# CSE 5243 INTRO. TO DATA MINING

## Word Embedding

Yu Su, CSE@The Ohio State University

# How to let a computer understand meaning?

A cat sits on a mat.

#_$@^_&*^&_()_@_+@^=

# Distributional semantics

□  You can get a lot of value by representing a word by means of its neighbors (context)

<span style="color:blue">"You shall know a word by the company it keeps"</span>

(J. R. Firth 1957: 11)

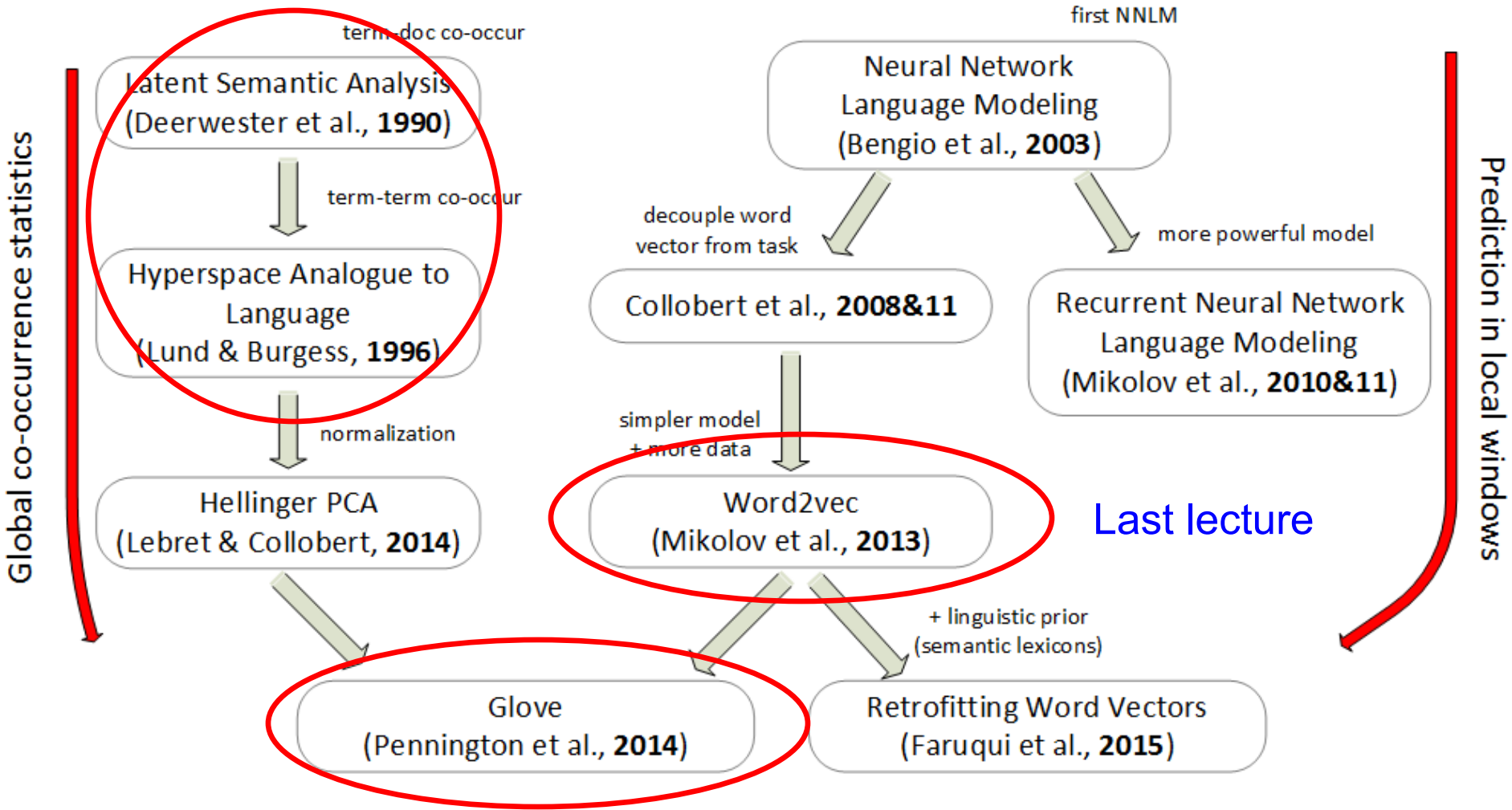□  One of the most successful ideas of modern statistical NLP

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* ↗

# History of word embedding



COUNT!

PREDICT!

term-doc co-occur

first NNLM

Latent Semantic Analysis
(Deerwester et al., **1990**)

Neural Network
Language Modeling
(Bengio et al., **2003**)

term-term co-occur

decouple word
vector from task

more powerful model

Hyperspace Analogue to
Language
(Lund & Burgess, **1996**)

Collobert et al., **2008&11**

Recurrent Neural Network
Language Modeling
(Mikolov et al., **2010&11**)

normalization

simpler model
+ more data

Hellinger PCA
(Lebret & Collobert, **2014**)

Word2vec
(Mikolov et al., **2013**)

Last lecture

+ linguistic prior
(semantic lexicons)

Glove
(Pennington et al., **2014**)

Retrofitting Word Vectors
(Faruqui et al., **2015**)

Global co-occurrence statistics

Prediction in local windows

# History of word embedding

**COUNT!**

**PREDICT!**

Global co-occurrence statistics

term-doc co-occur

Latent Semantic Analysis
(Deerwester et al., **1990**)

term-term co-occur

Hyperspace Analogue to
Language
(Lund & Burgess, **1996**)

normalization

Hellinger PCA
(Lebret & Collobert, **2014**)

first NNLM

Neural Network
Language Modeling
(Bengio et al., **2003**)

decouple word
vector from task

more powerful model

Collobert et al., **2008&11**

Recurrent Neural Network
Language Modeling
(Mikolov et al., **2010&11**)

simpler model
+ more data

Word2vec
(Mikolov et al., **2013**)

This lecture

Prediction in local windows

+ linguistic prior
(semantic lexicons)

Glove
(Pennington et al., **2014**)

Retrofitting Word Vectors
(Faruqui et al., **2015**)

5

# Different embeddings are based on different priors

Latent semantic analysis ⟹

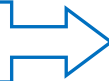"Words occur in same documents should be similar"

Word2vec ⟹

"Words occur in similar contexts should be similar"

Neural Network Language Modeling ⟹

"Word vectors should give plausible sentences high probability"

Collabert et al., 2008 & 2011 ⟹

"Word vectors should facilitate downstream classification tasks"

Faruqui et al., 2015 ⟹

"Words should follow linguistic constraints from semantic lexicons"

# Latent semantic analysis: word-doc occurrence matrix

| Terms | Docs | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
| data | 1 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 1 | 2 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| examples | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| introduction | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| mining | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| network | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 |
| package | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

- Word-doc occurrence matrix will give general topics,
  e.g., all sports words will have similar entries
- Apply SVD for dimensionality reduction

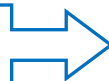# Different embeddings are based on different priors

Latent semantic analysis ⟹

"Words occur in same documents should be similar"

Word2vec ⟹

"Words occur in similar contexts should be similar"

Neural Network Language Modeling ⟹

"Word vectors should give plausible sentences high probability"

Collabert et al., 2008 & 2011 ⟹

"Word vectors should facilitate downstream classification tasks"

Faruqui et al., 2015 ⟹

"Words should follow linguistic constraints from semantic lexicons"

# Word2vec: "Words occur in similar contexts should be similar"

I just played with my dog.
I just played with my cat.
My dog likes to sleep on my bed.
My cat likes to sleep on my bed.

- Word2vec will adjust the vector of a word to be similar to the vectors of its context words

- Words with similar contexts thus end up with similar vectors

# Different embeddings are based on different priors

Latent semantic analysis ⟹

"Words occur in same documents should be similar"

Word2vec ⟹

"Words occur in similar contexts should be similar"

Neural Network Language Modeling ⟹

"Word vectors should assign high probability to plausible sentences"

Collabert et al., 2008 & 2011 ⟹

"Word vectors should facilitate downstream classification tasks"

Faruqui et al., 2015 ⟹

"Words should follow linguistic constraints from semantic lexicons"

# Probabilistic Language Modeling

- Goal: assign a probability to a sentence
  - Machine Translation:
    - Source sentence: 今晚大风
    - P(**large** winds tonight) < P(**strong** winds tonight)
  - Spell Correction
    - The office is about fifteen **minuets** from my house
      - P(about fifteen **minutes** from) > P(about fifteen **minuets** from)
  - Speech Recognition
    - P(I saw a van) >> P(eyes awe of an)
  - +Summarization, question answering, etc.

# Probabilistic Language Modeling

- Goal: compute the probability of a sentence or a sequence of words:

$$P(w_1^m) = P(w_1, w_2, ..., w_m)$$

- How to compute the joint probability?

$$P(a, dog, is, running, in, a, room)$$

- Chain rule:

$$P(w_1, w_2, ..., w_m) = P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1, w_2)...P(w_m \mid w_1, ...w_{m-1})$$

$$P(a, dog, is, running) =$$

$$P(a)P(dog \mid a)P(is \mid a, dog)P(running \mid a, dog, is)$$

# Probabilistic Language Modeling

$$P(w_1, w_2, ..., w_m) = \prod_{t}^{m-1} P(w_t \mid w_1, ... w_{t-1})$$

- Key: $P(w_t \mid w_1, ... w_{t-1})$

- Just count? Exponential number of entries and sparsity.

- Markov assumption:

$$P(w_t \mid w_1, ... w_{t-1}) \approx P(w_t \mid w_{t-n+1}, ... w_{t-1})$$

# Probabilistic Language Modeling

□ N-gram (bigram)

$$P(running \mid a, dog, is) \approx P(running \mid is) = \frac{count(is, running)}{count(is)}$$

□ What's the problem?
  ▪ Small context window (typically bigram or trigram)
  ▪ Not utilizing word similarity
    ▪ Seeing "A dog is running in a room" should increase probability of
    ▪ "The dog is walking in a room" and
    ▪ "A cat is running in the room" and
    ▪ "Some cats are running in the room"

□ Solution: Neural Network Language Modeling!

# Neural Network Language Model

Learn $P(w_t \mid w_{t-n+1}, \ldots w_{t-1})$



$i$-th output $= P(w_t = i \mid context)$

Softmax

softmax

most computation here

Fully connected
non-linear layer

tanh

Projection

$C(w_{t-n+1})$        $C(w_{t-2})$    $C(w_{t-1})$

Table
look−up
in $C$

Matrix $C$

shared parameters
across words

index for $w_{t-n+1}$        index for $w_{t-2}$        index for $w_{t-1}$

# The Lookup Table

- Each word in vocabulary maps to a vector in $\mathbb{R}^d$

- LookupTable: input of the i[th] word is

$$x = (0, 0, \ldots, 1, 0, \ldots, 0) \quad 1 \text{ at position } i$$

In the original space words are orthogonal.

cat = (0,0,0,0,0,0,0,0,0,1,0,0,0,0, …)
dog = (0,0,1,0,0,0,0,0,0,0,0,0,0,0, …)

To get the $\mathbb{R}^d$ embedding vector for the word we multiply *Cx* where *C* is a *d x D* matrix with *D* words in the vocabulary

*C* contains the word vectors!

# Neural Network Language Model

$$P(w_t = i) = \frac{\exp(y_i)}{\sum_{j=1}^{D} \exp(y_j)}$$



*i*-th output = $P(w_t = i \mid context)$

softmax

most computation here

tanh

$C(w_{t-n+1})$   $C(w_{t-2})$   $C(w_{t-1})$

Table
look-up
in $C$

Matrix $C$
shared parameters
across words

index for $w_{t-n+1}$   index for $w_{t-2}$   index for $w_{t-1}$

softmax

$$y = Uz + b_2$$

output

$$z = \tanh(Hx + b_1)$$

non-linearity

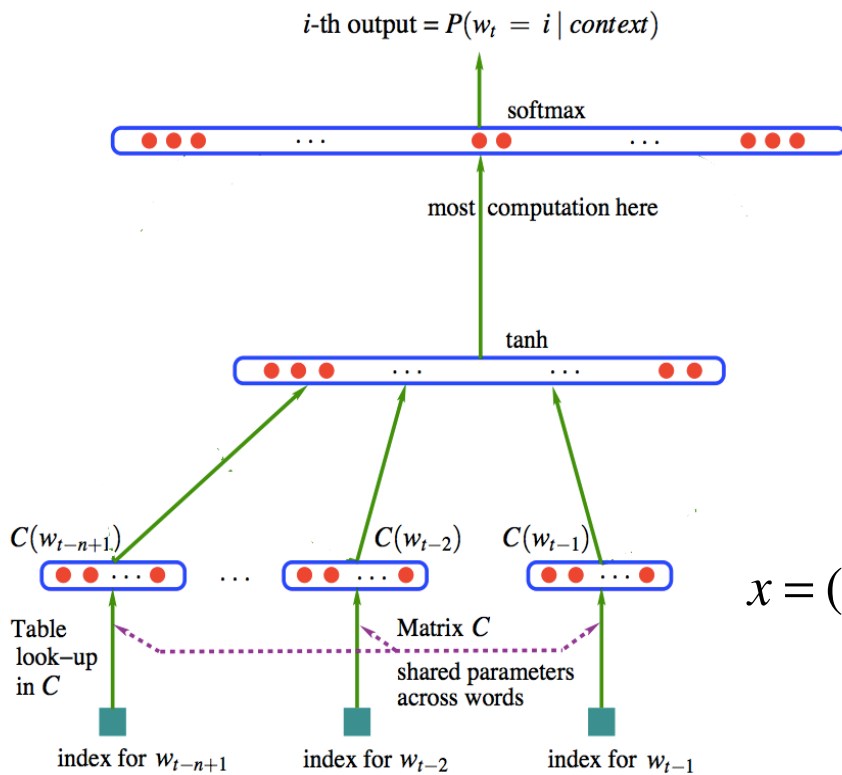$$x = (Cw_{t-n+1}, Cw_{t-n+2}, ..., Cw_{t-1})^T$$

projection

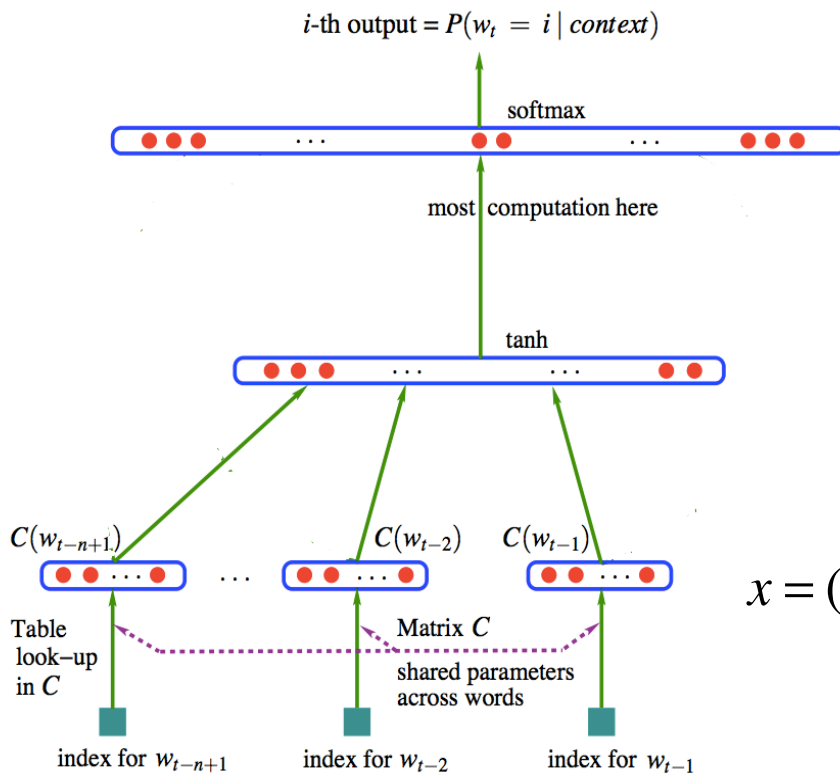$$w_{t-n+1}, w_{t-n+2}, ..., w_{t-1}$$

$d$ : word vector dimensionality

n: window size

D: vocabulary size

h: # of hidden units

Dimensionality of each layer?

$$P(w_t = i) = \frac{\exp(y_i)}{\sum_{j=1}^{D} \exp(y_j)}$$



i-th output = $P(w_t = i \mid context)$

softmax

most computation here

tanh

$C(w_{t-n+1})$    $C(w_{t-2})$    $C(w_{t-1})$

Table
look-up
in $C$

Matrix $C$

shared parameters
across words

index for $w_{t-n+1}$      index for $w_{t-2}$      index for $w_{t-1}$

softmax

$y = Uz + b_2$    $D$

output

$z = \tanh(Hx + b_1)$    $h$

non-linearity

$x = (Cw_{t-n+1}, Cw_{t-n+2}, ..., Cw_{t-1})^T$    $n * d$

projection

$w_{t-n+1}, w_{t-n+2}, ..., w_{t-1}$

$d$ : word vector dimensionality

n: window size

D: vocabulary size

h: # of hidden units

# of parameters in each layer?

$$P(w_t = i) = \frac{\exp(y_i)}{\sum_{j=1}^{D} \exp(y_j)}$$

softmax

$y = Uz + b_2$

$h * D + D$

output

$z = \tanh(Hx + b_1)$

$n * d * h + h$

non-linearity

$x = (Cw_{t-n+1}, Cw_{t-n+2}, ..., Cw_{t-1})^T$

$n * d$

projection

$w_{t-n+1}, w_{t-n+2}, ..., w_{t-1}$

# Training

- All free parameters

$$\boldsymbol{\theta} = (C, H, U, b_1, b_2)$$

- Backpropagation + Stochastic Gradient Ascent:

Costly!

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \varepsilon \frac{\partial \log P(w_t \mid w_{t-n+1}, \ldots, w_{t-1})}{\partial \boldsymbol{\theta}}$$

# Speed up training

- Most computations are at the output layer
  - In order to compute the normalization term of softmax, we have to compute the $y_i$ for every word!
  - Cost (almost) linear to vocabulary size.
  - Same problem in Skip-gram

- Solutions: Approximate the normalized probability
  - Negative sampling
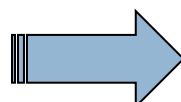  - Noise contrastive estimation
  - Hierarchical softmax
  - …

# Speed up training

- Most computations are at the output layer
  - In order to compute the normalization term of softmax, we have to compute the $y_i$ for every word!
  - Cost (almost) linear to vocabulary size.
  - Same problem in Skip-gram

- Solutions: Approximate the normalized probability
  - **Negative sampling**
  - Noise contrastive estimation
  - Hierarchical softmax
  - …

# Refresher: Skip-gram

- Given the central word, predict surrounding words in a window of length c

- Objective function:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^{T} \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j}|w_t)$$

- Softmax:

$$p(O|I) = \frac{\exp(v_O'^T v_I)}{\sum_{w \in V} \exp(v_w'^T v_I)}$$

$$\frac{\partial \log p(O|I)}{\partial v_I} = v'_O - \sum_w p(w|I)v'_w$$

# Negative sampling

- I: central word. O: a context word
- Original: Maximize $p(O|I,\theta)$
- We will derive an alternative which is less costly to compute
- Does pair (I,O) really come from the training data?

$$\theta = \arg\max_{\theta} p(D=1|I,O,\theta)$$

$$\text{where } p(D=1|I,O,\theta) = \sigma(v_I^T v_O') = \frac{1}{1+e^{-v_I^T v_O'}}$$

- Trivial solution: same (long enough) vector for all words
- Contrast with negative words!

# Negative sampling

- Solution: randomly sample k negative words $w_i$ from a noise distribution, assume $(I, w_i)$ are incorrect pairs

- $I =$ "is", $O =$ "running", $w_1 =$ "walk", $w_2 =$ "do", etc.

$$\text{maximize } p(D = 1 | I, O, \theta) \bullet \prod_{i=1}^{k} p(D = 0 | I, w_i, \theta)$$

$$\text{or } \log \sigma(v_I^T v_O') + \sum_{i=1}^{k} \mathbb{E}_{w_i \sim P_n(w)} [\log \sigma(-v_I^T v_{w_i}')]$$

$$\text{where } P_n(w) = \frac{U(w)^{3/4}}{Z}, U(w) \text{ the unigram distribution}$$

# Different embeddings are based on different priors

Latent semantic analysis →

"Words occur in same documents should be similar"

Word2vec →

"Words occur in similar contexts should be similar"

Neural Network Language Modeling →

"Word vectors should give plausible sentences high probability"

Collobert et al., 2008 & 2011 →

"Word vectors should facilitate downstream classification tasks"

Faruqui et al., 2015 →

"Words should follow linguistic constraints from semantic lexicons"

# What to get from this work

- How to supervise the learning of word embedding using external classification tasks

- How to do semi-supervised learning of word embedding

- How to apply word vectors and neural networks in other traditional NLP tasks

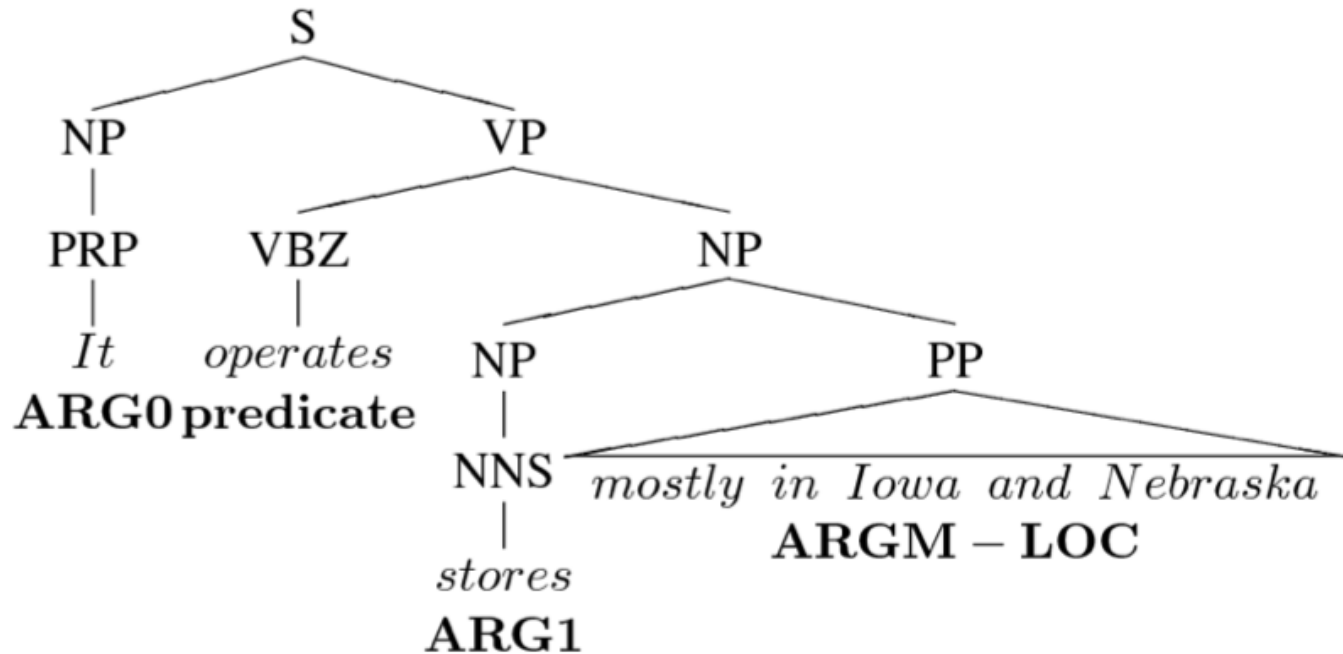# Embedding for other NLP tasks (Collobert et al., 2008&11)

- Part-Of-Speech Tagging (POS): syntactic roles (noun, adverb...)
- Chunking: syntactic constituents (noun phrase, verb phrase...)
- Name Entity Recognition (NER): person/company/location...
- Semantic Role Labeling (SRL):

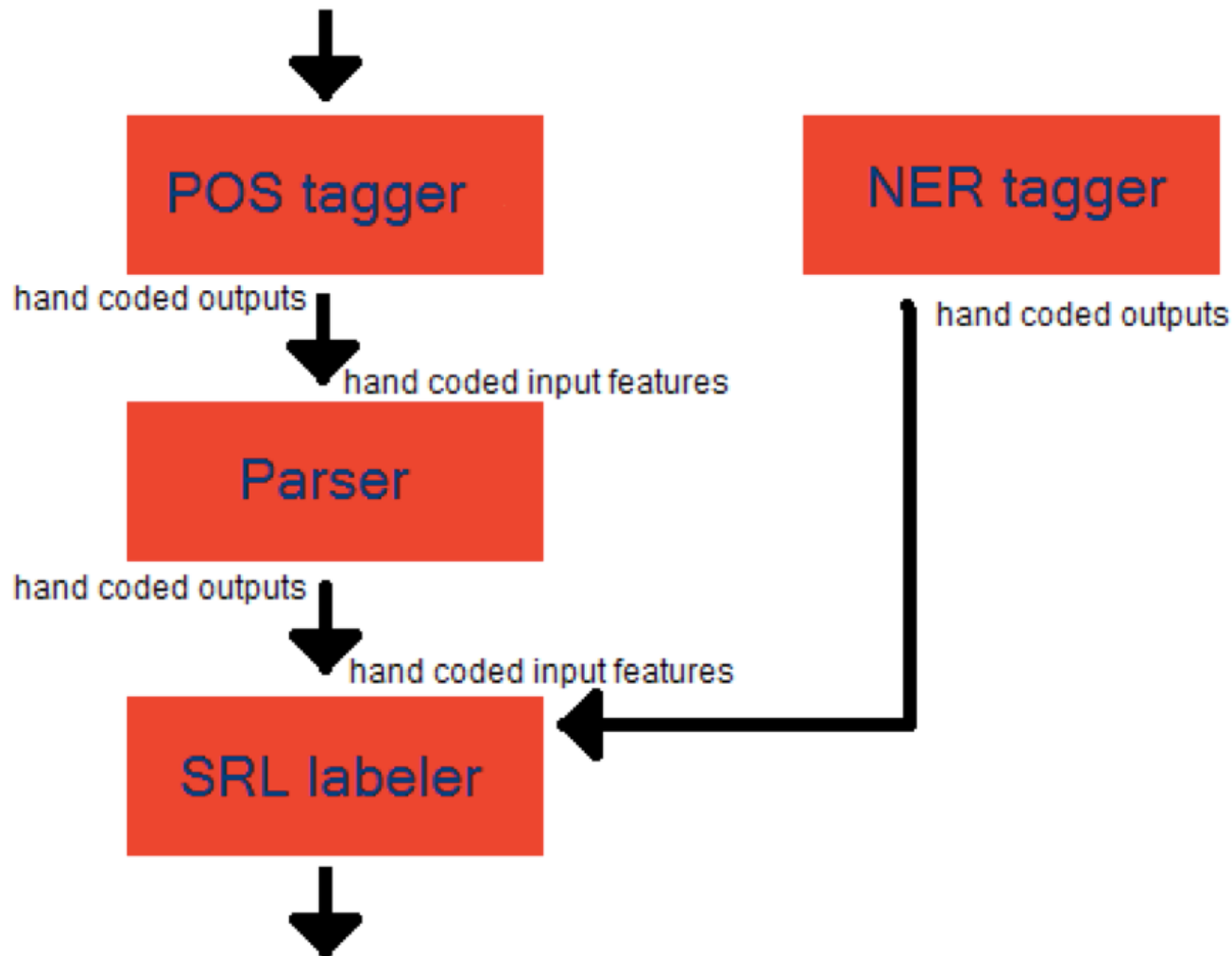$[John]_{ARG0}$ $[ate]_{REL}$ $[the\ apple]_{ARG1}$ $[in\ the\ garden]_{ARGM-LOC}$

# The Large-scale Feature Engineering Way



- Extract **hand-made features** e.g. from the parse tree
- Disjoint: all tasks trained separately,  Cascade features
- Feed these features to a shallow classifier like SVM

# The sub-optimal cascade

# NLP: Large scale machine learning

## Goals

- Task-specific engineering limits NLP scope
- Can we find unified hidden representations?
- Can we build unified NLP architecture?

## Means

- Start from scratch: forget (most of) NLP knowledge
- Compare against classical NLP benchmarks
- Our dogma: avoid task-specific engineering

# The big picture
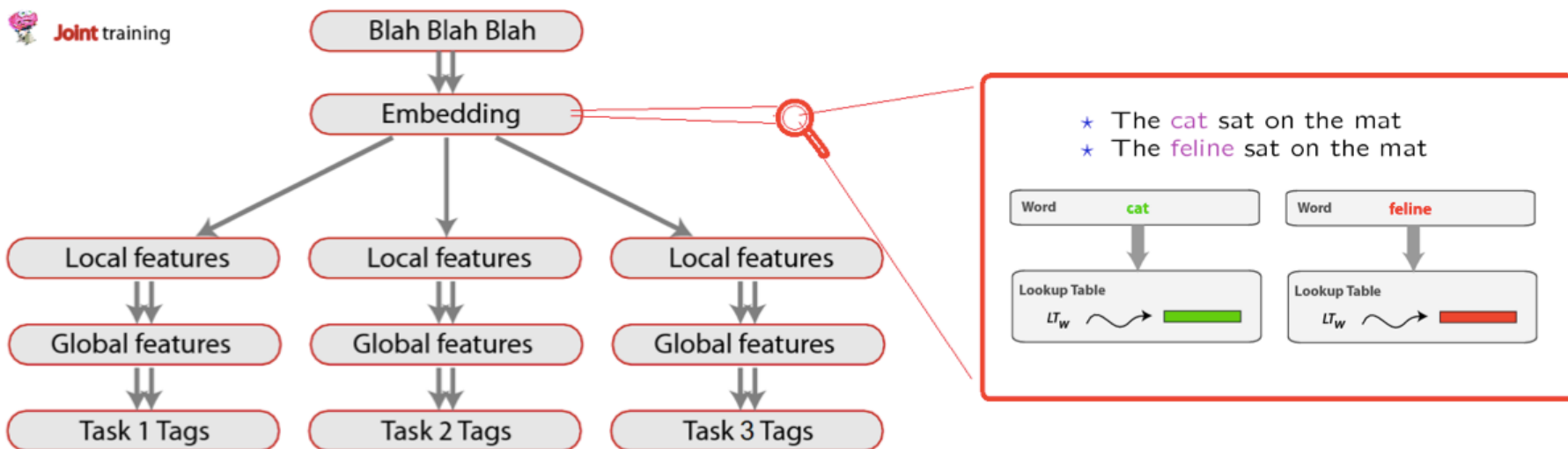
A unified architecture for all NLP (labeling) tasks:

| Sentence: | Felix | sat | on | the | mat | . |
|---|---|---|---|---|---|---|
| POS: | NNP | VBD | IN | DT | NN | . |
| CHUNK: | NP | VP | PP | NP | NP-I | . |
| NER: | PER | - | - | - | - | - |
| SRL: | ARG1 | REL | ARG2 | ARG2-I | ARG2-I | - |

# Different embeddings are based on different priors

Latent semantic analysis ⟹

"Words occur in same documents should be similar"

Word2vec ⟹

"Words occur in similar contexts should be similar"

Neural Network Language Modeling ⟹

"Word vectors should give plausible sentences high probability"
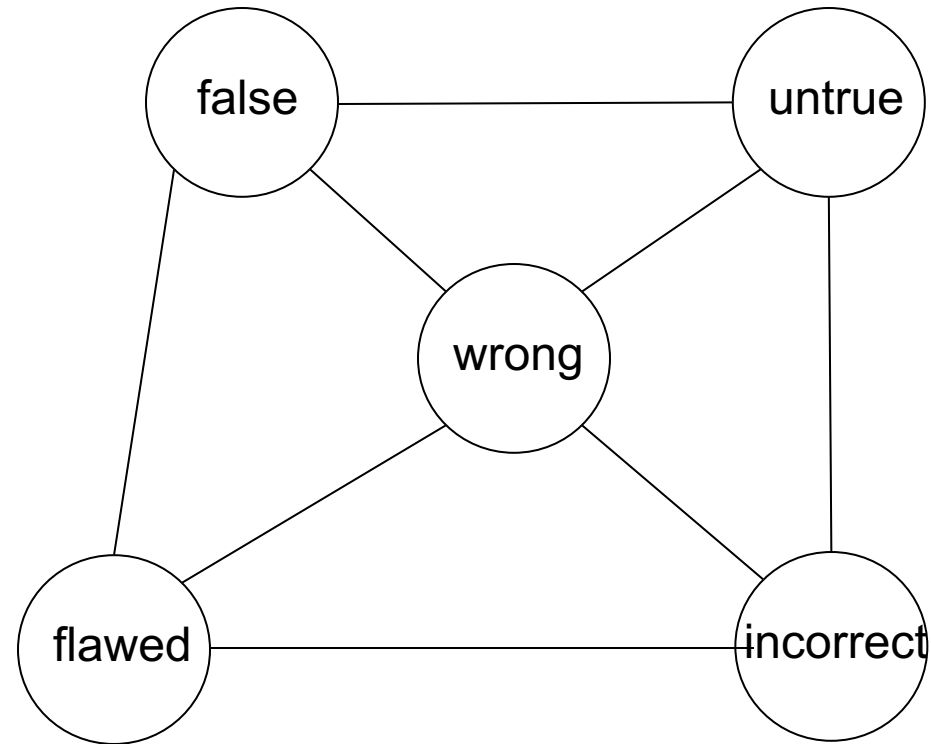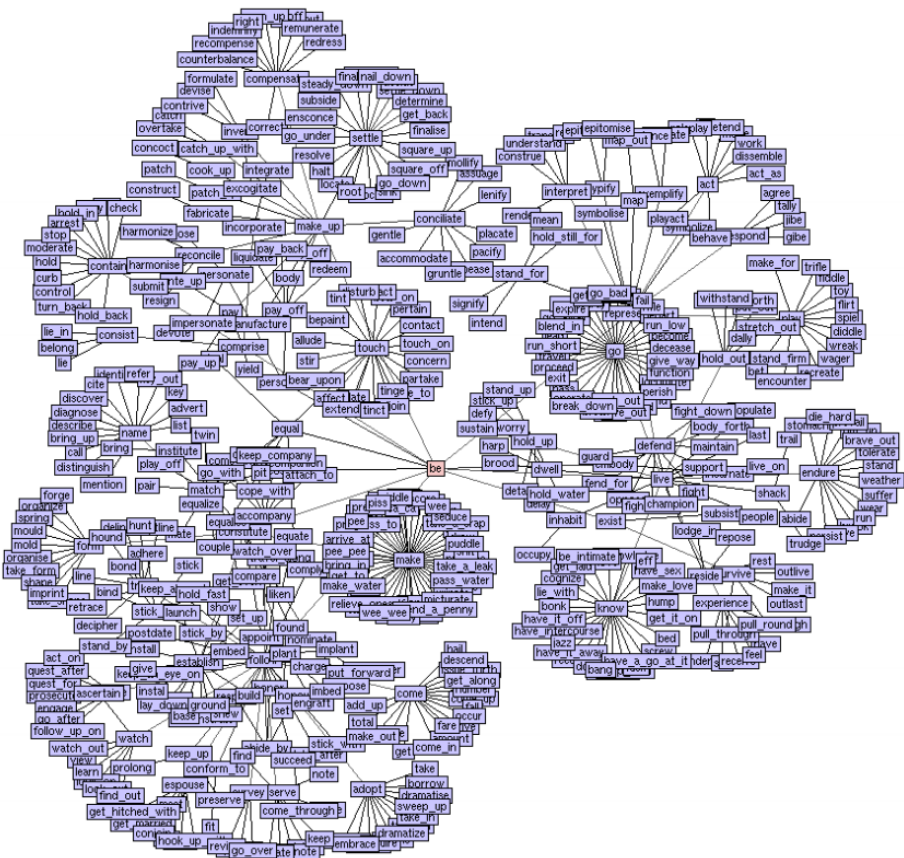
Collabert et al., 2008 & 2011 ⟹

"Word vectors should facilitate downstream classification tasks"

Faruqui et al., 2015 ⟹

"Words should follow linguistic constraints from semantic lexicons"
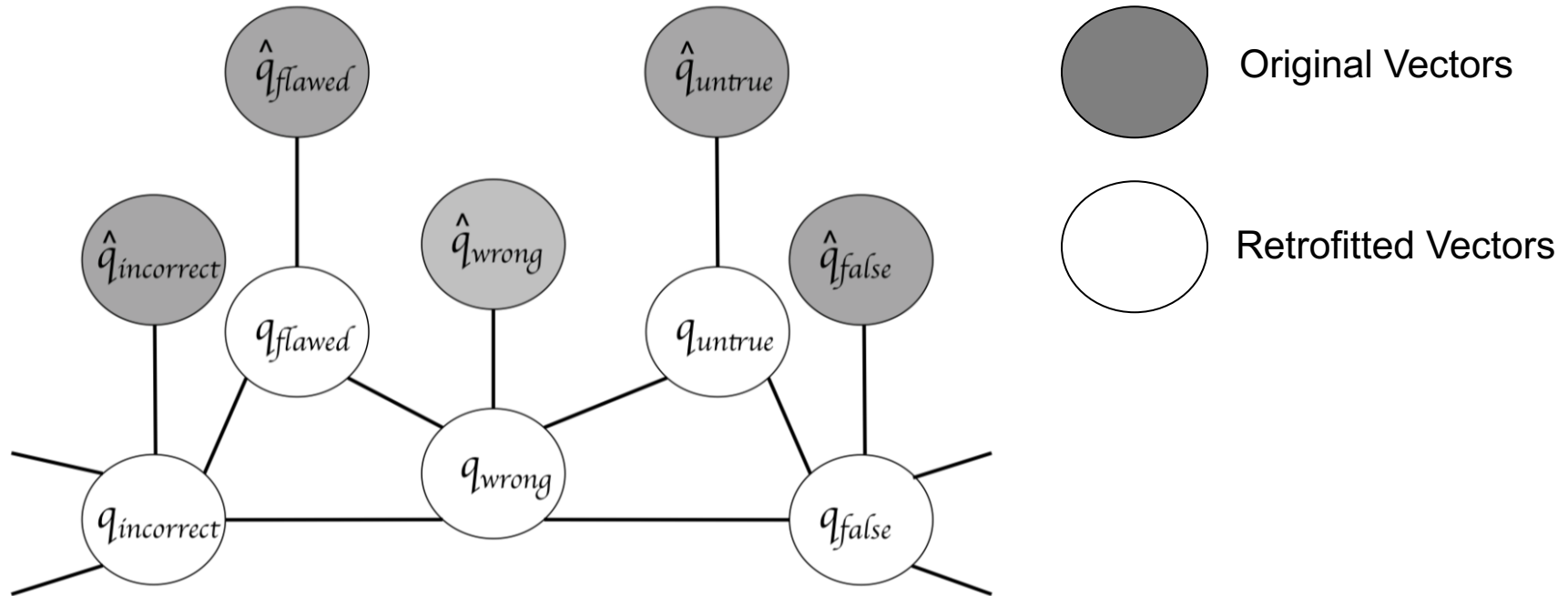
# Semantic lexicon: WordNet

# Retrofitting word vectors to semantic lexicons (NAACL'15)

□ Incorporates information from lexicons in word vectors

□ Post-processing approach

□ Applicable to **any** word embedding method

□ Applicable to **any** lexicon

# Retrofitting

Original Vectors

Retrofitted Vectors

$$\Psi(Q) = \sum_{i=1}^{n} \left[ \alpha_i \| q_i - \hat{q}_i \|^2 + \sum_{(i,j) \in E} \beta_{ij} \| q_i - q_j \|^2 \right]$$

# Semantic lexicons used in this work

- ☐ PPDB: Lexical paraphrases obtained from parallel texts
- ☐ WordNet: synonyms, hypernyms and hyponyms
- ☐ FrameNet: Cause_change_of_position -> push=raise=growth

| Lexicon | Words | Edges |
|---|---|---|
| PPDB | 102,902 | 374,555 |
| WordNet$_{syn}$ | 148,730 | 304,856 |
| WordNet$_{all}$ | 148,730 | 934,705 |
| FrameNet | 10,822 | 417,456 |

Table 1. Approximate size of the graphs obtained from different lexicons

# Experiment results

|  | Word Similarity | | | Synonym Selection | Syntactic Analysis | Sentiment Analysis |
|---|---|---|---|---|---|---|
| Lexicon | MEN-3k | RG-65 | WS-353 | TOEFL | SYN-REL | SA |
| Glove | 73.7 | 76.7 | 60.5 | 89.7 | 67.0 | 79.6 |
| +PPDB | 1.4 | 2.9 | −1.2 | **5.1** | −0.4 | **1.6** |
| +WN$_{syn}$ | 0.0 | 2.7 | 0.5 | **5.1** | −12.4 | 0.7 |
| +WN$_{all}$ | **2.2** | **7.5** | **0.7** | 2.6 | −8.4 | 0.5 |
| +FN | −3.6 | −1.0 | −5.3 | 2.6 | −7.0 | 0.0 |
| SG | 67.8 | 72.8 | 65.6 | 85.3 | 73.9 | 81.2 |
| +PPDB | **5.4** | 3.5 | **4.4** | **10.7** | −2.3 | **0.9** |
| +WN$_{syn}$ | 0.7 | 3.9 | 0.0 | 9.3 | −13.6 | 0.7 |
| +WN$_{all}$ | 2.5 | **5.0** | 1.9 | 9.3 | −10.7 | −0.3 |
| +FN | −3.2 | 2.6 | −4.9 | 1.3 | −7.3 | 0.5 |

Original Embedding

Semantic Lexicons

# In this lecture…

- More types of supervision used in training word embedding
  - Language modeling
  - NLP labeling tasks
  - Semantic lexicons

- Ways to speed up
  - E.g., negative sampling
  - Necessary for training on huge text corpora
  - Scale up from hundreds of millions to hundreds of billions

- How word embeddings help other NLP tasks