

Modeling Topics and Knowledge Bases with Embeddings

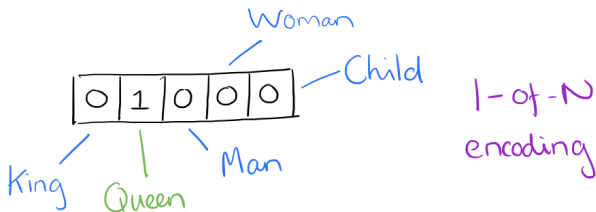
Dat Quoc Nguyen and Mark Johnson

Department of Computing
Macquarie University
Sydney, Australia

December 2016

Vector representations/embeddings of words

- One-hot representation: high-dimensional and sparse vector
 - ▶ Vocabulary size: N
 - ▶ N -dimensional vector: filled with 0s, except for a 1 at the position associated with word index



Vector representations/embeddings of words

- Deep learning evolution: most neural network toolkits do not play well with one-hot representations
 - ▶ Dense vector: low-dimensional distributed vector representation
 - ▶ Vector size $k \ll N$, for example: $k = 100$, Vocabulary size $N = 100000$

King

0.99
0.99
0.05
0.7
⋮

Queen

0.99
0.05
0.93
0.6

Woman

0.02
0.01
0.999
0.5

Princess

0.98
0.02
0.94
0.1

Vector representations/embeddings of words

- Unsupervised learning models are proposed to learn low-dimensional vectors of words efficiently, e.g W2V Skip-gram
- Word embeddings learned from large external corpora capture various aspects of word meanings
 - ▶ Possibly assign topics (i.e. labels) to clusters of “similar” meaning words
 - ▶ Topic **1** for {*banking*, *bank*, *transaction*, *finance*, *money*, *laundering*}

⇒ Can we incorporate word embeddings to topic models?



Improving topic models with word embeddings

- *Topic models* take a corpus of documents as input, and
 - ▶ Learn a set of latent *topics* for the corpus
 - ▶ Infer *document-to-topic* and *topic-to-word* distributions from co-occurrence of words within documents

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

- If the corpus is small and/or the documents are short, the topics will be noisy due to the limited information of word co-occurrence
- IDEA: Use the word embeddings learned on a large external corpus to improve the topic-word distributions in a topic model

Improving topic models with word embeddings

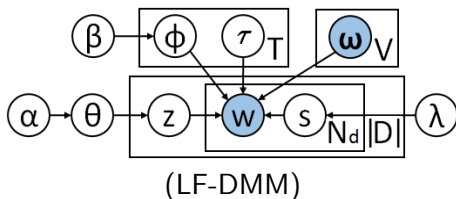
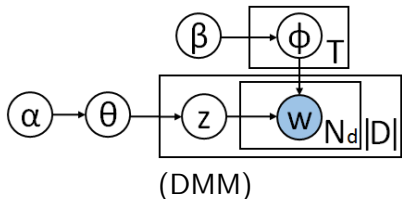
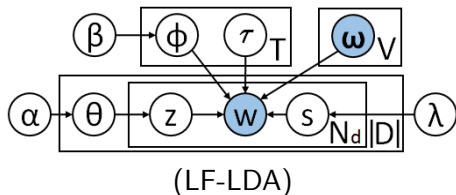
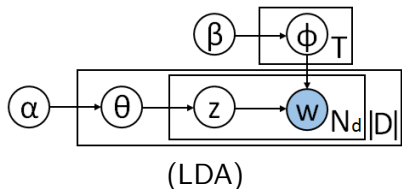
- Each word w is associated with a pre-trained *word embedding* ω_w
- Each topic t is associated with a *topic embedding* τ_t
- We define a latent feature topic-to-word distribution $\text{CatE}(w)$ over words:

$$\text{CatE}(w \mid t) = \frac{\exp(\omega_w \cdot \tau_t)}{\sum_{w' \in V} \exp(\omega_{w'} \cdot \tau_t)}$$

- ▶ Optimize the log-loss to learn τ_t
- Our new topic models *mix the CatE distribution* with a multinomial distribution over words
 - ▶ Combine information from a large, general corpus (via the CatE distribution) and a smaller but more specific corpus (via the multinomial distribution)

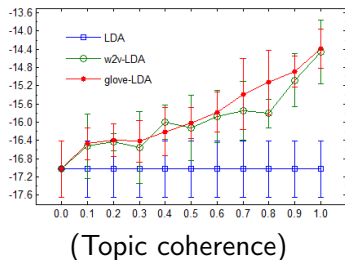
Improving topic models with word embeddings

- Combine Latent Dirichlet Allocation (LDA) and Dirichlet Multinomial Mixture (DMM) with the word embeddings: LF-LDA & LF-DMM



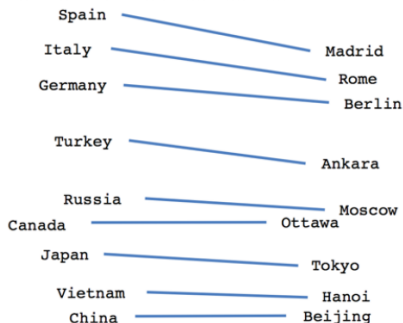
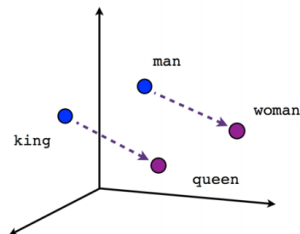
Improving topic models with word embeddings

Topic 1		Topic 3		Topic 4	
DMM	LF-DMM	DMM	LF-DMM	DMM	LF-DMM
japan	japan	u.s.	prices	egypt	libya
nuclear	nuclear	oil	sales	<u>china</u>	egypt
u.s.	u.s.	japan	oil	u.s	iran
crisis	plant	prices	u.s.	mubarak	mideast
plant	quake	stocks	profit	<u>bin</u>	opposition
<u>china</u>	radiation	sales	stocks	libya	protests
<u>libya</u>	earthquake	profit	japan	<u>laden</u>	leader
radiation	tsunami	<u>fed</u>	rise	<u>france</u>	syria
<u>u.n.</u>	nuke	rise	gas	bahrain	u.n.
<u>vote</u>	crisis	growth	growth	<u>air</u>	tunisia
<u>korea</u>	disaster	<u>wall</u>	shares	<u>report</u>	chief
<u>europe</u>	power	<u>street</u>	price	<u>rights</u>	protesters
<u>government</u>	oil	<u>china</u>	profits	<u>court</u>	mubarak
<u>election</u>	japanese	<u>fall</u>	rises	u.n.	crackdown
<u>deal</u>	plants	shares	earnings	<u>war</u>	bahrain



- Significant improvement of topic coherence scores on all models and experimental corpora
- Obtain **5+%** absolute improvements in clustering and classification evaluation scores on the small or short datasets
- No reliable difference between pre-trained Word2Vec and Glove vectors

Vector representations/embeddings: “distance” results



$$\mathbf{v}_{king} - \mathbf{v}_{queen} \approx \mathbf{v}_{man} - \mathbf{v}_{woman}$$

$$\mathbf{v}_{Vietnam} - \mathbf{v}_{Hanoi}$$

$$\approx \mathbf{v}_{Japan} - \mathbf{v}_{Tokyo}$$

$$\approx \mathbf{v}_{Germany} - \mathbf{v}_{Berlin}$$

$$\approx \mathbf{v}_{some_relationship}, \text{ saying: } is_capital_of$$

$$\mathbf{v}_{Hanoi} + \mathbf{v}_{is_capital_of} \approx \mathbf{v}_{Vietnam}$$

$$\mathbf{v}_{Tokyo} + \mathbf{v}_{is_capital_of} \approx \mathbf{v}_{Japan}$$

$$\mathbf{v}_{Berlin} + \mathbf{v}_{is_capital_of} \approx \mathbf{v}_{Germany}$$

Triple (*head entity, relation, tail entity*)

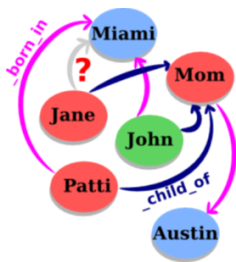
$$\mathbf{v}_h + \mathbf{v}_r \approx \mathbf{v}_t$$

$$\Rightarrow \|\mathbf{v}_h + \mathbf{v}_r - \mathbf{v}_t\|_{\ell_{1/2}} \approx 0$$

Link prediction in knowledge bases?

STransE: a new embedding model for link prediction

- Knowledge bases (KBs) of real-world triple facts (head entity, relation, tail entity) are useful resources for NLP tasks
- **Issue:** large KBs are still far from complete
- So it is useful to perform *link prediction in KBs* or *knowledge base completion*: predict which triples not in a knowledge base are likely to be true
- **Embedding models** for link prediction in KBs:
 - ▶ Associate entities and/or relations with dense feature vectors or matrices
 - ▶ Obtain SOTA performance and generalize to large KBs



STransE: a new embedding model for link prediction

- The TransE model (Bordes et al., 2013) represents each relation r by a translation vector \mathbf{v}_r , which is chosen so that $\|\mathbf{v}_h + \mathbf{v}_r - \mathbf{v}_t\|_{\ell_{1/2}} \approx 0$
 - ▶ Good for 1-to-1 relationships, e.g: *is_capital_of*
 - ▶ Not good for 1-to-Many, Many-to-1 and Many-to-Many, e.g: *gender*
- STransE: a novel embedding model of entities and relationships in KBs
 - ▶ Our STransE represents each entity as a low dimensional vector, and each relation by two matrices and a translation vector
 - ▶ STransE choose matrices $\mathbf{W}_{r,1}$ and $\mathbf{W}_{r,2}$, and vector \mathbf{v}_r so that:
$$\|\mathbf{W}_{r,1}\mathbf{v}_h + \mathbf{v}_r - \mathbf{W}_{r,2}\mathbf{v}_t\|_{\ell_{1/2}} \approx 0$$
 - ▶ Optimize a margin-based objective function to learn the vectors and matrices

STransE results for link prediction in KBs

- We conducted experiments on two benchmark datasets WN18 and FB15k (Bordes et al., 2013)

Dataset	#E	#R	#Train	#Valid	#Test
WN18	40,943	18	141,442	5,000	5,000
FB15k	14,951	1,345	483,142	50,000	59,071

- Link prediction task:
 - Predict h given $(?, r, t)$ or predict t given $(h, r, ?)$ where $?$ denotes the missing element
 - Evaluation metrics: mean rank (MR) and Hits@10 (H10)

Method	WN18		FB15k	
	MR	H10	MR	H10
TransE	251	89.2	125	47.1
TransH	303	86.7	87	64.4
TransR	225	92.0	77	68.7
CTransR	218	92.3	75	70.2
KG2E	348	93.2	59	74.0
TransD	212	92.2	91	77.3
TATEC	-	-	58	76.7
STransE	206	93.4	69	79.7
RTransE	-	-	50	76.2
PTransE	-	-	58	84.6

- Find a new relation-path based embedding model in our CoNLL paper!

STransE results for search personalization

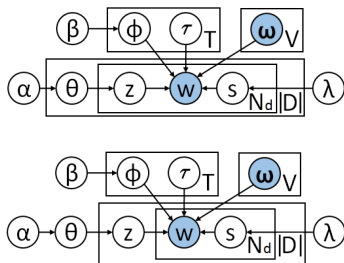
- Two users search using the same keywords, they are often looking for different information (i.e. difference due to the users' interests)
- Personalized search customizes results based on user's search history (i.e. submitted queries and clicked documents)
- Let (q, u, d) represent a triple (query, user, document)
- Represent the user u by two matrices $\mathbf{W}_{u,1}$ and $\mathbf{W}_{u,2}$ and a vector \mathbf{v}_u , which represents the user's topical interests, so that:

$$\|\mathbf{W}_{u,1}\mathbf{v}_q + \mathbf{v}_u - \mathbf{W}_{u,2}\mathbf{v}_d\|_{\ell_{1/2}} \approx 0$$

- \mathbf{v}_d and \mathbf{v}_q are pre-determined by employing the LDA topic model
- Optimize a margin-based objective function to learn the user embeddings and matrices

Metric	Search Eng.	L2R SP	STransE	TransE
MRR	0.559	0.631 _{+12.9%}	0.656 _{+17.3%}	0.645 _{+15.4%}
P@1	0.385	0.452 _{+17.4%}	0.501 _{+30.3%}	0.481 _{+24.9%}

Conclusions



$$\|\mathbf{W}_{r,1}\mathbf{v}_h + \mathbf{v}_r - \mathbf{W}_{r,2}\mathbf{v}_t\|_{\ell_{1/2}}$$

- Latent feature vector representations induced from large external corpora can be used to improve topic modeling on smaller datasets
- Our new embedding model STransE for link prediction in KBs
 - ▶ Applying to a search personalization task, STransE helps to significantly improve the ranking quality
- <https://github.com/datquocnguyen>

Thank you for your attention!

- Dat Quoc Nguyen, Richard Billingsley, Lan Du and Mark Johnson. Improving Topic Models with Latent Feature Word Representations. *Transactions of the ACL*, 2015.
- Dat Quoc Nguyen, Kairit Sirts, Lizhen Qu and Mark Johnson. STansE: a novel embedding model of entities and relationships in knowledge bases. In *Proceedings of NAACL-HLT*, 2016.
- Dat Quoc Nguyen, Kairit Sirts, Lizhen Qu and Mark Johnson. Neighborhood Mixture Model for Knowledge Base Completion. In *Proceedings of CoNLL*, 2016.
- Thanh Vu*, Dat Quoc Nguyen*, Mark Johnson, Dawei Song and Alistair Willis. Search Personalization with Embeddings. In *Proceedings of ECIR 2017*, to appear.