

# Deep Learning for NLP

## Part 2



CS224N

Christopher Manning

(Many slides borrowed from ACL 2012/NAACL 2013  
Tutorials by me, Richard Socher and Yoshua Bengio)

# 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] = 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

# Distributional word vectors

Distributional counts give same dimension, denser representation

	motel	hotel	bank
debt	3	9	17
crises	1	0	11
unified	0	0	5
Hodgepodge	0	1	2
the	122	147	183
pillow	21	25	1
reception	25	37	3
internet	8	19	8

# Two traditional word representations: Class-based and soft clustering

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

- Brown clustering (Brown et al. 1992, Liang 2005)
- Exchange clustering (Martin et al. 1998, Clark 2003)
  1. Clinton, Jiang, Bush, Wilensky, Suharto, Reagan, ...
  5. also, still, already, currently, actually, typically, ...
  6. recovery, strength, expansion, freedom, resistance, ...

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:

Word meaning is represented as a  
(dense) vector – a point in a  
(medium-dimensional) vector space

Neural word embeddings combine  
vector space semantics with the  
prediction of probabilistic models  
(Bengio et al. 2003, Collobert &  
Weston 2008, Huang et al. 2012)

*linguistics* =

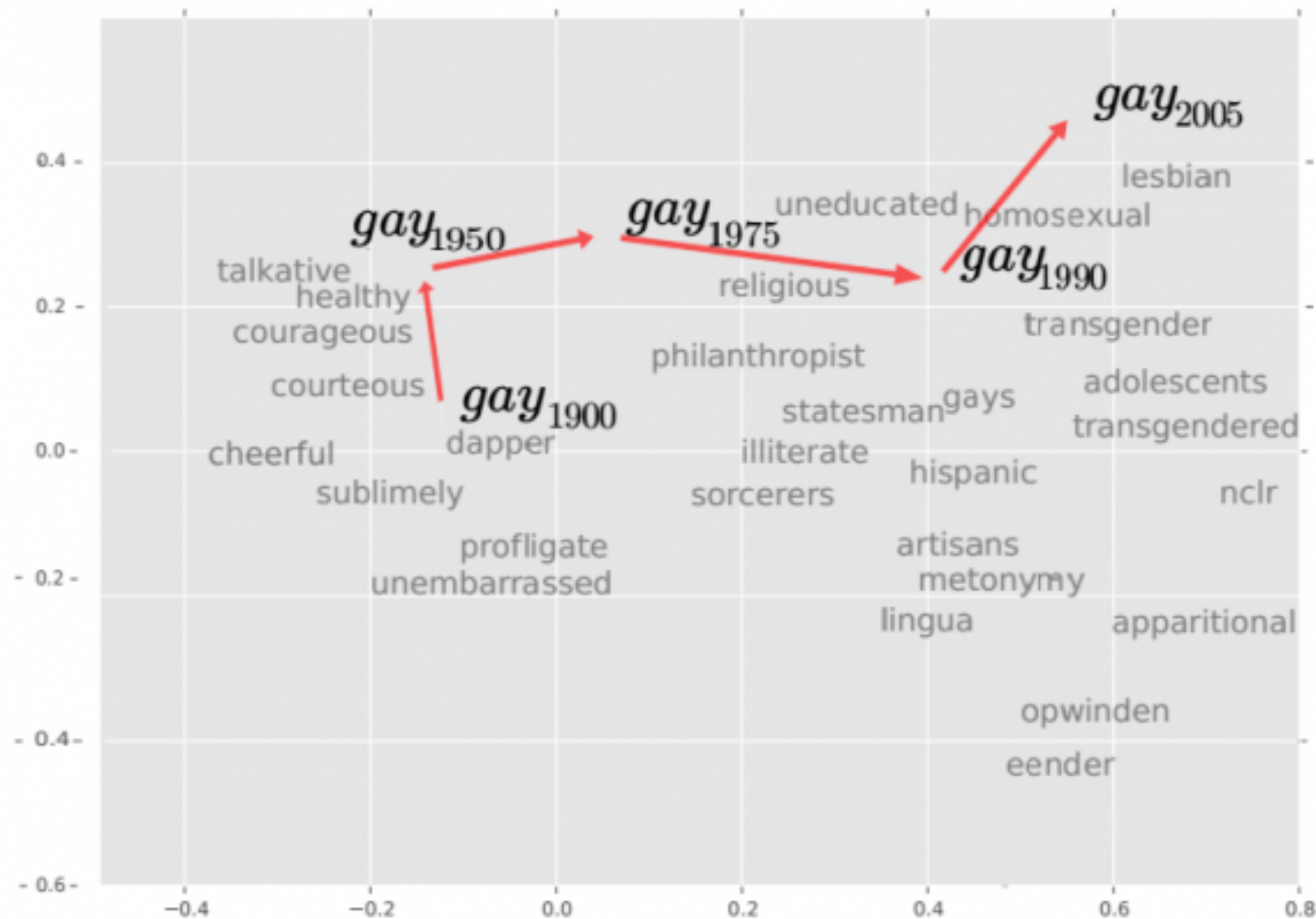
$$\begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}$$

# Neural word embeddings - visualization





# Kulkarni, Al-Rfou, Perozzi, & Skiena (2015)



# Distributed vs. distributional representations

**Distributed:** A concept is represented as continuous activation levels in a number of elements

[Contrast: **local**]

**Distributional:** Meaning is represented by contexts of use

[Contrast: **denotational**]



# Advantages of the neural word embedding approach

Compared to other methods, 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 compositional vector representations for longer phrases (next lecture)

# Unsupervised word vector Learning

# The mainstream methods for neural word vector learning in 2015

## 1. word2vec

- <https://code.google.com/p/word2vec/>
- Code for two different algorithms (Skipgram with negative sampling – SGNS, and Continuous Bag of Words – CBOW)
- Uses simple, fast bilinear models to learn very good word representations
- Mikolov, Sutskever, Chen, Corrado, and Dean. NIPS 2013


## 2. GloVe

- <http://nlp.stanford.edu/projects/glove/>
- A non-linear matrix factorization: Starts with count matrix and factorizes it with an explicit loss function
- Sort of similar to SGNS, really, but from Stanford 😊
- Pennington, Socher, and Manning. EMNLP 2014

# But we won't look at either, but ... Contrastive Estimation of Word Vectors

**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 Ohio a mat

# A neural network for learning word vectors

How do we formalize this idea? Ask that

score(cat chills on a mat) > score(cat chills Ohio 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  $L \in \mathbb{R}^{n \times |V|}$

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

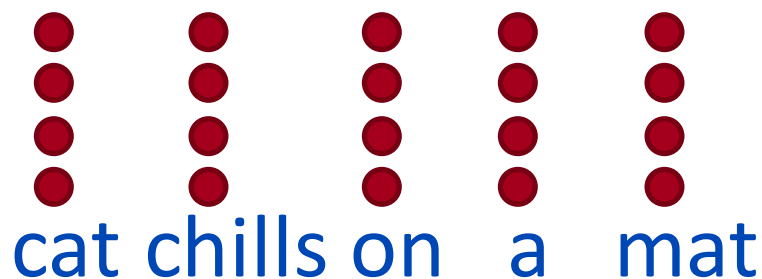
the   cat   mat ...

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



# 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 form a  $5n$  vector:
- $x = [ \text{●●●●} \text{●●●●} \text{●●●●} \text{●●●●} \text{●●●●} ]$
- How do we then compute  $\text{score}(x)$ ?

# Scoring 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:

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

# Summary: Feed-forward Computation

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

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

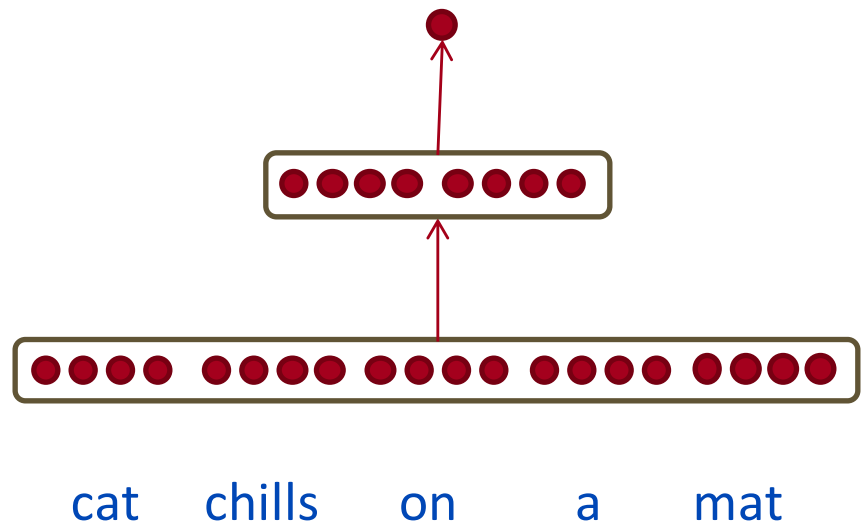
$$s = U^T a$$

$$a = f(z)$$

$$z = Wx + b$$

$$x = [x_{cat} \ x_{chills} \ x_{on} \ x_a \ x_{mat}]$$

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

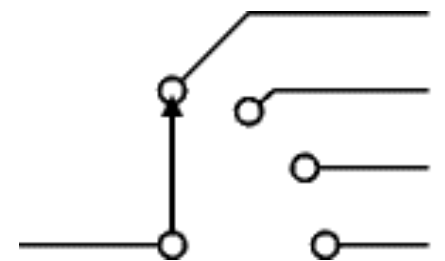


## Summary: Feed-forward Computation

- $s = \text{score}(\text{cat chills on a mat})$
- $s_c = \text{score}(\text{cat chills Ohio a mat})$
- Idea for training objective: make score of true window larger and corrupt window's score lower (until they're sufficiently separated). Minimize

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

- This is continuous, can perform SGD
  - Look at a few examples, nudge weights to make  $J$  smaller



# The Backpropagation Algorithm

- **Backpropagation** is a way of computing gradients of expressions **efficiently** through recursive application of **chain rule**.
  - Either Rumelhart, Hinton & McClelland 1986 or Werbos 1974 or Linnainmaa 1970
- The derivative of a loss (an objective function) on each variable tells you the sensitivity of the loss to its value.
- An individual step of the chain rule is local. It tells you how the sensitivity of the objective function is modulated by that function in the network
  - The input changes at some rate for a variable
  - The function at the node scales that rate

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

# Training the parameters of the model

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

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

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

If cost  $J$  is  $< 0$  (or  $= 0!$ ), the derivative is 0. Do nothing.

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:  $U$ ,  $W$ ,  $b$ ,  $x$ . Then take differences for  $J$ .

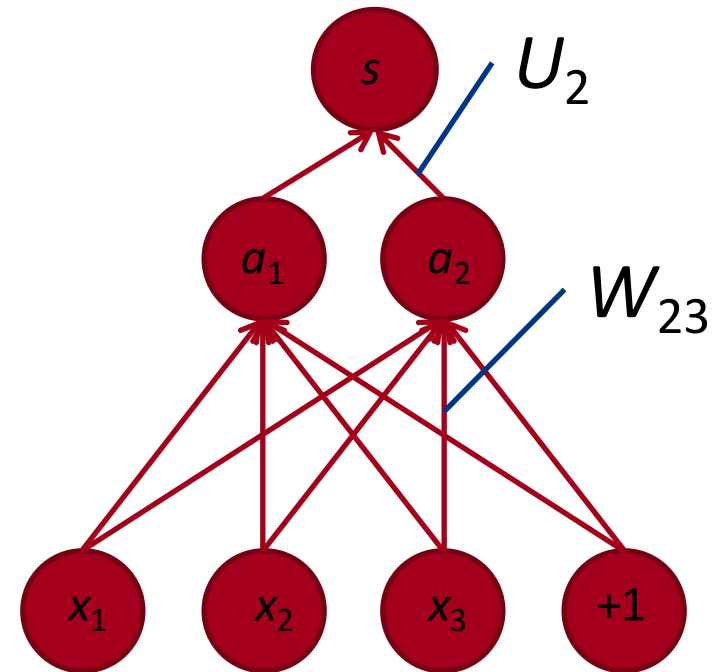
$$\frac{\partial s}{\partial U} = \frac{\partial}{\partial U} U^T a \qquad \frac{\partial s}{\partial U} = a$$

# 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  $W_{ij}$ :

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

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

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

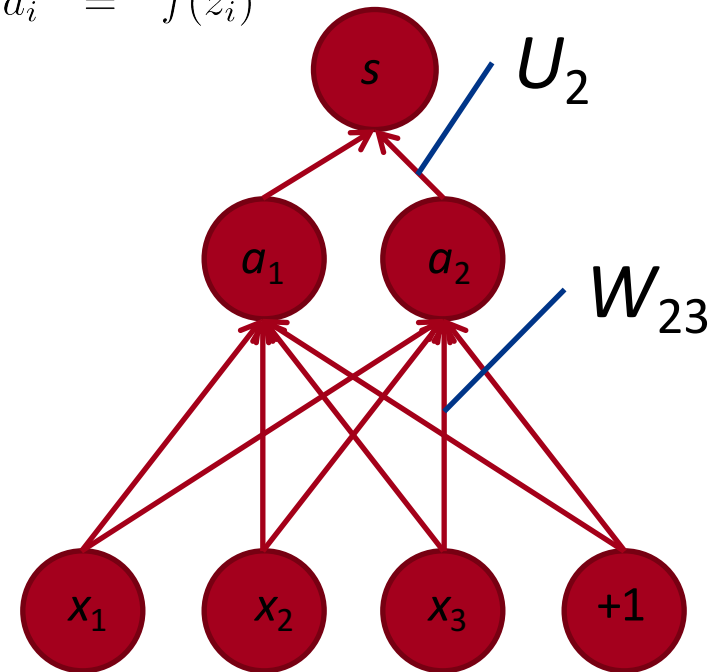
$$= U_i \frac{\partial f(z_i)}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}}$$

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

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

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

$$a_i = f(z_i)$$





# Training with Backpropagation

Derivative of single weight  $W_{ij}$  :

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

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

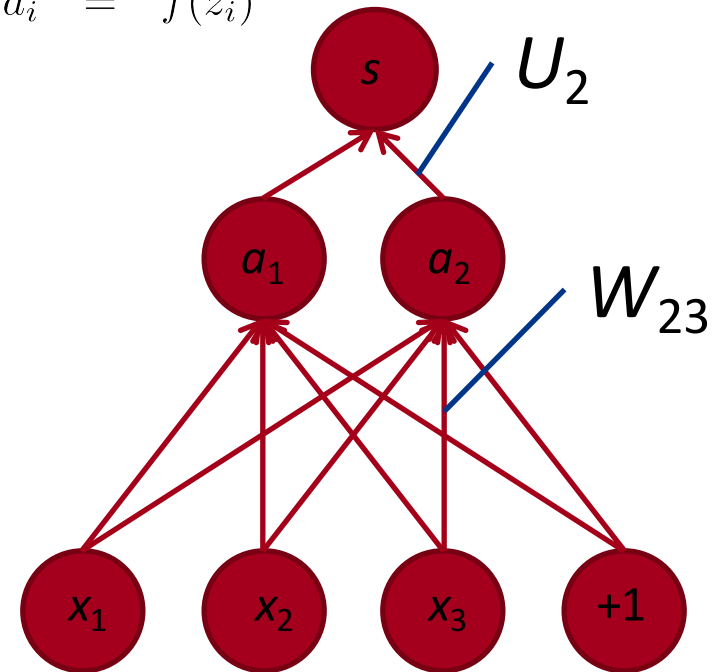
$$= U_i f'(z_i) \frac{\partial}{\partial W_{ij}} \sum_k W_{ik} x_k$$

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

$\delta_i$ 
Local error signal
 $x_j$ 
Local input signal

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

$$a_i = f(z_i)$$



25 where  $f'(z) = f(z)(1 - f(z))$  for logistic  $f$

# Training with Backpropagation

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

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

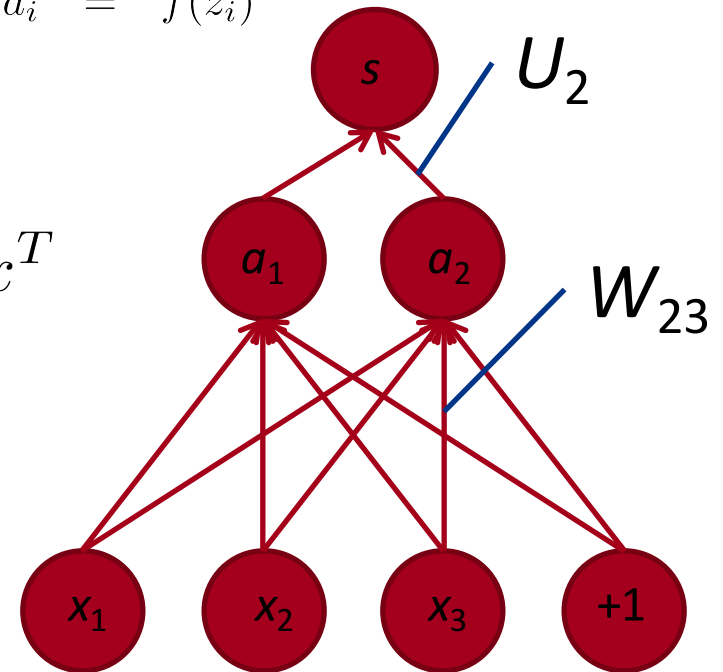
$$z_i = W_i \cdot x + b_i = \sum_{j=1}^3 W_{ij} x_j + b_i$$

$$a_i = f(z_i)$$

- 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  $a$

$$\begin{pmatrix} \delta_1 x_1 & \delta_1 x_2 & \delta_1 x_3 \\ \delta_2 x_1 & \delta_2 x_2 & \delta_2 x_3 \end{pmatrix}$$



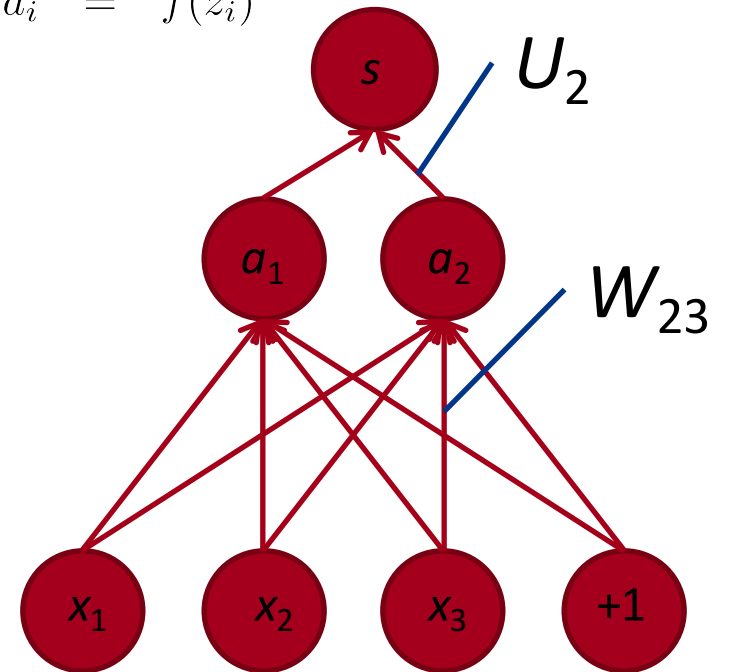
# Training with Backpropagation

- 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

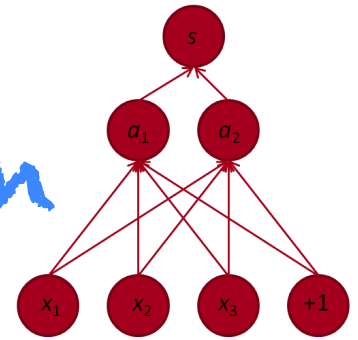
It's simply taking derivatives and using the chain rule!

Remaining trick: **Efficiency**

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 were 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:

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

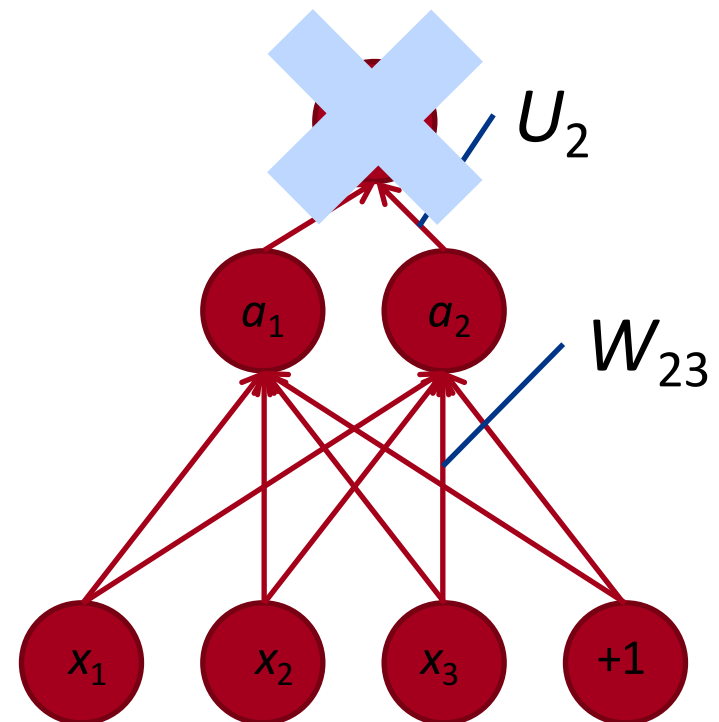


# Learning word-level classifiers: POS and NER

# The Model

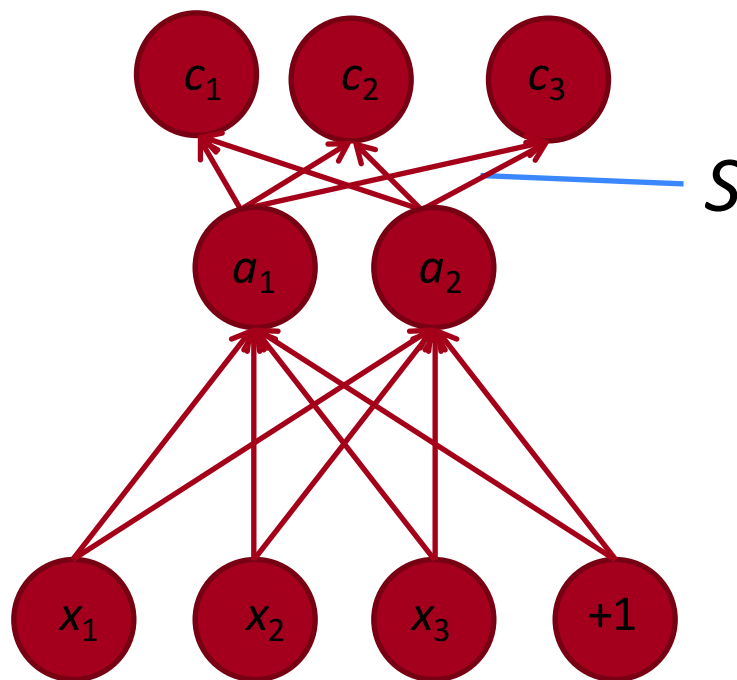
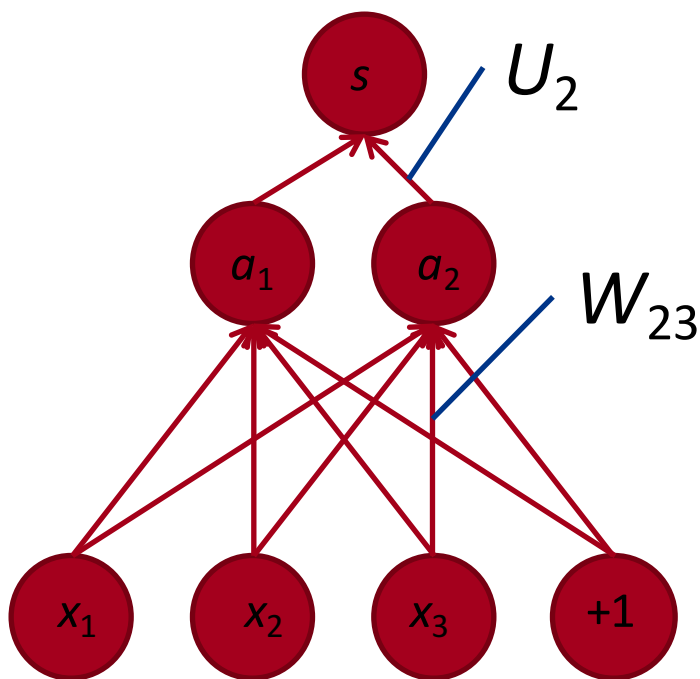
(Collobert & Weston 2008;  
Collobert et al. 2011)

- 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 (*but not the same as!*) the score in the unsupervised word vector learning model



# The Model - Training

- We already know softmax/MaxEnt 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:





# Training with Backpropagation: softmax

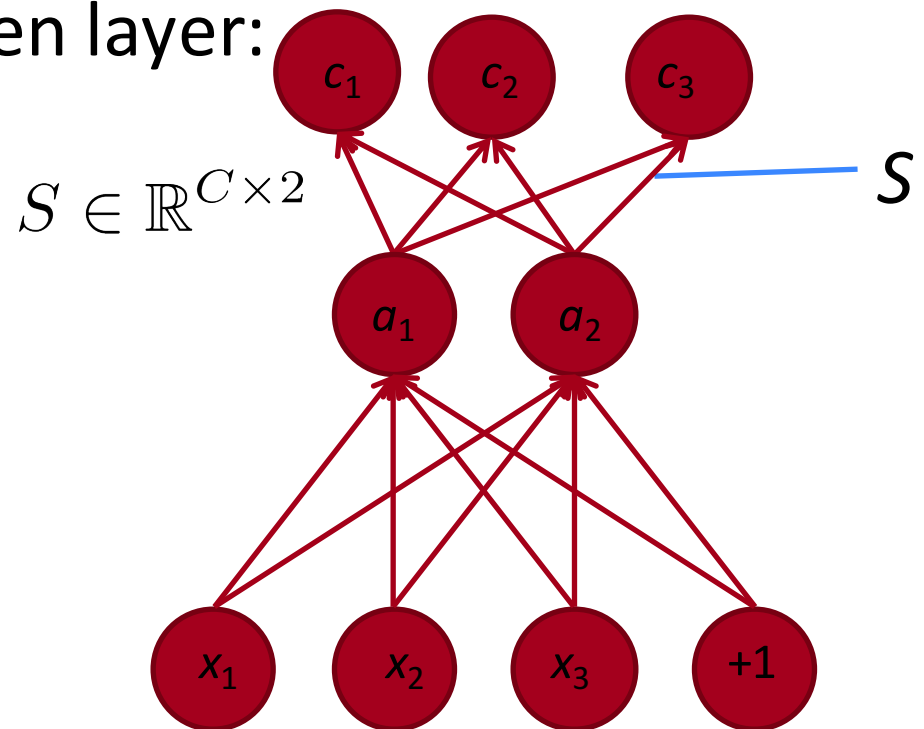
What is the major benefit of learned word vectors?

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

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



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



For small supervised data sets,  
unsupervised pre-training helps a lot

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

\* Results used in [Collobert & Weston \(2011\)](#).

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*	<b>97.10</b>	<b>88.87</b>
<b>C&amp;W 2011**</b>	<b>97.29</b>	<b>89.59</b>

\* 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

Part 1.5: The Basics

# Backpropagation Training

# Back-Prop

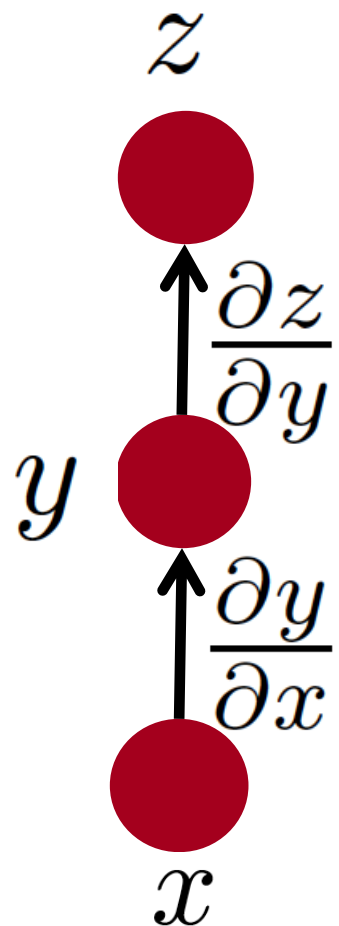
- Compute gradient of example-wise loss wrt parameters

- 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



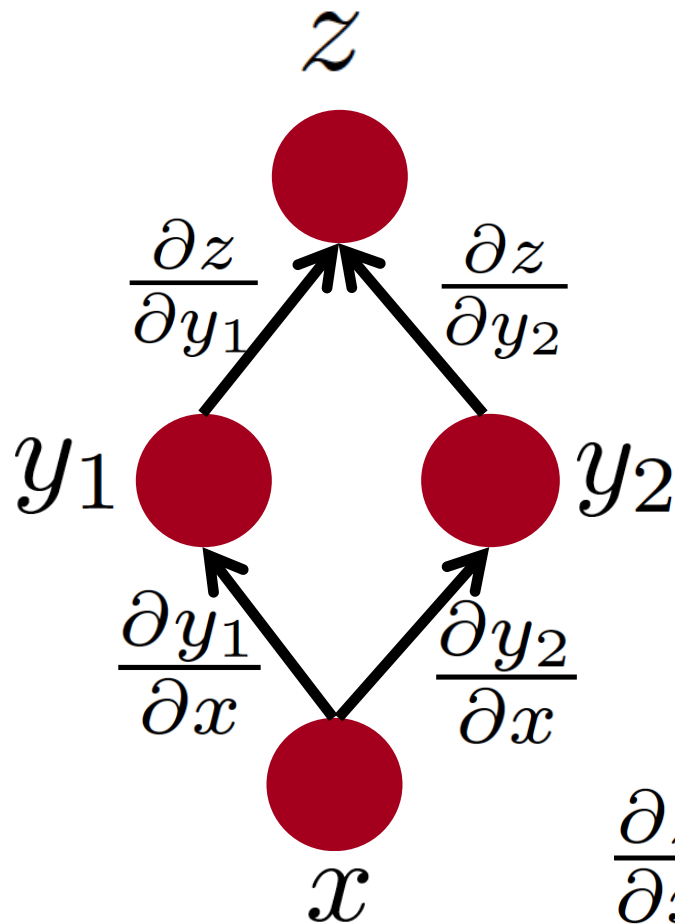
$$\Delta z = \frac{\partial z}{\partial y} \Delta y$$

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

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

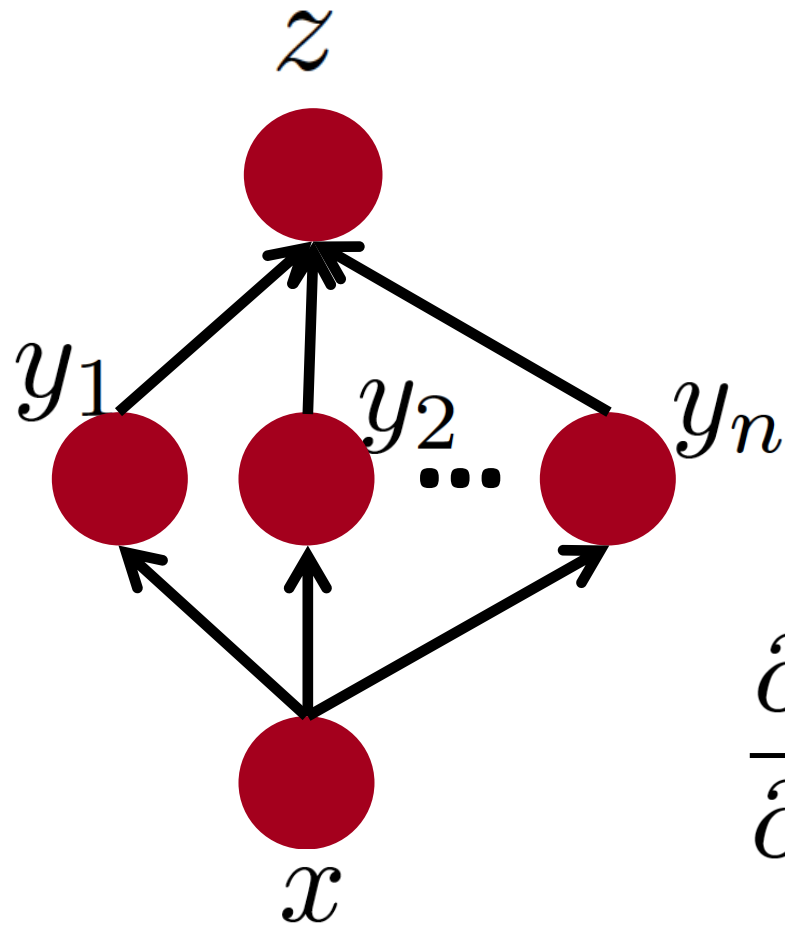
$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

# Multiple Paths Chain Rule



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

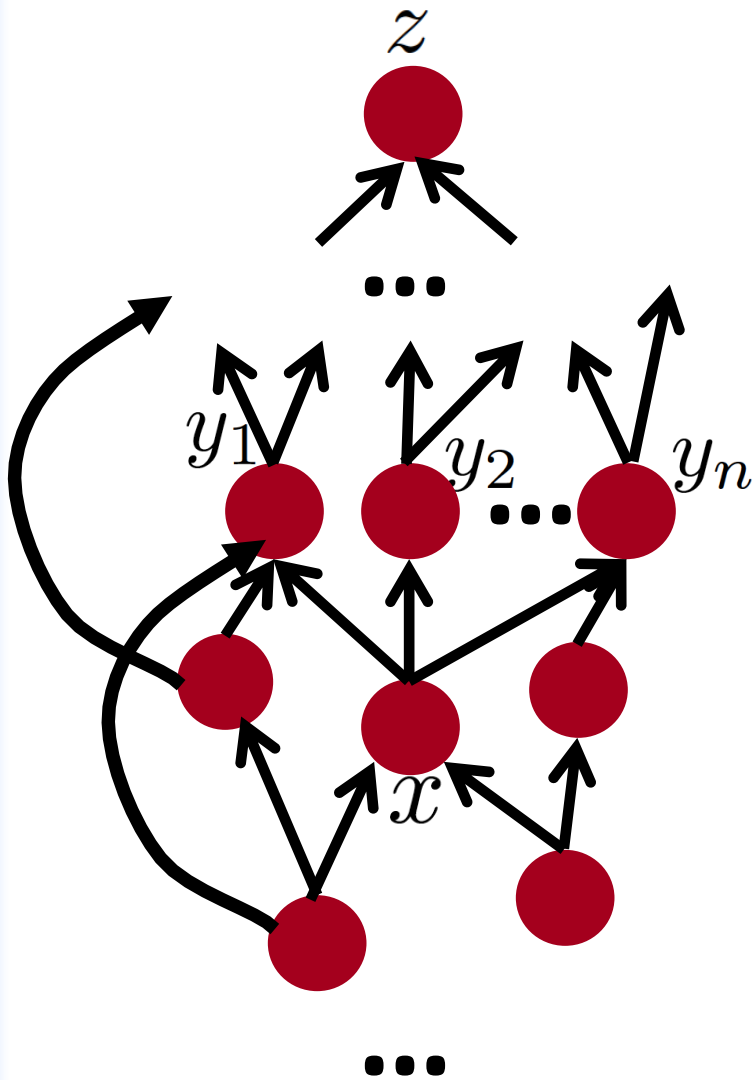
# Multiple Paths Chain Rule - General



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



# Chain Rule in Flow Graph

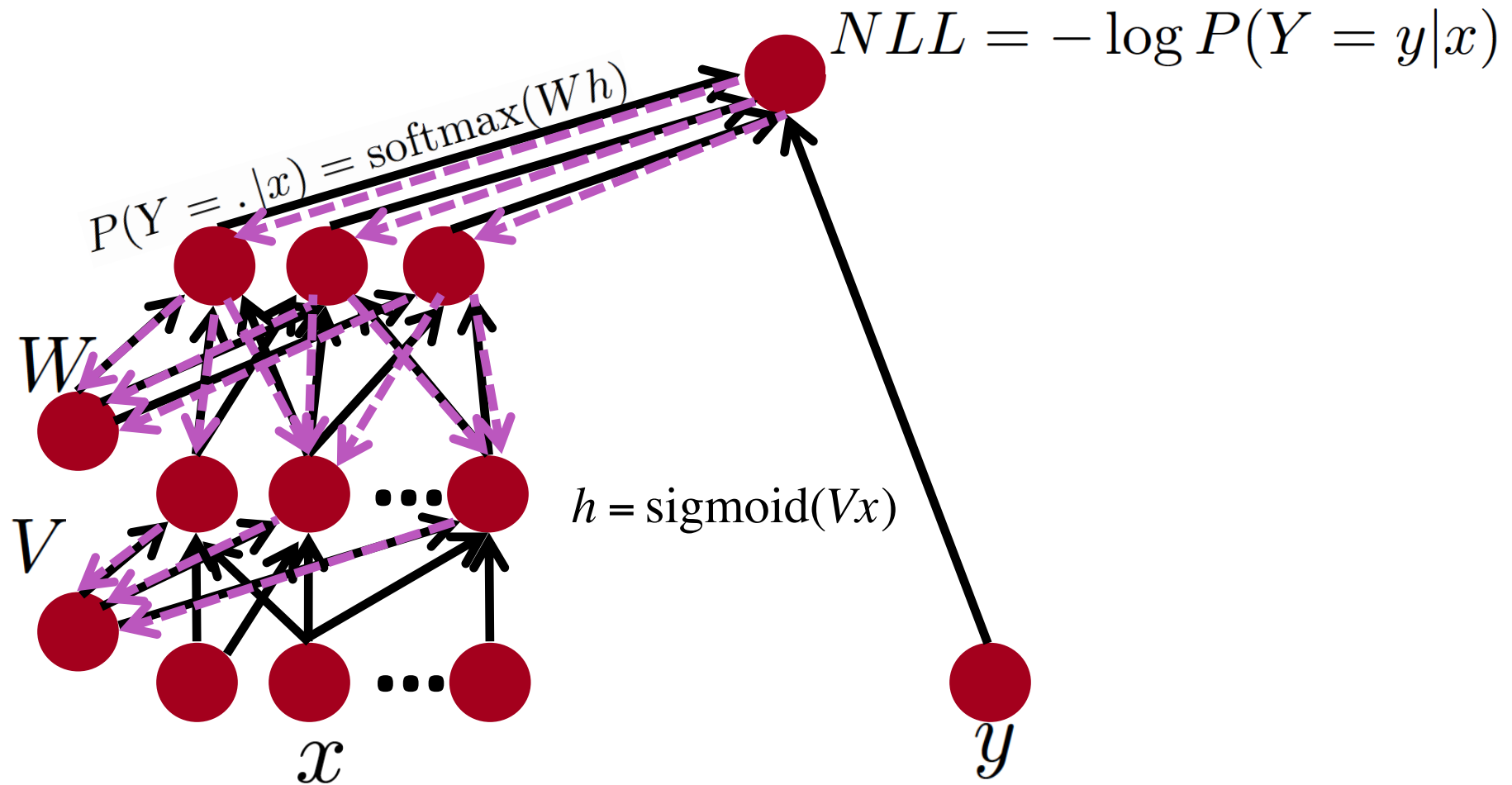


Flow graph: any directed acyclic graph  
node = computation result  
arc = computation dependency

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

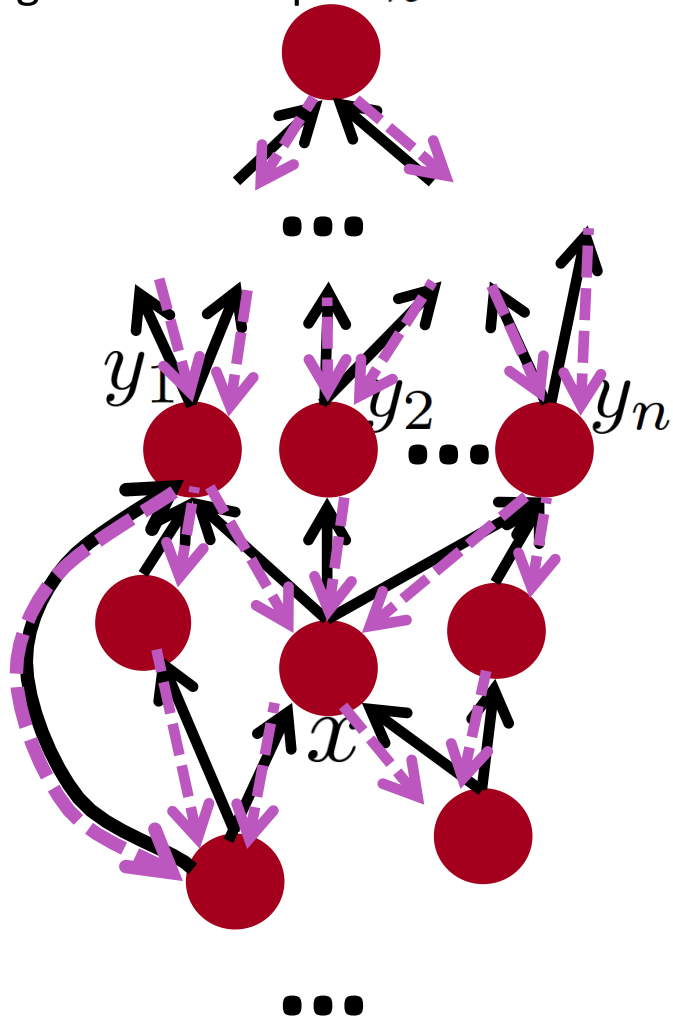
$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$$

# Back-Prop in Multi-Layer Net



# Back-Prop in General Flow Graph

Single scalar output  $z$

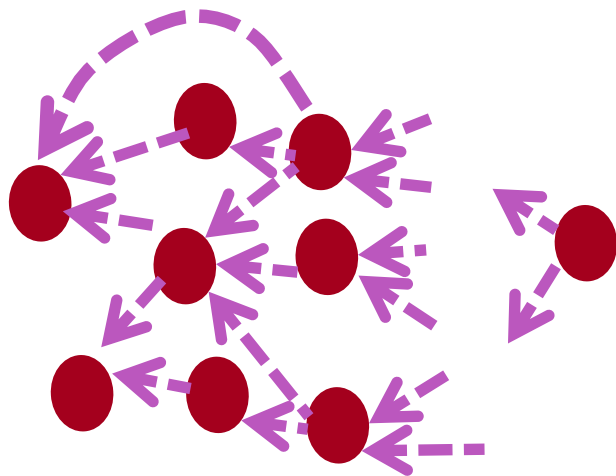
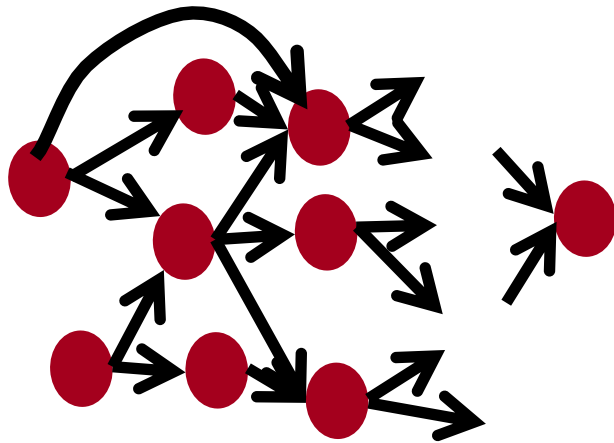


1. Fprop: visit nodes in topo-sort order
  - Compute value of node given predecessors
2. Bprop:
  - initialize output gradient = 1
  - visit nodes in reverse order:
    - Compute gradient wrt each node using gradient wrt successors

$\{y_1, y_2, \dots, y_n\}$  = 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
  - See Theano (Python)

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

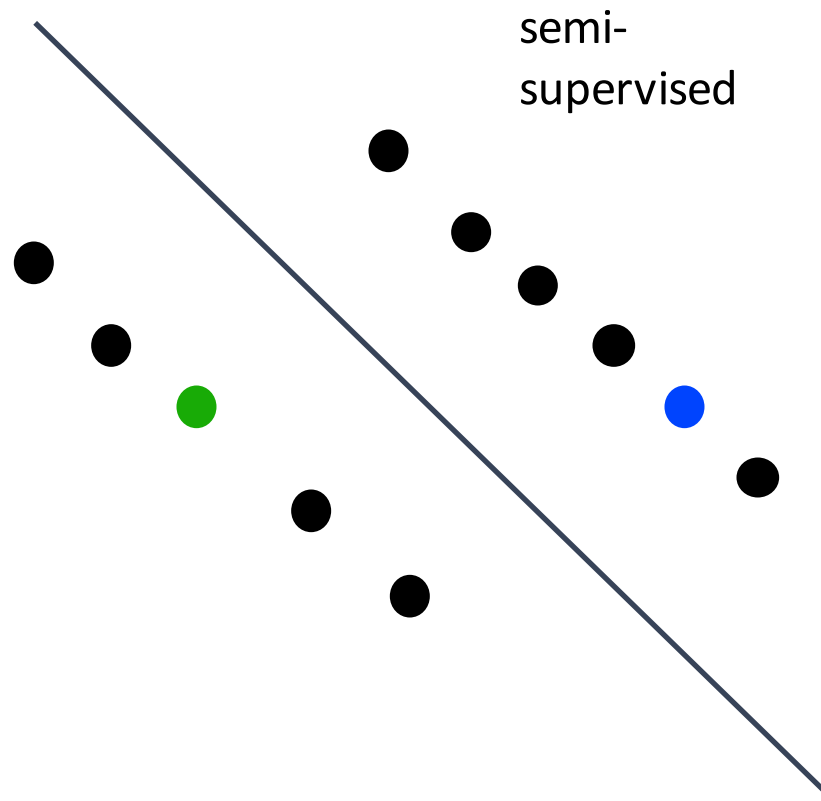
- Hypothesis:  $P(c|x)$  can be more accurately computed using shared structure with  $P(x)$

purely  
supervised



# Semi-Supervised Learning

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





# 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

