

# Modeling Topics and Knowledge Bases with Embeddings

Dat Quoc Nguyen

Department of Computing  
Macquarie University  
Sydney, Australia

August 2016

- Mark Johnson
- Kairit Sirts
- Lizhen Qu
- Richard Billingsley
- Lan Du

- 1 General introduction
- 2 Improving topic models with word embeddings
  - Introduction
  - Latent-feature topic models
  - Experimental evaluation
  - Summary
- 3 A new embedding model for knowledge base completion
  - Introduction
  - Our neighborhood mixture model
  - Experimental evaluation
  - Summary

Use vector representations for improving topic models as well as for improving link prediction in knowledge bases (i.e. knowledge base completion)

- Incorporate word embeddings trained on large external corpora to improve topic modeling on smaller datasets

Nguyen et al. "Improving Topic Models with Latent Feature Word Representations." *Transactions of ACL*, 2015, vol. 3, pp. 299-313.

- Predict the missing relationships between entities in knowledge bases

Nguyen et al. "Neighborhood Mixture Model for Knowledge Base Completion." In *Proceedings of CoNLL 2016*, pp. 40-50.

Nguyen et al. "STransE: a novel embedding model of entities and relationships in knowledge bases." In *Proceedings of NAACL-HLT 2016*, pp. 460-466.

- 1 General introduction
- 2 Improving topic models with word embeddings
  - Introduction
  - Latent-feature topic models
  - Experimental evaluation
  - Summary
- 3 A new embedding model for knowledge base completion
  - Introduction
  - Our neighborhood mixture model
  - Experimental evaluation
  - Summary

- 1 General introduction
- 2 Improving topic models with word embeddings
  - Introduction
  - Latent-feature topic models
  - Experimental evaluation
  - Summary
- 3 A new embedding model for knowledge base completion

- *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
- 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
- *Latent word representations* learnt from large external corpora capture various aspects of word meanings
  - ▶ We used the pre-trained Word2Vec (Mikolov et al., 2013) and Glove (Pennington et al., 2014) word representations

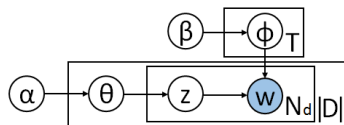
- Use the word representations learnt on a large external corpus to improve the topic-word distributions in a topic model
  - ▶ Combine Latent Dirichlet Allocation (Blei et al., 2003) and Dirichlet Multinomial Mixture (Nigam et al., 2000) with the word representations
  - ▶ Improvement is greatest on small corpora with short documents



# LDA and DMM

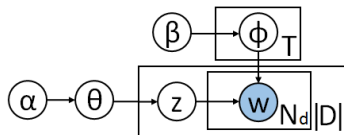
- Latent Dirichlet Allocation (LDA)

$$\begin{aligned}\theta_d &\sim \text{Dir}(\alpha) & z_{d_i} &\sim \text{Cat}(\theta_d) \\ \phi_z &\sim \text{Dir}(\beta) & w_{d_i} &\sim \text{Cat}(\phi_{z_{d_i}})\end{aligned}$$



- Dirichlet Multinomial Mixture (DMM) model: one-topic-per-document

$$\begin{aligned}\theta &\sim \text{Dir}(\alpha) & z_d &\sim \text{Cat}(\theta) \\ \phi_z &\sim \text{Dir}(\beta) & w_{d_i} &\sim \text{Cat}(\phi_{z_d})\end{aligned}$$



- Inference is typically performed with a *Gibbs sampler*, integrating out  $\theta$  and  $\phi$  (Griffiths et al., 2004; Yin and Wang, 2014)

- 1 General introduction
- 2 Improving topic models with word embeddings
  - Introduction
  - Latent-feature topic models
  - Experimental evaluation
  - Summary
- 3 A new embedding model for knowledge base completion

# Latent-feature topic-to-word distributions

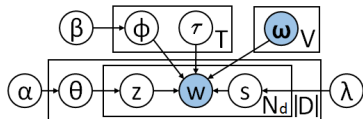
- We assume that each word  $w$  is associated with a *word vector*  $\omega_w$
- We learn a *topic vector*  $\tau_t$  for each topic  $t$
- We use these to define a latent feature topic-to-word distribution  $\text{CatE}(w)$  over words:

$$\text{CatE}(w \mid \tau_t \omega_w^\top) \propto \exp(\tau_t \cdot \omega_w)$$

- ▶  $\tau_t \omega_w^\top$  is a vector of unnormalized scores, one per word
- In our topic models, we *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)
  - ▶ Use a Boolean *indicator variable* that records whether a word is generated from CatE or the multinomial distribution

# The Latent Feature LDA (LF-LDA) model

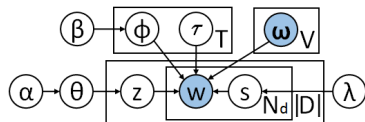
$$\begin{aligned}\theta_d &\sim \text{Dir}(\alpha) & z_{d_i} &\sim \text{Cat}(\theta_d) \\ \phi_z &\sim \text{Dir}(\beta) & s_{d_i} &\sim \text{Ber}(\lambda) \\ w_{d_i} &\sim (1 - s_{d_i})\text{Cat}(\phi_{z_{d_i}}) + s_{d_i}\text{Cat}(\tau_{z_{d_i}} \omega^\top)\end{aligned}$$



- Replace the topic-to-word Dirichlet multinomial component in LDA with a two-component mixture of a topic-to-word Dirichlet multinomial component and a latent feature topic-to-word component
- $s_{d_i}$  is the Boolean indicator variable indicating whether word  $w_{d_i}$  is generated from the latent feature component
- $\lambda$  is a user-specified hyper-parameter determining how often words are generated from the latent feature component
  - ▶ If we estimated  $\lambda$  from data, we expect it would never generate through the latent feature component

# The Latent Feature DMM (LF-DMM) model

$$\begin{aligned}\theta &\sim \text{Dir}(\alpha) & z_d &\sim \text{Cat}(\theta) \\ \phi_z &\sim \text{Dir}(\beta) & s_{d_i} &\sim \text{Ber}(\lambda) \\ w_{d_i} &\sim (1 - s_{d_i})\text{Cat}(\phi_{z_d}) + s_{d_i}\text{CatE}(\tau_{z_d} \omega^\top)\end{aligned}$$



- Replace the topic-to-word Dirichlet multinomial component in DMM with a two-component mixture of a topic-to-word Dirichlet multinomial component and a latent feature topic-to-word component
- $s_{d_i}$  is the Boolean indicator variable indicating whether word  $w_{d_i}$  is generated from the latent feature component
- $\lambda$  is a user-specified hyper-parameter determining how often words are generated from the latent feature component

# Inference for the LF-LDA model

- We integrate out  $\theta$  and  $\phi$  as in the Griffiths et al. (2004) sampler, and *interleave MAP estimation for  $\tau$  with Gibbs sweeps for the other variables*
- Algorithm outline:

initialize the word-topic variables  $z_{d_i}$  using the LDA sampler  
repeat:

for each topic  $t$ :

use LBFGS to optimize the L2-regularized log-loss

$$\tau_t = \arg \max_{\tau_t} P(\tau_t \mid \mathbf{z}, \mathbf{s})$$

for each document  $d$  and each word location  $i$ :

sample  $z_{d_i}$  from  $P(z_{d_i} \mid \mathbf{z}_{\neg d_i}, \mathbf{s}_{\neg d_i}, \boldsymbol{\tau})$

sample  $s_{d_i}$  from  $P(s_{d_i} \mid \mathbf{z}, \mathbf{s}_{\neg d_i}, \boldsymbol{\tau})$

- We integrate out  $\theta$  and  $\phi$  as in the Yin and Wang (2014) sampler, and *interleave MAP estimation for  $\tau$  with Gibbs sweeps*
- Algorithm outline:

initialize the word-topic variables  $z_d$ ; using the DMM sampler

repeat:

for each topic  $t$ :

use LBFGS to optimize the L2-regularized log-loss

$$\tau_t = \arg \max_{\tau_t} P(\tau_t | \mathbf{z}, \mathbf{s})$$

for each document  $d$ :

$$\text{sample } z_d \text{ and } s_d \text{ from } P(z_d, s_d | \mathbf{z}_{-d}, \mathbf{s}_{-d}, \tau)$$

- 1 General introduction
- 2 Improving topic models with word embeddings
  - Introduction
  - Latent-feature topic models
  - **Experimental evaluation**
  - Summary
- 3 A new embedding model for knowledge base completion



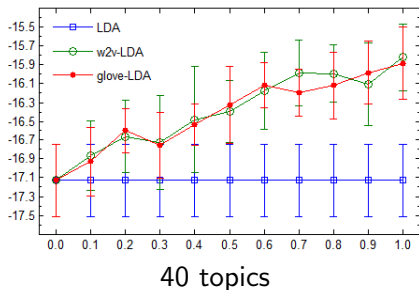
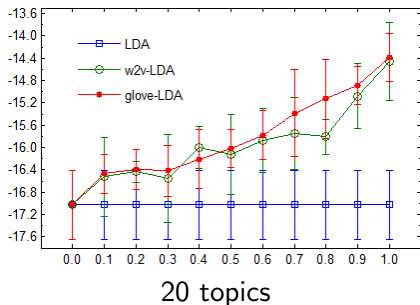
# Goals of evaluation

- A topic model learns document-topic and topic-word distributions:
  - ▶ *Topic coherence* evaluates the topic-word distributions
  - ▶ *Document clustering* and *document classification* evaluate the document-topic distribution
- Do the Word2Vec and Glove word vectors behave differently in topic modelling? (w2v-LDA, glove-LDA, w2v-DMM, glove-DMM)
- We expect that the latent feature component will have *the greatest impact on small corpora*, so our evaluation focuses on them:

| Dataset  |                  | # labels | # docs | words/doc | # types |
|----------|------------------|----------|--------|-----------|---------|
| N20      | 20 newsgroups    | 20       | 18,820 | 103.3     | 19,572  |
| N20short | $\leq 20$ words  | 20       | 1,794  | 13.6      | 6,377   |
| N20small | 400 docs         | 20       | 400    | 88.0      | 8,157   |
| TMN      | TagMyNews        | 7        | 32,597 | 18.3      | 13,428  |
| TMNtitle | TagMyNews titles | 7        | 32,503 | 4.9       | 6,347   |
| Twitter  |                  | 4        | 2,520  | 5.0       | 1,390   |

# Topic coherence evaluation

- Lau et al. (2014) showed that *human scores on a word intrusion task* are highly correlated with the *normalized pointwise mutual information (NPMI)*
- We found latent feature vectors produced a *significant improvement of NPMI scores on all models and corpora*
  - ▶ Greatest improvement when  $\lambda = 1$  (unsurprisingly)



NPMI scores on the N20short dataset, varying the mixture weight  $\lambda$  from 0.0 to 1.0.

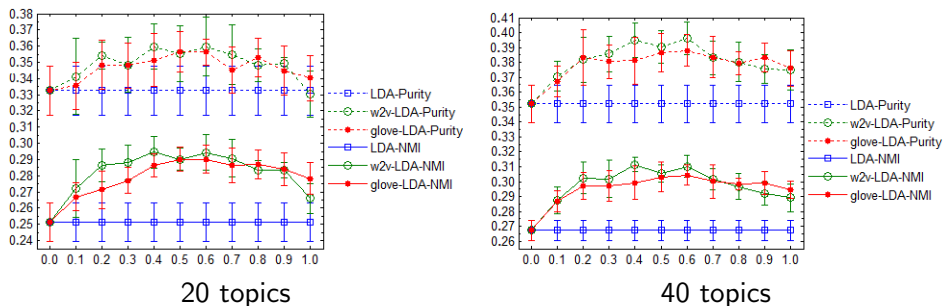
# w2v-DMM on TagMyNews titles corpus

| Topic 1           |                   | Topic 3       |                 | Topic 4       |                   |
|-------------------|-------------------|---------------|-----------------|---------------|-------------------|
| DMM               | w2v-DMM           | DMM           | w2v-DMM         | DMM           | w2v-DMM           |
| japan             | japan             | u.s.          | prices          | egypt         | libya             |
| nuclear           | nuclear           | oil           | sales           | <u>china</u>  | egypt             |
| u.s.              | u.s.              | japan         | oil             | u.s           | <b>iran</b>       |
| crisis            | plant             | prices        | u.s.            | mubarak       | <b>mideast</b>    |
| plant             | <b>quake</b>      | stocks        | profit          | <u>bin</u>    | <b>opposition</b> |
| <u>china</u>      | radiation         | sales         | stocks          | libya         | <b>protests</b>   |
| <u>libya</u>      | <b>earthquake</b> | profit        | japan           | <u>laden</u>  | <b>leader</b>     |
| radiation         | <b>tsunami</b>    | <u>fed</u>    | rise            | <u>france</u> | <b>syria</b>      |
| <u>u.n.</u>       | <b>nuke</b>       | rise          | <b>gas</b>      | bahrain       | u.n.              |
| <u>vote</u>       | crisis            | growth        | growth          | <u>air</u>    | <b>tunisia</b>    |
| <u>korea</u>      | <b>disaster</b>   | <u>wall</u>   | shares          | <u>report</u> | <b>chief</b>      |
| europa            | <b>power</b>      | <u>street</u> | <b>price</b>    | <u>rights</u> | <b>protesters</b> |
| <u>government</u> | <b>oil</b>        | <u>china</u>  | <b>profits</b>  | <u>court</u>  | mubarak           |
| <u>election</u>   | <b>japanese</b>   | <u>fall</u>   | <b>rises</b>    | u.n.          | <b>crackdown</b>  |
| <u>deal</u>       | <b>plants</b>     | shares        | <b>earnings</b> | <u>war</u>    | bahrain           |

- Table shows the 15 most probable topical words found by 20-topic w2v-DMM on the TMNtitle corpus
- Words found by DMM but not by w2v-DMM are underlined
- Words found by w2v-DMM but not DMM are in bold

# Document clustering evaluation (1)

- Cluster documents by assigning them to the *highest probability topic*
- Evaluate clusterings by *purity* and *normalized mutual information (NMI)*



Purity and NMI results on the N20short dataset, varying the mixture weight  $\lambda$  from 0.0 to 1.0.

- In general, best results with  $\lambda = 0.6$
- ⇒ Set  $\lambda = 0.6$  in all further experiments

## Document clustering evaluation (2)

| Data    | Method    | Purity               |                      | NMI                  |                      |
|---------|-----------|----------------------|----------------------|----------------------|----------------------|
|         |           | T=4                  | T=20                 | T=4                  | T=20                 |
| Twitter | LDA       | 0.559 ± 0.020        | 0.614 ± 0.016        | 0.196 ± 0.018        | 0.174 ± 0.008        |
|         | w2v-LDA   | <b>0.598</b> ± 0.023 | <b>0.635</b> ± 0.016 | <b>0.249</b> ± 0.021 | <b>0.191</b> ± 0.011 |
|         | glove-LDA | 0.597 ± 0.016        | <b>0.635</b> ± 0.014 | 0.242 ± 0.013        | <b>0.191</b> ± 0.007 |
|         | Improve.  | <b>0.039</b>         | <b>0.021</b>         | <b>0.053</b>         | <b>0.017</b>         |
| Twitter | DMM       | 0.523 ± 0.011        | 0.619 ± 0.015        | 0.222 ± 0.013        | 0.213 ± 0.011        |
|         | w2v-DMM   | <b>0.589</b> ± 0.017 | 0.655 ± 0.015        | 0.243 ± 0.014        | 0.215 ± 0.009        |
|         | glove-DMM | 0.583 ± 0.023        | <b>0.661</b> ± 0.019 | <b>0.250</b> ± 0.020 | <b>0.223</b> ± 0.014 |
|         | Improve.  | <b>0.066</b>         | <b>0.042</b>         | <b>0.028</b>         | <b>0.01</b>          |

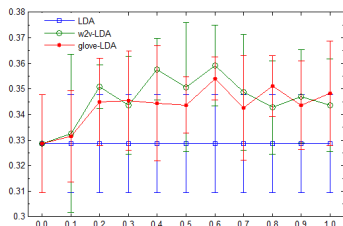
- On the short, our models obtain better clustering results than the baseline models:
  - ▶ on N20small, we get 6.0% improvement on NMI at  $T = 6$
  - ▶ on TMN and TMNtitle, we obtain 6.1% and 2.5% higher Purity at  $T = 80$

## Document clustering evaluation (3)

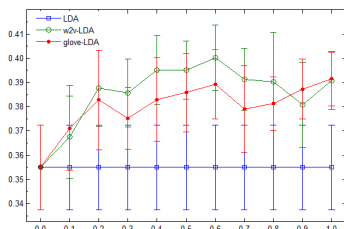
- For small  $T \leq 7$ , on the large datasets of N20, TMN and TMNtitle, our models and baseline models obtain similar clustering results
- With larger  $T$ , our models perform better than baselines on the short TMN and TMNtitle datasets. On the N20 dataset, the baseline LDA model obtains slightly higher clustering results than ours
- No reliable difference between Word2Vec and Glove vectors

# Document classification (1)

- Use SVM to predict the ground truth label from the topic-proportion vector of each document



20 topics



40 topics

$F_1$  scores on N20short dataset, varying the mixture weight  $\lambda$  from 0.0 to 1.0.

| Data     | Method    | $\lambda = 0.6$          |                          |                          |                          |
|----------|-----------|--------------------------|--------------------------|--------------------------|--------------------------|
|          |           | T=6                      | T=20                     | T=40                     | T=80                     |
| N20small | LDA       | 0.204 $\pm$ 0.020        | 0.392 $\pm$ 0.029        | 0.459 $\pm$ 0.030        | 0.477 $\pm$ 0.025        |
|          | w2v-LDA   | <b>0.213</b> $\pm$ 0.018 | <b>0.442</b> $\pm$ 0.025 | <b>0.502</b> $\pm$ 0.031 | <b>0.509</b> $\pm$ 0.022 |
|          | glove-LDA | 0.181 $\pm$ 0.011        | 0.420 $\pm$ 0.025        | 0.474 $\pm$ 0.029        | 0.498 $\pm$ 0.012        |
|          | Improve.  | <b>0.009</b>             | <b>0.05</b>              | <b>0.043</b>             | <b>0.032</b>             |

# Document classification (2)

| Data     | Method    | $\lambda = 0.6$          |                          |                          |                          |
|----------|-----------|--------------------------|--------------------------|--------------------------|--------------------------|
|          |           | T=7                      | T=20                     | T=40                     | T=80                     |
| TMN      | DMM       | 0.607 $\pm$ 0.040        | 0.694 $\pm$ 0.026        | 0.712 $\pm$ 0.014        | 0.721 $\pm$ 0.008        |
|          | w2v-DMM   | 0.607 $\pm$ 0.019        | 0.736 $\pm$ 0.025        | <b>0.760</b> $\pm$ 0.011 | 0.771 $\pm$ 0.005        |
|          | glove-DMM | <b>0.621</b> $\pm$ 0.042 | <b>0.750</b> $\pm$ 0.011 | 0.759 $\pm$ 0.006        | <b>0.775</b> $\pm$ 0.006 |
|          | Improve.  | <b>0.014</b>             | <b>0.056</b>             | <b>0.048</b>             | <b>0.054</b>             |
| TMNtitle | DMM       | 0.500 $\pm$ 0.021        | 0.600 $\pm$ 0.015        | 0.630 $\pm$ 0.016        | 0.652 $\pm$ 0.005        |
|          | w2v-DMM   | 0.528 $\pm$ 0.028        | 0.663 $\pm$ 0.008        | 0.682 $\pm$ 0.008        | <b>0.681</b> $\pm$ 0.006 |
|          | glove-DMM | <b>0.565</b> $\pm$ 0.022 | <b>0.680</b> $\pm$ 0.011 | <b>0.684</b> $\pm$ 0.009 | <b>0.681</b> $\pm$ 0.004 |
|          | Improve.  | <b>0.065</b>             | <b>0.08</b>              | <b>0.054</b>             | <b>0.029</b>             |
| Data     | Method    | $\lambda = 0.6$          |                          |                          |                          |
|          |           | T=4                      | T=20                     | T=40                     | T=80                     |
| Twitter  | LDA       | 0.526 $\pm$ 0.021        | 0.636 $\pm$ 0.011        | 0.650 $\pm$ 0.014        | 0.653 $\pm$ 0.008        |
|          | w2v-LDA   | <b>0.578</b> $\pm$ 0.047 | 0.651 $\pm$ 0.015        | 0.661 $\pm$ 0.011        | <b>0.664</b> $\pm$ 0.010 |
|          | glove-LDA | 0.569 $\pm$ 0.037        | <b>0.656</b> $\pm$ 0.011 | <b>0.662</b> $\pm$ 0.008 | 0.662 $\pm$ 0.006        |
|          | Improve.  | <b>0.052</b>             | <b>0.02</b>              | <b>0.012</b>             | <b>0.011</b>             |
| Twitter  | DMM       | 0.469 $\pm$ 0.014        | 0.600 $\pm$ 0.021        | 0.645 $\pm$ 0.009        | 0.665 $\pm$ 0.014        |
|          | w2v-DMM   | <b>0.539</b> $\pm$ 0.016 | 0.649 $\pm$ 0.016        | 0.656 $\pm$ 0.007        | 0.676 $\pm$ 0.012        |
|          | glove-DMM | 0.536 $\pm$ 0.027        | <b>0.654</b> $\pm$ 0.019 | <b>0.657</b> $\pm$ 0.008 | <b>0.680</b> $\pm$ 0.009 |
|          | Improve.  | <b>0.07</b>              | <b>0.054</b>             | <b>0.012</b>             | <b>0.015</b>             |



- 1 General introduction
- 2 Improving topic models with word embeddings
  - Introduction
  - Latent-feature topic models
  - Experimental evaluation
  - Summary
- 3 A new embedding model for knowledge base completion

# Conclusions and future directions

- Latent feature vectors induced from large external corpora can be used to improve topic modeling
  - ▶ Latent features significantly improve topic coherence across a range of corpora with both the LDA and DMM models
  - ▶ Document clustering and document classification also significantly improve, even though these depend directly only on the document-topic distribution
- The improvements were greatest for small document collections and/or for short documents
- We did not detect any reliable difference between Word2Vec and Glove vectors
- Retrain the word vectors to fit the topic-modeling corpus
- More sophisticated latent-feature models of topic-word distributions
- More efficient training procedures

- 1 General introduction
- 2 Improving topic models with word embeddings
  - Introduction
  - Latent-feature topic models
  - Experimental evaluation
  - Summary
- 3 A new embedding model for knowledge base completion
  - Introduction
  - Our neighborhood mixture model
  - Experimental evaluation
  - Summary

- 1 General introduction
- 2 Improving topic models with word embeddings
- 3 A new embedding model for knowledge base completion
  - **Introduction**
  - Our neighborhood mixture model
  - Experimental evaluation
  - Summary

# Introduction

- 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* (KBC): predict which triples not in a knowledge base are likely to be true

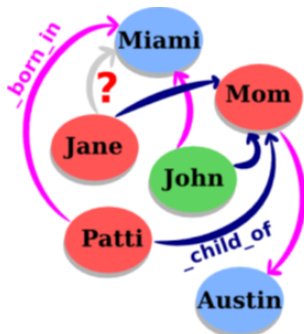


Figure extracted from "Jason Weston and Antoine Bordes. 2014. Embedding Methods for NLP. *EMNLP 2014 tutorial*."

- **Embedding models** for KBC:
  - ▶ Associate entities and/or relations with dense feature vectors or matrices
  - ▶ Obtain SOTA performance and generalize to large KBs
- Most embedding models for KBC learn only from triples
- Recent works show that the relation paths between entities in KBs provide useful information and improve KBC

(Harrison Ford, **born\_in\_hospital**/ $r_1$ , Swedish Covenant Hospital)

$\Rightarrow$ (Swedish Covenant Hospital, **located\_in\_city**/ $r_2$ , Chicago)

$\Rightarrow$ (Chicago, **city\_in\_country**/ $r_3$ , United States)

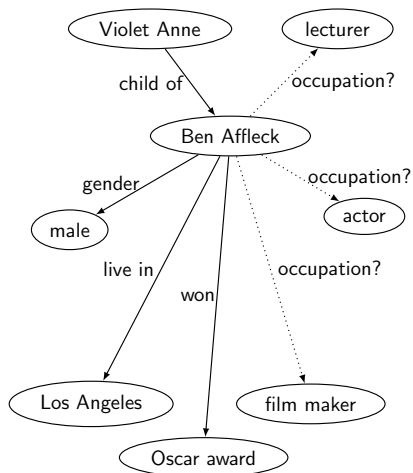
Relation path  $p = \{r_1, r_2, r_3\}$  is useful for predicting the relationship “*nationality*” between the head and tail entities

# Introduction

- **Our motivation:** neighborhoods could provide lots of useful information for predicting the relationship between the entities

Ben\_Affleck

$$\begin{aligned} &= \omega_{r,1}(\text{Violet\_Anne}, \text{child\_of}) \\ &+ \omega_{r,2}(\text{male}, \text{gender}^{-1}) \\ &+ \omega_{r,3}(\text{Los\_Angeles}, \text{live\_in}^{-1}) \\ &+ \omega_{r,4}(\text{Oscar\_award}, \text{won}^{-1}) \end{aligned}$$



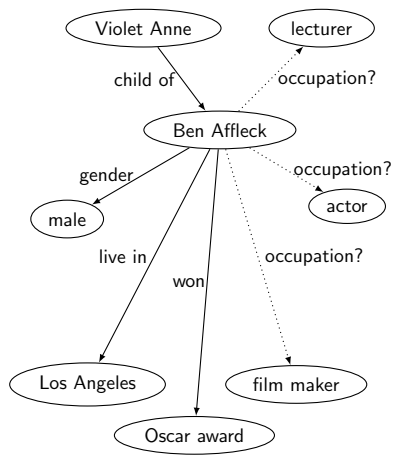
- 1 General introduction
- 2 Improving topic models with word embeddings
- 3 A new embedding model for knowledge base completion
  - Introduction
  - **Our neighborhood mixture model**
  - Experimental evaluation
  - Summary



# Our neighbor-based entity representation

$$\mathcal{E} = \{\text{Ben\_Affleck}, \text{Los\_Angeles}, \dots\}$$
$$\mathcal{R} = \{\text{live\_in}, \text{won}, \text{child\_of}, \text{gender}, \dots\}$$
$$\mathcal{G} = \{(\text{Violet\_Anne}, \text{child\_of}, \text{Ben\_Affleck}),$$
  
 $(\text{Ben\_Affleck}, \text{won}, \text{Oscar\_award}),$   
 $(\text{Ben\_Affleck}, \text{live\_in}, \text{Los\_Angeles}), \dots\}$ 

$\mathcal{N}_e$  is the set of all entity and relation pairs that are neighbors for entity  $e$

$$\mathcal{N}_{\text{Ben\_Affleck}} = \{(\text{Violet\_Anne}, \text{child\_of}),$$
  
 $(\text{male}, \text{gender}^{-1}),$   
 $(\text{Los\_Angeles}, \text{live\_in}^{-1}),$   
 $(\text{Oscar\_award}, \text{won}^{-1})\}$ 

# Our neighbor-based entity representation

- $\mathbf{v}_e \in \mathbb{R}^k$ :  $k$ -dimensional “base” vector associated with entity  $e$
- $\mathbf{u}_{e,r} \in \mathbb{R}^k$ : relation-specific entity vector,  $e \in \mathcal{E}$ ,  $r \in \mathcal{R} \cup \mathcal{R}^{-1}$
- The neighborhood-based entity representation  $\mathbf{v}_{e,r}$  for an entity  $e$  for predicting the relation  $r$  is defined as follows:

$$\mathbf{v}_{e,r} = a_e \mathbf{v}_e + \sum_{(e',r') \in \mathcal{N}_e} b_{r,r'} \mathbf{u}_{e',r'} \quad (1)$$

$a_e$  and  $b_{r,r'}$  are the mixture weights that are constrained to sum to 1:

$$a_e \propto \delta + \exp \alpha_e \quad (2)$$

$$b_{r,r'} \propto \exp \beta_{r,r'} \quad (3)$$

$\delta \geq 0$ : hyper-parameter

$\alpha_e, \beta_{r,r'}$ : learnable exponential mixture parameters

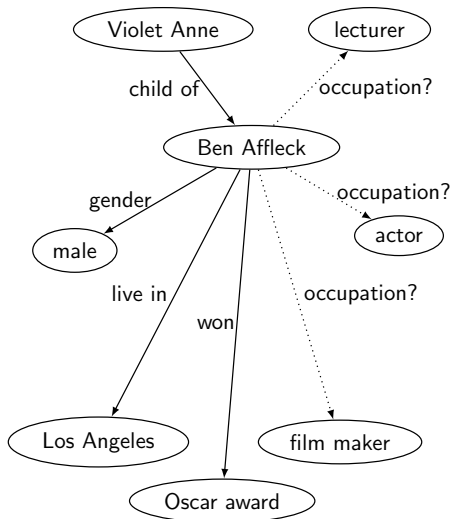
# Our neighbor-based entity representation

$$\mathbf{v}_{e,r} = a_e \mathbf{v}_e + \sum_{(e',r') \in \mathcal{N}_e} b_{r,r'} \mathbf{u}_{e',r'}$$

$e = \text{Ben\_Affleck}$

$r = \text{occupation}$

$\mathcal{N}_e = \{(\text{Violet\_Anne}, \text{child\_of}),$   
 $(\text{male}, \text{gender}^{-1}),$   
 $(\text{Los\_Angeles}, \text{live\_in}^{-1}),$   
 $(\text{Oscar\_award}, \text{won}^{-1})\}$



# Our new embedding model TransE-NMM for KBC

- Embedding models define for each triple  $(h, r, t) \in \mathcal{G}$ , a *score function*  $f(h, r, t)$  that measures its implausibility
- **Goal:** choose  $f$  such that the score  $f(h, r, t)$  of a plausible triple  $(h, r, t)$  is smaller than the score  $f(h', r', t')$  of an implausible triple  $(h', r', t')$ .
- Entity  $e$  and relation  $r$  are represented with vectors  $\mathbf{v}_e \in \mathbb{R}^k$  and  $\mathbf{v}_r \in \mathbb{R}^k$

$$f(h, r, t)_{\text{TransE}} = \|\mathbf{v}_h + \mathbf{v}_r - \mathbf{v}_t\|_{\ell_{1/2}}$$

- The score function of **our new model TransE-NMM** is defined as follows:

$$f(h, r, t) = \|\vartheta_{h,r} + \mathbf{v}_r - \vartheta_{t,r^{-1}}\|_{\ell_{1/2}} \quad (4)$$

$$\vartheta_{e,r} = a_e \mathbf{v}_e + \sum_{(e',r') \in \mathcal{N}_e} b_{r,r'} \mathbf{u}_{e',r'}$$

$$\mathbf{u}_{e,r} = \mathbf{v}_e + \mathbf{v}_r \quad (5)$$

$$\mathbf{v}_{r^{-1}} = -\mathbf{v}_r \quad (6)$$

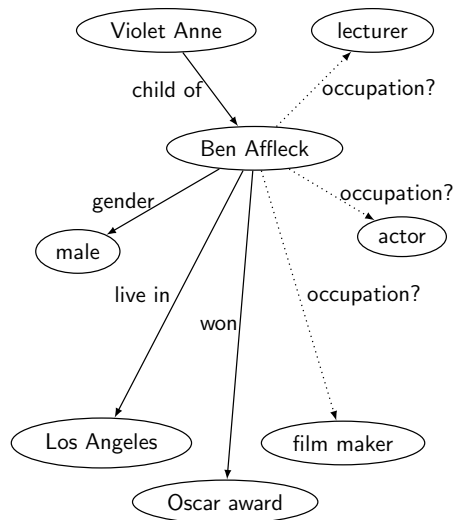
# Our new embedding model TransE-NMM for KBC

$$\mathbf{v}_{e,r} = a_e \mathbf{v}_e + \sum_{(e',r') \in \mathcal{N}_e} b_{r,r'} (\mathbf{v}_{e'} + \mathbf{v}_{r'})$$

$e = \text{Ben\_Affleck}$

$r = \text{occupation}$

$\mathcal{N}_e = \{(\text{Violet\_Anne}, \text{child\_of}),$   
 $(\text{male}, \text{gender}^{-1}),$   
 $(\text{Los\_Angeles}, \text{live\_in}^{-1}),$   
 $(\text{Oscar\_award}, \text{won}^{-1})\}$



# Parameter optimization

- Model parameters:
  - ▶ Entity vectors  $\mathbf{v}_e$
  - ▶ Relation type vectors  $\mathbf{v}_r$
  - ▶  $\boldsymbol{\alpha} = \{\alpha_e | e \in \mathcal{E}\}$ : entity-specific weights
  - ▶  $\boldsymbol{\beta} = \{\beta_{r,r'} | r, r' \in \mathcal{R} \cup \mathcal{R}^{-1}\}$ : relation-specific weights
- Minimize the  $L_2$ -regularized margin-based objective function:

$$\mathcal{L} = \sum_{\substack{(h,r,t) \in \mathcal{G} \\ (h',r,t') \in \mathcal{G}'_{(h,r,t)}}} [\gamma + f(h,r,t) - f(h',r,t')]_+ + \frac{\lambda}{2} (\|\boldsymbol{\alpha}\|_2^2 + \|\boldsymbol{\beta}\|_2^2)$$

$$\mathcal{G}'_{(h,r,t)} = \{(h',r,t) \mid h' \in \mathcal{E}, (h',r,t) \notin \mathcal{G}\} \\ \cup \{(h,r,t') \mid t' \in \mathcal{E}, (h,r,t') \notin \mathcal{G}\}$$

- ▶  $[x]_+ = \max(0, x)$
- ▶  $\gamma$ : the margin hyper-parameter
- ▶  $\lambda$ : the  $L_2$  regularization parameter
- ▶ Impose constraints during training with RMSProp:  $\|\mathbf{v}_e\|_2 \leq 1, \|\mathbf{v}_r\|_2 \leq 1$

# Related work

| Model         | Score function $f(h, r, t)$   |
|---------------|---|
| STransE       | $\ \mathbf{W}_{r,1}\mathbf{v}_h + \mathbf{v}_r - \mathbf{W}_{r,2}\mathbf{v}_t\ _{\ell_{1/2}}; \mathbf{W}_{r,1}, \mathbf{W}_{r,2} \in \mathbb{R}^{k \times k}; \mathbf{v}_r \in \mathbb{R}^k$  |
| TransE        | $\ \mathbf{v}_h + \mathbf{v}_r - \mathbf{v}_t\ _{\ell_{1/2}}; \mathbf{v}_r \in \mathbb{R}^k$  |
| TransH        | $\ (\mathbf{I} - \mathbf{r}_p \mathbf{r}_p^\top)\mathbf{v}_h + \mathbf{v}_r - (\mathbf{I} - \mathbf{r}_p \mathbf{r}_p^\top)\mathbf{v}_t\ _{\ell_{1/2}}$<br>$\mathbf{r}_p, \mathbf{v}_r \in \mathbb{R}^k; \mathbf{I}$ : Identity matrix size $k \times k$  |
| TransD        | $\ (\mathbf{I} + \mathbf{r}_p \mathbf{h}_p^\top)\mathbf{v}_h + \mathbf{v}_r - (\mathbf{I} + \mathbf{r}_p \mathbf{t}_p^\top)\mathbf{v}_t\ _{\ell_{1/2}}$<br>$\mathbf{r}_p, \mathbf{v}_r \in \mathbb{R}^n; \mathbf{h}_p, \mathbf{t}_p \in \mathbb{R}^k; \mathbf{I}$ : Identity matrix size $n \times k$                   |
| TransR        | $\ \mathbf{W}_r \mathbf{v}_h + \mathbf{v}_r - \mathbf{W}_r \mathbf{v}_t\ _{\ell_{1/2}}; \mathbf{W}_r \in \mathbb{R}^{n \times k}; \mathbf{v}_r \in \mathbb{R}^n$  |
| NTN           | $\mathbf{v}_r^\top \tanh(\mathbf{v}_h^\top \mathbf{M}_r \mathbf{v}_t + \mathbf{W}_{r,1}\mathbf{v}_h + \mathbf{W}_{r,2}\mathbf{v}_t + \mathbf{b}_r)$<br>$\mathbf{v}_r, \mathbf{b}_r \in \mathbb{R}^n; \mathbf{M}_r \in \mathbb{R}^{k \times k \times n}; \mathbf{W}_{r,1}, \mathbf{W}_{r,2} \in \mathbb{R}^{n \times k}$ |
| DISTMULT      | $\mathbf{v}_h^\top \mathbf{W}_r \mathbf{v}_t; \mathbf{W}_r$ is a diagonal matrix $\in \mathbb{R}^{k \times k}$  |
| Bilinear-COMP | $\mathbf{v}_h^\top \mathbf{W}_{r_1} \mathbf{W}_{r_2} \dots \mathbf{W}_{r_m} \mathbf{v}_t; \mathbf{W}_{r_1}, \mathbf{W}_{r_2}, \dots, \mathbf{W}_{r_m} \in \mathbb{R}^{k \times k}$  |
| TransE-COMP   | $\ \mathbf{v}_h + \mathbf{v}_{r_1} + \mathbf{v}_{r_2} + \dots + \mathbf{v}_{r_m} - \mathbf{v}_t\ _{\ell_{1/2}}; \mathbf{v}_{r_1}, \mathbf{v}_{r_2}, \dots, \mathbf{v}_{r_m} \in \mathbb{R}^k$   |
| TransE-NMM    | $\ \vartheta_{h,r} + \mathbf{v}_r - \vartheta_{t,r-1}\ _{\ell_{1/2}}$   |

- 1 General introduction
- 2 Improving topic models with word embeddings
- 3 A new embedding model for knowledge base completion
  - Introduction
  - Our neighborhood mixture model
  - **Experimental evaluation**
  - Summary



# Evaluation: experimental setup

| Dataset: | WN11    | FB13    | NELL186 |
|----------|---------|---------|---------|
| #R       | 11      | 13      | 186     |
| #E       | 38,696  | 75,043  | 14,463  |
| #Train   | 112,581 | 316,232 | 31,134  |
| #Valid   | 2,609   | 5,908   | 5,000   |
| #Test    | 10,544  | 23,733  | 5,000   |

- #E: number of entities
- #R: number of relation types
- #Train, #Valid and #Test are the numbers of correct triples in the training, validation, and test sets, respectively
- Each validation and test set also contains the same number of incorrect triples as the number of correct triples

## Triple classification task:

- Predict whether a triple  $(h, r, t)$  is correct or not
- Set a relation-specific threshold  $\theta_r$  for each relation type  $r$
- For an unseen test triple  $(h, r, t)$ , if  $f(h, r, t)$  is smaller than  $\theta_r$  then the triple will be classified as correct, otherwise incorrect
- The relation-specific thresholds are determined by maximizing the micro-averaged accuracy on the validation set

# Evaluation: experimental setup

- **Entity prediction task:**

- ▶ Predict  $h$  given  $(?, r, t)$  or predict  $t$  given  $(h, r, ?)$  where  $?$  denotes the missing element
- ▶ Corrupt each correct test triple  $(h, r, t)$  by replacing either  $h$  or  $t$  by each of the possible entities in turn
- ▶ Rank these candidates in ascending order of their implausibility value computed by the score function
- ▶ “Raw” and “Filtered” setting protocols in which “Filtered” setting is to filter out before ranking any corrupted triples that appear in the KB
- ▶ Metrics: mean rank (MR), mean reciprocal rank (MRR) and Hits@10 (H10)

- **Relation prediction task:**

- ▶ Predict  $r$  given  $(h, ?, t)$  where  $?$  denotes the missing element
- ▶ Corrupt each correct test triple  $(h, r, t)$  by replacing  $r$  by each of the possible relations in turn

# Evaluation: quantitative results

| Data    | Method |            | Triple class. |              | Entity prediction |              |              | Relation prediction |              |              |
|---------|--------|------------|---------------|--------------|-------------------|--------------|--------------|---------------------|--------------|--------------|
|         |        |            | Mic.          | Mac.         | MR                | MRR          | H@10         | MR                  | MRR          | H@10         |
| WN11    | R      | TransE     | 85.21         | 82.53        | 4324              | 0.102        | 19.21        | 2.37                | 0.679        | <b>99.93</b> |
|         |        | TransE-NMM | <b>86.82</b>  | <b>84.37</b> | <b>3466</b>       | <b>0.123</b> | <b>20.59</b> | <b>2.14</b>         | <b>0.687</b> | 99.92        |
|         | F      | TransE     |               |              | 4304              | 0.122        | 21.86        | 2.37                | 0.679        | <b>99.93</b> |
|         |        | TransE-NMM |               |              | <b>3447</b>       | <b>0.137</b> | <b>23.03</b> | <b>2.14</b>         | <b>0.687</b> | 99.92        |
| FB13    | R      | TransE     | 87.57         | 86.66        | 9037              | 0.204        | 35.39        | 1.01                | 0.996        | 99.99        |
|         |        | TransE-NMM | <b>88.58</b>  | <b>87.99</b> | <b>8289</b>       | <b>0.258</b> | <b>35.53</b> | 1.01                | 0.996        | <b>100.0</b> |
|         | F      | TransE     |               |              | 5600              | 0.213        | 36.28        | 1.01                | 0.996        | 99.99        |
|         |        | TransE-NMM |               |              | <b>5018</b>       | <b>0.267</b> | <b>36.36</b> | 1.01                | 0.996        | <b>100.0</b> |
| NELL186 | R      | TransE     | 92.13         | 88.96        | 309               | 0.192        | 36.55        | 8.43                | 0.580        | 77.18        |
|         |        | TransE-NMM | <b>94.57</b>  | <b>90.95</b> | <b>238</b>        | <b>0.221</b> | <b>37.55</b> | <b>6.15</b>         | <b>0.677</b> | <b>82.16</b> |
|         | F      | TransE     |               |              | 279               | 0.268        | 47.13        | 8.32                | 0.602        | 77.26        |
|         |        | TransE-NMM |               |              | <b>214</b>        | <b>0.292</b> | <b>47.82</b> | <b>6.08</b>         | <b>0.690</b> | <b>82.20</b> |

- **Mic.:** Micro-averaged accuracy; **Mac.:** Macro-averaged accuracy
- “R” and “F” denote the “Raw” and “Filtered” settings used in the entity prediction and relation prediction tasks, respectively
- Better results are in **bold**

# Evaluation: quantitative results

| Method        | W11         | F13         |
|---------------|-------------|-------------|
| TransR        | 85.9        | 82.5        |
| CTransR       | 85.7        | -           |
| TransD        | <u>86.4</u> | <b>89.1</b> |
| TranSparse-S  | <u>86.4</u> | 88.2        |
| TranSparse-US | <b>86.8</b> | 87.5        |
| NTN           | 70.6        | 87.2        |
| TransH        | 78.8        | 83.3        |
| SLogAn        | 75.3        | 85.3        |
| KG2E          | 85.4        | 85.3        |
| Bilinear-COMP | 77.6        | 86.1        |
| TransE-COMP   | 80.3        | 87.6        |
| TransE        | 85.2        | 87.6        |
| TransE-NMM    | <b>86.8</b> | <u>88.6</u> |

Micro-averaged accuracy for triple classification on WN11 and FB13

Results on the NELL186 test set:

| Method       | Triple class. |              | Entity pred. |              |
|--------------|---------------|--------------|--------------|--------------|
|              | Mic.          | Mac.         | MR           | H@10         |
| TransE-LLE   | 90.08         | 84.50        | 535          | 20.02        |
| SME-LLE      | 93.64         | 89.39        | <u>253</u>   | 37.14        |
| SE-LLE       | <u>93.95</u>  | 88.54        | 447          | 31.55        |
| TransE-SkipG | 85.33         | 80.06        | 385          | 30.52        |
| SME-SkipG    | 92.86         | <u>89.65</u> | 293          | <b>39.70</b> |
| SE-SkipG     | 93.07         | 87.98        | 412          | 31.12        |
| TransE       | 92.13         | 88.96        | 309          | 36.55        |
| TransE-NMM   | <b>94.57</b>  | <b>90.95</b> | <b>238</b>   | <u>37.55</u> |

The entity prediction results are in the “Raw” setting

## Evaluation: qualitative results

- Take the relation-specific mixture weights from the learned TransE-NMM
- Extract neighbor relations with the largest mixture weights given a relation

| <b>Relation</b>        | <b>Top 3-neighbor relations</b>                                |
|------------------------|--|
| has_instance<br>(WN11) | type_of<br>subordinate_instance_of<br>domain_topic             |
| nationality<br>(FB13)  | place_of_birth<br>place_of_death<br>location                   |
| CEOof<br>(NELL186)     | WorksFor<br>TopMemberOfOrganization<br>PersonLeadsOrganization |

- 1 General introduction
- 2 Improving topic models with word embeddings
- 3 A new embedding model for knowledge base completion
  - Introduction
  - Our neighborhood mixture model
  - Experimental evaluation
  - **Summary**

## Conclusions and future work

- We introduced a neighborhood mixture model for knowledge base completion by constructing neighbor-based vector representations for entities
- We demonstrated its effect by extending the state-of-the-art embedding model TransE with our neighborhood mixture model
- Our model significantly improves TransE and obtains better results than the other state-of-the-art embedding models on three evaluation tasks
- We plan to apply the neighborhood mixture model to the relation path models to combine the useful information from both relation paths and entity neighborhoods

Thank you for your attention!