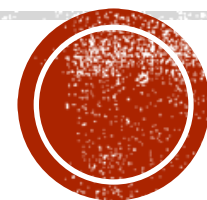


# THE LINEAR ALGEBRAIC STRUCTURE OF WORD MEANINGS

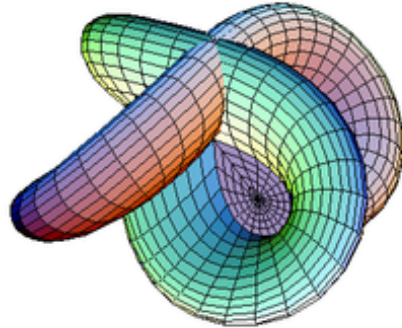
Tengyu Ma



Joint works with Sanjeev Arora, Yuanzhi Li, Yingyu  
Liang, and Andrej Risteski

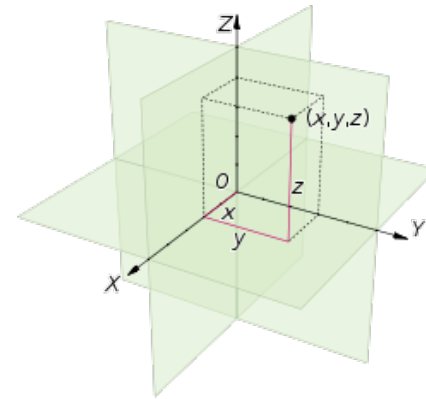
Princeton University

# EMBEDDINGS (IN MACHINE LEARNING)



$$x \in \mathcal{X}$$

complicated space



$$v_x \in \mathbb{R}^d$$

Euclidean space with  
**meaningful** inner products

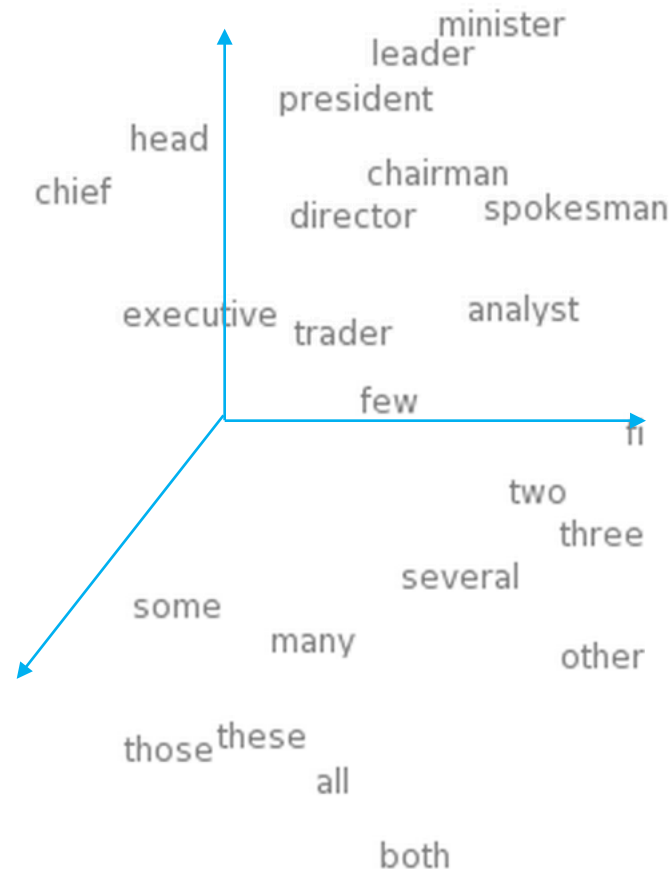
- Kernel methods  $\xrightarrow{\text{handcrafted lifting}}$  Linearly separable
- Neural nets  $\xrightarrow{\text{trained neural nets}}$  Multi-class linear classifier



# WORD EMBEDDING

Vocabulary =  
{ 60k most frequent words }

→  $\mathbb{R}^{300}$



Goal: Embedding captures semantics information  
(via linear algebraic operations)

- inner products characterize similarity
  - similar words have large inner products
- differences characterize relationship
  - analogous pairs have similar differences
- more?

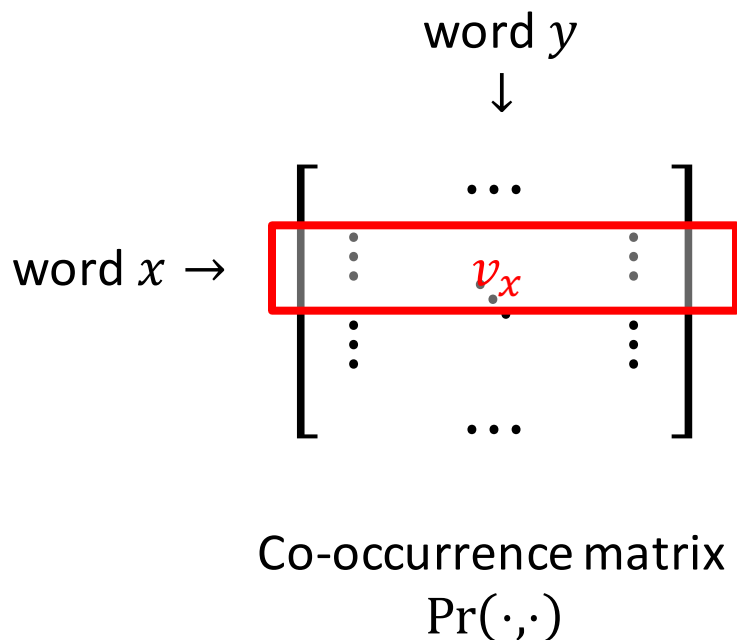
picture: [Chris Olah's blog](#)



# WORD EMBEDDING, AN OLD IDEA

Meaning of a word is determined by words it co-occurs with.

(*Distributional hypothesis of meaning*, [Harris'54], [Firth'57])



- $\Pr(x, y) \triangleq$  prob. of co-occurrences of  $x, y$  in a window of size 5
- $\langle v_x, v_y \rangle$  - a good measure of similarity of  $(x, y)$  [Lund-Burgess'96]
- $v_x$  = row of **entry-wise** square-root of co-occurrence matrix [Rohde et al'05]
- $v_x$  = row of PMI( $x, y$ ) =  $\log \frac{\Pr[x, y]}{\Pr[x] \Pr[y]}$  matrix [Church-Hanks'90]



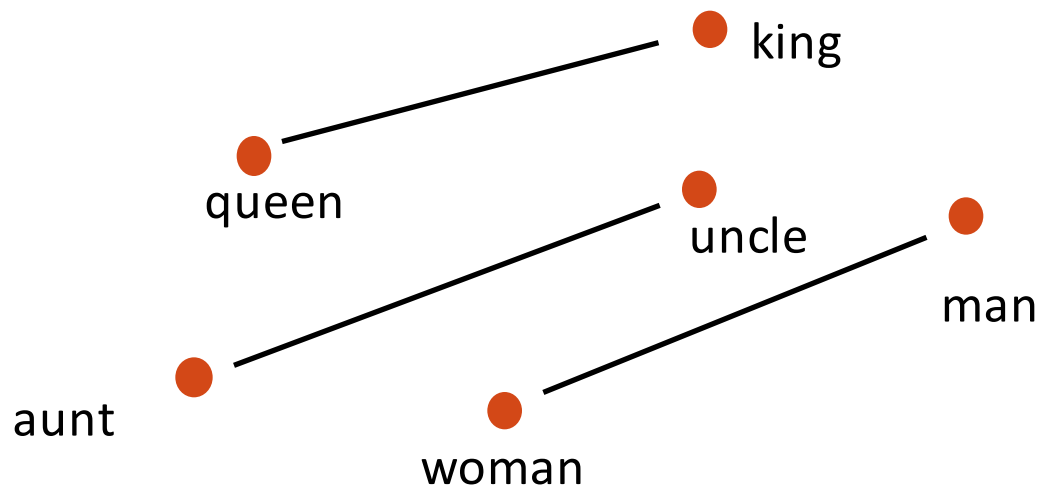
# LINEAR STRUCTURE AFTER NON-LINEAR EMBEDDING

Algorithm [Levy-Goldberg]: (dimension-reduction version of [Church-Hanks'90])

- Compute  $\text{PMI}(x, y) = \log \frac{\text{Pr}[x, y]}{\text{Pr}[x] \text{Pr}[y]}$
- Take rank-300 SVD (best rank-300 approximation) of PMI
  - $\Leftrightarrow$  Fit  $\text{PMI}(x, y) \approx \langle v_x, v_y \rangle$  (with squared loss), where  $v_x \in \mathbb{R}^{300}$

- “Linear structure” in the found  $v_x$ 's :

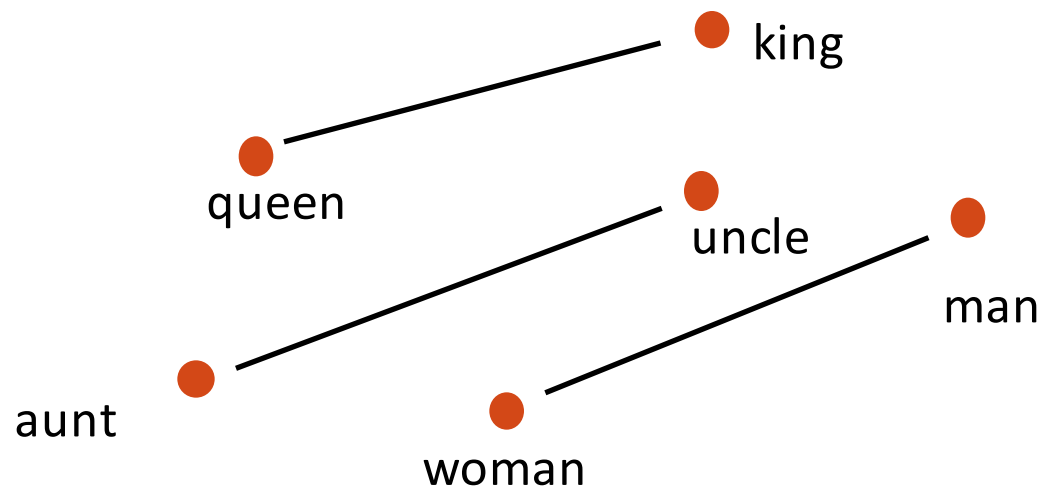
$$v_{\text{woman}} - v_{\text{man}} \approx v_{\text{queen}} - v_{\text{king}} \approx v_{\text{uncle}} - v_{\text{aunt}} \approx \dots$$



# APPLICATIONS/TESTS : SOLVING ANALOGY TASKS

- Questions: woman: man  
queen: ?  
aunt: ?

- Answers:  $king = \operatorname{argmin}_w \|(v_{queen} - v_w) - (v_{woman} - v_{man})\|$   
 $aunt = \operatorname{argmin}_w \|(v_{uncle} - v_w) - (v_{woman} - v_{man})\|$



# NON-LINEAR EMBEDDING METHODS

➤ recurrent neural network based model [Mikolov et al'12]

➤ word2vec [Mikolov et al'13] :

$$\Pr[x_{i+6} | x_{i+1}, \dots, x_{i+5}] \propto \exp\langle v_{x_{i+6}}, \frac{1}{5} (v_{x_{i+1}} + \dots + v_{x_{i+5}}) \rangle$$

➤ GloVe [Pennington et al'14] :

$$\log \Pr[x, y] \approx \langle v_x, v_y \rangle + s_x + s_y + C$$

➤ [Levy-Goldberg'14] (Previous slide)

$$\text{PMI}(x, y) = \log \frac{\Pr[x, y]}{\Pr[x] \Pr[y]} \approx \langle v_x, v_y \rangle + C$$

Logarithm (or exponential) seems to exclude linear algebra!



# Why co-occurrence statistics + log $\rightarrow$ linear structure

[Levy-Goldberg'13, Pennington et al'14, rephrased]

➤ For most of the words  $\chi$ :

$$\frac{\Pr[\chi \mid \textit{king}]}{\Pr[\chi \mid \textit{queen}]} \approx \frac{\Pr[\chi \mid \textit{man}]}{\Pr[\chi \mid \textit{woman}]}$$

- For  $\chi$  unrelated to gender: LHS, RHS  $\approx 1$
- for  $\chi = \textit{dress}$ , LHS, RHS  $\ll 1$  ; for  $\chi = \textit{John}$ , LHS, RHS  $\gg 1$

➤ It suggests

$$\sum_{\chi} \left( \log \frac{\Pr[\chi \mid \textit{king}]}{\Pr[\chi \mid \textit{queen}]} - \log \frac{\Pr[\chi \mid \textit{man}]}{\Pr[\chi \mid \textit{woman}]} \right)^2 \approx 0$$
$$= \sum_{\chi} \left( (\text{PMI}(\chi, \textit{king}) - \text{PMI}(\chi, \textit{queen})) - (\text{PMI}(\chi, \textit{man}) - \text{PMI}(\chi, \textit{woman})) \right)^2 \approx 0$$

➤ Rows of PMI matrix has “linear structure”

➤ Empirically one can find  $v_w$ 's s.t.  $\text{PMI}(\chi, w) \approx \langle v_{\chi}, v_w \rangle$

➤ Suggestion:  $v_w$ 's also have linear structure





# WHY THESE METHODS CAN WORK?

M1: Why do low-dim vectors capture essence of huge co-occurrence statistics? That is, why is a low-dim fit of PMI matrix even possible?

$$\text{PMI}(x, y) \approx \langle v_x, v_y \rangle \quad (*)$$

➤ NB: PMI matrix is not necessarily PSD.

M2: Why low-dim vectors solves analogy when (\*) is only roughly true?

↑  
empirical fit has 17% error

➤ NB: solving analogy task requires inner products of 6 pairs of word vectors, and that “king” survives against all other words – noise is potentially an issue!

$$\textit{king} = \operatorname{argmax}_w \|(v_{\textit{queen}} - v_w) - (v_{\textit{woman}} - v_{\textit{man}})\|^2$$

➤ Fact: low-dim word vectors have **more accurate** linear structure than the rows of PMI (therefore better analogy task performance).



# OUR INSIGHTS

M1: Why do low-dim vectors capture essence of huge co-occurrence statistics? That is, why is a low-dim fit of PMI matrix even possible?

$$\text{PMI}(x, y) \approx \langle v_x, v_y \rangle \quad (*)$$

A1: Under a generative model (named RAND-WALK) , (\*) provably holds

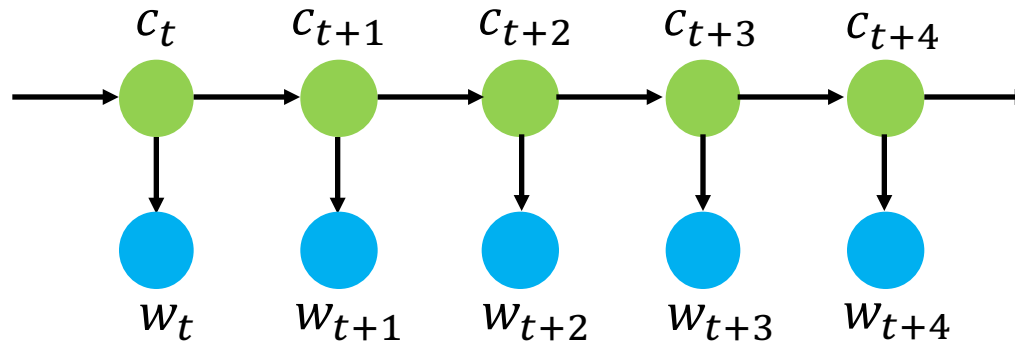
M2: Why low-dim vectors solves analogy when (\*) is only roughly true?

A2: (\*) + isotropy of word vectors  $\Rightarrow$  low-dim fitting reduces noise

(Quite intuitive, though doesn't follow Occam's bound for PAC-learning)



# RAND-WALK: A GENERATIVE MODEL FOR LANGUAGE



➤ Hidden Markov Model:

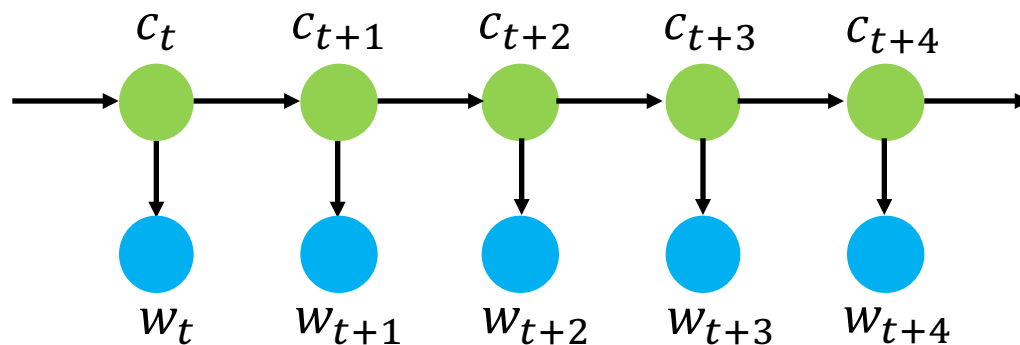
- discourse vector  $c_t \in \mathbb{R}^d$  governs the discourse/theme/context of time  $t$
- words  $w_t$  (observable); embedding  $v_{w_t} \in \mathbb{R}^d$  (parameters to learn)
- log-linear observation model

$$\Pr[w_t | c_t] \propto \exp\langle v_{w_t}, c_t \rangle$$

➤ Closely related to [\[Mnih-Hinton'07\]](#)

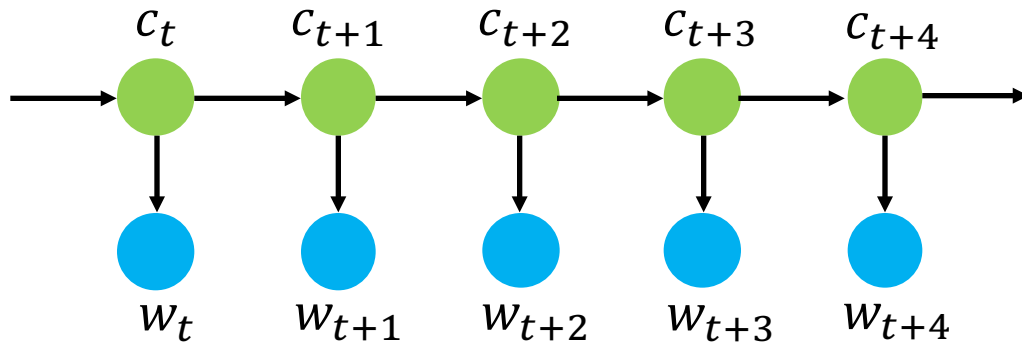


# RAND-WALK: A GENERATIVE MODEL FOR LANGUAGE (CONT'D)



- Ideally,  $c_t, v_w \in \mathbb{R}^d$  should contain semantic information in its coordinates
  - E.g. (0.5, -0.3, ...) could mean “0.5 gender, -0.3 age,..”
- But, the whole system is rotational invariant:  $\langle c_t, v_w \rangle = \langle R c_t, R v_w \rangle$
- There should exist a rotation so that the coordinates are meaningful (back to this later)





➤ Assumptions:

- $\{v_w\}$  consists of vectors drawn from  $s \cdot \mathcal{N}(0, \text{Id})$ ;  $s$  is bounded scalar r.v.
- $c_t$  does a slow random walk (doesn't change much in a window of 5)
- log-linear observation model:  $\Pr[w_t | c_t] \propto \exp\langle v_{w_t}, c_t \rangle$

➤ Main Theorem:

(1)  $\log \Pr[w, w'] = \|v_w + v_{w'}\|^2 / d - 2 \log Z \pm \epsilon$

(2)  $\log \Pr[w] = \|v_w\|^2 / d - \log Z \pm \epsilon$

(3)  $\text{PMI}(w, w') = \langle v_w, v_{w'} \rangle / d \pm \epsilon$

Fact: (2) implies that the words have power law dist.

➤ Norm determines frequency; spatial orientation determines “meaning”



# EXPLAINING EXISTING METHODS

➤ word2vec [Mikolov et al'13] :

$$\Pr[ w_{i+6} \mid w_{i+1}, \dots, w_{i+5} ] \propto \exp\langle v_{w_{i+6}}, \frac{1}{5} (v_{w_{i+1}} + \dots + v_{w_{i+5}}) \rangle$$

➤ GloVe [Pennington et al'14] :

$$\log \Pr[w, w'] \approx \langle v_w, v_{w'} \rangle + s_w + s_{w'} + C$$

Eq. (1)  $\log \Pr[w, w'] = \|v_w + v_{w'}\|^2 / d - 2 \log Z \pm \epsilon$

➤ [Levy-Goldberg'14]

$$\text{PMI}(w, w') \approx \langle v_w, v_{w'} \rangle + C$$

Eq. (3)  $\text{PMI}(w, w') = \langle v_w, v_{w'} \rangle / d \pm \epsilon$



# EXPLAINING EXISTING METHODS CONT'D

➤ word2vec [Mikolov et al'13]:

$$\Pr[w_{i+6} | w_{i+1}, \dots, w_{i+5}] \propto \exp\langle v_{w_{i+6}}, \frac{1}{5}(v_{w_{i+1}} + \dots + v_{w_{i+5}}) \rangle$$

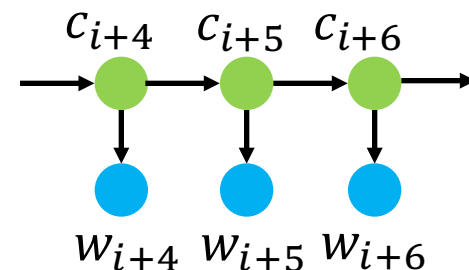
↑

max-likelihood  
estimate of  $c_{i+6}$

➤ Under our model,

- Random walk is slow:  $c_{i+1} \approx c_{i+2} \approx \dots \approx c_{i+6} \approx c$
- Best estimate for current discourse  $c_{i+6}$ :

$$\operatorname{argmax}_{c, \|c\|=1} \Pr[c | w_{i+1}, \dots, w_{i+5}] = \alpha(v_{w_{i+1}} + \dots + v_{w_{i+5}})$$



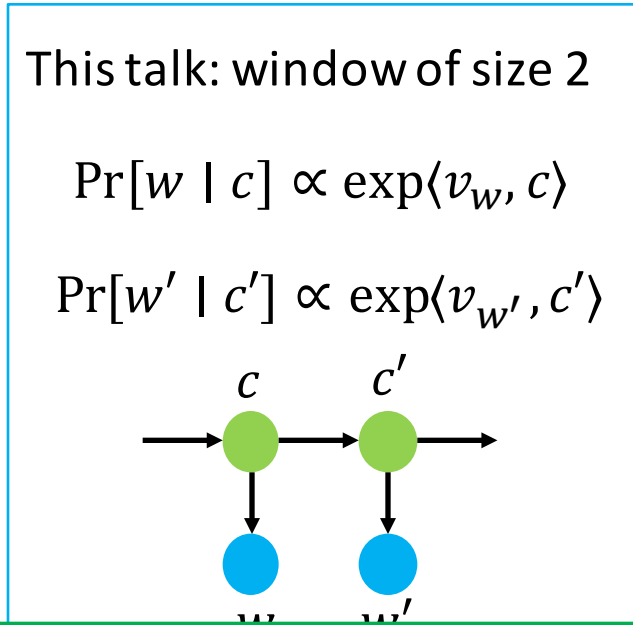
- Prob. distribution of next word given the best guess  $c$ :

$$\Pr[w_{i+6} | c_{i+6} = \alpha(v_{w_{i+1}} + \dots + v_{w_{i+5}})] \propto \exp\langle v_{w_{i+6}}, \alpha(v_{w_{i+1}} + \dots + v_{w_{i+5}}) \rangle$$



# PROOF SKETCH OF MAIN THM.

- $\Pr[w | c] = \frac{1}{Z_c} \cdot \exp\langle v_w, c \rangle$
- $Z_c = \sum_w \exp\langle v_w, c \rangle$  partition function



$$\Pr[w, w'] = \int \Pr[w | c] \Pr[w' | c'] p(c, c') dc dc'$$

$$= \int \underbrace{\frac{1}{Z_c Z_{c'}}}_{??} \cdot \underbrace{\exp\langle v_w, c \rangle \exp\langle v_{w'}, c' \rangle p(c, c')}_{\text{Assume } c = c' \text{ with probability 1,}}$$

spherical Gaussian vector  $c$

- $\mathbb{E}[\exp\langle v, c \rangle] = \exp\|v\|^2 / d$

??

- Assume  $c = c'$  with probability 1,

$$= \int \exp\langle v_w + v_{w'}, c \rangle p(c) dc = \exp\|v_w + v_{w'}\|^2 / d$$

**Eq. (1)**  $\log \Pr[w, w'] = \|v_w + v_{w'}\|^2 / d - 2 \log Z \pm \epsilon$





# PROOF SKETCH OF MAIN THM CONT'D

- $\Pr[w \mid c] = \frac{1}{Z_c} \cdot \exp\langle v_w, c \rangle$
- $Z_c = \sum_w \exp\langle v_w, c \rangle$  partition function

Lemma 1: for almost all  $c$ , almost all  $\{v_w\}$ ,  
 $Z_c = (1 + o(1))Z$

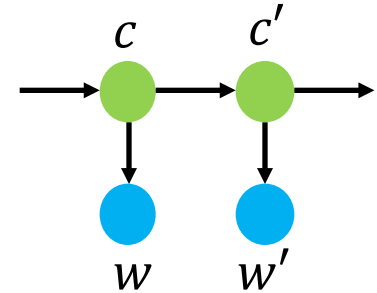
- Proof (sketch) :
  - for most  $c$ ,  $Z_c$  concentrates around its mean
  - mean of  $Z_c$  is determined by  $\|c\|$ , which in turn concentrates
  - caveat:  $\exp\langle v, c \rangle$  for  $v \sim \mathcal{N}(0, \text{Id})$  is not subgaussian, nor sub-exponential. ( $\alpha$ -Orlicz norm is not bounded for any  $\alpha > 0$ )

$$\text{Eq. (1)} \quad \log \Pr[w, w'] = \|v_w + v_{w'}\|^2 / d - 2 \log Z \pm \epsilon$$

This talk: window of size 2

$$\Pr[w \mid c] \propto \exp\langle v_w, c \rangle$$

$$\Pr[w' \mid c'] \propto \exp\langle v_{w'}, c' \rangle$$



# A HEAVY TAIL PHENOMENON

Lemma 1: for almost all  $c$ , almost all  $\{v_w\}$ ,  
$$Z_c = (1 + o(1))Z$$

- Proof Sketch:
- Fixing  $c$ , to show high probability over choices of  $v_w$ 's

$$Z_c = \sum_w \exp\langle v_w, c \rangle = (1 + o(1))\mathbb{E}[Z_c]$$

- $z_w = \langle v_w, c \rangle$  scalar Gaussian random variable
- $\|c\|$  governs the mean and variance of  $z_w$ .
- $\|c\|$  in turns is concentrated



# A HEAVY TAIL PHENOMENON

- Question:  $z_1, \dots, z_n \sim \mathcal{N}(0,1)$

$$Z = \sum_{i=1}^n \exp(z_i)$$

- How is  $Z$  concentrated?

- $\mathbb{E}[Z_c] = \Theta(n)$ , and  $\text{Var}[Z_c] = O(n)$

- The tail of  $\exp(z_i)$  is bad!

- $\Pr[\exp z_i > t] \approx t^{-\log t}$

- Claim:

$$\Pr[Z > \mathbb{E}Z + C\sqrt{n} \cdot \log n] \leq \exp(-\log^2 n)$$

- Trick: truncate  $z_i$  at  $\log n$  and deal with the tail by union bound

- (sub)-Gaussian tail

$$\Pr[X > t] \leq \exp(-t^2/2)$$

- (sub)-exponential tail

$$\Pr[X > t] \leq \exp(-t/2)$$



# A HEAVY TAIL PHENOMENON

Lemma 1: for almost all  $c$ , almost all  $\{v_w\}$ ,  
$$Z_c = (1 + o(1))Z$$



- Proof Sketch:
- Fixing  $c$ , we have with high probability over choices of  $v_w$ 's

$$Z_c = \sum_w \exp\langle v_w, c \rangle = (1 + o(1))\mathbb{E}[Z_c]$$

- $z_w = \langle v_w, c \rangle$  scalar Gaussian random variable
- $\|c\|$  governs the mean and variance of  $z_w$ .
- $\|c\|$  in turns is concentrated



# PROOF SKETCH OF MAIN THM CONT'D

- $\Pr[w | c] = \frac{1}{Z_c} \cdot \exp\langle v_w, c \rangle$
- $Z_c = \sum_w \exp\langle v_w, c \rangle$  partition function

Lemma 1: for almost all  $c$ , almost all  $\{v_w\}$ ,  
 $Z_c = (1 + o(1))Z$

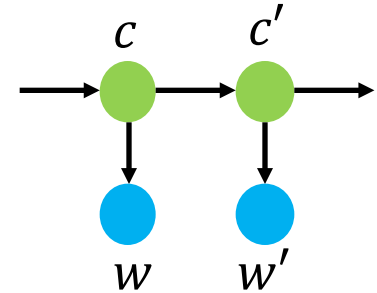
$$\begin{aligned}\Pr[w, w'] &= \int \frac{1}{Z_c Z_{c'}} \cdot \exp\langle v_w + v_{w'}, c \rangle p(c) dc \\ &= (1 \pm o(1)) \frac{1}{Z^2} \int \exp\langle v_w + v_{w'}, c \rangle p(c) dc \\ &= (1 \pm o(1)) \frac{1}{Z^2} \exp(\|v_w + v_{w'}\|^2 / d)\end{aligned}$$

**Eq. (1)**  $\log \Pr[w, w'] = \|v_w + v_{w'}\|^2 / d - 2 \log Z \pm \epsilon$

This talk: window of size 2

$$\Pr[w | c] \propto \exp\langle v_w, c \rangle$$

$$\Pr[w' | c'] \propto \exp\langle v_{w'}, c' \rangle$$



# MODEL VERIFICATION

- Our theory predicts

$$\text{Eq. (1)} \quad \log \Pr[w, w'] = \|v_w + v_{w'}\|^2 / d - 2 \log Z \pm \epsilon$$

- (Approximate) maximum likelihood objective (**SN**)

$$\min_{\{v_w\}, Y} \sum_{w, w'} \widehat{\Pr}[w, w'] (\log \widehat{\Pr}[w, w'] - \|v_w + v_{w'}\|^2 - Y)^2$$

Simplest word embedding method yet (**fewest “knobs”** to turn)

Comparable performance on analogy test

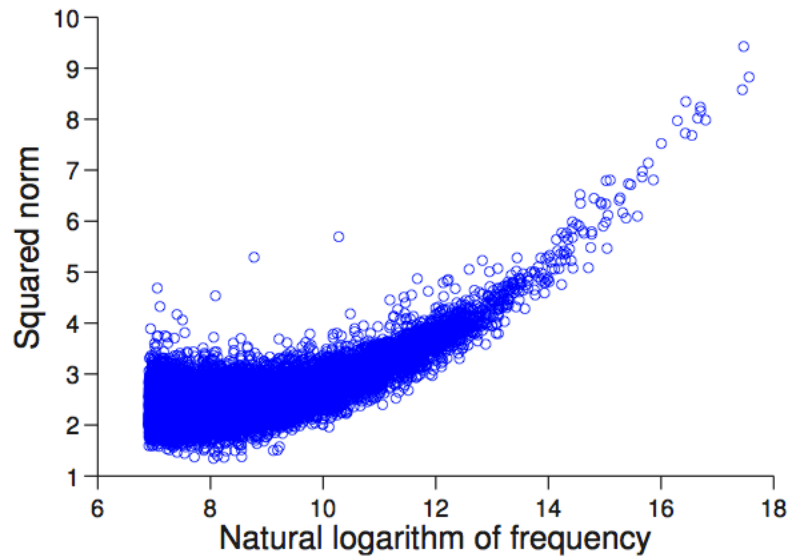
	Relations	SN	GloVe	CBOw	skip-gram
G	semantic	0.84	0.85	0.79	0.73
	syntactic	0.61	0.65	0.71	0.68
	total	0.71	0.73	0.74	0.70
M	adjective	0.50	0.56	0.58	0.58
	noun	0.69	0.70	0.56	0.58
	verb	0.48	0.53	0.64	0.56
	total	0.53	0.57	0.62	0.57



# MODEL VERIFICATION CONT'D

➤ Our theory predicts

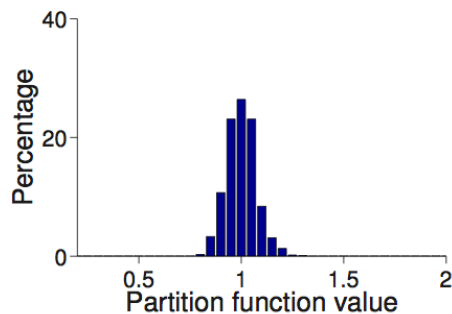
$$\text{Eq. (2)} \quad \log \Pr[w] = \|v_w\|^2 / d - \log Z \pm \epsilon$$



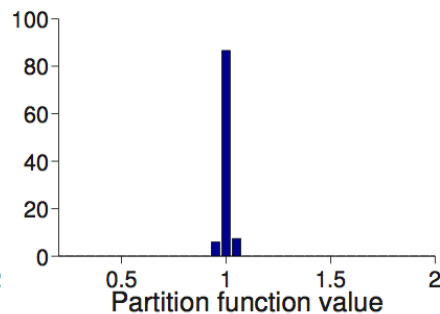
# MODEL VERIFICATION CONT'D

➤ Our theory predicts

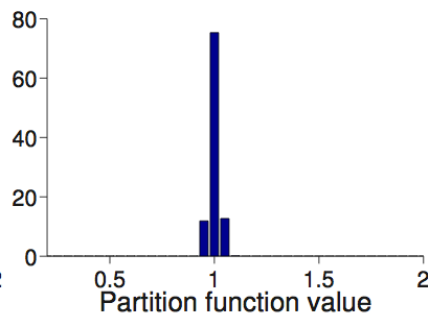
$$Z_c = (1 \pm o(1))Z$$



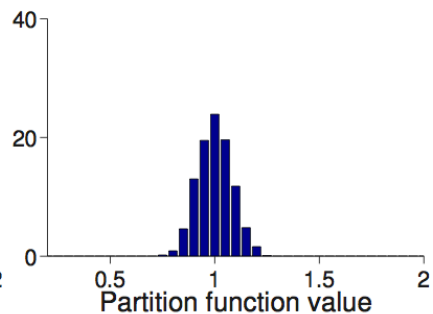
(a) SN



(b) GloVe



(c) CBOW



(d) skip-gram





# WRAP UP

- Under generative model RANK-WALK

For most of the words  $\chi$ :

$$\frac{\Pr[\chi \mid a]}{\Pr[\chi \mid b]} \approx \frac{\Pr[\chi \mid c]}{\Pr[\chi \mid d]} \iff v_a - v_b \approx v_c - v_d$$

↑

semantic def. of analogy

↑

algebraic def. of analogy

- Beyond only solving analogy task?
- Extracting more information from analogy/embeddings?





Some recent work:

Extracting different meanings from word embeddings

(same team: Arora, Li, Liang, M., Risteski)



# POLYSEMY

➤ “Tie” can mean article of clothing, or physical act

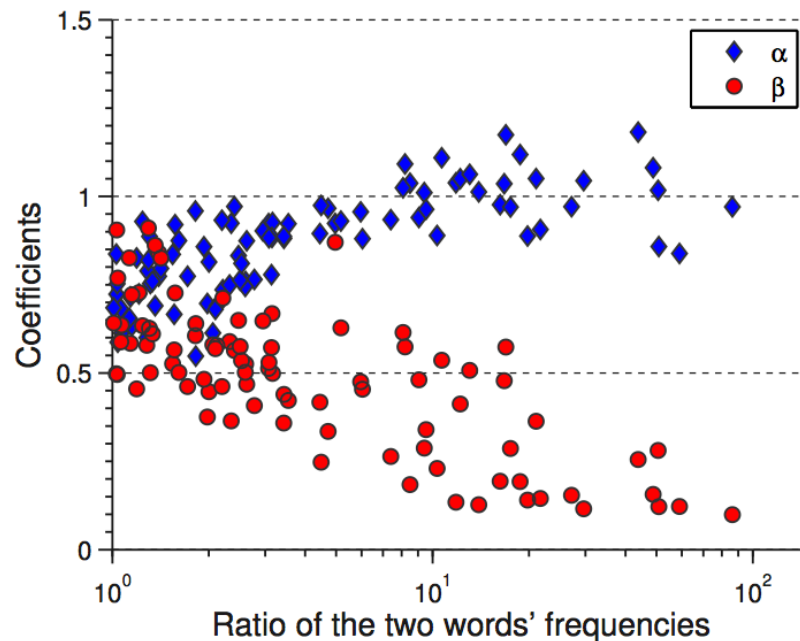


➤ *Tie* represents unrelated words  $tie_1, tie_2$ , etc.

Quick experiment: Take two **random/unrelated** words  $w_1, w_2$  where  $w_1$  is  $\sim 100$  times **more frequent** than  $w_2$ . Declare these to be a single word and compute its embedding in our model.

Result: close to something like  $0.8v_{w_1} + 0.2v_{w_2}$





- Mathematical explanation
- Merge  $w_1, w_2$  as  $w$ . Let  $r = \frac{\Pr[w_1]}{\Pr[w_2]} > 1$
- Then  $v_w \approx \alpha v_{w_1} + \beta v_{w_2}$ , where
  - $\alpha = 1 - c_1 \log\left(1 + \frac{1}{r}\right) \approx 1$
  - $\beta = 1 - c_2 \log r$
- $\beta > .1$  even if  $r = 100$  !
- Rare meaning is not swamped, thanks to the **log** !



# EXTRACTING DIFFERENT MEANINGS

- “Tie” can mean article of clothing, or physical act



- *Tie* represents unrelated words

which correspond to different discourses

$$v_{tie} = 0.8a_1 + 0.2a_2$$

↑

discourse discourse  
for *tie*<sub>1</sub> for *tie*<sub>2</sub>

- Sparse coding for extracting different meanings

- Find  $m = 2000$  “discourses”  $a_1, a_2, \dots \in \mathbb{R}^n$  such that each word vector  $v_w$  is expressed as weighted sum of **at most 5** of them, plus “noise vector.”

$$v_w = x_{w,1}a_1 + x_{w,2}a_2 + \dots + noise$$

$x_w$  has only 5 non-zeros

- Training objective:

$$\min_{\substack{A=[a_1, \dots, a_m] \\ \text{sparse } x'_w\text{'s}}} \sum_w \|v_w - Ax_w\|^2$$

- local search algo. [EAB’05], provable algo. [SWW’12, AGM’14, AGMM’15..]



Representative subset of 2000 discourses (represented using their nearest words)

Atom 1978	825	231	616	1638	149	330
drowning	instagram	stakes	membrane	slapping	orchestra	conferences
suicides	twitter	thoroughbred	mitochondria	pulling	philharmonic	meetings
overdose	facebook	guineas	cytosol	plucking	philharmonia	seminars
murder	tumblr	preakness	cytoplasm	squeezing	conductor	workshops
poisoning	vimeo	filly	membranes	twisting	symphony	exhibitions
commits	linkedin	fillies	organelles	bowing	orchestras	organizes
stabbing	reddit	epsom	endoplasmic	slamming	toscanini	concerts
strangulation	myspace	racecourse	proteins	tossing	concertgebouw	lectures
gunshot	tweets	sired	vesicles	grabbing	solti	presentations



closest words to  $a_{231}$



5 atoms that express  $v_{tie}$

Atom 1005	31	1561	2060	1563
trousers	season	scoreline	wires	operatic
blouse	teams	goalless	cables	soprano
waistcoat	winning	equaliser	wiring	mezzo
skirt	league	clinching	electrical	contralto
sleeved	finished	scoreless	wire	baritone
pants	championship	replay	cable	coloratura



# MULTI-LAYER SPARSE CODING

- Atoms of discourse found are fairly fine-grained
- Maybe  $a_{biochemistry} = \alpha \cdot b_{biology} + \beta \cdot b_{chemistry}$ ?
- Another layer:

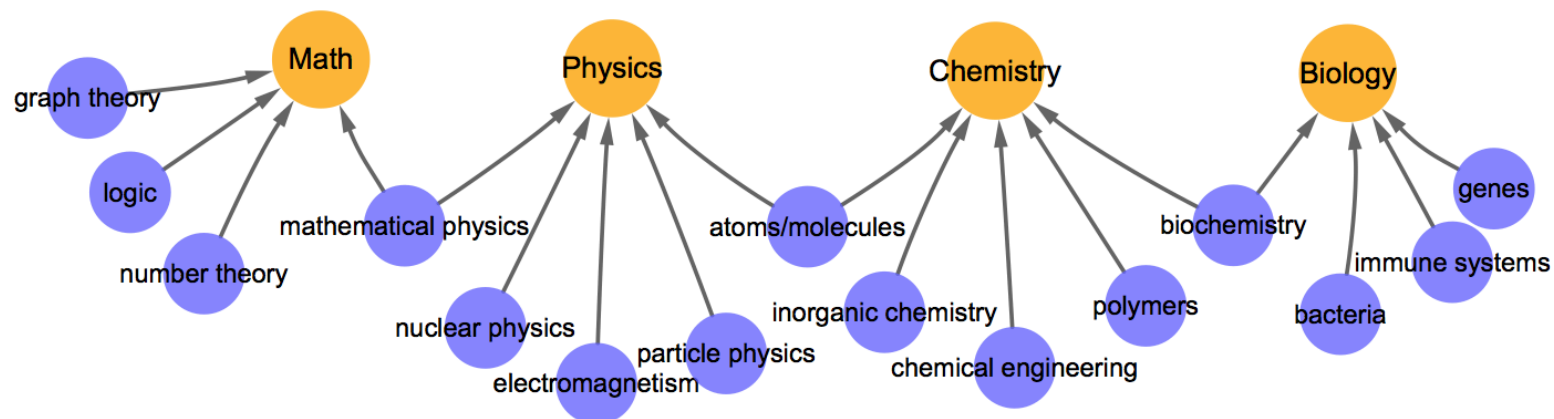
$$\min_{B, Y \text{ sparse}} \|A - BY\|^2$$

411
acids
amino
biosynthesis
peptide
<i>biochemistry</i>





# MULTI-LAYER SPARSE CODING CONT'D



Atom	28	2016	468	1318	411
	logic deductive propositional semantics	graph subgraph bipartite vertex	boson massless particle higgs	polyester polypropylene resins epoxy	acids amino biosynthesis peptide
tag	<i>logic</i>	<i>graph theory</i>	<i>particle physics</i>	<i>polymer</i>	<i>biochemistry</i>



# CONCLUSIONS

- Part I: new generative model that captures semantics.
- Provable guarantee:
  - log of co-occurrence matrix has low rank structure
  - semantic analogy  $\Leftrightarrow$  linear algebraic structure for word vectors
- Simplistic assumptions, but good fit to reality
  
- Part II: automatic way of detect word meanings
  - Hierarchical basis in the embedding space
  
- Other applications of our model/method?





# AN IDEAL SCENARIO

- Each ordinate of  $v_w$  means something:

	currency	country	American	Chinese
	↓	↓	↓	↓
$v_{USA}$	[... .., 0, ... ..]	[... .., 1, ... ..]	[... .., 1, ... ..]	[... .., 0, ... ..]
$v_{dollar}$	[... .., 1, ... ..]	[... .., 0, ... ..]	[... .., 1, ... ..]	[... .., 0, ... ..]
$v_{China}$	[... .., 0, ... ..]	[... .., 1, ... ..]	[... .., 0, ... ..]	[... .., 1, ... ..]
$v_{RMB}$	[... .., 1, ... ..]	[... .., 0, ... ..]	[... .., 0, ... ..]	[... .., 1, ... ..]

- On other coordinates, the values are either very small or the supports are non-overlapping

$$v_{USA} - v_{dollar} = [\dots, -1, \dots, 1, \dots, 0, \dots, 0, \dots]$$

$$v_{China} - v_{RMB} = [\dots, -1, \dots, 1, \dots, 0, \dots, 0, \dots]$$

- Problem: rotational invariance – rotation of word vectors doesn't change the model.



# REVISED IDEA: SPARSE CODING

currency country American Chinese

↓ ↓ ↓ ↓

$$v_{USA} = [\dots, 0, \dots, 1, \dots, 1, \dots, 0, \dots]$$

$$v_{dollar} = [\dots, 1, \dots, 0, \dots, 1, \dots, 0, \dots]$$

$$v_{China} = [\dots, 0, \dots, 1, \dots, 0, \dots, 1, \dots]$$

$$v_{RMB} = [\dots, 1, \dots, 0, \dots, 0, \dots, 1, \dots]$$

$$\cdot \begin{bmatrix} R \end{bmatrix}$$

↑

sparse coefficients

↑

basis vectors

- With sparsity, the model is identifiable; allows overcomplete basis; is tractable under mild assumptions. [SWW'12] [AGM'13][AAJNT'13][AGMM'14]



# EXPERIMENTS

$$\min_{X \text{ sparse}, R} \|V - X \cdot R\|_F^2$$

- $V$  contains word vectors as rows (obtained from any embedding method)
- Sparsity of rows of  $X$  is chosen to be 5
- $R$  contains 2000 basis vectors (as rows), each of which is 300-dim



# RESOLVING MYSTERY 2

Assuming M1 was answered,

$$\text{PMI}(w, w') = \langle v_w, v_{w'} \rangle + \xi \quad (*)$$

with large  $\xi$

M2: Why low-dim vectors solves analogy when (\*) is only roughly true?

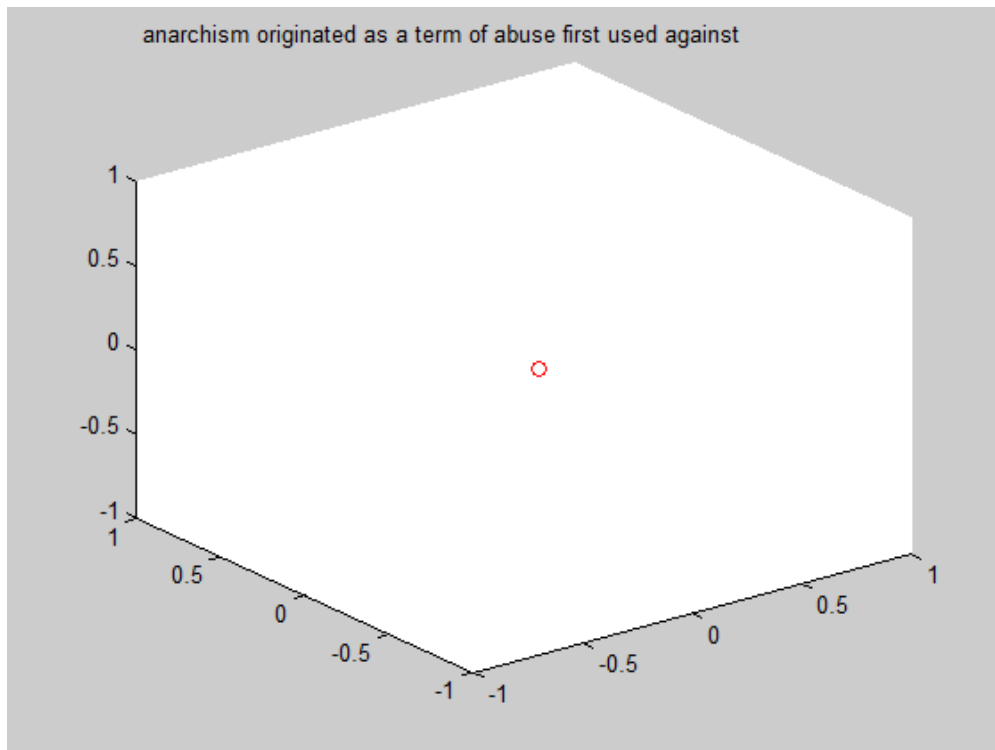
A2: (\*) + isotropy of word vectors  $\Rightarrow$  low-dim fitting reduces noise

(Quite intuitive, though doesn't follow Occam's bound for PAC-learning)



# SLOW RANDOM WALK ILLUSTRATION

- Our theory assumes that  $c_t$  does a slow random walk



- red dot: the estimate hidden variable  $c_t$  at time  $t$
- sentence at top: the window of size 10 at time  $t$





# RESOLVING MYSTERY 2

Assuming M1 was answered,

$$\text{PMI}(w, w') = \langle v_w, v_{w'} \rangle + \xi \quad (*)$$

with large  $\xi$

M2: Why low-dim vectors solves analogy when (\*) is only roughly true?

A2: (\*) + isotropy of word vectors  $\Rightarrow$  low-dim fitting reduces noise

(Quite intuitive, though doesn't follow Occam's bound for PAC-learning)

