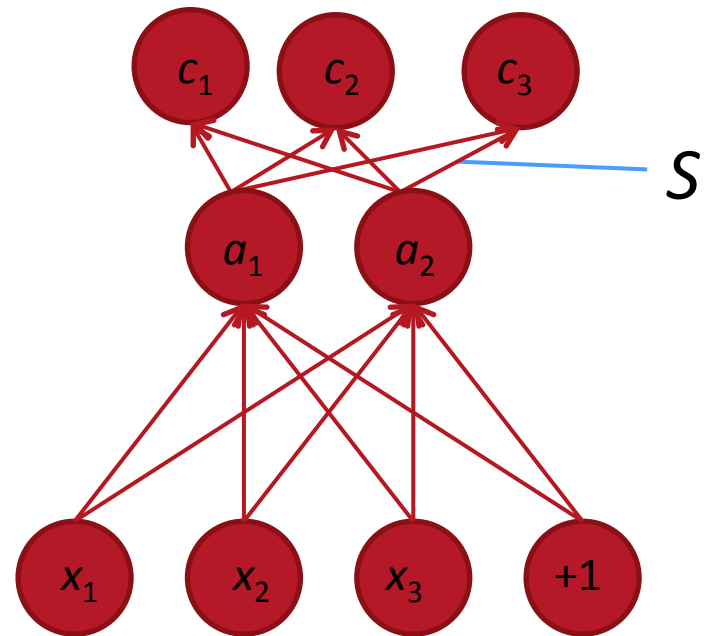
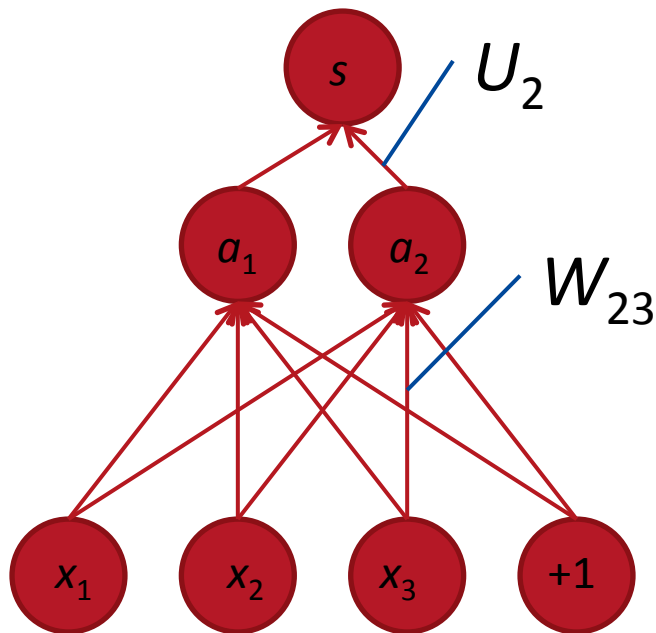
The Model - Training

- We already know the softmax classifier and how to optimize it
- The interesting twist in deep learning is that the input features are also learned, similar to learning word vectors with a score:



The Model - Training

- All derivatives of layers beneath the score were multiplied by U , for the *softmax*, the error vector becomes the difference between predicted and gold standard distributions



The secret sauce is the unsupervised pre-training on a large text collection

	POS WSJ (acc.)	NER CoNLL (F1)
State-of-the-art*	97.24	89.31
Supervised NN	96.37	81.47
Unsupervised pre-training followed by supervised NN**	97.20	88.87
+ hand-crafted features***	97.29	89.59

* Representative systems: POS: (Toutanova et al. 2003), NER: (Ando & Zhang 2005)

** 130,000-word embedding trained on Wikipedia and Reuters with 11 word window, 100 unit hidden layer – **for 7 weeks!** – then supervised task training

*** Features are character suffixes for POS and a gazetteer for NER

Supervised refinement of the unsupervised word representation helps

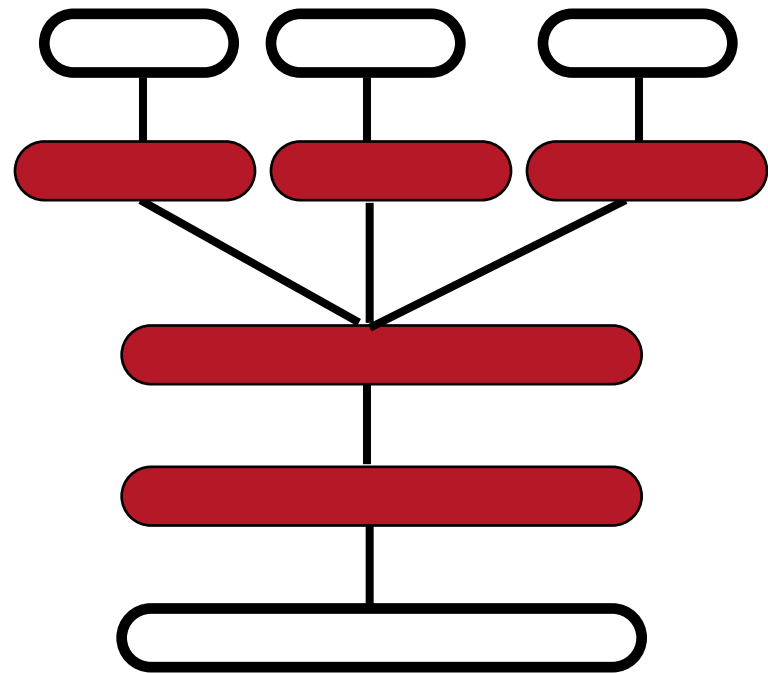
	POS WSJ (acc.)	NER CoNLL (F1)
Supervised NN	96.37	81.47
NN with Brown clusters	96.92	87.15
Fixed embeddings*	97.10	88.87
C&W 2011**	97.29	89.59

* Same architecture as C&W 2011, but word embeddings are kept constant during the supervised training phase

** C&W is unsupervised pre-train + supervised NN + features model of last slide

Multi-Task Learning

- Generalizing better to new tasks is crucial to approach AI
- Deep architectures learn good intermediate representations that can be shared across tasks
- Good representations make sense for many tasks



Combining Multiple Sources of Evidence with Shared Embeddings

- Relational learning
- Multiple sources of information / relations
- Some symbols (e.g. words, wikipedia entries) shared
- Shared embeddings help propagate information among data sources: e.g., WordNet, XWN, Wikipedia, FreeBase, ...

Part 1.7

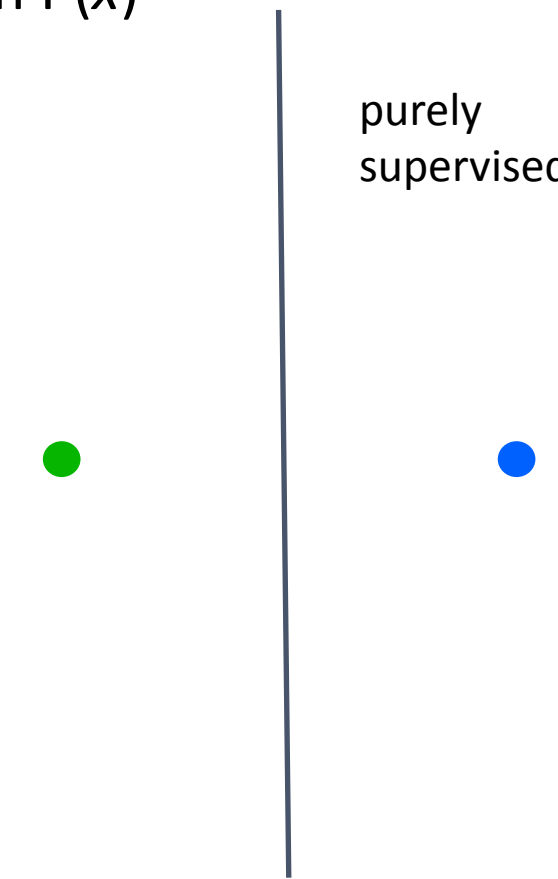
Sharing statistical strength

Sharing Statistical Strength

- Besides very fast prediction, the main advantage of deep learning is **statistical**
- Potential to learn from less labeled examples because of sharing of statistical strength:
 - Unsupervised pre-training & Multi-task learning
 - Semi-supervised learning →

Semi-Supervised Learning

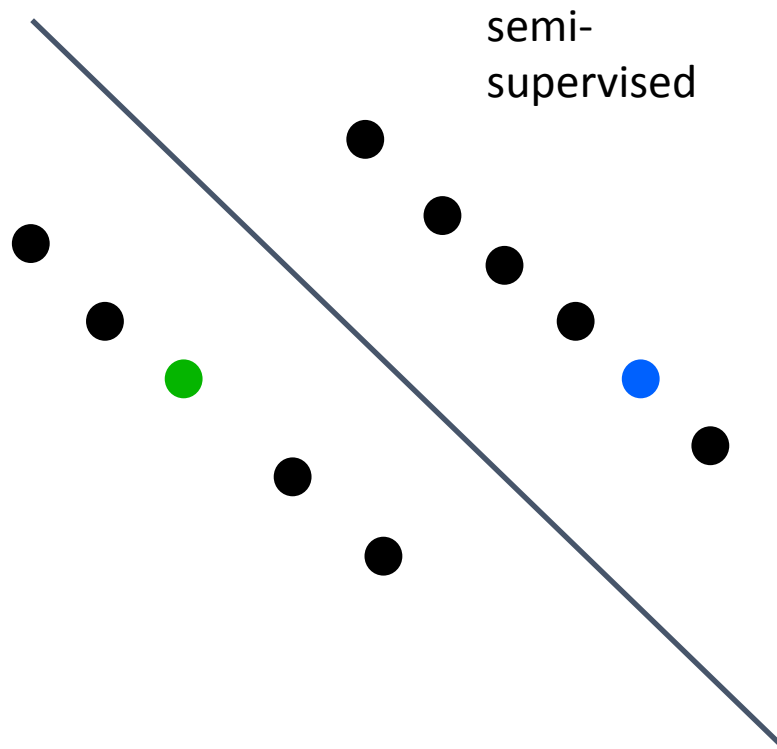
- Hypothesis: $P(c|x)$ can be more accurately computed using shared structure with $P(x)$



purely
supervised

Semi-Supervised Learning

- Hypothesis: $P(c|x)$ can be more accurately computed using shared structure with $P(x)$



Deep autoencoders

- Alternative to contrastive unsupervised word learning
 - Another is RBMs ([Hinton et al. 2006](#)), which we don't cover today
- Works well for fixed input representations (word vectors are not but bag of word representations are)
 1. Definition, intuition and variants of autoencoders
 2. Stacking for deep autoencoders
 3. Why do autoencoders improve deep neural nets so much?

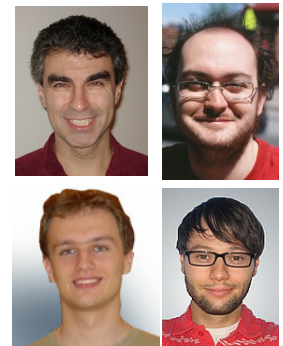
Auto-Encoders

- Multilayer neural net with target output = input
- Reconstruction=decoder(encoder(input))

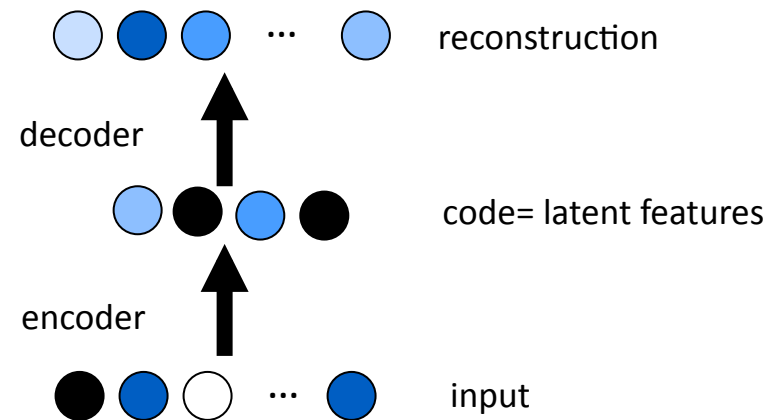
$$a = \tanh(Wx + b)$$

$$x' = \tanh(W^T a + c)$$

$$cost = \|x' - x\|^2$$



- Probable inputs have small reconstruction error



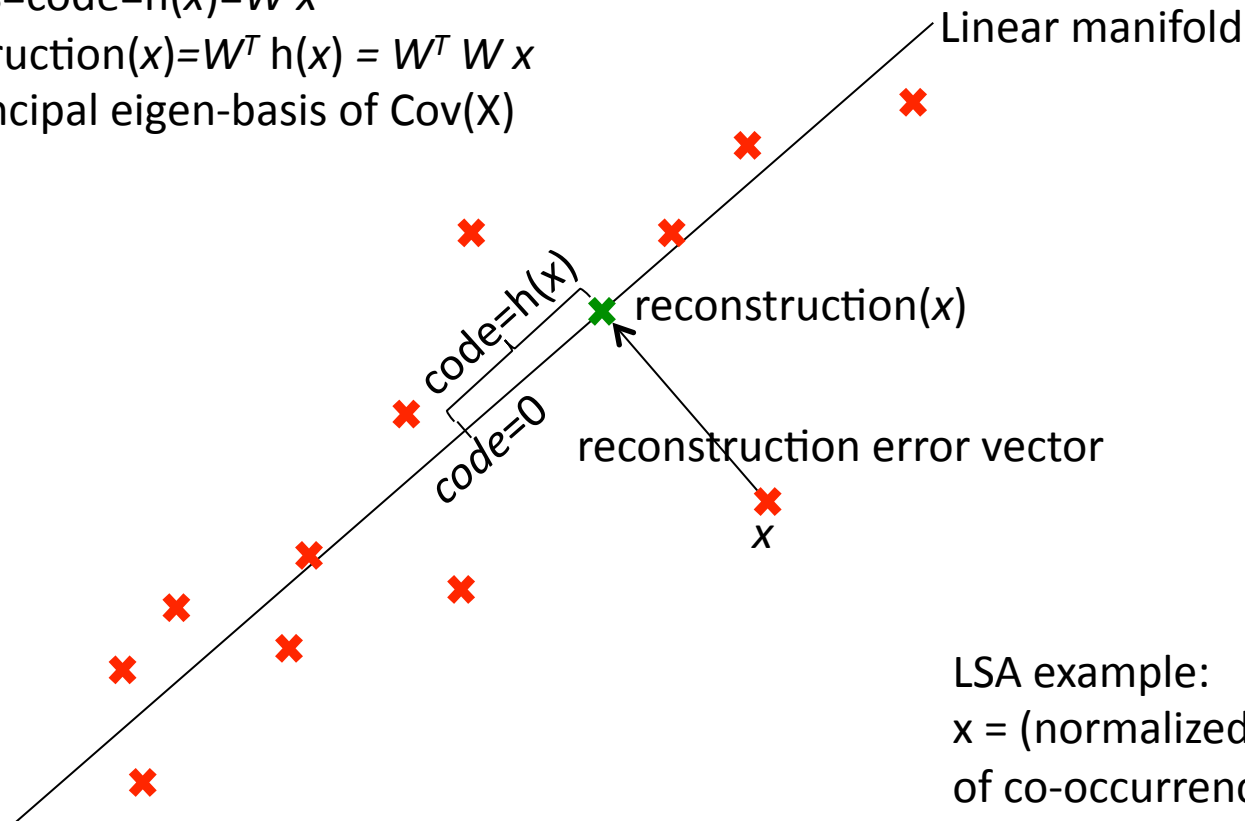
PCA = Linear Manifold = Linear Auto-Encoder

input x , 0-mean

features=code= $h(x)=W x$

reconstruction(x)= $W^T h(x) = W^T W x$

W = principal eigen-basis of $\text{Cov}(X)$

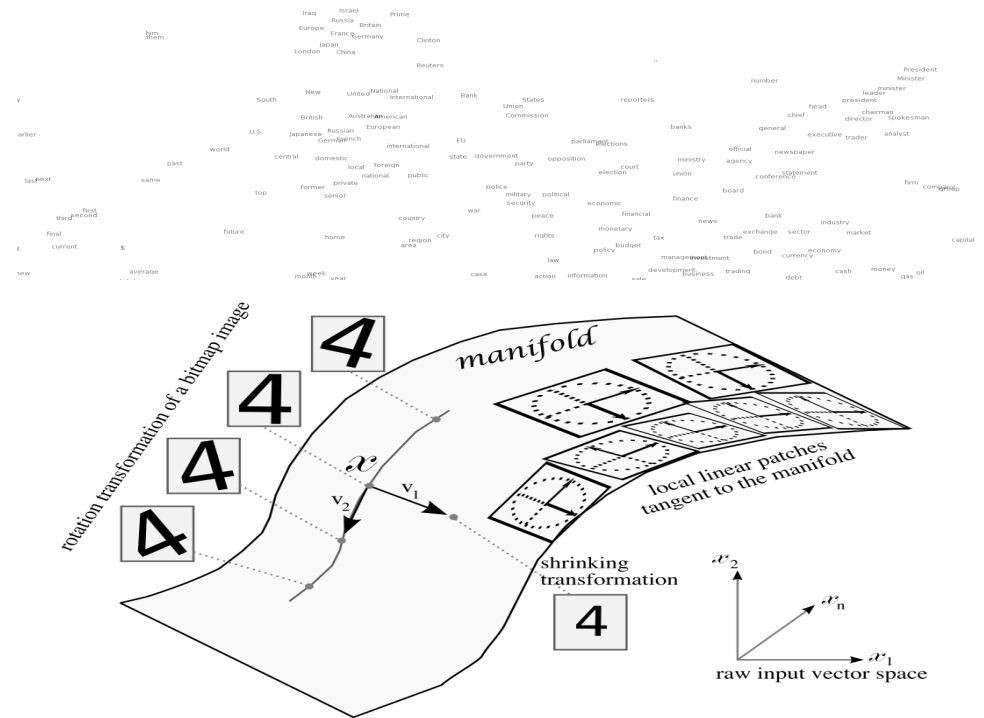
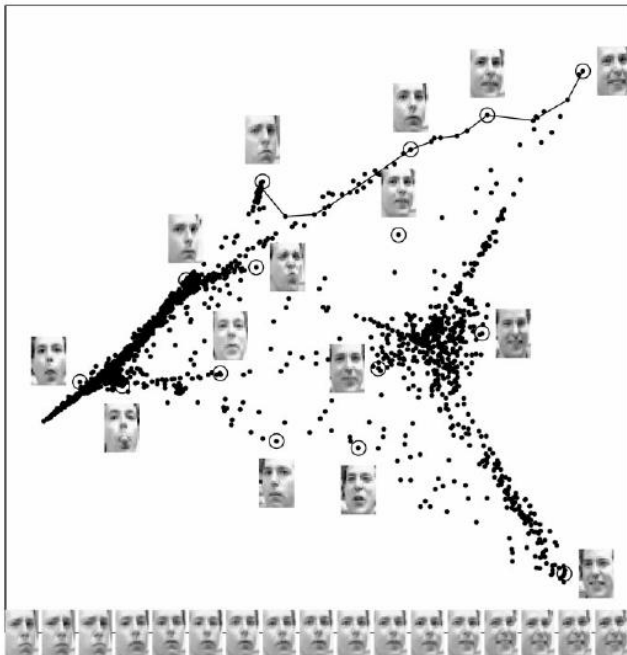


LSA example:

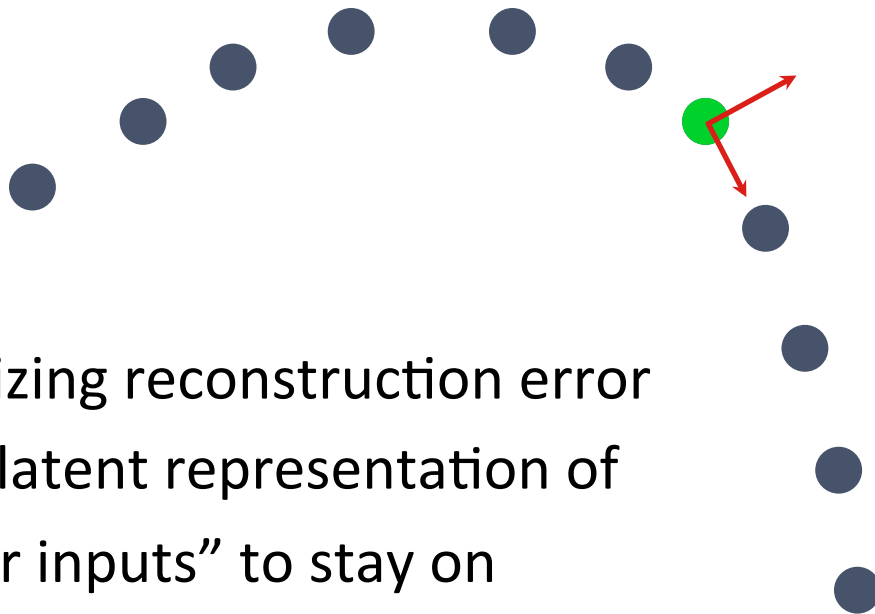
x = (normalized) distribution
of co-occurrence frequencies

The Manifold Learning Hypothesis

- Examples concentrate near a lower dimensional “manifold” (region of high density where small changes are only allowed in certain directions)



Auto-Encoders Learn Salient Variations, like a non-linear PCA



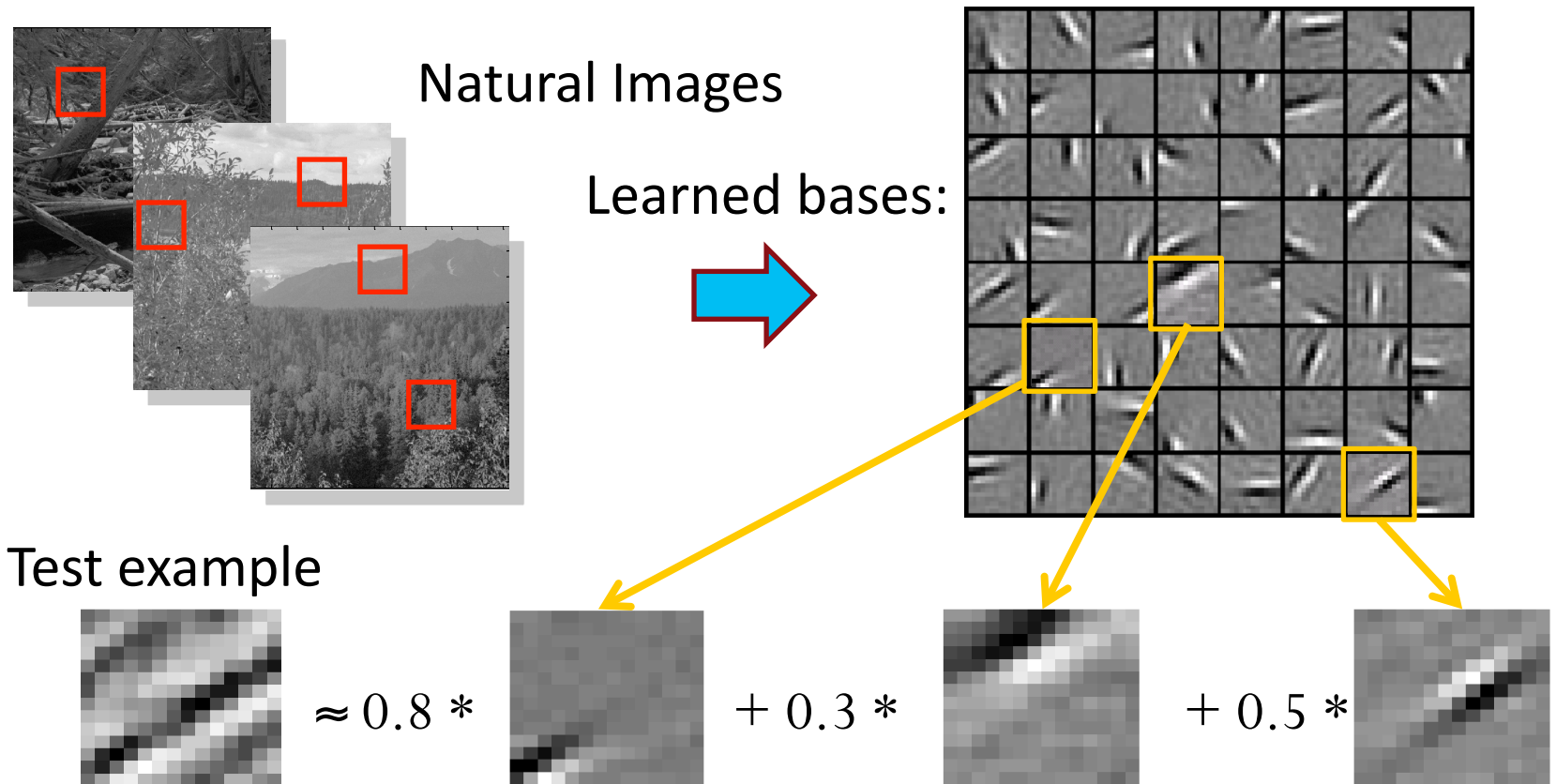
Minimizing reconstruction error forces latent representation of “similar inputs” to stay on manifold

Auto-Encoder Variants

- Discrete inputs: cross-entropy or log-likelihood reconstruction criterion (similar to used for discrete targets for MLPs)
- Preventing them to learn the identity everywhere:
 - Undercomplete (eg PCA): bottleneck code smaller than input
 - Sparsity: penalize hidden unit activations so at or near 0
[Goodfellow et al 2009]
 - Denoising: predict true input from corrupted input
[Vincent et al 2008]
 - Contractive: force encoder to have small derivatives
[Rifai et al 2011]



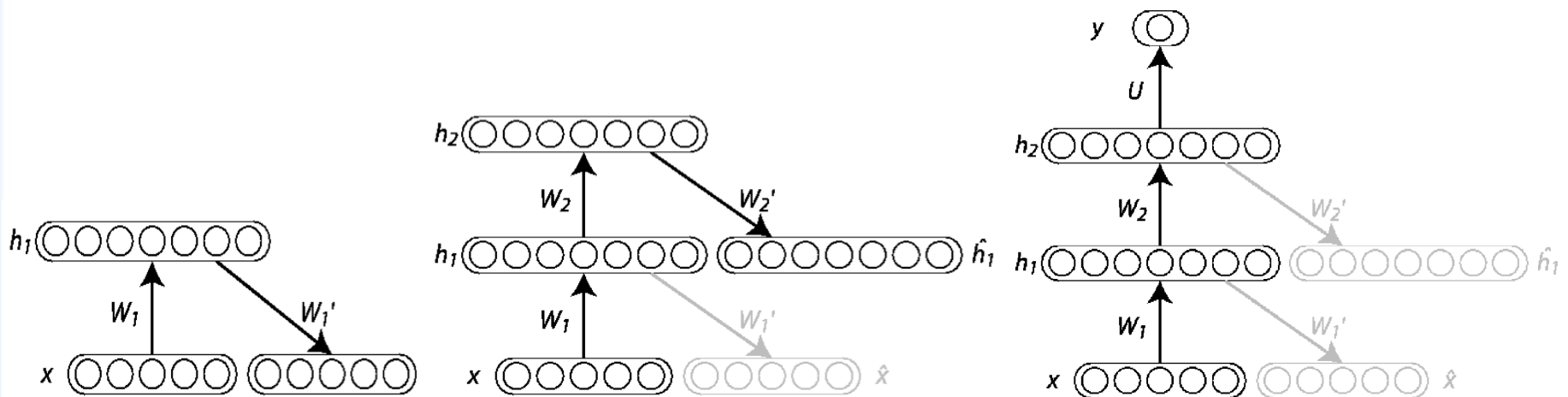
Sparse autoencoder illustration for images



$$[a_1, \dots, a_{64}] = [0, 0, \dots, 0, \mathbf{0.8}, 0, \dots, 0, \mathbf{0.3}, 0, \dots, 0, \mathbf{0.5}, 0]$$

Stacking Auto-Encoders

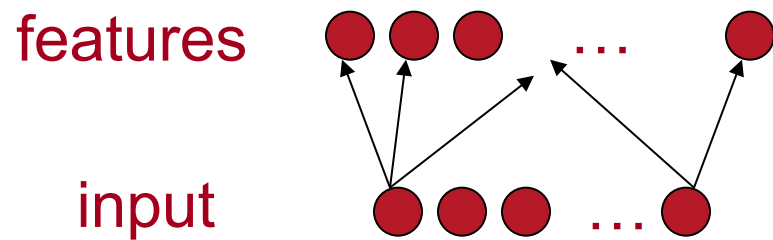
- Can be stacked successfully (Bengio et al NIPS'2006) to form highly non-linear representations



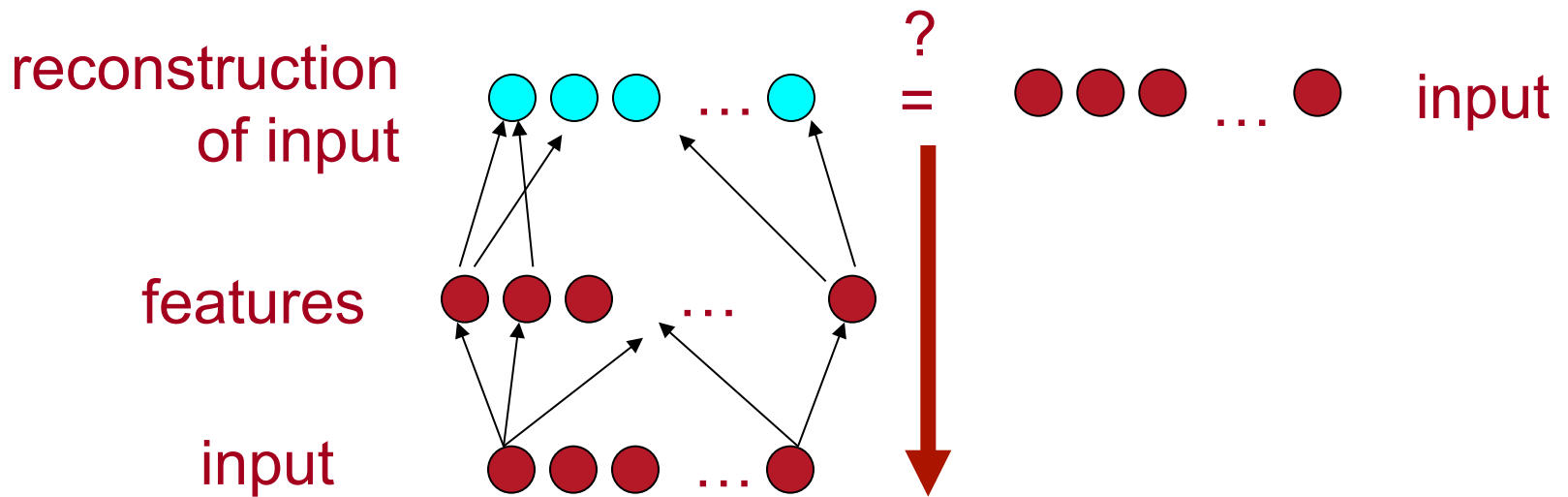
Layer-wise Unsupervised Learning

input ● ● ● ... ●

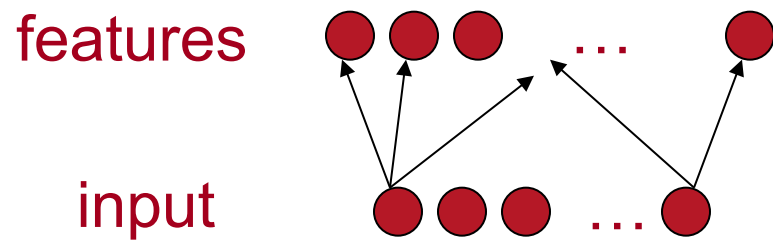
Layer-wise Unsupervised Pre-training



Layer-wise Unsupervised Pre-training



Layer-wise Unsupervised Pre-training

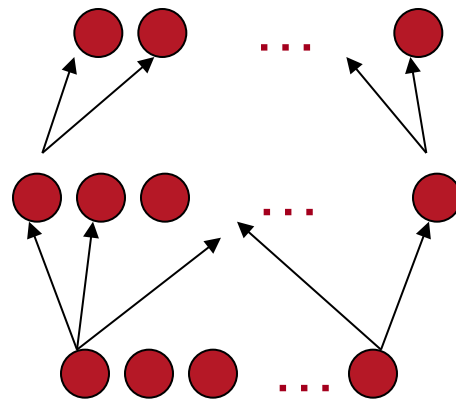


Layer-wise Unsupervised Pre-training

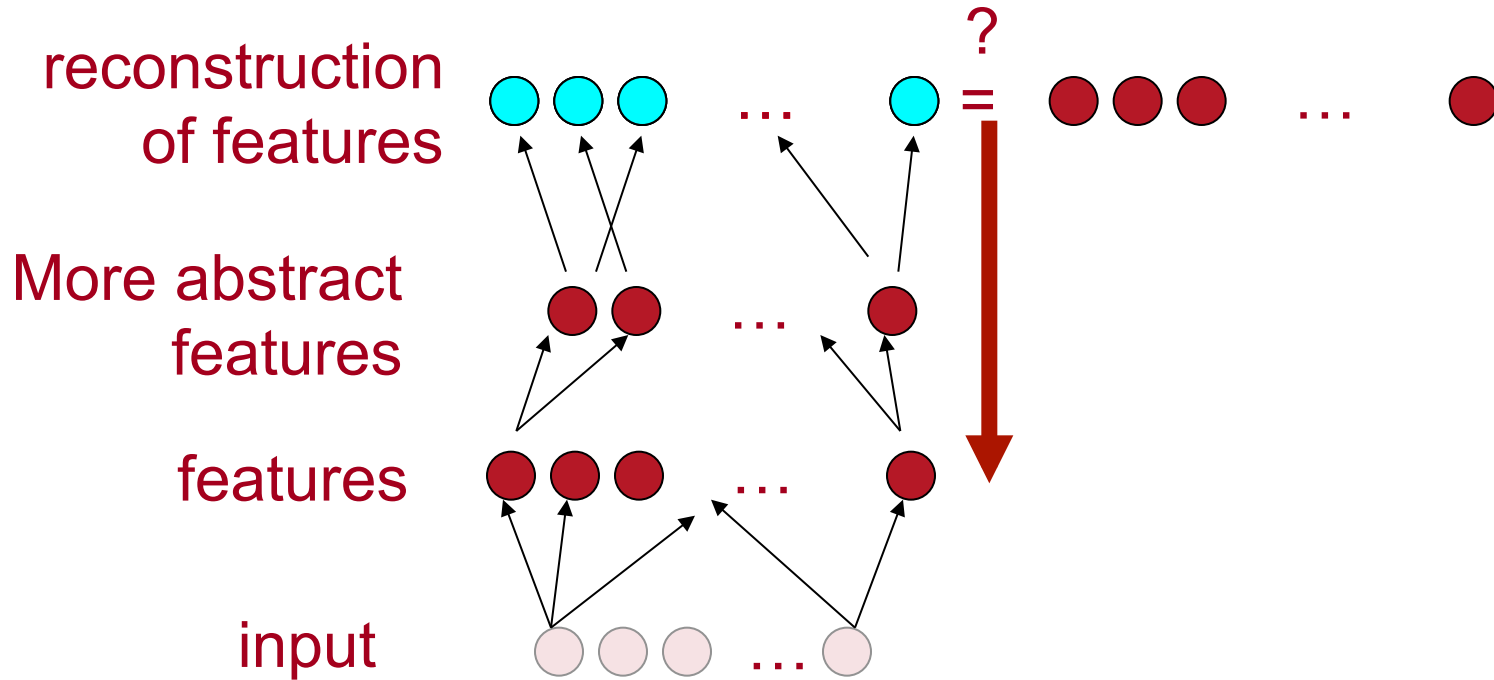
More abstract
features

features

input



Layer-wise Unsupervised Learning

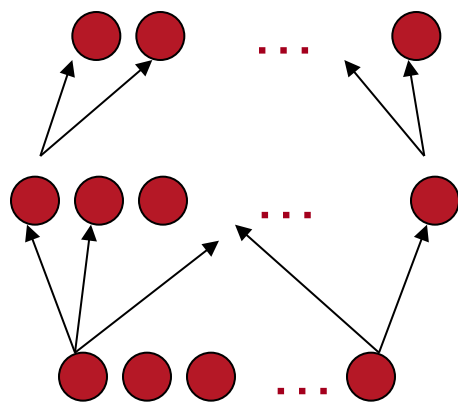


Layer-wise Unsupervised Pre-training

More abstract
features

features

input



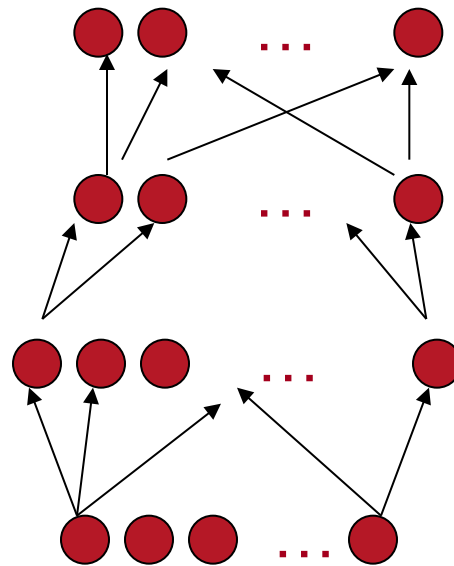
Layer-wise Unsupervised Learning

Even more abstract
features

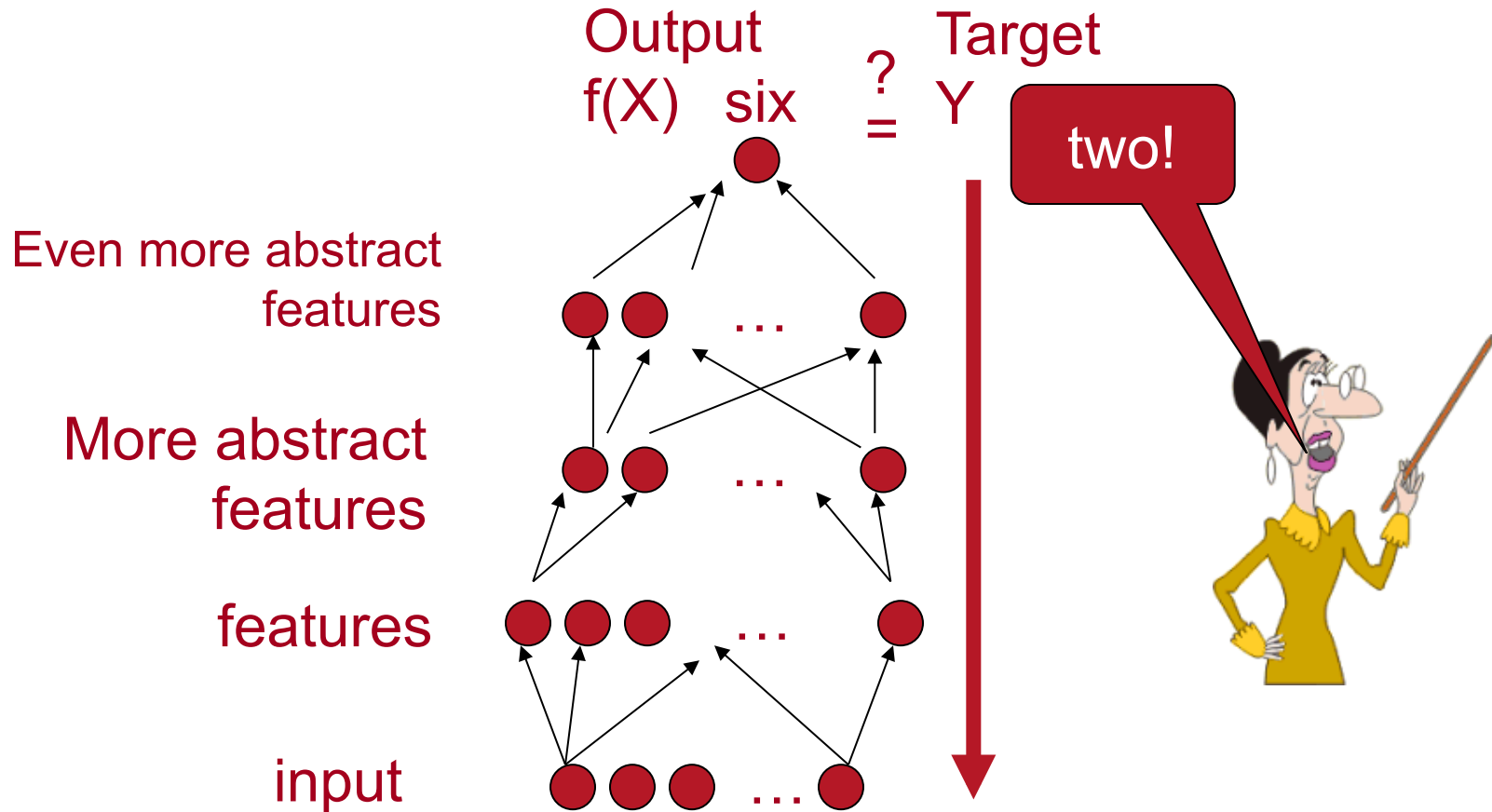
More abstract
features

features

input

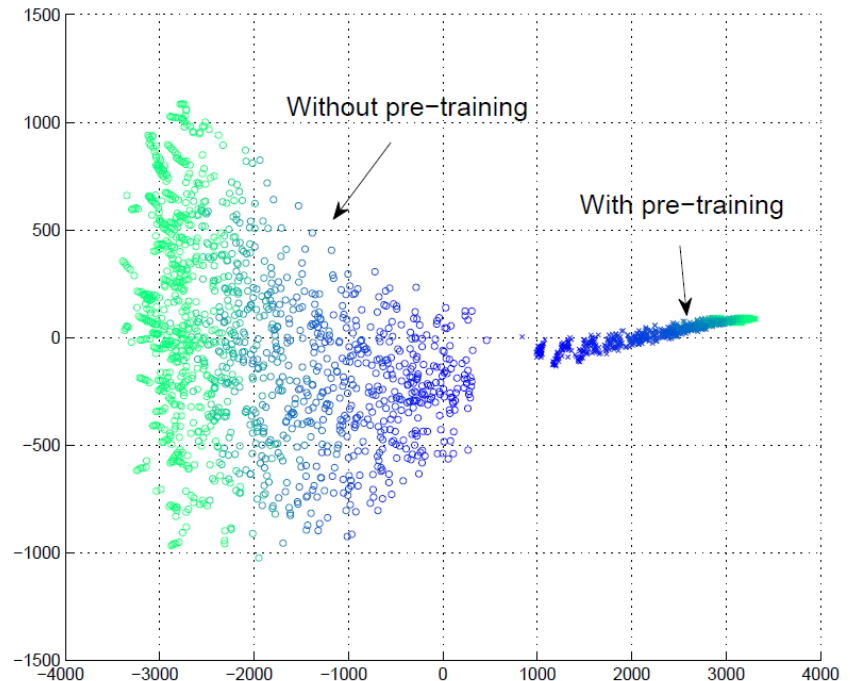


Supervised Fine-Tuning



Why is unsupervised pre-training working so well?

- Regularization hypothesis:
 - Representations good for $P(x)$ are good for $P(y|x)$
- Optimization hypothesis:
 - Unsupervised initializations start near better local minimum of supervised training error
 - Minima otherwise not achievable by random initialization



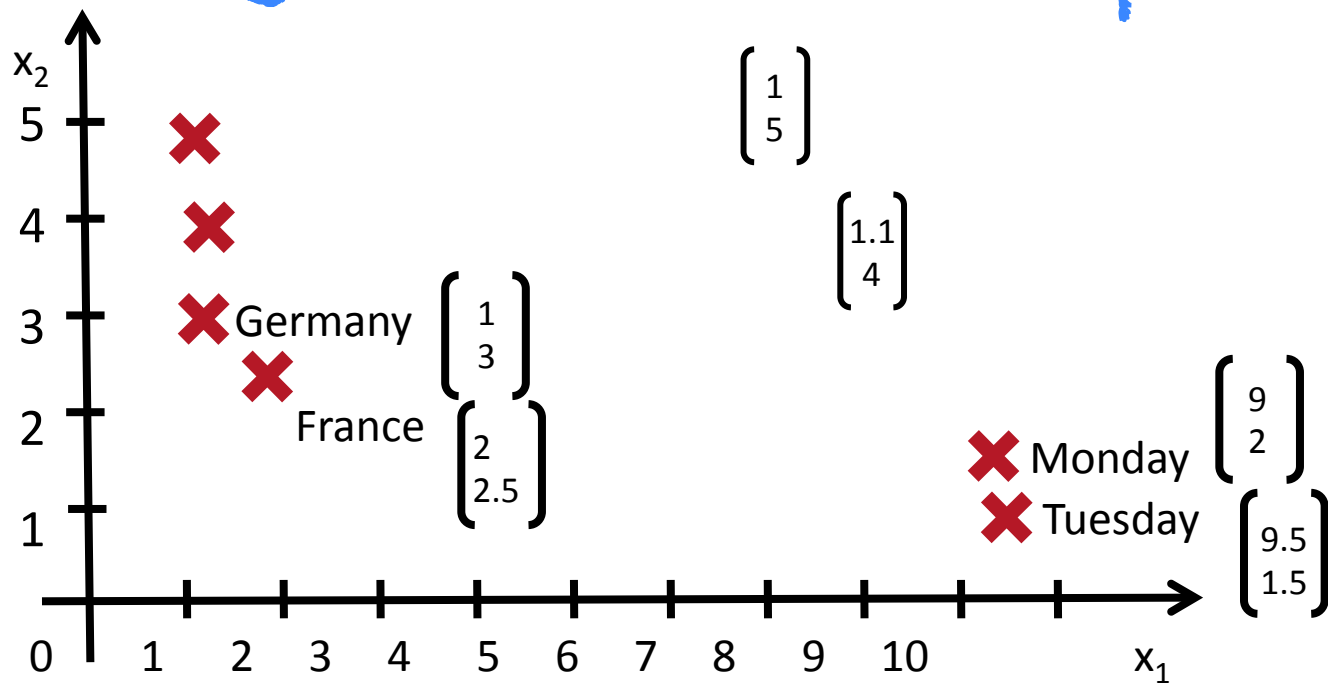
Erhan, Courville, Manzagol,
Vincent, Bengio (JMLR, 2010)



Part 2

Recursive Neural Networks

Building on Word Vector Space Models



the country of my birth
the place where I was born

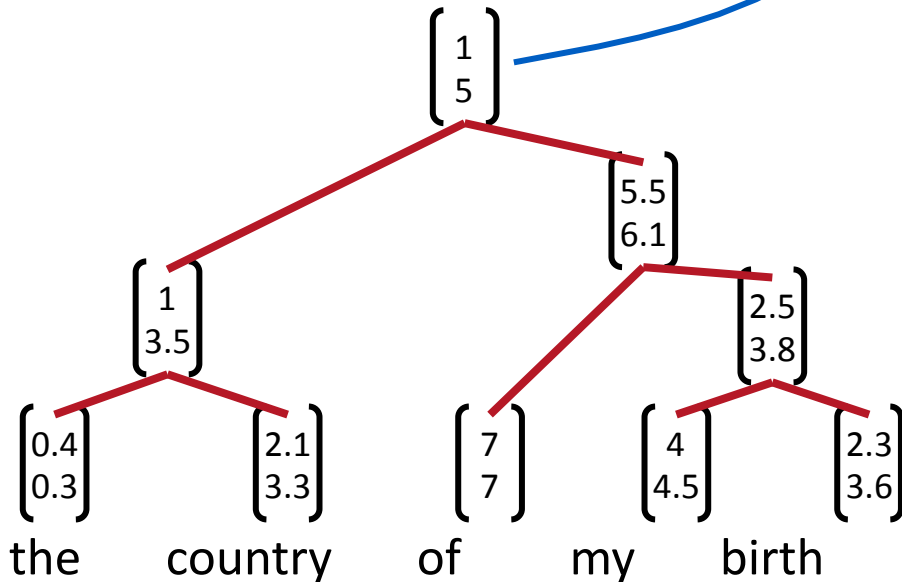
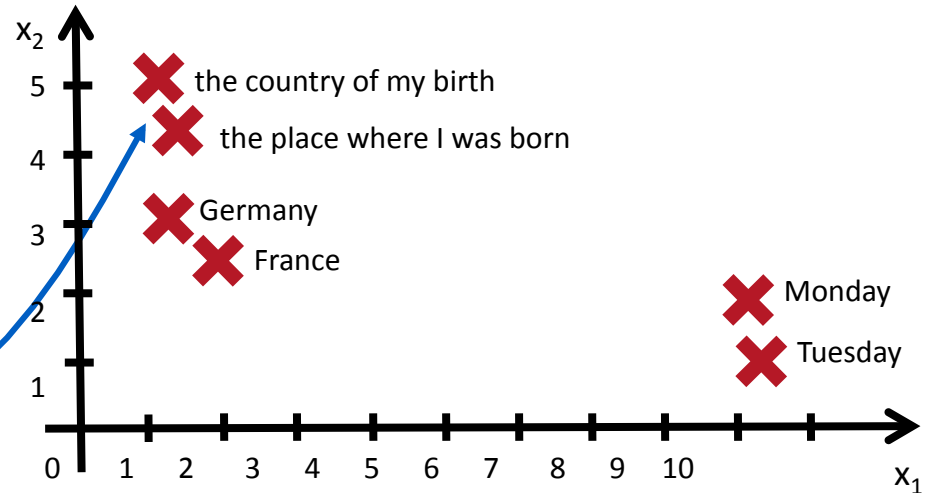
But how can we represent the meaning of longer phrases?

98 By mapping them into the same vector space!

How should we map phrases into a vector space?

Use principle of compositionality

The meaning (vector) of a sentence is determined by
(1) the meanings of its words and
(2) the rules that combine them.

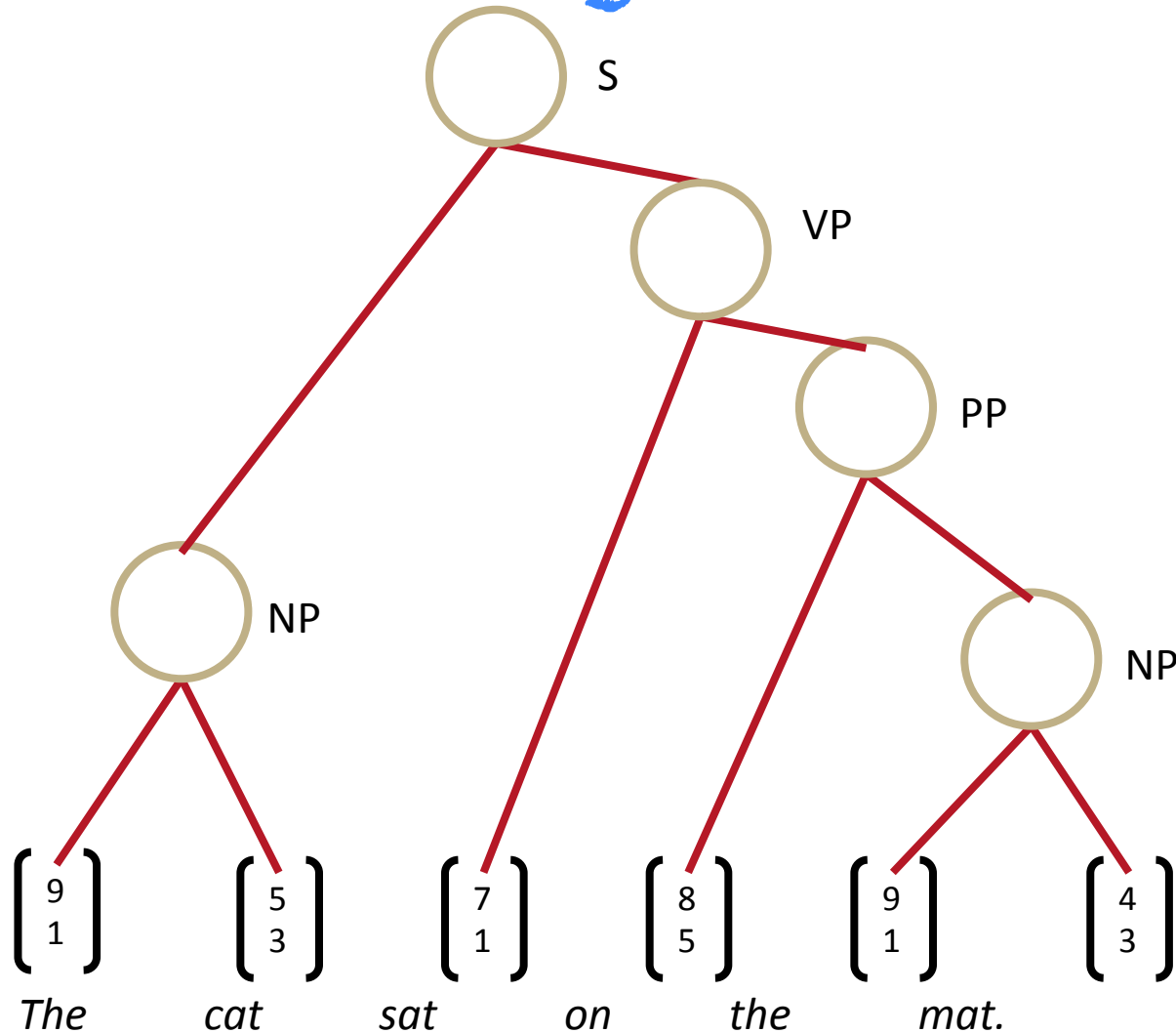


Recursive Neural Nets
can jointly learn
compositional vector
representations and
parse trees

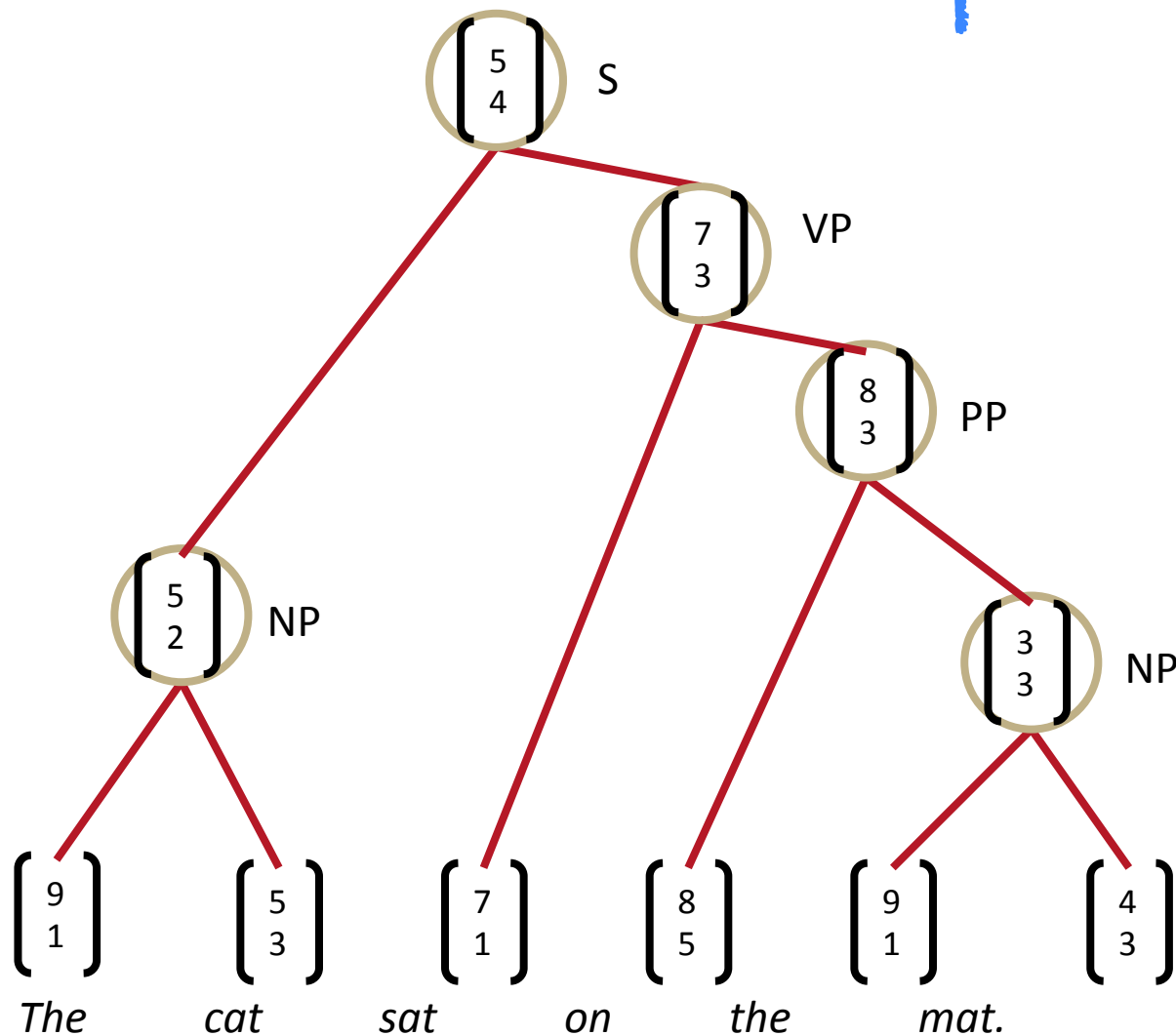
Recursive Neural Networks

1. Motivation
2. Recursive Neural Networks for Parsing
3. Theory: Backpropagation Through Structure
4. Recursive Autoencoders
5. Application to Sentiment Analysis and Paraphrase Detection
6. Compositionality Through Recursive Matrix-Vector Spaces
7. Relation classification

Sentence Parsing: What we want



Learn Structure and Representation

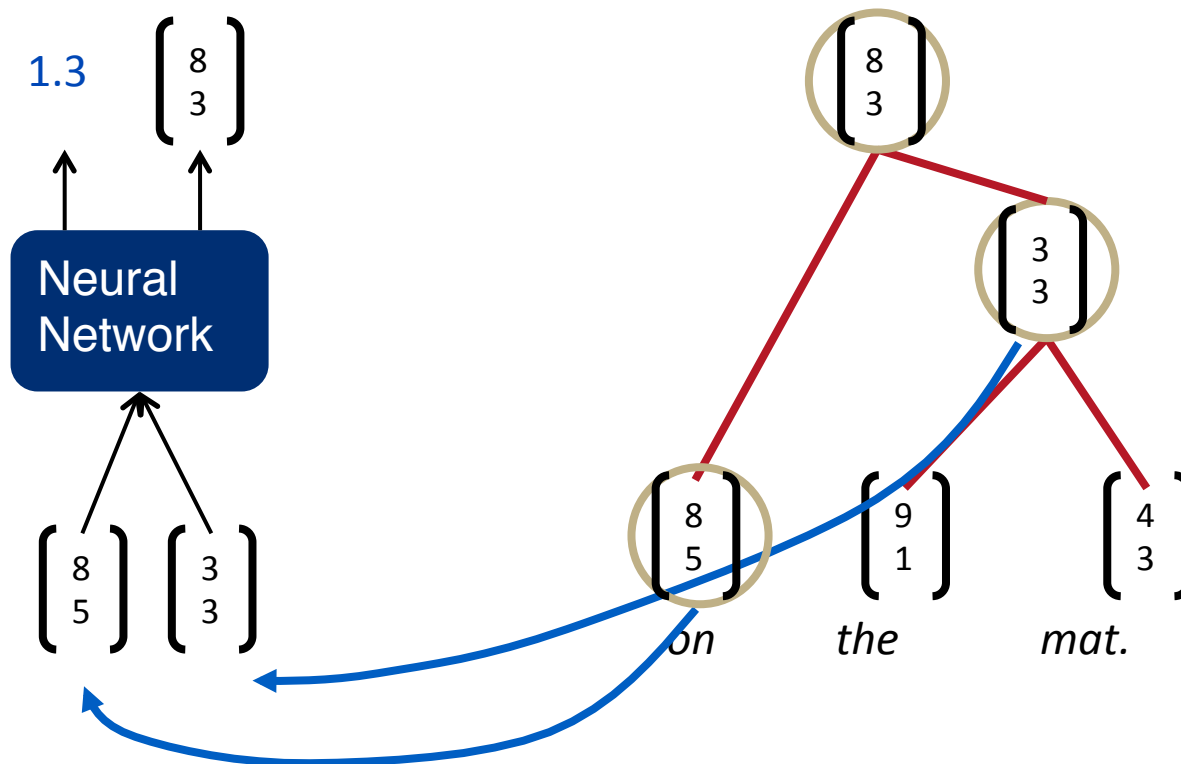


Recursive Neural Networks for Structure Prediction

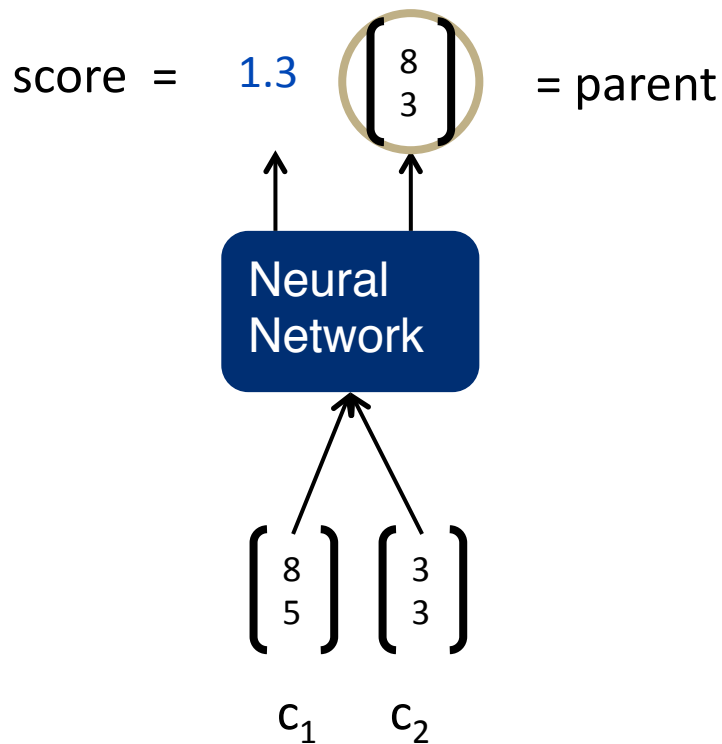
Inputs: two candidate children's representations

Outputs:

1. The semantic representation if the two nodes are merged.
2. Score of how plausible the new node would be.



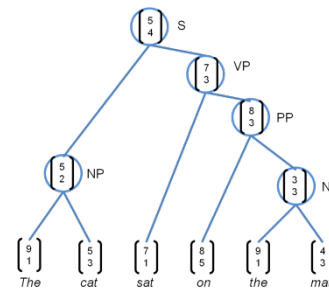
Recursive Neural Network Definition



$$\text{score} = U^T p$$

$$p = \tanh\left(W \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} + b\right),$$

Same W parameters at all nodes of the tree



Related Work to Socher et al. (ICML 2011)

- Pollack (1990): Recursive auto-associative memories



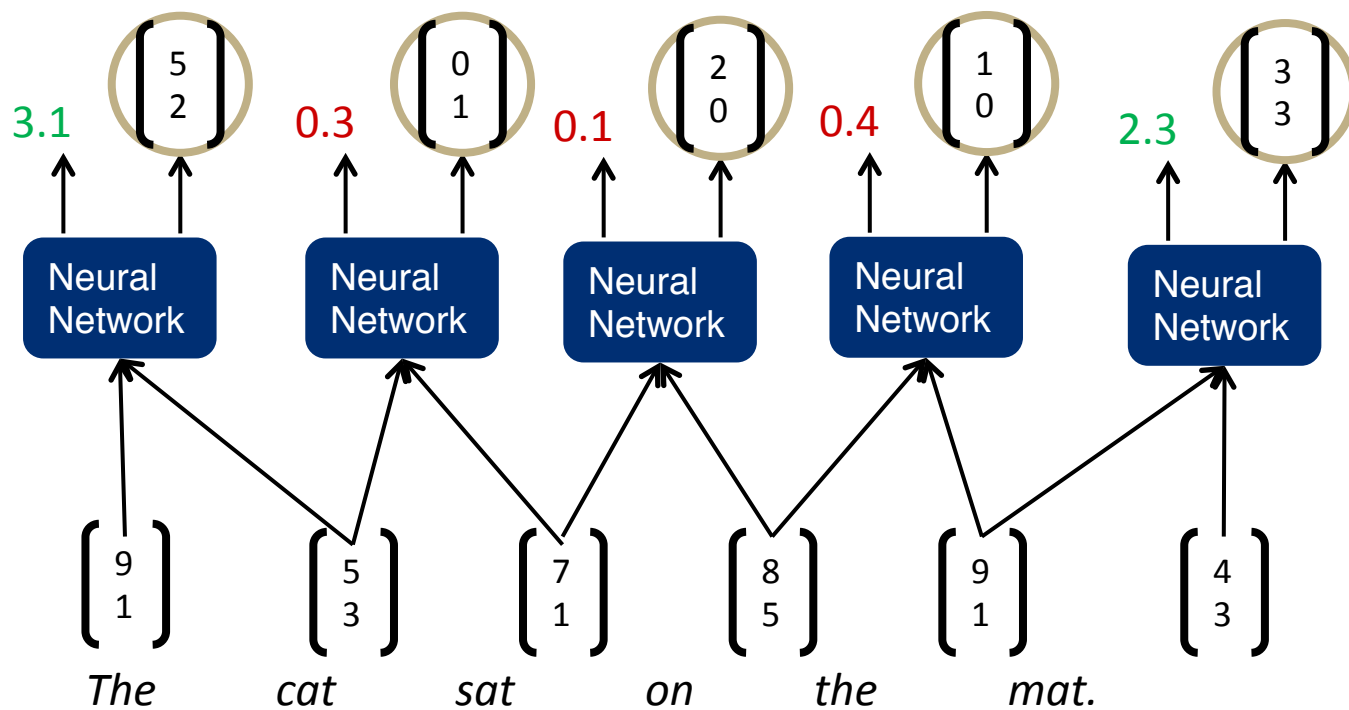
- Previous Recursive Neural Networks work by Goller & Küchler (1996), Costa et al. (2003) assumed fixed tree structure and used one hot vectors.



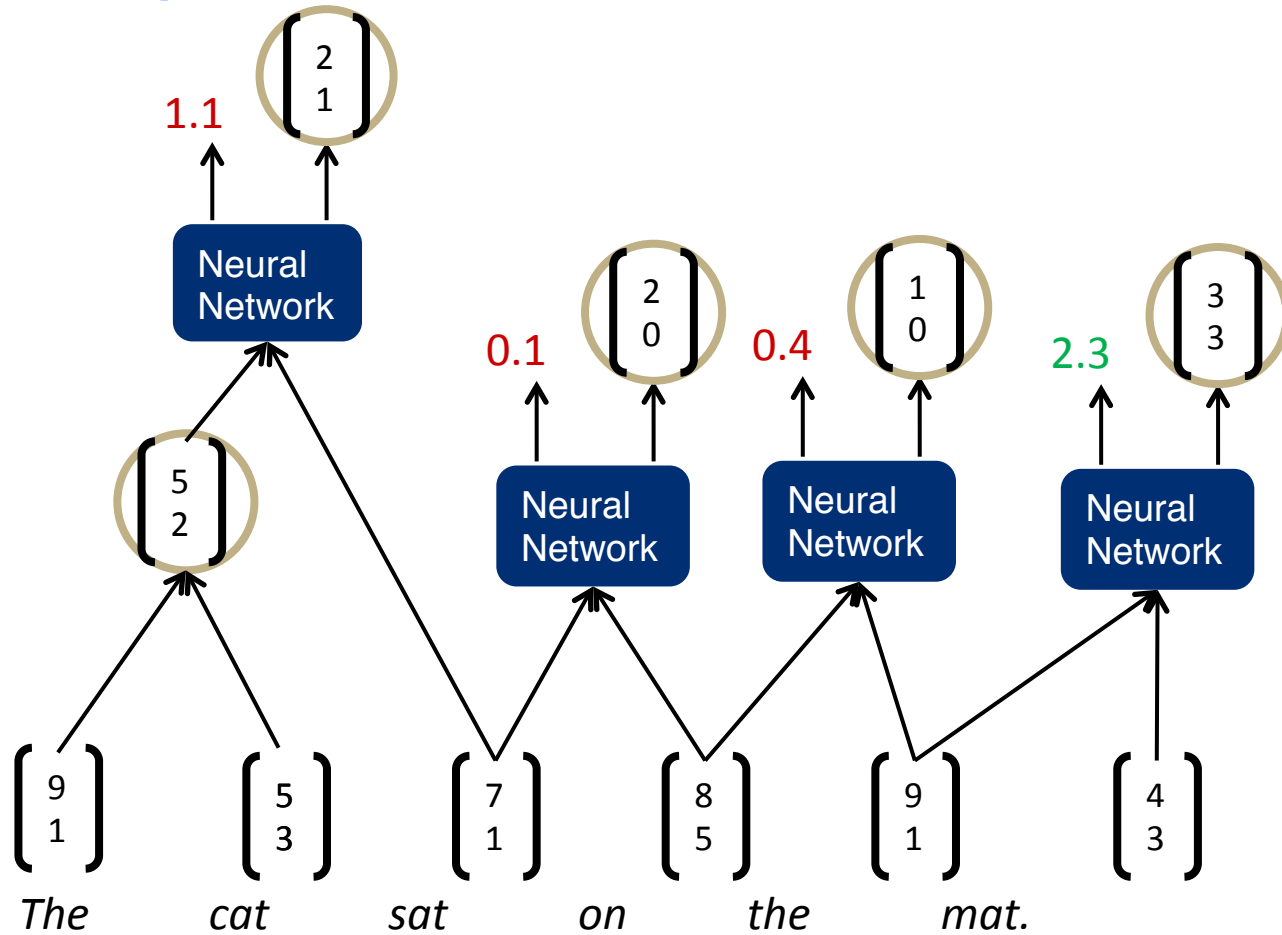
- Hinton (1990) and Bottou (2011): Related ideas about recursive models and recursive operators as smooth versions of logic operations



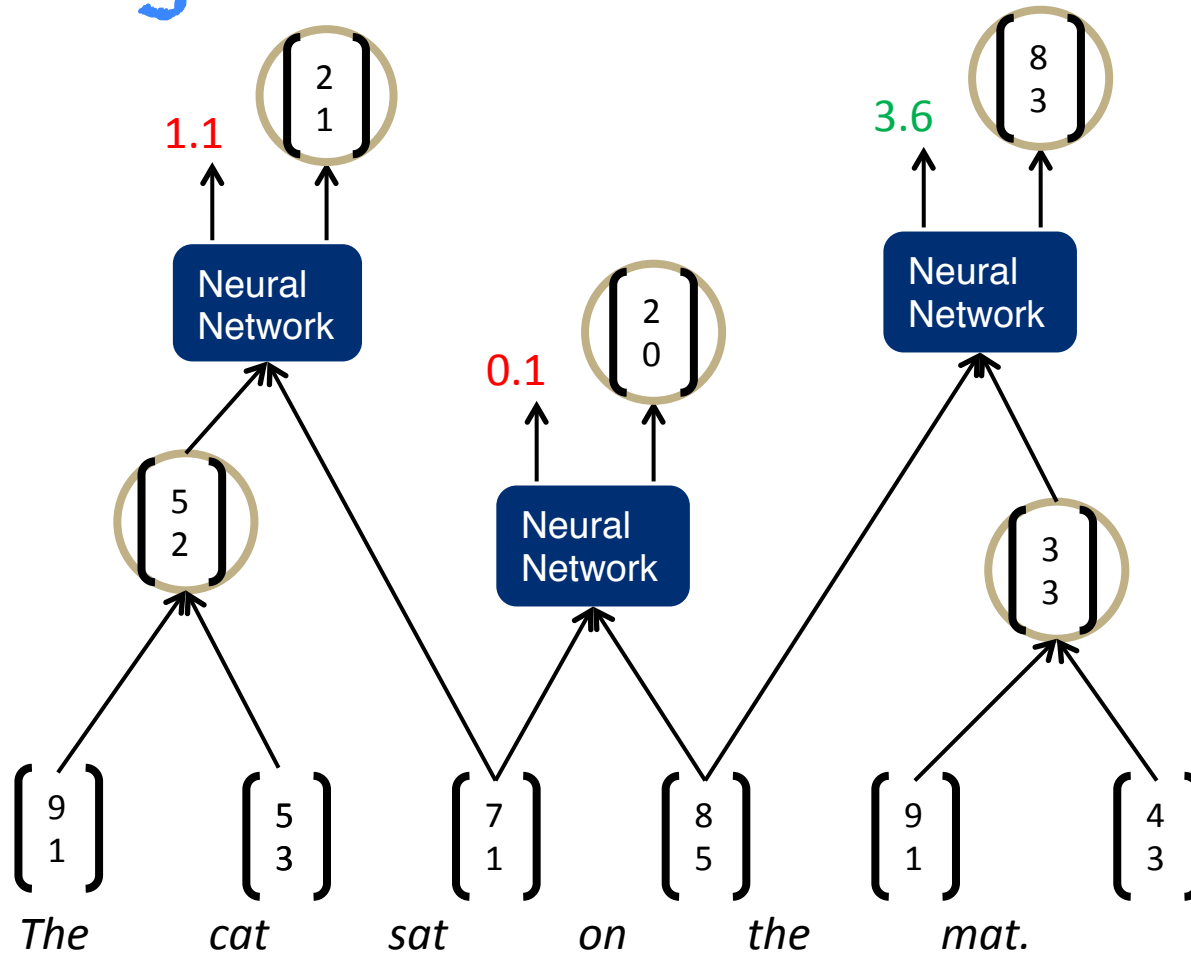
Parsing a sentence with an RNN



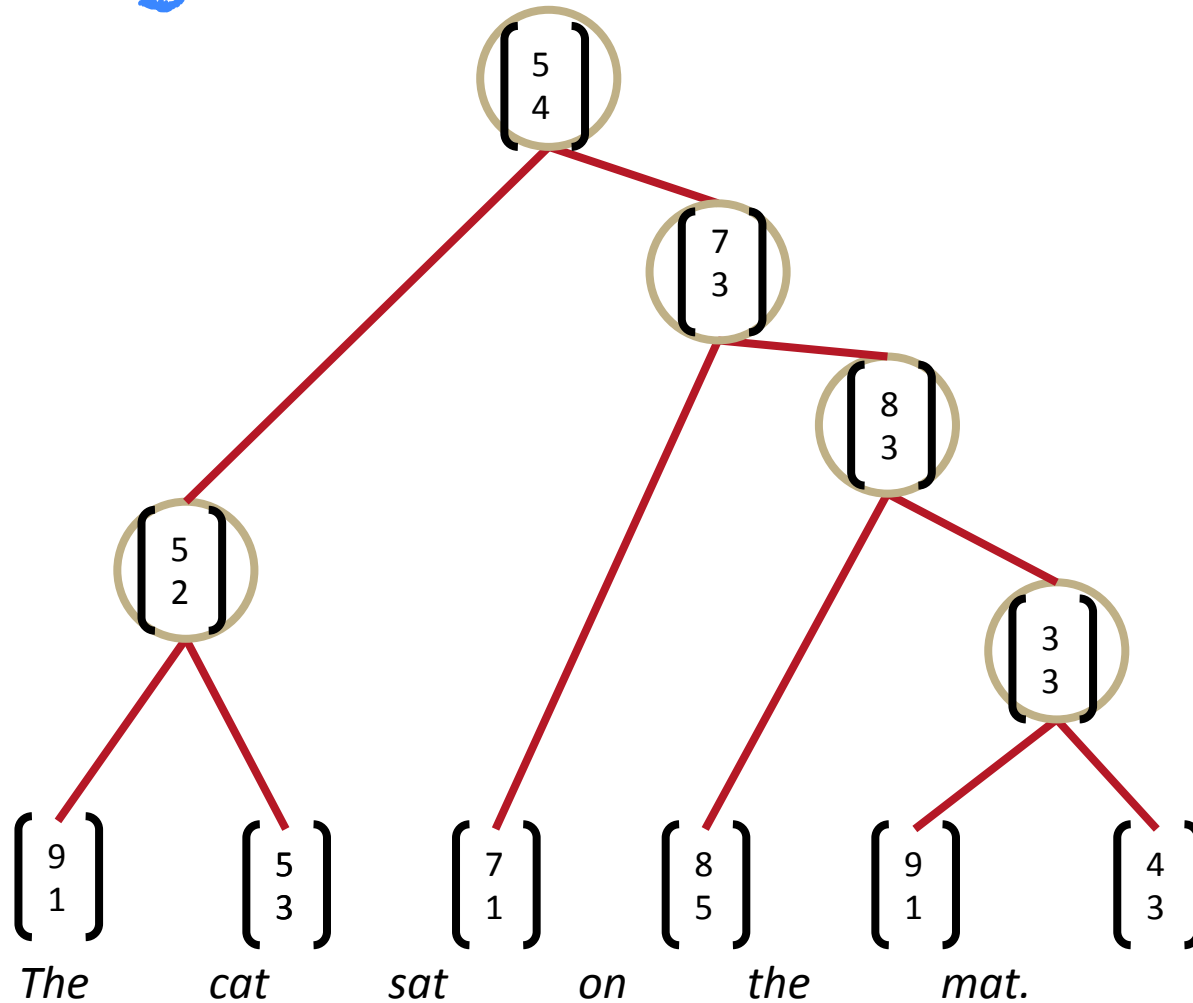
Parsing a sentence



Parsing a sentence



Parsing a sentence



Max-Margin Framework - Details

- The score of a tree is computed by the sum of the parsing decision scores at each node.



- Similar to max-margin parsing (Taskar et al. 2004), a supervised max-margin objective

$$J = \sum_i s(x_i, y_i) - \max_{y \in A(x_i)} (s(x_i, y) + \Delta(y, y_i))$$

- The loss $\Delta(y, y_i)$ penalizes all incorrect decisions

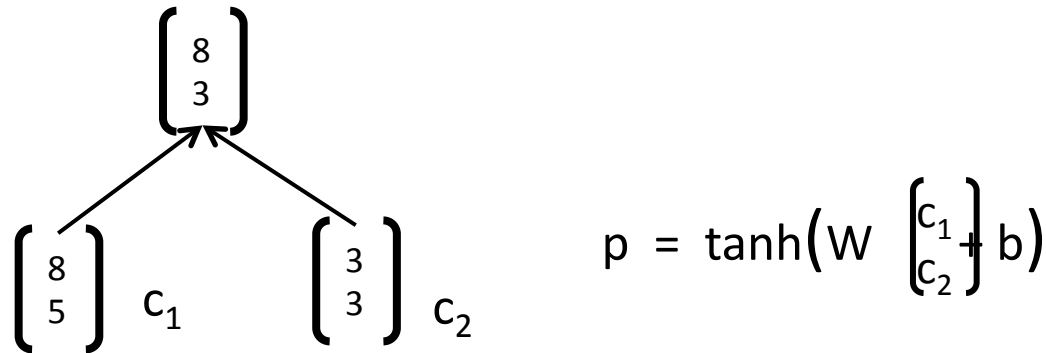
Backpropagation Through Structure (BTS)



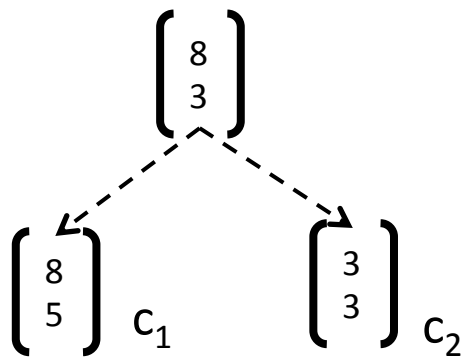
- Introduced by Goller & Küchler (1996)
- Principally the same as general backpropagation
- Two differences resulting from the tree structure:
- Split derivatives at each node
- Sum derivatives of W from all nodes

BTS: Split derivatives at each node

- During forward prop, the parent is computed using 2 children



- Hence, the errors need to be computed wrt each of them:



where each child's error is n-dimensional

BTS: Sum derivatives of all nodes

- You can actually assume it's a different W at each node
- Intuition via example:

$$\begin{aligned} & \frac{\partial}{\partial W} f(W(f(Wx))) \\ = & f'(W(f(Wx))) \left(\left(\frac{\partial}{\partial W} W \right) f(Wx) + W \frac{\partial}{\partial W} f(Wx) \right) \\ = & f'(W(f(Wx))) (f(Wx) + W f'(Wx)x) \end{aligned}$$

- If take separate derivatives of each occurrence, we get same:

$$\begin{aligned} & \frac{\partial}{\partial W_2} f(W_2(f(W_1x))) + \frac{\partial}{\partial W_1} f(W_2(f(W_1x))) \\ = & f'(W_2(f(W_1x))) (f(W_1x)) + f'(W_2(f(W_1x))) (W_2 f'(W_1x)x) \\ = & f'(W_2(f(W_1x))) (f(W_1x) + W_2 f'(W_1x)x) \\ \stackrel{=}{=} & f'(W(f(Wx))) (f(Wx) + W f'(Wx)x) \end{aligned}$$

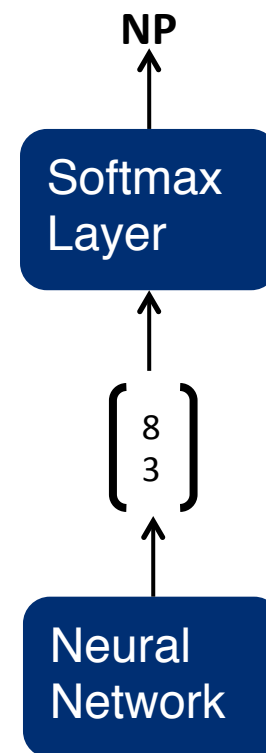
BTS: Optimization

- As before, we can plug the gradients into a standard off-the-shelf L-BFGS optimizer
- For non-continuous objective use subgradient *method* (Ratliff et al. 2007)

Labeling in Recursive Neural Networks

- We can use each node's representation as features for a *softmax* classifier:

$$p(c|p) = \textit{softmax}(Sp)$$



Experiments: Parsing Short Sentences

- Standard *WSJ* train/test
- Good results on short sentences
- More work is needed for longer sentences

Model	L15 Dev	L15 Test
Recursive Neural Network	92.1	90.3
Sigmoid NN (Titov & Henderson 2007)	89.5	89.3
Berkeley Parser (Petrov & Klein 2006)	92.1	91.6

All the figures are adjusted for seasonal variations

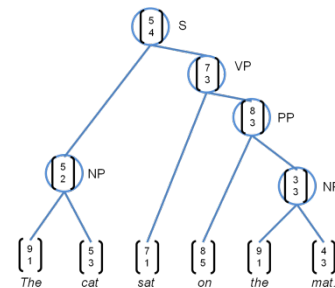
1. All the numbers are adjusted for seasonal fluctuations
2. All the figures are adjusted to remove usual seasonal patterns

Knight-Ridder wouldn't comment on the offer

1. Harsco declined to say what country placed the order
2. Coastal wouldn't disclose the terms

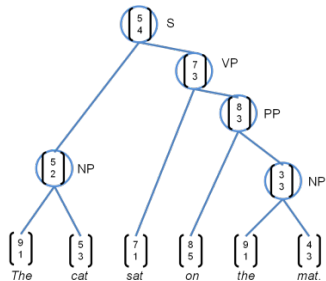
Sales grew almost 7% to \$UNK m. from \$UNK m.

1. Sales rose more than 7% to \$94.9 m. from \$88.3 m.
2. Sales surged 40% to UNK b. yen from UNK b.

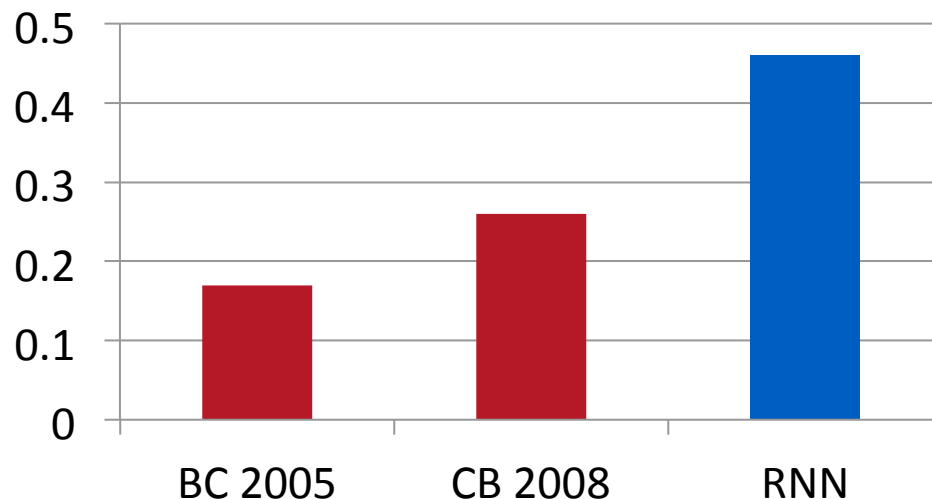


Short Paraphrase Detection

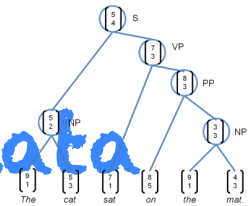
- Goal is to say which of candidate phrases are a good paraphrase of a given phrase
 - Motivated by Machine Translation
 - Initial algorithms: **Bannard & Callison-Burch 2005** (BC 2005), **Callison-Burch 2008** (CB 2008) exploit bilingual sentence-aligned corpora and hand-built linguistic constraints
 - We simply re-use our system learned on parsing the WSJ



F1 of Paraphrase Detection



Paraphrase detection task, CCB data

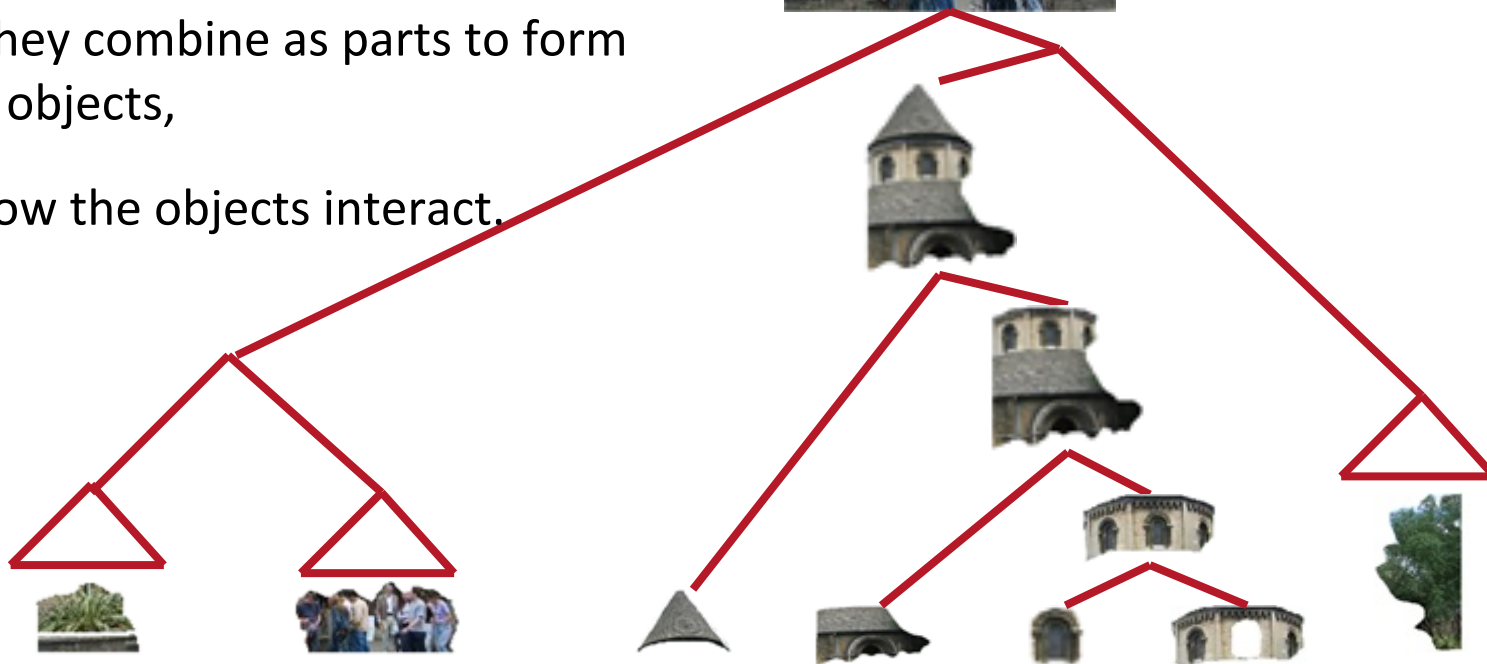


Target	Candidates with human goodness label (1–5) ordered by our system
the united states	the usa (5) the us (5) united states (5) north america (4) united (1) the (1) of the united states (3) america (5) nations (2) we (3)
around the world	around the globe(5) throughout the world(5) across the world(5) over the world(2) in the world(5) of the budget(2) of the world(5)
it would be	it would represent (5) there will be (2) that would be (3) it would be ideal (2) it would be appropriate (2) it is (3) it would (2)
of capital punishment	of the death penalty (5) to death (2) the death penalty (2) of (1)
in the long run	in the long term (5) in the short term (2) for the longer term (5) in the future (5) in the end (3) in the long-term (5) in time (5) of the (1)

Scene Parsing

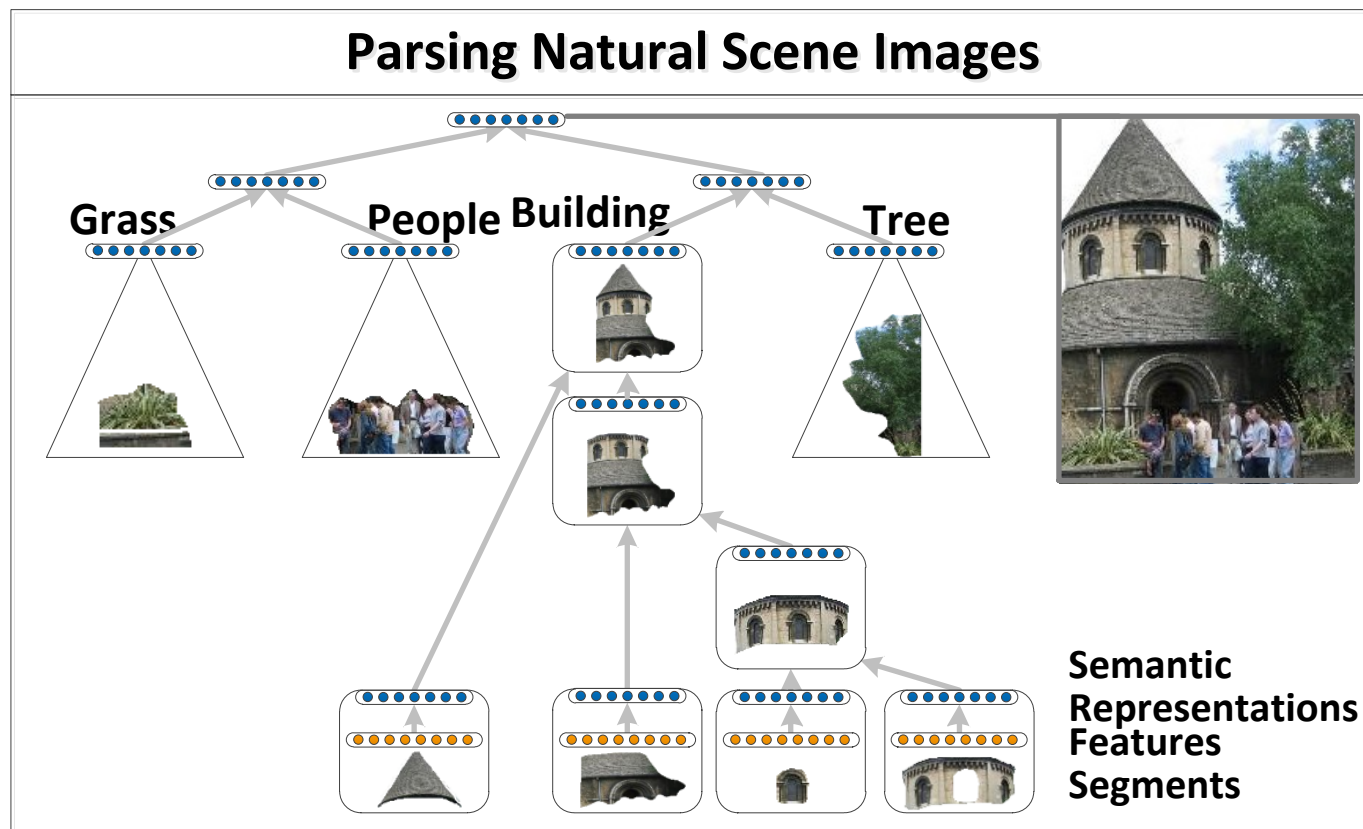
Similar principle of compositionality.

- The meaning of a scene image is also a function of smaller regions,
- how they combine as parts to form larger objects,
- and how the objects interact.



Algorithm for Parsing Images

Same Recursive Neural Network as for natural language parsing!
(Socher et al. ICML 2011)



Multi-class segmentation



■ sky ■ tree ■ road ■ grass ■ water ■ bldg ■ mntn ■ fg obj.

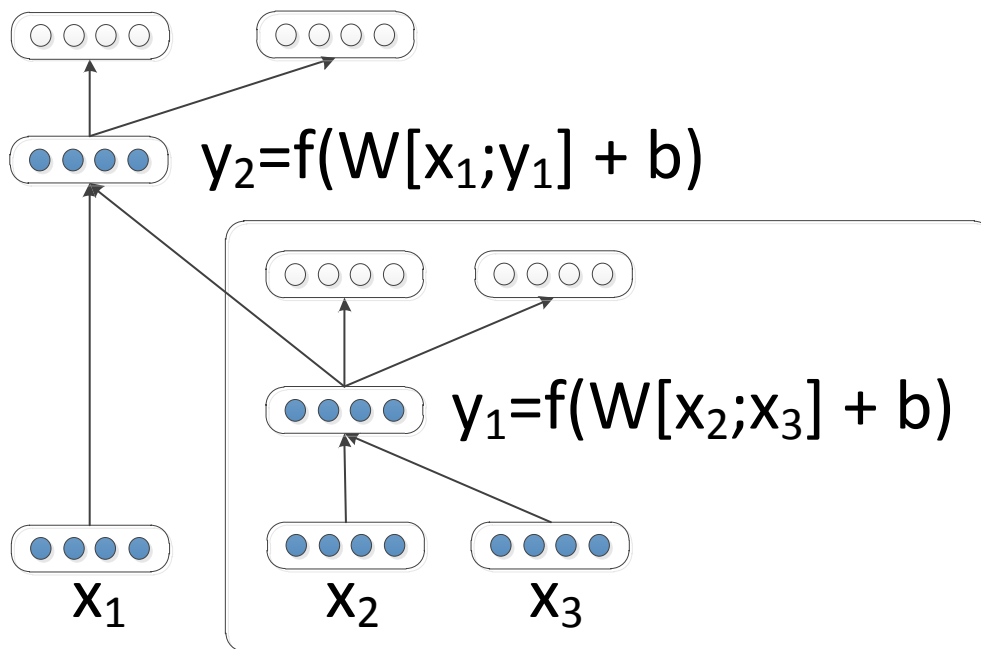
Method	Accuracy
Pixel CRF (Gould et al., ICCV 2009)	74.3
Classifier on superpixel features	75.9
Region-based energy (Gould et al., ICCV 2009)	76.4
Local labelling (Tighe & Lazebnik, ECCV 2010)	76.9
Superpixel MRF (Tighe & Lazebnik, ECCV 2010)	77.5
Simultaneous MRF (Tighe & Lazebnik, ECCV 2010)	77.5
Recursive Neural Network	78.1

Recursive Neural Networks

1. Motivation
2. Recursive Neural Networks for Parsing
3. Theory: Backpropagation Through Structure
4. Recursive Autoencoders
5. Application to Sentiment Analysis and Paraphrase Detection
6. Compositionality Through Recursive Matrix-Vector Spaces
7. Relation classification

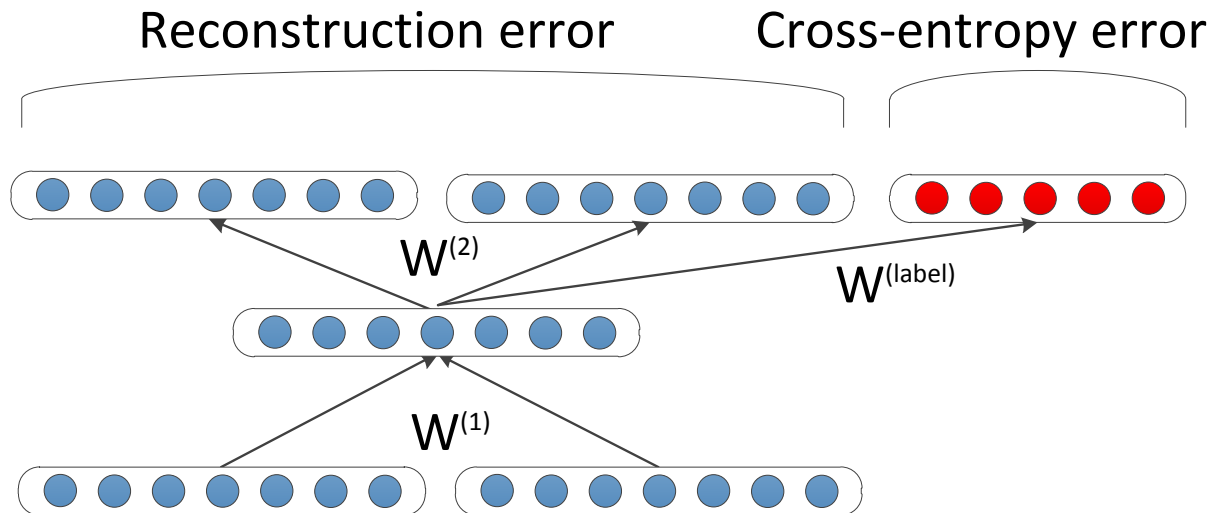
Recursive Autoencoders

- Similar to Recursive Neural Net but instead of a supervised score we compute a reconstruction error at each node. $E_{rec}([c_1; c_2]) = \frac{1}{2} ||[c_1; c_2] - [c'_1; c'_2]||^2$



Semi-supervised Recursive Autoencoder

- To capture sentiment and solve antonym problem, add a softmax classifier
- Error is a weighted combination of reconstruction error and cross-entropy
- Socher et al. (EMNLP 2011)



Sentiment Detection

- Sentiment detection is crucial to business intelligence, stock trading, ...



Maybe she'll change her name to Halliburton. Just to see.

3/18/11 at 4:00 PM | 17 Comments
Mentions of the Name 'Anne Hathaway' May Drive Berkshire Hathaway Stock

By Patrick Huguenin



The Huffington Post recently [pointed out](#) that whenever Anne Hathaway is in the news, the stock price for Warren Buffett's Berkshire Hathaway goes up. Really. When *Bride Wars* opened, the stock rose 2.61 percent. (*Rachel*

Getting Married only kicked it up 0.44 percent, but, you know, that one was so light on plot compared to *Bride Wars*.)

Sentiment Detection and Bag-of-Words Models

- Most methods start with a bag of words + linguistic features/processing/lexica
- But such methods (including tf-idf) can't distinguish:
 - + white blood cells destroying an infection
 - an infection destroying white blood cells

Single Scale Experiments: Movies

Stealing Harvard doesn't care about cleverness, wit or any other kind of intelligent humor.

a film of ideas and wry comic mayhem.

Accuracy of Positive/Negative Sentiment Classification

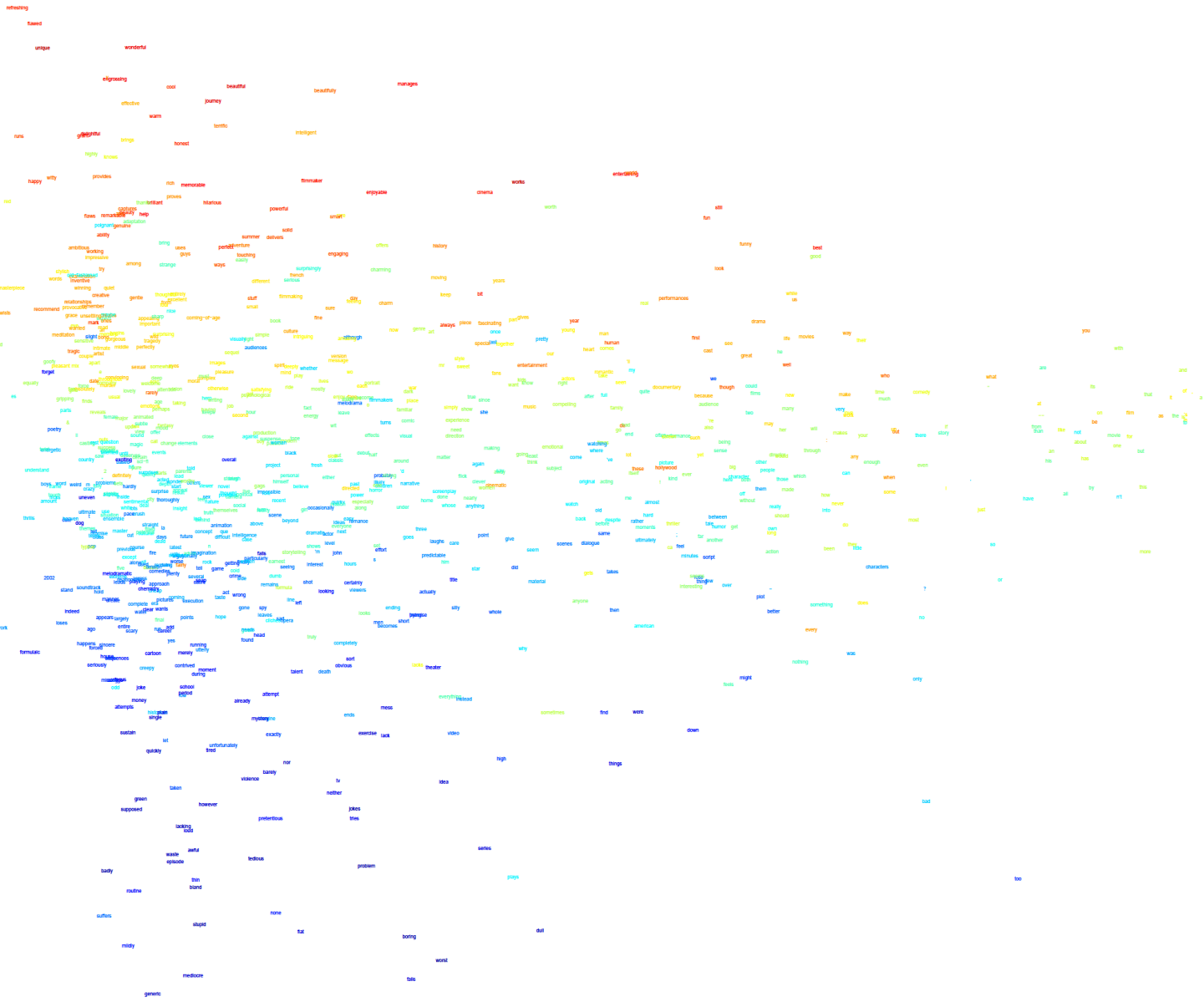
- Results on movie reviews (MR) and opinions (MPQA).
- All other methods use hand-designed polarity shifting rules or sentiment lexica.
- RAE: no hand-designed features, learns vector

Method	MR	MPQA
Phrase voting with lexicons	63.1	81.7
Bag of features with lexicons	76.4	84.1
Tree-CRF (Nakagawa et al. 2010)	77.3	86.1
RAE (this work)	77.7	86.4



Sorted Negative and Positive N-grams

Most Negative N-grams	Most Positive N-grams
bad; boring; dull; flat; pointless	touching; enjoyable; powerful
that bad; abysmally pathetic	the beautiful; with dazzling
is more boring; manipulative and contrived	funny and touching; a small gem
boring than anything else.; a major waste ... generic	cute, funny, heartwarming; with wry humor and genuine
loud, silly, stupid and pointless. ; dull, dumb and derivative horror film.	, deeply absorbing piece that works as a; ... one of the most ingenious and entertaining;

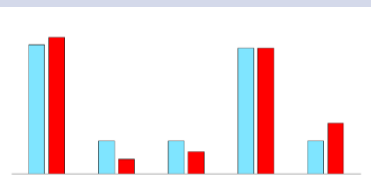
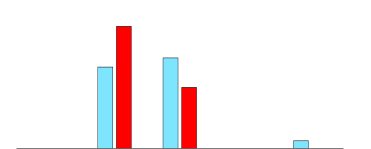
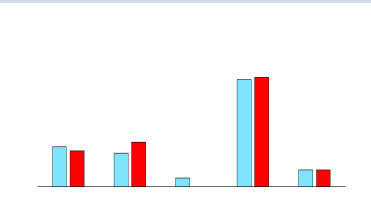


Sentiment Distribution Experiments

- Learn distributions over multiple complex sentiments → New dataset and task
- Experience Project
 - <http://www.experienceproject.com>
 - “I walked into a parked car”
 - Sorry, Hugs; You rock; Tee-hee ; I understand; Wow just wow
 - Over 31,000 entries with 113 words on average

Sentiment distributions

- Sorry, Hugs; You rock; Tee-hee ; I understand; Wow just wow

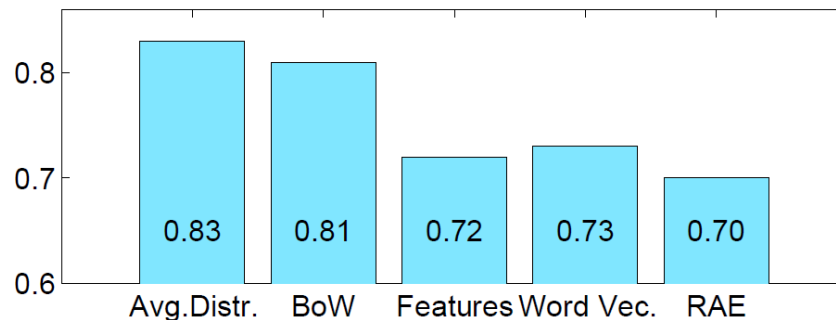
Predicted and Gold Distribution	Anonymous Confession
	i am a very succesfull business man. i make good money but i have been addicted to crack for 13 years. i moved 1 hour away from my dealers 10 years ago to stop using now i dont use daily but ...
	well i think hairy women are attractive
	Dear Love, I just want to say that I am looking for you. Tonight I felt the urge to write, and I am becoming more and more frustrated that I have not found you yet. I'm also tired of spending so much heart on an old dream. ...

Experience Project most votes results

Method	Accuracy %
Random	20
Most frequent class	38
Bag of words; MaxEnt classifier	46
Spellchecker, sentiment lexica, SVM	47
SVM on neural net word features	46
RAE (this work)	50

Average KL between gold and predicted label distributions:

133



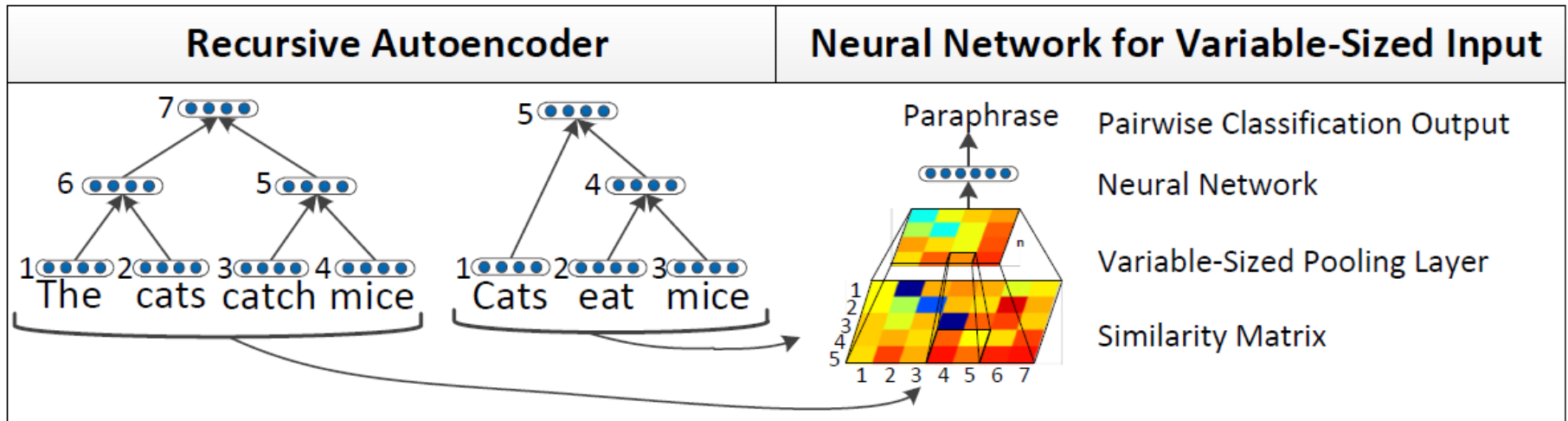
Paraphrase Detection

- Pollack said the plaintiffs failed to show that Merrill and Blodget directly caused their losses
- Basically , the plaintiffs did not show that omissions in Merrill's research caused the claimed losses
- The initial report was made to Modesto Police December 28
- It stems from a Modesto police report

How to compare the
meaning of two
sentences?

Recursive Autoencoders for Full Sentence Paraphrase Detection

- Unsupervised Unfolding RAE and a pair-wise sentence comparison of nodes in parsed trees
- Socher et al. (NIPS 2011)



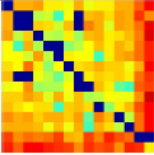
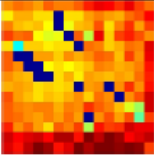
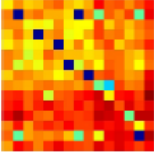
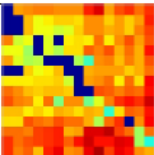
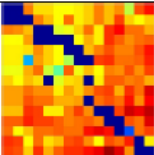
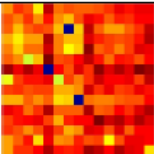
Recursive Autoencoders for Full Sentence Paraphrase Detection

- Experiments on Microsoft Research Paraphrase Corpus
- (Dolan et al. 2004)

Method	Acc.	F1
Rus et al.(2008)	70.6	80.5
Mihalcea et al.(2006)	70.3	81.3
Islam et al.(2007)	72.6	81.3
Qiu et al.(2006)	72.0	81.6
Fernando et al.(2008)	74.1	82.4
Wan et al.(2006)	75.6	83.0
Das and Smith (2009)	73.9	82.3
Das and Smith (2009) + 18 Surface Features	76.1	82.7
F. Bu et al. (ACL 2012): String Re-writing Kernel	76.3	--
Unfolding Recursive Autoencoder (NIPS 2011)	76.8	83.6



Recursive Autoencoders for Full Sentence Paraphrase Detection

L	Pr	Sentences	Sim.Mat.
P	0.95	(1) LLEYTON Hewitt yesterday traded his tennis racquet for his first sporting passion - Australian football - as the world champion relaxed before his Wimbledon title defence (2) LLEYTON Hewitt yesterday traded his tennis racquet for his first sporting passion-Australian rules football-as the world champion relaxed ahead of his Wimbledon defence	
P	0.82	(1) The lies and deceptions from Saddam have been well documented over 12 years (2) It has been well documented over 12 years of lies and deception from Saddam	
P	0.67	(1) Pollack said the plaintiffs failed to show that Merrill and Blodget directly caused their losses (2) Basically , the plaintiffs did not show that omissions in Merrill's research caused the claimed losses	
N	0.49	(1) Prof Sally Baldwin, 63, from York, fell into a cavity which opened up when the structure collapsed at Tiburtina station, Italian railway officials said (2) Sally Baldwin, from York, was killed instantly when a walkway collapsed and she fell into the machinery at Tiburtina station	
N	0.44	(1) Bremer, 61, is a onetime assistant to former Secretaries of State William P. Rogers and Henry Kissinger and was ambassador-at-large for counterterrorism from 1986 to 1989 (2) Bremer, 61, is a former assistant to former Secretaries of State William P. Rogers and Henry Kissinger	
N	0.11	(1) The initial report was made to Modesto Police December 28 (2) It stems from a Modesto police report	

Recursive Neural Networks

1. Motivation
2. Recursive Neural Networks for Parsing
3. Theory: Backpropagation Through Structure
4. Recursive Autoencoders
5. Application to Sentiment Analysis and Paraphrase Detection
6. Compositionality Through Recursive Matrix-Vector Spaces
7. Relation classification

Compositionality Through Recursive Matrix-Vector Spaces

$$p = \tanh\left(W \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} b\right)$$

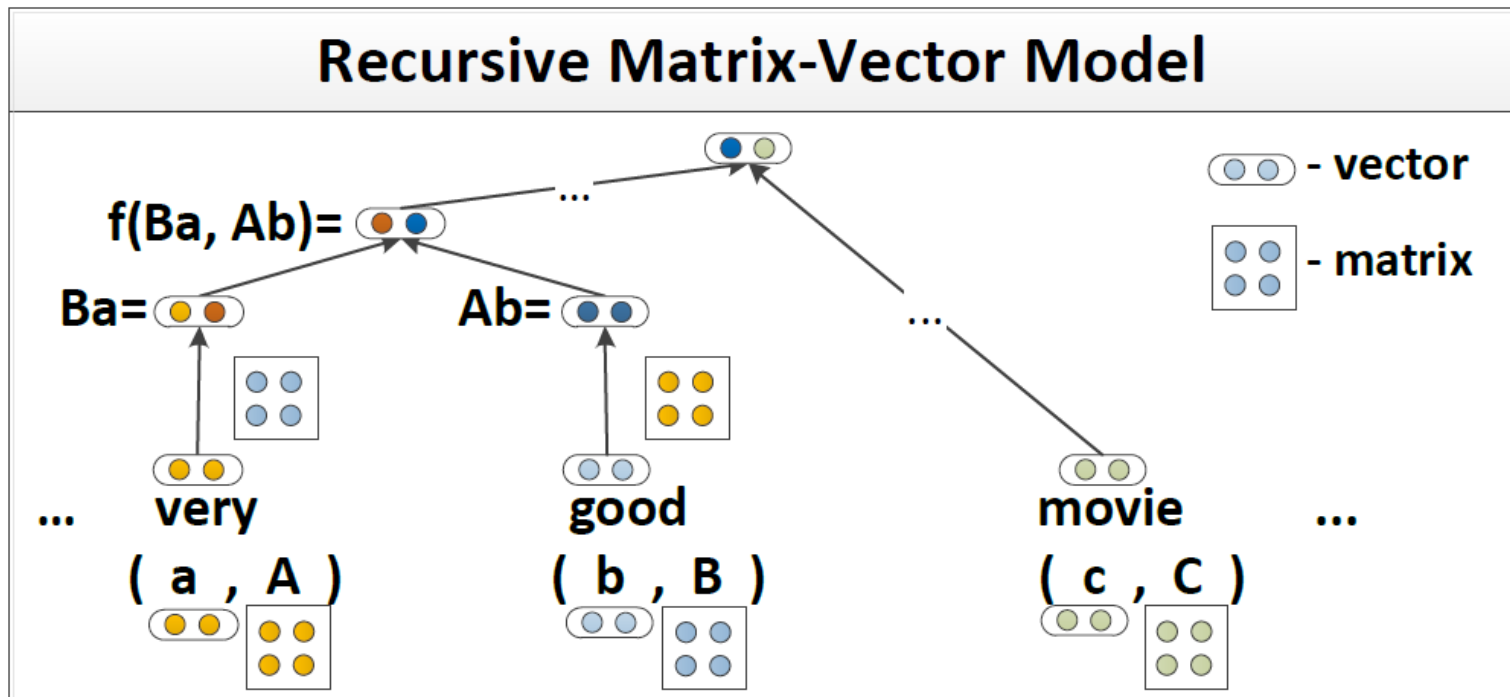
- But what if words act mostly as an operator, e.g. “very” in

very good

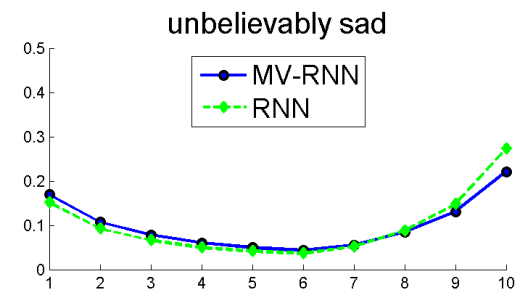
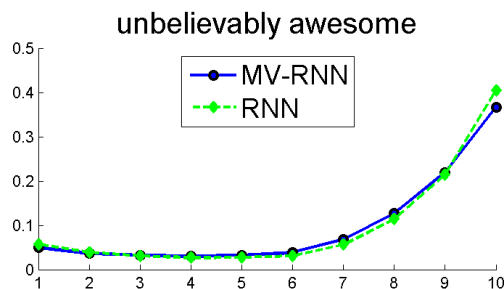
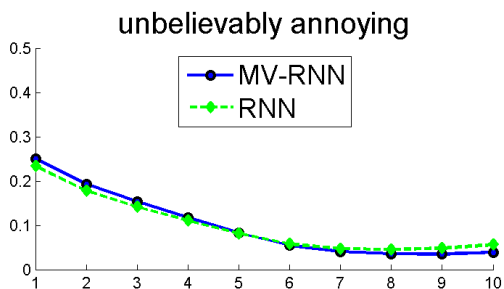
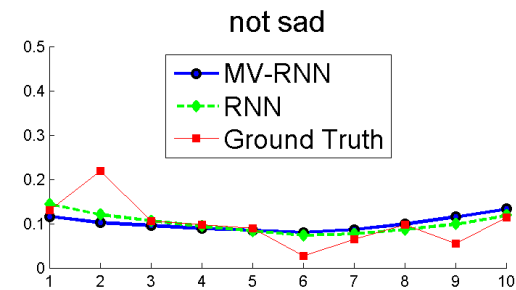
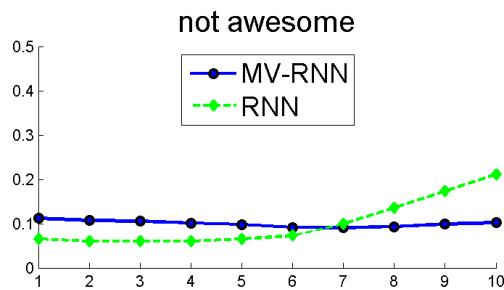
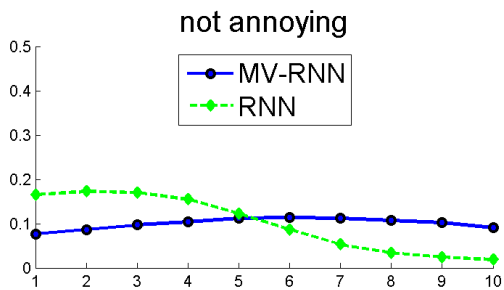
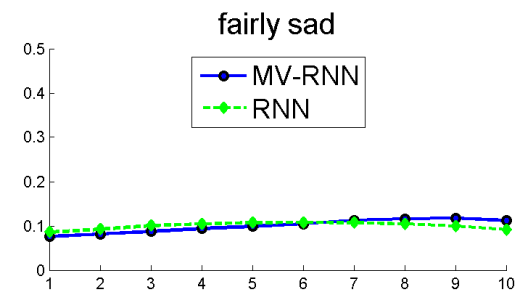
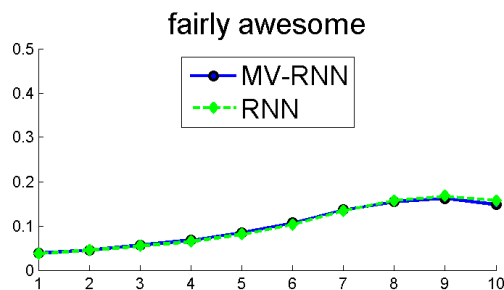
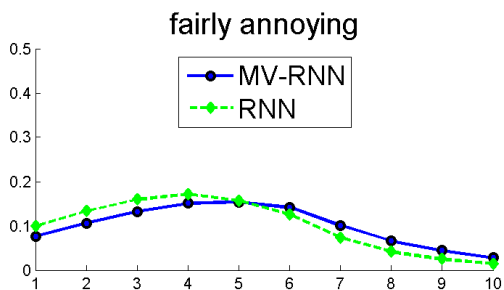
Compositionality Through Recursive Matrix-Vector Recursive Neural Networks

$$p = \tanh\left(W \begin{pmatrix} c_1 \\ + \\ c_2 \end{pmatrix} b\right)$$

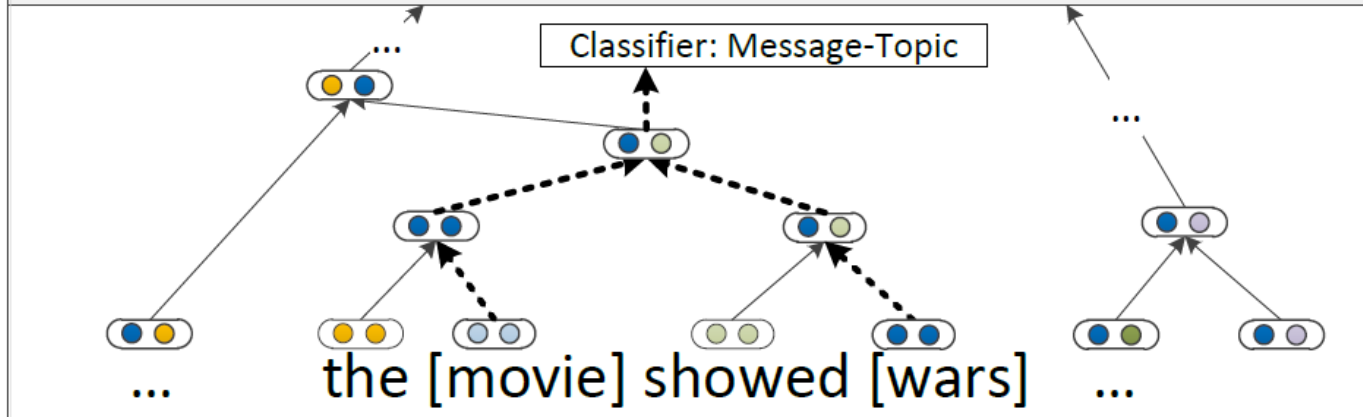
$$p = \tanh\left(W \begin{pmatrix} c_2 c_1 \\ + \\ c_1 c_2 \end{pmatrix} b\right)$$



Predicting Sentiment Distributions



MV-RNN for Relationship Classification



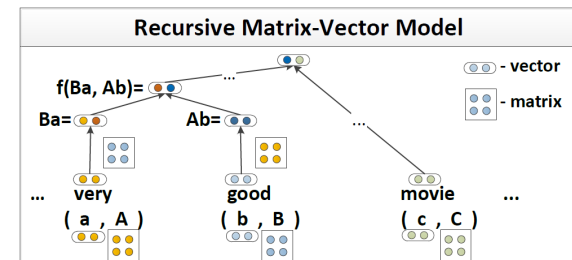
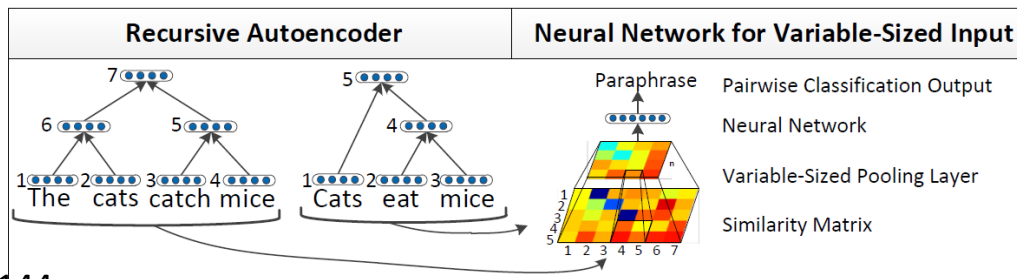
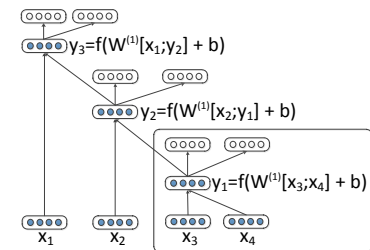
More info at
EMNLP talk
on July 14th

Relationship	Sentence with labeled nouns for which to predict relationships
Cause-Effect(e2,e1)	Avian [influenza] _{e1} is an infectious disease caused by type a strains of the influenza [virus] _{e2} .
Entity-Origin(e1,e2)	The [mother] _{e1} left her native [land] _{e2} about the same time and they were married in that city.
Message-Topic(e2,e1)	Roadside [attractions] _{e1} are frequently advertised with [billboards] _{e2} to attract tourists.

Classifier	Feature Sets	F1
SVM	POS, stemming, syntactic patterns	60.1
SVM	word pair, words in between	72.5
SVM	POS, WordNet, stemming, syntactic patterns	74.8
SVM	POS, WordNet, morphological features, thesauri, Google <i>n</i> -grams	77.6
MaxEnt	POS, WordNet, morphological features, noun compound system, thesauri, Google <i>n</i> -grams	77.6
SVM	POS, WordNet, prefixes and other morphological features, POS, dependency parse features, Levin classes, PropBank, FrameNet, NomLex-Plus, Google <i>n</i> -grams, paraphrases, TextRunner	82.2
RNN	-	74.8
Lin.MVR	-	73.0
MV-RNN	-	79.1
RNN	POS, WordNet, NER	77.6
Lin.MVR	POS, WordNet, NER	78.7
MV-RNN	POS, WordNet, NER	82.4

Summary: Recursive Deep Learning

- Recursive Deep Learning can predict hierarchical structure and classify the structured output using compositional vectors
- State-of-the-art performance on
 - Sentiment Analysis on multiple corpora
 - Paraphrase detection on the MSRP dataset
 - Relation Classification on SemEval 2011, Task8
 - Vision modality (Stanford background dataset)
- Code on www.socher.org



Part 3

1. Applications
 1. Neural language models
 2. Structured embedding of knowledge bases
 3. Assorted other speech and NLP applications
2. Resources (readings, code, ...)
3. Tricks of the trade
4. Discussion: Limitations, advantages, future directions

Existing NLP Applications

- Language Modeling
 - Speech Recognition
 - Machine Translation
- Part-Of-Speech Tagging
- Chunking
- Named Entity Recognition
- Semantic Role Labeling
- Sentiment Analysis
- Paraphrasing
- Question-Answering
- Word-Sense Disambiguation

Part 3.1: Applications

Neural Language Models

Language Modeling

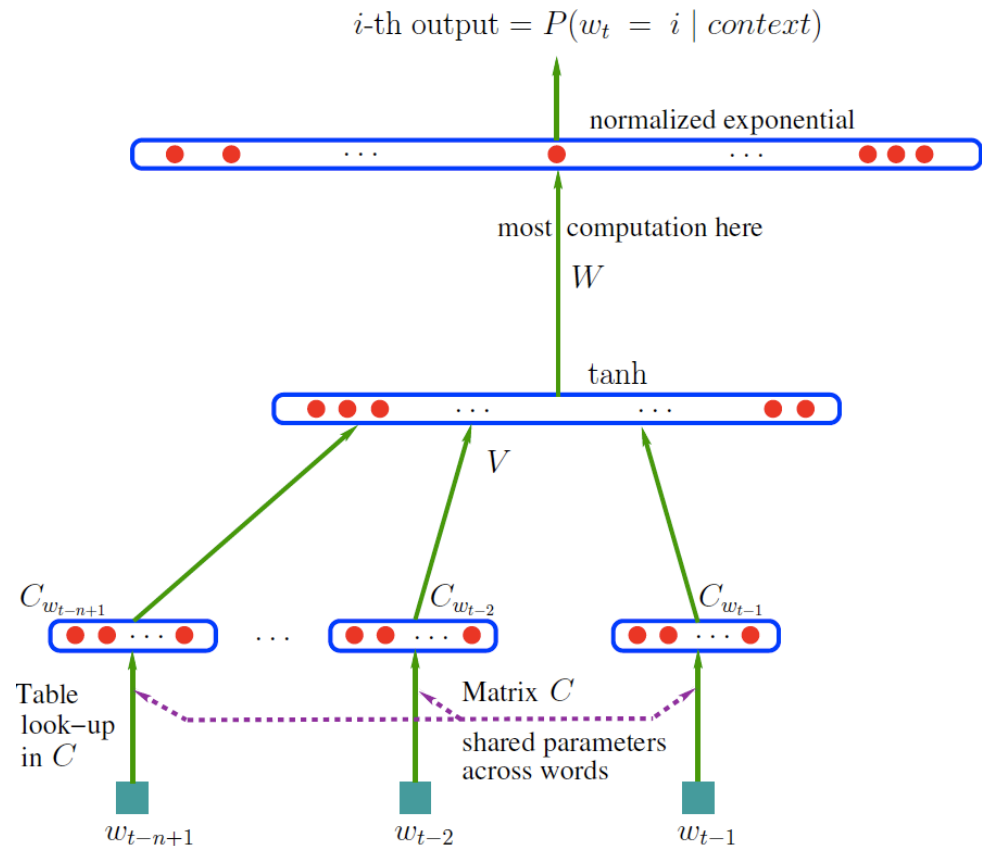
- Predict $P(\text{next word} \mid \text{previous word})$
- Gives a probability for a longer sequence
- Applications to Speech, Translation and Compression
- Computational bottleneck: large vocabulary V means that computing the output costs $\# \text{hidden units} \times |V|$.

Neural Language Model

- *Bengio et al NIPS'2000 and JMLR 2003 "A Neural Probabilistic Language Model"*



- Each word represented by a distributed continuous-valued code
- Generalizes to sequences of words that are semantically similar to training sequences



Recurrent Neural Net Language Modeling for ASR

- [Mikolov et al 2011]

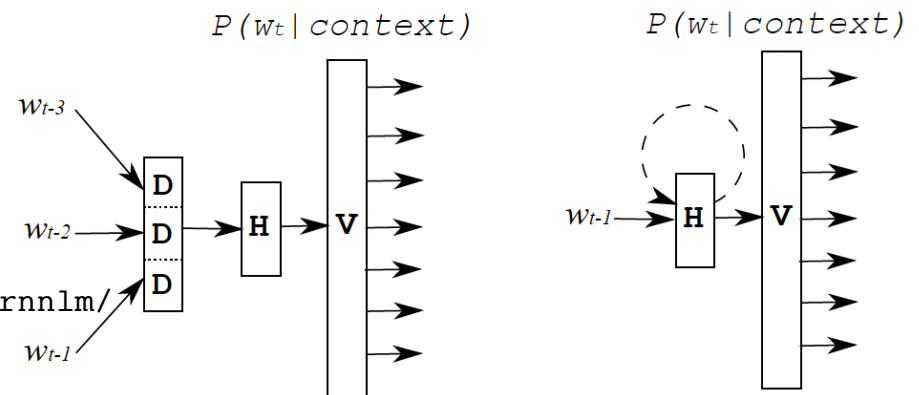
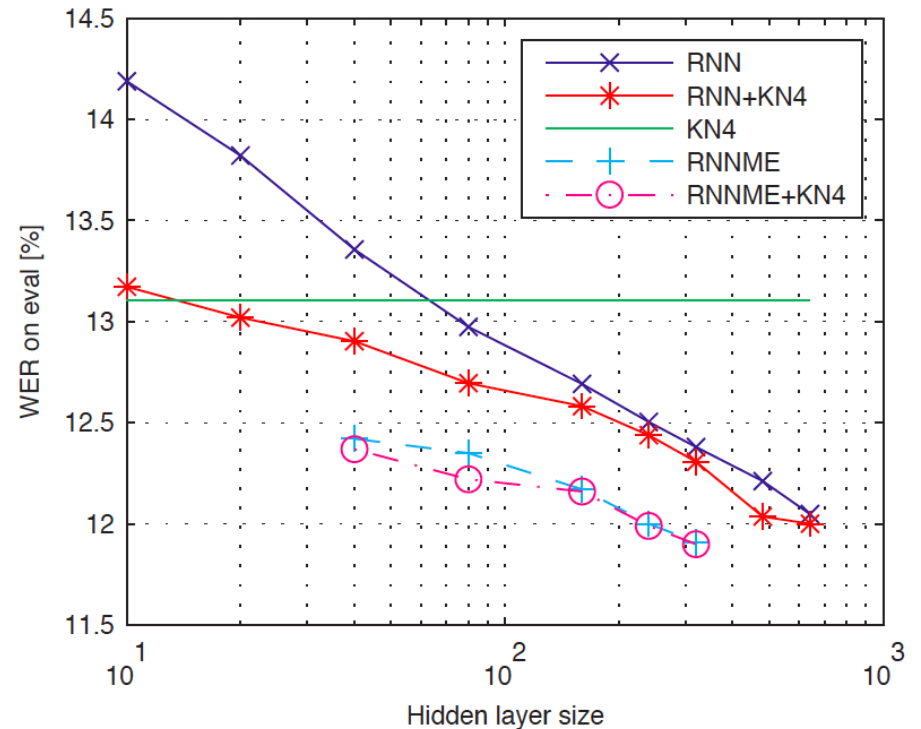


Bigger is better...
experiments on Broadcast
News NIST-RT04

perplexity goes from
140 to 102

Paper shows how to
train a recurrent neural net
with a single core in a few
days, with > 1% absolute
improvement in WER

Code: <http://www.fit.vutbr.cz/~imikolov/rnnlm/>



Language Modeling Output Bottleneck

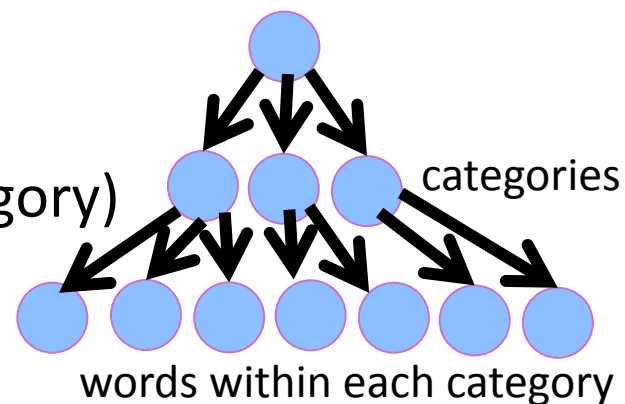
- [Schwenk et al 2002]: only predict most frequent words (short list) and use n-gram for the others



- [Morin & Bengio 2005; Blitzer et al 2005; Mnih & Hinton 2007,2009; Mikolov et al 2011]: **hierarchical representations**, multiple output groups, conditionally computed, predict

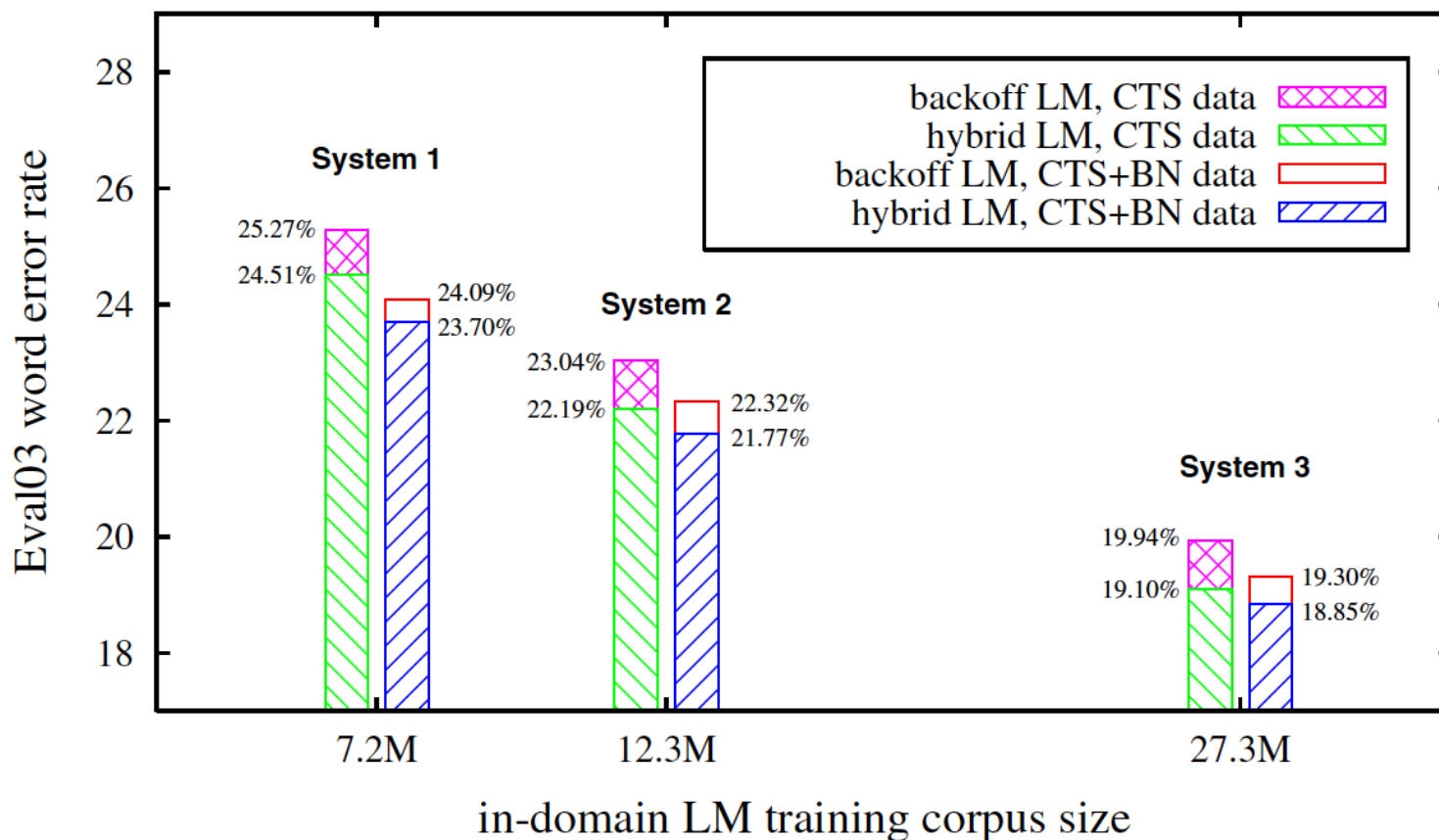
- $P(\text{word category} \mid \text{context})$
- $P(\text{sub-category} \mid \text{context, category})$
- $P(\text{word} \mid \text{context, sub-category, category})$

- Hard categories, can be arbitrary [Mikolov et al 2011]



Neural Net Language Modeling for ASR

- [Schwenk 2007], real-time ASR, perplexity AND word error rate improve (CTS evaluation set 2003), perplexities go from 50.1 to 45.5



Application to Statistical Machine Translation



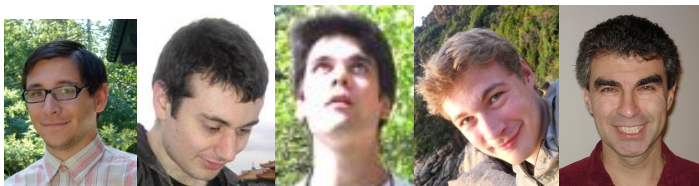
- Schwenk (NAACL 2012 workshop on the future of LM)
 - 41M words, Arabic/English bitexts + 151M English from LDC
- Perplexity down from 71.1 (6 Gig back-off) to 56.9 (neural model, 500M memory)
- +1.8 BLEU score (50.75 to 52.28)
- Can take advantage of longer contexts
- Code: <http://lium.univ-lemans.fr/cs1m/>

Part 3.1: Applications

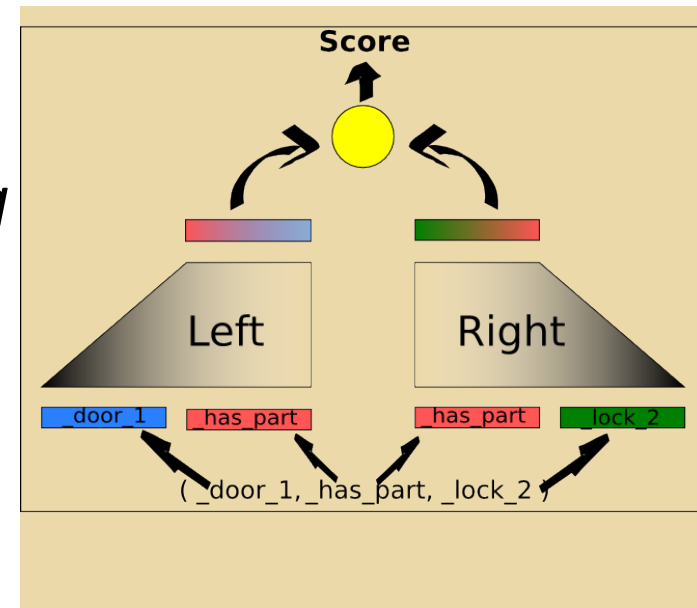
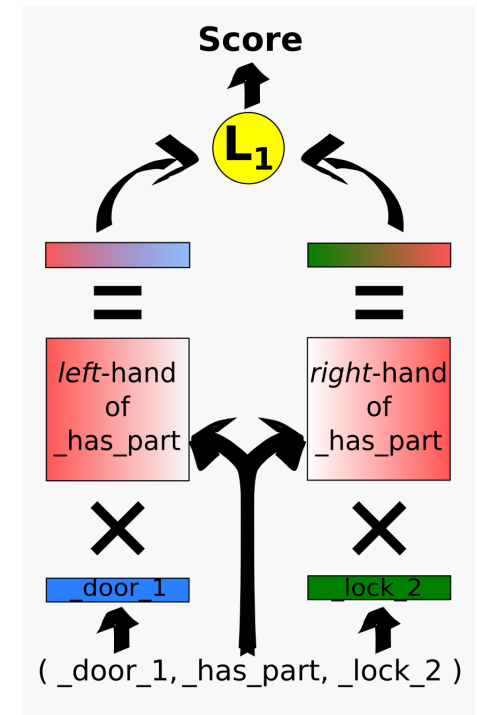
Structured embedding of knowledge bases

Modeling Semantics

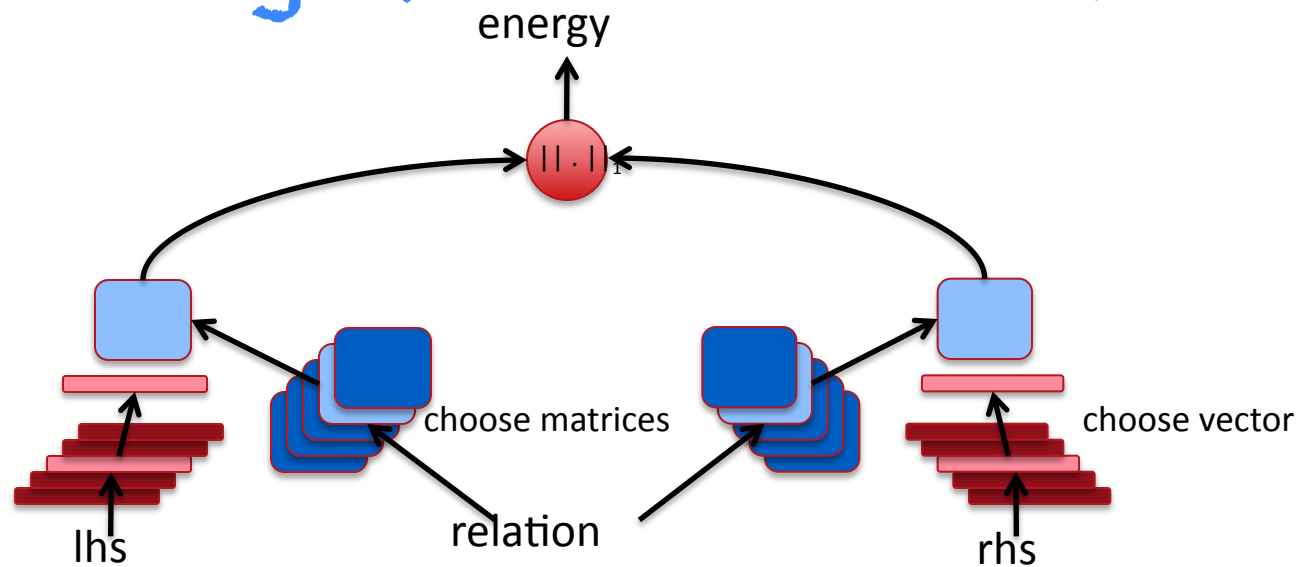
Learning Structured Embeddings of Knowledge Bases, (Bordes, Weston, Collobert & Bengio, AAAI 2011)



Joint Learning of Words and Meaning Representations for Open-Text Semantic Parsing, (Bordes, Glorot, Weston & Bengio, AISTATS 2012)



Modeling Relations with Matrices



Model (lhs, relation, rhs)

Each concept = 1 embedding vector

Each relation = 2 matrices. **Matrix acts like an operator.**

Ranking criterion

Energy = low for training examples, high o/w

Question Answering: implicitly adding new relations to WN or FB

	Model (All)	<i>TextRunner</i>
<i>lhs</i>	_army_NN_1	<i>army</i>
<i>rel</i>	_attack_VB_1	<i>attacked</i>
top ranked <i>rhs</i>	_troop_NN_4 _armed_service_NN_1 _ship_NN_1 _territory_NN_1 _military_unit_NN_1	<i>Israel</i> <i>the village</i> <i>another army</i> <i>the city</i> <i>the fort</i>
top ranked <i>lhs</i>	_business_firm_NN_1 _person_NN_1 _family_NN_1 _payoff_NN_3 _card_game_NN_1	<i>People</i> <i>Players</i> <i>one</i> <i>Students</i> <i>business</i>
<i>rel</i>	_earn_VB_1	<i>earn</i>
<i>rhs</i>	_money_NN_1	<i>money</i>

MRs inferred from text define triplets between WordNet synsets.

Model captures knowledge about relations between nouns and verbs.

→ Implicit addition of new relations to WordNet!

→ Generalize Freebase!

Embedding Nearest Neighbors of Words & Senses

<p>_mark_NN</p> <p>_indication_NN _print_NN_3 _print_NN _roll_NN _pointer_NN</p>	<p>_mark_NN_1</p> <p>_score_NN_1 _number_NN_2 _gradation_NN _evaluation_NN_1 _tier_NN_1</p>	<p>_mark_NN_2</p> <p>_marking_NN_1 _symbolizing_NN_1 _naming_NN_1 _marking_NN _punctuation_NN_3</p>
<p>_take_VB</p> <p>_bring_VB _put_VB _ask_VB _hold_VB _provide_VB</p>	<p>_canary_NN</p> <p>_sea_mew_NN_1 _yellowbird_NN_2 _canary_bird_NN_1 _larus_marinus_NN_1 _mew_NN</p>	<p>_different_JJ_1</p> <p>_eccentric_NN _dissimilar_JJ _same_JJ_2 _similarity_NN_1 _common_JJ_1</p>

Word Sense Disambiguation

- Senseval-3 results
(only sentences with Subject-Verb-Object structure)

MFS=most frequent sense

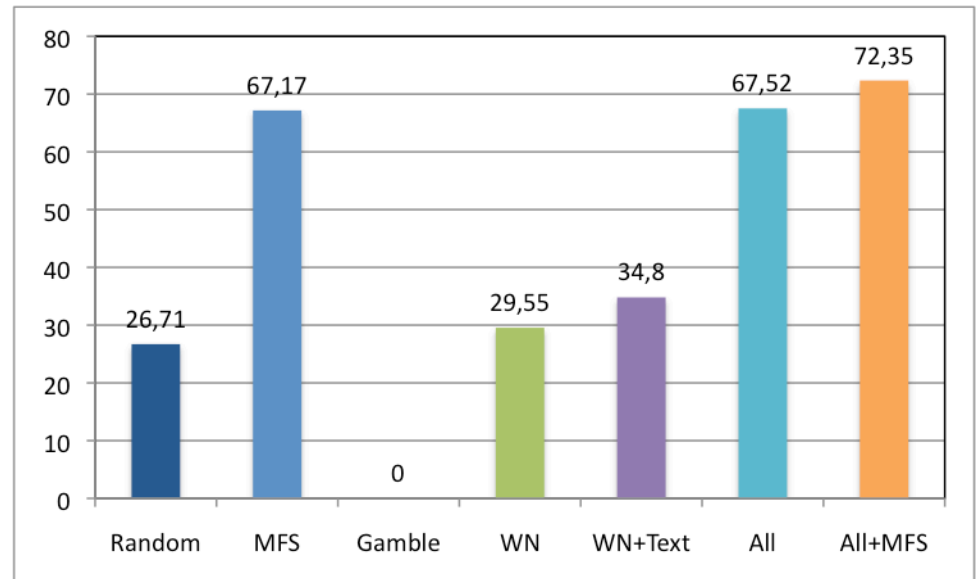
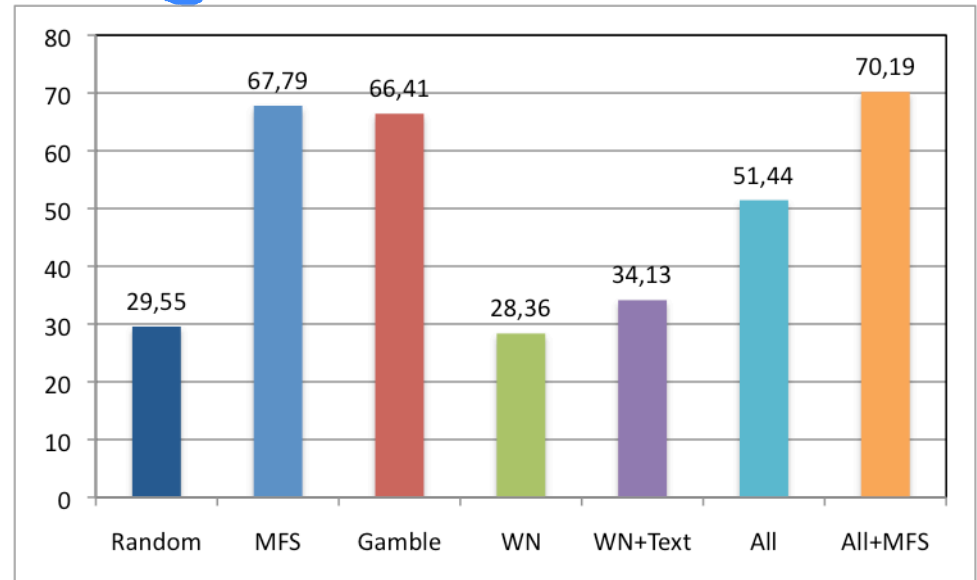
All=training from all sources

Gamble=Decadt et al 2004

(Senseval-3 SOA)

- XWN results

XWN = eXtended WN

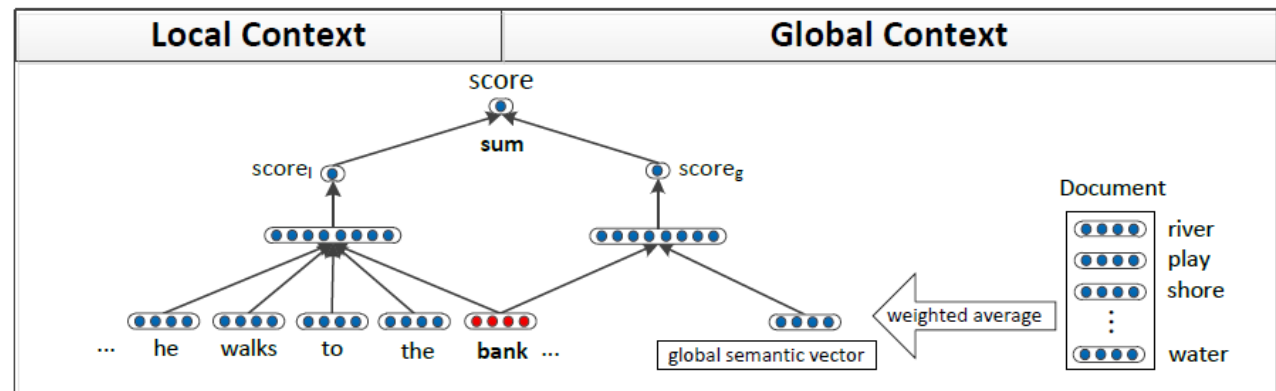


Part 3.1: Applications

Assorted Speech and NLP Applications

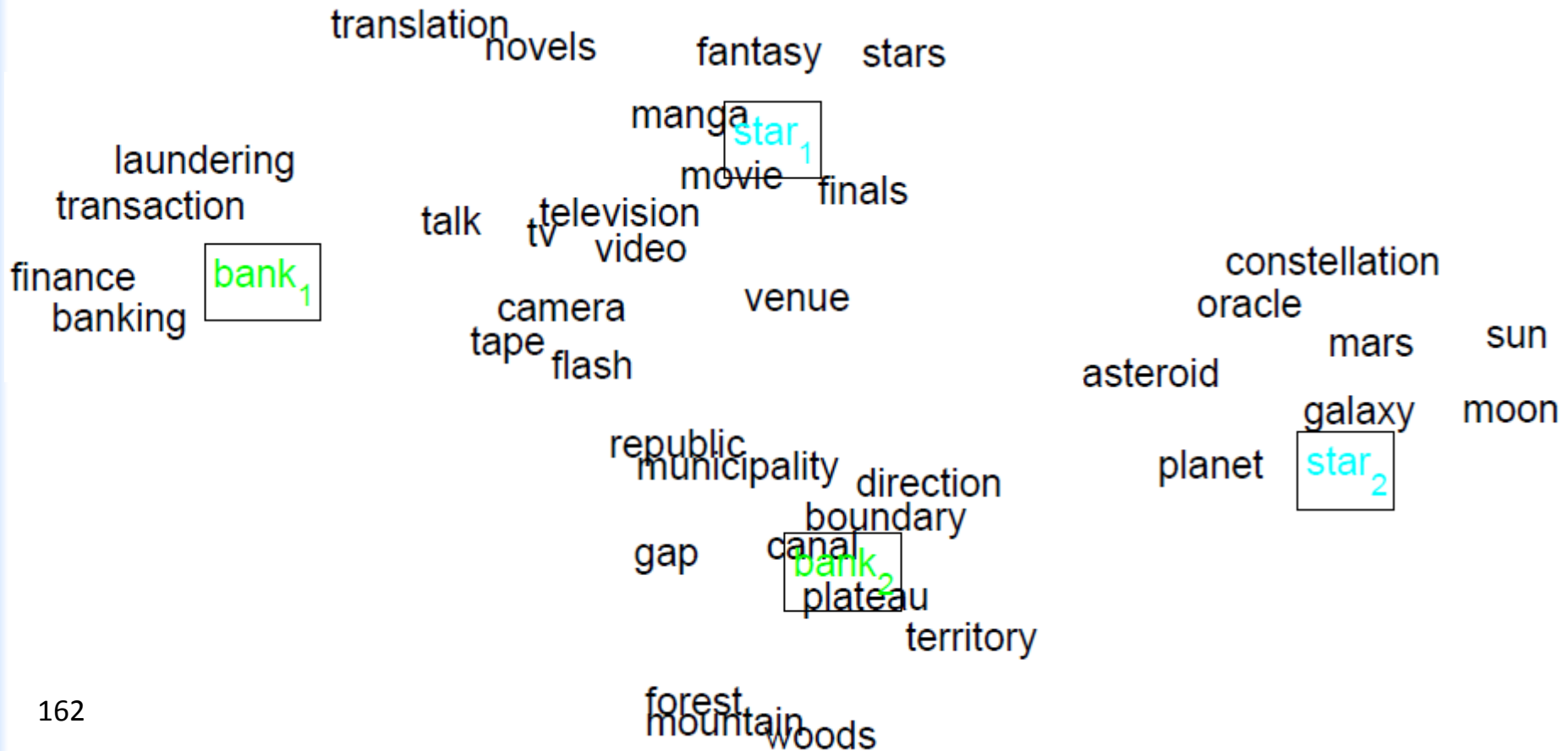
Learning Multiple Word Vectors

- Tackles problems with polysemous words
- Can be done with both standard tf-idf based methods [Reisinger and Mooney, NAACL 2010]
- Recent neural word vector model by [Huang et al. ACL 2012] learns multiple prototypes using both local and global context
- State of the art correlations with human similarity judgments



Learning Multiple Word Vectors

- Visualization of learned word vectors from Huang et al. (ACL 2012)



Phoneme-Level Acoustic Models

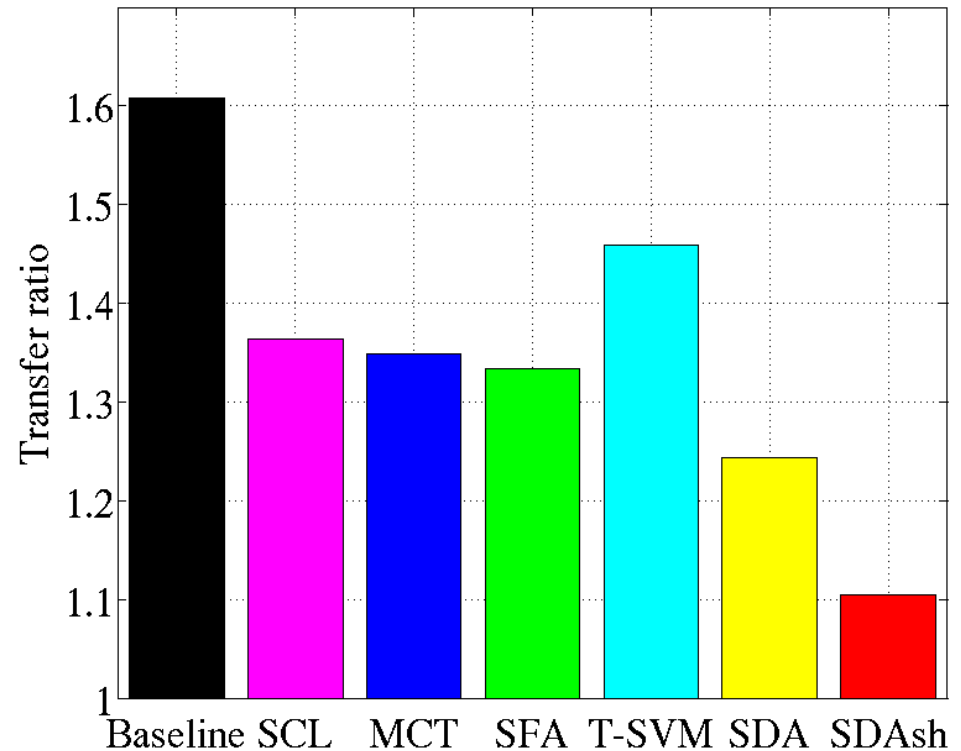


- [Mohamed et al, 2011, IEEE Tr.ASLP]
- Unsupervised pre-training as Deep Belief Nets (a stack of RBMs), supervised fine-tuning to predict phonemes
- Phoneme classification on TIMIT:
 - CD-HMM: 27.3% error
 - CRFs: 26.6%
 - Triphone HMMs w. BMMI: 22.7%
 - Unsupervised DBNs: 24.5%
 - Fine-tuned DBNs: 20.7%
- Improved version by Dong Yu is **RELEASED IN MICROSOFT'S ASR** system for Audio Video Indexing Service

Domain Adaptation for Sentiment Analysis



- [Glorot et al, ICML 2011] beats SOTA on Amazon benchmark, 25 domains
- Embeddings pre-trained in denoising auto-encoder
- Disentangling effect (features specialize to domain or sentiment)



Part 3.2: Resources

Resources: Tutorials and Code

Related Tutorials

- See “*Neural Net Language Models*” **Scholarpedia** entry
- Deep Learning tutorials: <http://deeplearning.net/tutorials>
- Stanford deep learning tutorials with simple programming assignments and reading list
<http://deeplearning.stanford.edu/wiki/>
- Recursive Autoencoder class project
<http://cseweb.ucsd.edu/~elkan/250B/learningmeaning.pdf>
- Graduate Summer School: Deep Learning, Feature Learning
<http://www.ipam.ucla.edu/programs/gss2012/>
- ICML 2012 Representation Learning tutorial <http://www.iro.umontreal.ca/~bengioy/talks/deep-learning-tutorial-2012.html>
- Paper references in separate pdf

Software

- Theano (Python CPU/GPU) mathematical and deep learning library <http://deeplearning.net/software/theano>
 - Can do automatic, symbolic differentiation
- Senna: POS, Chunking, NER, SRL
 - by Collobert et al. <http://ronan.collobert.com/senna/>
 - State-of-the-art performance on many tasks
 - 3500 lines of C, extremely fast and using very little memory
- Recurrent Neural Network Language Model
<http://www.fit.vutbr.cz/~imikolov/rnnlm/>
- Recursive Neural Net and RAE models for paraphrase detection, sentiment analysis, relation classification www.socher.org


Software: what's next

- Off-the-shelf SVM packages are useful to researchers from a wide variety of fields (no need to understand RKHS).
- One of the goals of deep learning: Build off-the-shelf NLP classification packages that are using as input only raw text, possibly with a label.

Part 3.3: Deep Learning Tricks

Deep Learning Tricks

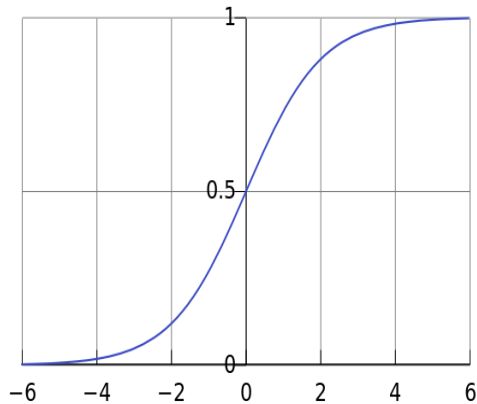
Deep Learning Tricks of the Trade

- Y. Bengio (2012), “Practical Recommendations for Gradient-Based Training of Deep Architectures”
 - Unsupervised pre-training
 - Stochastic gradient descent and setting learning rates
 - Main hyper-parameters
 - Learning rate schedule & Early stopping
 - Minibatches
 - Parameter initialization
 - Number of hidden units
 - L1 or L2 weight decay
 - Sparsity regularization
 - Debugging → Finite difference gradient check (Yay)
 - How to efficiently search for hyper-parameter configurations

Non-linearities: What's used

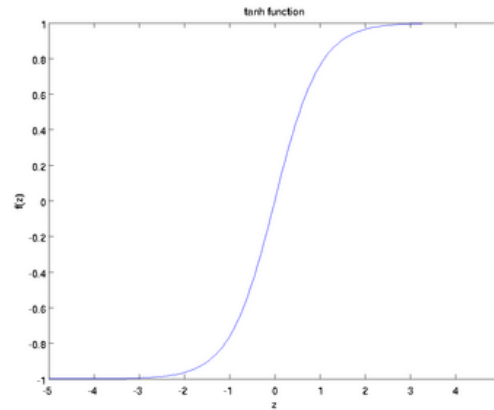
logistic (“sigmoid”)

$$f(z) = \frac{1}{1 + \exp(-z)}$$



tanh

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



tanh is just a rescaled and shifted sigmoid ($2 \times$ as steep, $[-1,1]$):

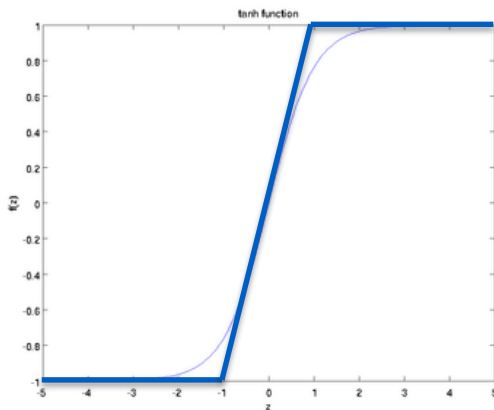
$$\tanh(z) = 2\text{logistic}(2z) - 1$$

tanh is what is most used and often performs best for deep nets

Non-linearities: There are various other choices

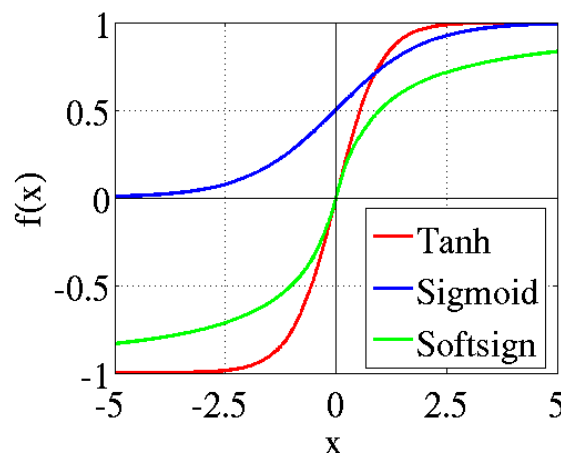
hard tanh

$$\text{HardTanh}(x) = \begin{cases} -1 & \text{if } x < -1 \\ x & \text{if } -1 \leq x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$$



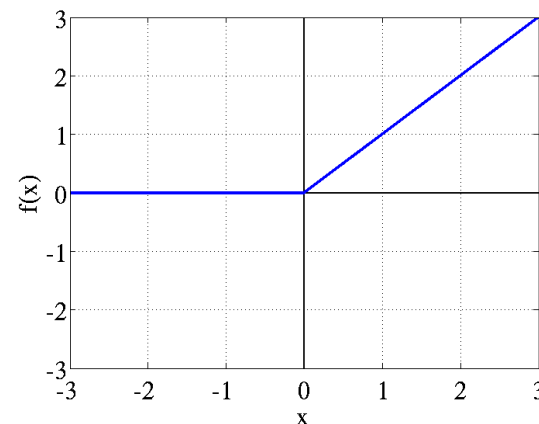
soft sign

$$\text{softsign}(z) = \frac{a}{1+|a|}$$



rectifier

$$\text{rect}(z) = \max(z, 0)$$



- hard tanh similar but computationally cheaper than tanh and saturates hard.
- [Glorot and Bengio *AISTATS* 2010, 2011] discuss softsign and rectifier

Stochastic Gradient Descent (SGD)

- Gradient descent uses total gradient over all examples per update, SGD updates after only 1 or few examples:

$$\theta^{(t)} \leftarrow \theta^{(t-1)} - \epsilon_t \frac{\partial L(z_t, \theta)}{\partial \theta}$$

- L = loss function, z_t = current example, θ = parameter vector, and ϵ_t = learning rate.
- Ordinary gradient descent is a batch method, very slow, **should never be used**. Use 2nd order batch method such as LBFGS. On large datasets, SGD usually wins over all batch methods. On smaller datasets LBFGS or Conjugate Gradients win. Large-batch LBFGS extends the reach of LBFGS [Le et al ICML'2011].

Learning Rates

- Simplest recipe: keep it fixed and use the same for all parameters.
- Collobert scales them by the inverse of square root of the fan-in of each neuron
- Better results can generally be obtained by allowing learning rates to decrease, typically in $O(1/t)$ because of theoretical convergence guarantees, e.g.,

$$\epsilon_t = \frac{\epsilon_0 \tau}{\max(t, \tau)}$$

with hyper-parameters ϵ_0 and τ .

Long-Term Dependencies and Clipping Trick



- In very deep networks such as recurrent networks (or possibly recursive ones), the gradient is a product of Jacobian matrices, each associated with a step in the forward computation. This can become very small or very large quickly [Bengio et al 1994], and the locality assumption of gradient descent breaks down.

$$L = L(s_T(s_{T-1}(\dots s_{t+1}(s_t, \dots))))$$

$$\frac{\partial L}{\partial s_t} = \frac{\partial L}{\partial s_T} \frac{\partial s_T}{\partial s_{T-1}} \dots \frac{\partial s_{t+1}}{\partial s_t}$$

- The solution first introduced by Mikolov is to clip gradients to a maximum value. Makes a big difference in RNNs



Parameter Initialization

- Initialize hidden layer biases to 0 and output (or reconstruction) biases to optimal value if weights were 0 (e.g. mean target or inverse sigmoid of mean target).
- Initialize weights \sim Uniform(-r,r), r inversely proportional to fan-in (previous layer size) and fan-out (next layer size):

$$\sqrt{6 / (\text{fan-in} + \text{fan-out})}$$

for tanh units, and 4x bigger for sigmoid units [Glorot AISTATS 2010]

Note: for embedding weights, fan-in=1 and we don't care about fan-out, Collobert uses Uniform(-1,1).

Part 3.4: Discussion

Discussion: Limitations, Advantages, Future Directions

Concerns

- Many algorithms and variants (burgeoning field)
- Hyper-parameters (layer size, regularization, possibly learning rate)
 - Use multi-core machines, clusters and random sampling for cross-validation (Bergstra & Bengio 2012)
 - Pretty common for powerful methods, e.g. BM25
 - Can use (mini-batch) L-BFGS instead of SGD

Concerns

- Not always obvious how to combine with existing NLP
 - Simple: Add word or phrase vectors as features. Gets close to state of the art for NER, [Turian et al, ACL 2010]
 - Integrate with known structures: Recursive and recurrent networks for trees and chains
 - Your research here

Concerns

- Slower to train than linear models
 - Only by a small constant factor, and much more compact than non-parametric (e.g. n-gram models)
 - Very fast during inference/test time (feed-forward pass is just a few matrix multiplies)
- Need more training data
 - Can *handle and benefit from* more training data, suitable for age of Big Data (Google trains neural nets with a billion connections, [Le et al, ICML 2012])

CONCERNS

- There aren't many good ways to encode prior knowledge about the structure of language into deep learning models
 - There is some truth to this. However:
 - You can choose architectures suitable for a problem domain, as we did for linguistic structure
 - You can include human-designed features in the first layer, just like for a linear model
 - And the goal is to get the machine doing the learning!

Concern:

Problems with model interpretability

- No discrete categories or words, everything is a continuous vector. We'd like have symbolic features like NP, VP, etc. and see why their combination makes sense.
 - True, but most of language is fuzzy and many words have soft relationships to each other. Also, many NLP features are already not human-understandable (e.g., concatenations/ combinations of different features).

Concern: non-convex optimization

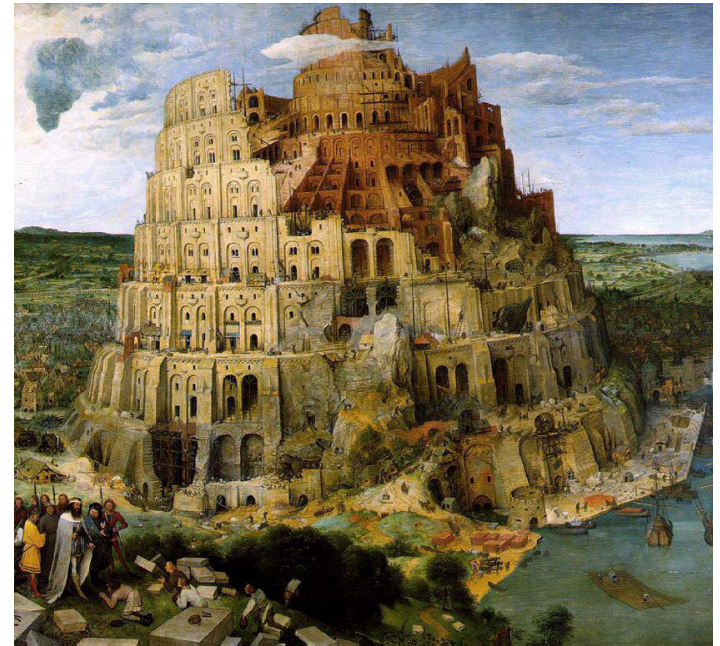
- Can initialize system with convex learner
 - Convex SVM
 - Fixed feature space
- Then optimize non-convex variant (add and tune learned features), can't be worse than convex learner

Advantages

- Despite a small community in the intersection of deep learning and NLP, already many state of the art results on a variety of language tasks
- Often very simple matrix derivatives (backprop) for training and matrix multiplications for testing → fast implementation
- Fast inference and well suited for multi-core CPUs/GPUs and parallelization across machines

Learning Multiple Levels of Abstraction

- The big payoff of deep learning is to allow learning higher levels of abstraction
- Higher-level abstractions disentangle the factors of variation, which allows much easier generalization and transfer
- More abstract representations
 - Successful transfer (domains, languages)



The End

