

# Deep Learning for Unsupervised Relation Extraction

---

Étienne Simon

5 july 2022

LTG

## Education background:

- Algorithmic, Automata Theory, Lambda Calculus...
- Machine Learning, Statistics...

## Research background:

- Machine-learning-oriented NLP
- Not used to work with POS, dependency trees...
- Machine Translation, Diachronic Topic Modeling, Multimodal Semantic Role Labeling, Relation Extraction



⇒ An "otter" entity exists.



⇒ An "otter" entity exists.



⇒ An "inside of" relation exists.



⇒ An "otter" entity exists.



⇒ An "inside of" relation exists.

**Structuralism:** interrelations are keys to our understanding of the world.



⇒ An "otter" entity exists.



⇒ An "inside of" relation exists.

**Structuralism:** interrelations are keys to our understanding of the world.

Realism or Nominalism – **Episteme**

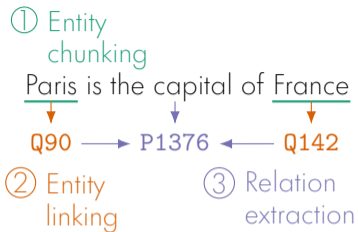
**Techne**

A way to abstract information for easier processing.

## Knowledge Base Population

Maps between two symbolic representations  
(text and knowledge bases).

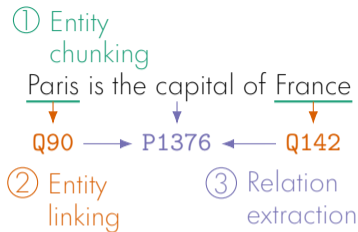
Knowledge bases are set of facts:  
(entity, *relation*, entity)



## Knowledge Base Population

Maps between two symbolic representations (text and knowledge bases).

Knowledge bases are set of facts:  
(entity, relation, entity)



## Symbolic Representations

symbol  $\leftrightarrow$  concept

e.g.: one-hot vector, text (Paris is the capital of France),  
knowledge base (Paris<sup>Q90</sup>, capital<sup>P1376</sup>, France<sup>Q142</sup>)

## Distributed Representations

concept  $\rightarrow$  several units; unit  $\rightarrow$  part of several concepts

e.g.: embeddings, neural network activations



Megrez <sub>$e_1$</sub> <sup>Q850779</sup> is a star in the northern circumpolar constellation of Ursa Major <sub>$e_2$</sub> <sup>Q10460</sup>.



$e_1$  part of constellation  $e_2$

Posidonius <sub>$e_1$</sub> <sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria <sub>$e_2$</sub> <sup>Q617550</sup>.



$e_1$  born in  $e_2$

Hipparchus <sub>$e_1$</sub> <sup>Q159905</sup> was born in Nicaea, Bithynia <sub>$e_2$</sub> <sup>Q739037</sup>, and probably died on the island of Rhodes, Greece.



$e_1$  born in  $e_2$

In an **unsupervised** fashion.

Two kind of approaches: clustering and similarity function.

Megrez <sub>$e_1$</sub> <sup>Q850779</sup> is a star in the northern circumpolar constellation of Ursa Major <sub>$e_2$</sub> <sup>Q10460</sup>.



$e_1$  part of constellation  $e_2$

Posidonius <sub>$e_1$</sub> <sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria <sub>$e_2$</sub> <sup>Q617550</sup>.



$e_1$  born in  $e_2$

Hipparchus <sub>$e_1$</sub> <sup>Q159905</sup> was born in Nicaea, Bithynia <sub>$e_2$</sub> <sup>Q739037</sup>, and probably died on the island of Rhodes, Greece.



$e_1$  born in  $e_2$

Megrez <sub>$e_1$</sub> <sup>Q850779</sup> is a star in the northern circumpolar constellation of Ursa Major <sub>$e_2$</sub> <sup>Q10460</sup>.



*$e_1$  part of constellation  $e_2$*

Posidonius <sub>$e_1$</sub> <sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria <sub>$e_2$</sub> <sup>Q617550</sup>.



*$e_1$  born in  $e_2$*

Hipparchus <sub>$e_1$</sub> <sup>Q159905</sup> was born in Nicaea, Bithynia <sub>$e_2$</sub> <sup>Q739037</sup>, and probably died on the island of Rhodes, Greece.



*$e_1$  born in  $e_2$*

Megrez <sub>$e_1$</sub> <sup>Q850779</sup> is a star in the northern circumpolar constellation of Ursa Major <sub>$e_2$</sub> <sup>Q10460</sup>.



*$e_1$  part of constellation  $e_2$*

Posidonius <sub>$e_1$</sub> <sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria <sub>$e_2$</sub> <sup>Q617550</sup>.



*$e_1$  born in  $e_2$*

Hipparchus <sub>$e_1$</sub> <sup>Q159905</sup> was born in Nicaea, Bithynia <sub>$e_2$</sub> <sup>Q739037</sup>, and probably died on the island of Rhodes, Greece.



*$e_1$  born in  $e_2$*

Same cluster  $\iff$  Same relation

Induced clusters need **not** be labeled with a relation.

Megrez <sub>$e_1$</sub> <sup>Q850779</sup> is a star in the northern circumpolar constellation of Ursa Major <sub>$e_2$</sub> <sup>Q10460</sup>.



*$e_1$  part of constellation  $e_2$*

Posidonius <sub>$e_1$</sub> <sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria <sub>$e_2$</sub> <sup>Q617550</sup>.



*$e_1$  born in  $e_2$*

Hipparchus <sub>$e_1$</sub> <sup>Q159905</sup> was born in Nicaea, Bithynia <sub>$e_2$</sub> <sup>Q739037</sup>, and probably died on the island of Rhodes, Greece.



*$e_1$  born in  $e_2$*

Same cluster  $\iff$  Same relation  
Induced clusters need **not** be labeled with a relation.

### Clustering Metrics

- B<sup>3</sup>** Similar to standard  $F_1$
- V-measure** Entropic  $F_1$
- ARI** Pair of samples consistency

Megrez <sub>$e_1$</sub> <sup>Q850779</sup> is a star in the northern circumpolar constellation of Ursa Major <sub>$e_2$</sub> <sup>Q10460</sup>.



$e_1$  part of constellation  $e_2$

Posidonius <sub>$e_1$</sub> <sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria <sub>$e_2$</sub> <sup>Q617550</sup>.



$e_1$  born in  $e_2$

Hipparchus <sub>$e_1$</sub> <sup>Q159905</sup> was born in Nicaea, Bithynia <sub>$e_2$</sub> <sup>Q739037</sup>, and probably died on the island of Rhodes, Greece.



$e_1$  born in  $e_2$

Megrez<sub>e<sub>1</sub></sub><sup>Q850779</sup> is a star in the northern circumpolar constellation of Ursa Major<sub>e<sub>2</sub></sub><sup>Q10460</sup>.

Posidonius<sub>e<sub>1</sub></sub><sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria<sub>e<sub>2</sub></sub><sup>Q617550</sup>.

Hipparchus<sub>e<sub>1</sub></sub><sup>Q159905</sup> was born in Nicaea, Bithynia<sub>e<sub>2</sub></sub><sup>Q739037</sup>, and probably died on the island of Rhodes, Greece.

Megrez <sub>$e_1$</sub> <sup>Q850779</sup> is a star in the northern circum-polar constellation of Ursa Major <sub>$e_2$</sub> <sup>Q10460</sup>. }  $x_1$

Posidonius <sub>$e_1$</sub> <sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria <sub>$e_2$</sub> <sup>Q617550</sup>. }  $x_2$

Hipparchus <sub>$e_1$</sub> <sup>Q159905</sup> was born in Nicaea, Bithynia <sub>$e_2$</sub> <sup>Q739037</sup>, and probably died on the island of Rhodes, Greece. }  $x_3$

Learn a similarity function  
 $\text{sim} : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$

$\text{sim}(x_1, x_2) < \text{sim}(x_2, x_3)$

$\text{sim}(x_1, x_3) < \text{sim}(x_2, x_3)$

5 way 1 shot: given 1 query and 5 candidates, which of the candidates is most similar to the query?

Evaluated using accuracy.



## Regularizing Discriminative Models

---

Megrez <sub>$e_1$</sub> <sup>Q850779</sup> is a star in the northern circumpolar constellation of Ursa Major <sub>$e_2$</sub> <sup>Q10460</sup>.



$e_1$  part of constellation  $e_2$

Posidonius <sub>$e_1$</sub> <sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria <sub>$e_2$</sub> <sup>Q617550</sup>.



$e_1$  born in  $e_2$

Hipparchus <sub>$e_1$</sub> <sup>Q159905</sup> was born in Nicaea, Bithynia <sub>$e_2$</sub> <sup>Q739037</sup>, and probably died on the island of Rhodes, Greece.



$e_1$  born in  $e_2$

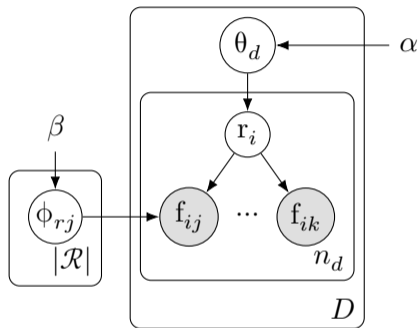
Same cluster  $\iff$  Same relation

Induced clusters need **not** be labeled with a relation.

Evaluated using clustering metrics similar to standard  $F_1$ /precision/recall.

1. Related work
2. Limitation: can't train deep classifier
3. Model details
4. Analysis of limitation
5. Proposed solution
6. Results

An LDA-like model:



$\theta_d$  distribution of relations in document  $d$

$r_i$  conveyed relation

$\phi_{rj}$  associate features to relations

$f_i$  features:

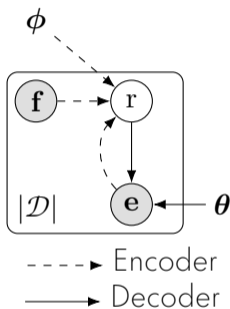
1. bag of words of the infix;
2. surface form of the entities;
3. lemma words on the dependency path;
4. POS of the infix words;

...

Assume  $\mathcal{H}_{\text{BICLIQUE}}$ :  $\forall r \in \mathcal{R} : \exists A, B \subseteq \mathcal{E} : r \bullet \check{r} = A^2 \wedge \check{r} \bullet r = B^2$

**Problem:** Makes large independance assumptions.

A conditional  $\beta$ -VAE:



Autoencode the entities  $\mathbf{e}$  given the sentence features  $\mathbf{f}$ .

$$\mathcal{L}_{\text{VAE}}(\boldsymbol{\theta}, \phi) = \mathcal{L}_{\text{reconstruction}}(\boldsymbol{\theta}, \phi) + \mathcal{L}_{\text{VAE REG}}(\phi)$$

$$\mathcal{L}_{\text{VAE REG}}(\phi) = \text{D}_{\text{KL}}(Q(r | \mathbf{e}; \phi) \| \mathcal{U}(\mathcal{R}))$$

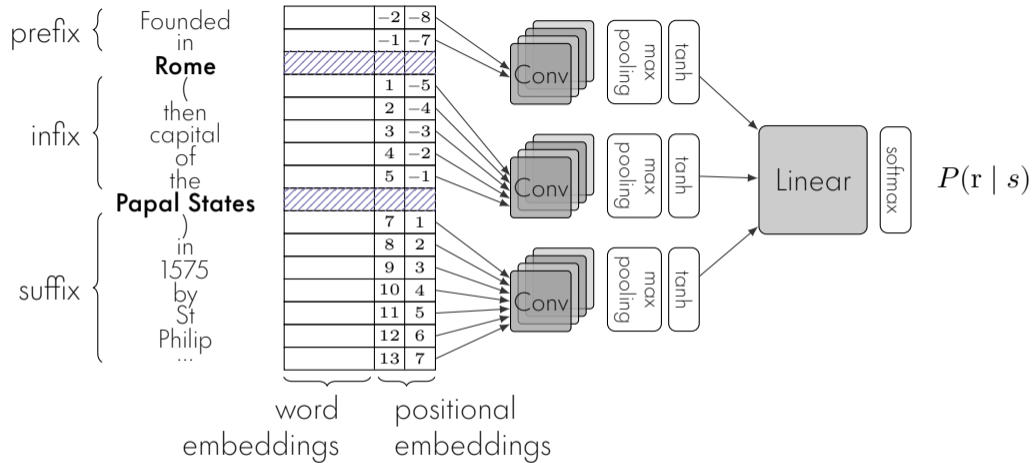
Assume  $\mathcal{H}_{\text{UNIFORM}}$ : All relations occur with equal frequency.

$$\forall r \in \mathcal{R}: P(r) = \frac{1}{|\mathcal{R}|}$$

Assume  $\mathcal{H}_{1 \rightarrow 1}$ : All relations are bijective.

$$\forall r \in \mathcal{R}: r \bullet \check{r} \cup \mathbf{I} = \check{r} \bullet r \cup \mathbf{I} = \mathbf{I}$$

**Problem:** Still uses hand designed features.



Zeng et al. "Distant Supervision for Relation Extraction via Piecewise Convolutional Neural Networks" EMNLP 2015

## Experimental Setup

We introduced:

- 2 metrics (V-measure, ARI)
- 2 datasets (T-RExes)

**B<sup>3</sup>** Similar to standard  $F_1$   
**V-measure** Entropic  $F_1$   
**ARI** Pair of samples consistency

Model		B <sup>3</sup>			V-measure			ARI
Classifier	Reg.	$F_1$	Prec.	Rec.	$F_1$	Hom.	Comp.	
rel-LDA		29.1	24.8	35.2	30.0	26.1	35.1	13.3
rel-LDA1		36.9	30.4	47.0	37.4	31.9	45.1	24.2
Linear	$\mathcal{L}_{\text{VAE REG}}$	35.2	23.8	67.1	27.0	18.6	49.6	18.7
PCNN	$\mathcal{L}_{\text{VAE REG}}$	27.6	24.3	31.9	24.7	21.2	29.6	15.7

## Experimental Setup

We introduced:

- 2 metrics (V-measure, ARI)
- 2 datasets (T-RExes)

**B<sup>3</sup>** Similar to standard  $F_1$   
**V-measure** Entropic  $F_1$   
**ARI** Pair of samples consistency

Model		B <sup>3</sup>			V-measure			ARI
Classifier	Reg.	$F_1$	Prec.	Rec.	$F_1$	Hom.	Comp.	
rel-LDA		29.1	24.8	35.2	30.0	26.1	35.1	13.3
rel-LDA1		36.9	30.4	47.0	37.4	31.9	45.1	24.2
Linear	$\mathcal{L}_{VAE REG}$	35.2	23.8	67.1	27.0	18.6	49.6	18.7
PCNN	$\mathcal{L}_{VAE REG}$	27.6	24.3	31.9	24.7	21.2	29.6	15.7

Yao et al. "Structured Relation Discovery using Generative Models" EMNLP 2011



## Experimental Setup

We introduced:

- 2 metrics (V-measure, ARI)
- 2 datasets (T-RExes)

**B<sup>3</sup>** Similar to standard  $F_1$   
**V-measure** Entropic  $F_1$   
**ARI** Pair of samples consistency

Model		B <sup>3</sup>			V-measure			ARI
Classifier	Reg.	$F_1$	Prec.	Rec.	$F_1$	Hom.	Comp.	
rel-LDA		29.1	24.8	35.2	30.0	26.1	35.1	13.3
rel-LDA1		36.9	30.4	47.0	37.4	31.9	45.1	24.2
Linear	$\mathcal{L}_{\text{VAE REG}}$	35.2	23.8	67.1	27.0	18.6	49.6	18.7
PCNN	$\mathcal{L}_{\text{VAE REG}}$	27.6	24.3	31.9	24.7	21.2	29.6	15.7

Marcheggiani and Titov “Discrete-State Variational Autoencoders for Joint Discovery and Factorization of Relations” TACL 2016

## Experimental Setup

We introduced:

- 2 metrics (V-measure, ARI)
- 2 datasets (T-RExes)

**B<sup>3</sup>** Similar to standard  $F_1$

**V-measure** Entropic  $F_1$

**ARI** Pair of samples consistency

Model		B <sup>3</sup>			V-measure			ARI
Classifier	Reg.	$F_1$	Prec.	Rec.	$F_1$	Hom.	Comp.	
rel-LDA		29.1	24.8	35.2	30.0	26.1	35.1	13.3
rel-LDA1		36.9	30.4	47.0	37.4	31.9	45.1	24.2
Linear	$\mathcal{L}_{\text{VAE REG}}$	35.2	23.8	67.1	27.0	18.6	49.6	18.7
PCNN	$\mathcal{L}_{\text{VAE REG}}$	27.6	24.3	31.9	24.7	21.2	29.6	15.7

**Problem:** Using a deep encoder does not work.

- We introduce a new formalism.
- The encoder and decoder are sub-models performing different tasks.
- The interaction between these two sub-models is problematic.

"The sol<sub>e<sub>1</sub></sub> was the currency of ?<sub>e<sub>2</sub></sub> between 1863 and 1985."

“The sol <sub>$e_1$</sub>  was the currency of ? <sub>$e_2$</sub>  between 1863 and 1985.”

$e_{-i}$  missing entity,  $e_i$  remaining entity,  $s$  conveying sentence

$$\text{for } i = 1, 2 : \quad \overbrace{P(e_{-i} \mid s, e_i)}^{\text{fill-in-the-blank}}$$

"The sol <sub>$e_1$</sub>  was the currency of ? <sub>$e_2$</sub>  between 1863 and 1985."

$e_{-i}$  missing entity,  $e_i$  remaining entity,  $s$  conveying sentence,  $r$  conveyed relation

$$\text{for } i = 1, 2 : \quad \overbrace{P(e_{-i} \mid s, e_i)}^{\text{fill-in-the-blank}} \qquad \overbrace{P(e_{-i} \mid r, e_i)}^{\text{entity predictor}}$$

“The sol <sub>$e_1$</sub>  was the currency of ? <sub>$e_2$</sub>  between 1863 and 1985.”

$e_{-i}$  missing entity,  $e_i$  remaining entity,  $s$  conveying sentence,  $r$  conveyed relation

$$\text{for } i = 1, 2 : \quad \overbrace{P(e_{-i} | s, e_i)}^{\text{fill-in-the-blank}} = \sum_{r \in \mathcal{R}} \overbrace{P(r | s)}^{\text{classifier}} \overbrace{P(e_{-i} | r, e_i)}^{\text{entity predictor}}$$

Assume  $\mathcal{H}_{\text{BLANKABLE}}$ : The relation can be predicted from the text surrounding the two entities alone.

“The sol <sub>$e_1$</sub>  was the currency of ? <sub>$e_2$</sub>  between 1863 and 1985.”

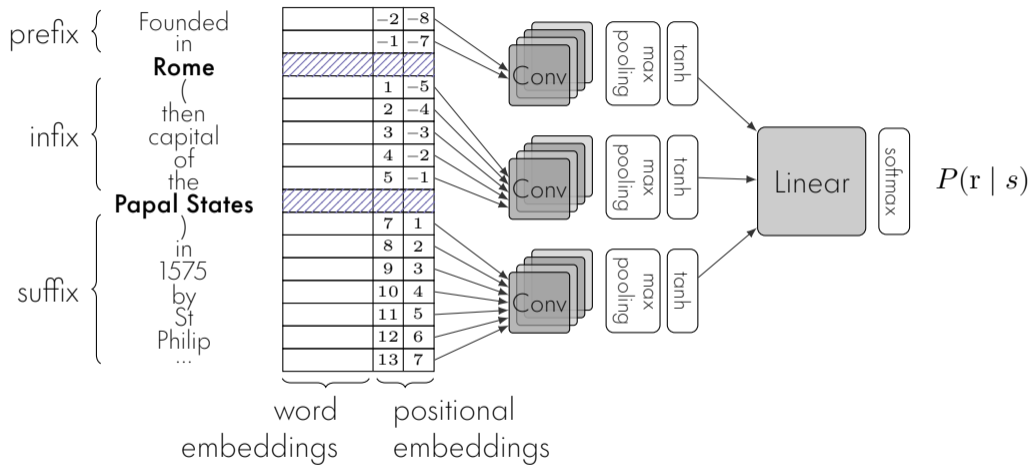
$e_{-i}$  missing entity,  $e_i$  remaining entity,  $s$  conveying sentence,  $r$  conveyed relation

$$\text{for } i = 1, 2 : \quad \overbrace{P(e_{-i} | s, e_i)}^{\text{fill-in-the-blank}} = \sum_{r \in \mathcal{R}} \overbrace{P(r | s)}^{\text{classifier}} \overbrace{P(e_{-i} | r, e_i)}^{\text{entity predictor}}$$

Assume  $\mathcal{H}_{\text{BLANKABLE}}$ : The relation can be predicted from the text surrounding the two entities alone.

1. Train a fill-in-the-blank model on an unsupervised dataset.
2. Throw away the entity predictor.
3. Use the classifier on new samples.





$$P(e_{-i} | s, e_i) = \sum_{r \in \mathcal{R}} \overbrace{P(r | s)}^{\text{fill-in-the-blank classifier}} \overbrace{P(e_{-i} | r, e_i)}^{\text{entity predictor}}$$

## Hybrid (Marcheggiani and Titov 2016)

$$\psi(e_1, r, e_2) = \psi_{\text{SP}}(e_1, r, e_2) + \psi_{\text{RESCAL}}(e_1, r, e_2)$$

$$P(e_1 | r, e_2) = \frac{\exp \psi(e_1, r, e_2)}{\sum_{e' \in \mathcal{E}} \exp \psi(e', r, e_2)}$$

## Selectional Preferences

$$\psi_{\text{SP}}(e_1, r, e_2) = \mathbf{u}_{e_1}^\top \mathbf{a}_r + \mathbf{u}_{e_2}^\top \mathbf{b}_r$$

$\mathbf{U} \in \mathbb{R}^{\mathcal{E} \times d}$  entity embeddings

$\mathbf{A}, \mathbf{B} \in \mathbb{R}^{\mathcal{R} \times d}$  relation embeddings

## RESCAL

$$\psi_{\text{RESCAL}}(e_1, r, e_2) = \mathbf{u}_{e_1}^\top \mathbf{C}_r \mathbf{u}_{e_2}$$

$\mathbf{U} \in \mathbb{R}^{\mathcal{E} \times d}$  entity embeddings

$\mathbf{C} \in \mathbb{R}^{\mathcal{R} \times d \times d}$  relation embeddings

$$\overbrace{P(e_{-i} | s, e_i)}^{\text{fill-in-the-blank}} = \sum_{r \in \mathcal{R}} \overbrace{P(r | s)}^{\text{classifier}} \overbrace{P(e_{-i} | r, e_i)}^{\text{entity predictor}}$$

$$\mathcal{L}_{\text{EP}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{\substack{(\mathbf{s}, \mathbf{e}_1, \mathbf{e}_2) \sim \mathcal{U}(\mathcal{D}) \\ \mathbf{r} \sim \text{PCNN}(\mathbf{s}; \boldsymbol{\phi})}} \left[ \begin{aligned} & -\log \sigma(\psi(\mathbf{e}_1, \mathbf{r}, \mathbf{e}_2; \boldsymbol{\theta})) \\ & - \sum_{j=1}^k \mathbb{E}_{\mathbf{e}' \sim \mathcal{U}_{\mathcal{D}}(\mathcal{E})} [\log \sigma(-\psi(\mathbf{e}_1, \mathbf{r}, \mathbf{e}'; \boldsymbol{\theta}))] \\ & - \sum_{j=1}^k \mathbb{E}_{\mathbf{e}' \sim \mathcal{U}_{\mathcal{D}}(\mathcal{E})} [\log \sigma(-\psi(\mathbf{e}', \mathbf{r}, \mathbf{e}_2; \boldsymbol{\theta}))] \end{aligned} \right]$$

$$\mathcal{L}_{\text{EP}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{\substack{(\mathbf{s}, \mathbf{e}_1, \mathbf{e}_2) \sim \mathcal{U}(\mathcal{D}) \\ \mathbf{r} \sim \text{PCNN}(\mathbf{s}; \boldsymbol{\phi})}} \left[ \begin{aligned} & -\log \sigma(\psi(\mathbf{e}_1, \mathbf{r}, \mathbf{e}_2; \boldsymbol{\theta})) \\ & - \sum_{j=1}^k \mathbb{E}_{\mathbf{e}' \sim \mathcal{U}_{\mathcal{D}}(\mathcal{E})} [\log \sigma(-\psi(\mathbf{e}_1, \mathbf{r}, \mathbf{e}'; \boldsymbol{\theta}))] \\ & - \sum_{j=1}^k \mathbb{E}_{\mathbf{e}' \sim \mathcal{U}_{\mathcal{D}}(\mathcal{E})} [\log \sigma(-\psi(\mathbf{e}', \mathbf{r}, \mathbf{e}_2; \boldsymbol{\theta}))] \end{aligned} \right]$$

1. Take a sample uniformly from the dataset.

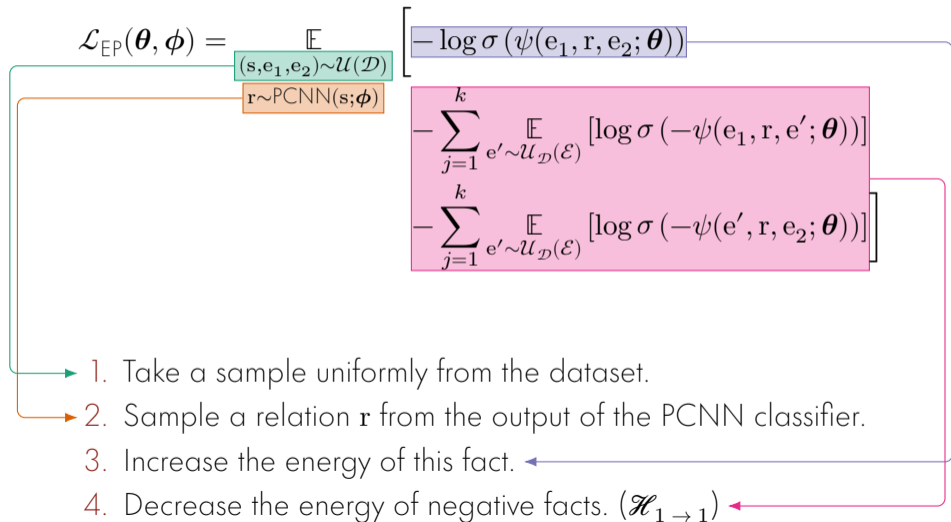
$$\mathcal{L}_{EP}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{\substack{(s, e_1, e_2) \sim \mathcal{U}(\mathcal{D}) \\ \mathbf{r} \sim \text{PCNN}(s; \boldsymbol{\phi})}} \left[ \begin{aligned} & -\log \sigma(\psi(e_1, \mathbf{r}, e_2; \boldsymbol{\theta})) \\ & - \sum_{j=1}^k \mathbb{E}_{e' \sim \mathcal{U}_{\mathcal{D}}(\mathcal{E})} [\log \sigma(-\psi(e_1, \mathbf{r}, e'; \boldsymbol{\theta}))] \\ & - \sum_{j=1}^k \mathbb{E}_{e' \sim \mathcal{U}_{\mathcal{D}}(\mathcal{E})} [\log \sigma(-\psi(e', \mathbf{r}, e_2; \boldsymbol{\theta}))] \end{aligned} \right]$$

1. Take a sample uniformly from the dataset.

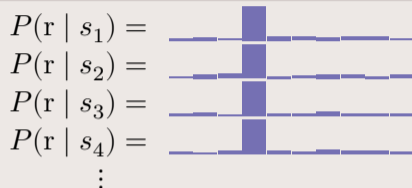
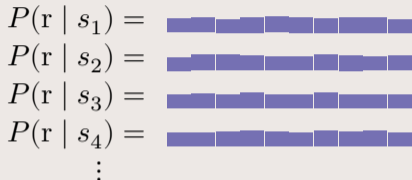
2. Sample a relation  $\mathbf{r}$  from the output of the PCNN classifier.

$$\mathcal{L}_{\text{EP}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{\substack{(\mathbf{s}, \mathbf{e}_1, \mathbf{e}_2) \sim \mathcal{U}(\mathcal{D}) \\ \mathbf{r} \sim \text{PCNN}(\mathbf{s}; \boldsymbol{\phi})}} \left[ \begin{aligned} & -\log \sigma(\psi(\mathbf{e}_1, \mathbf{r}, \mathbf{e}_2; \boldsymbol{\theta})) \\ & - \sum_{j=1}^k \mathbb{E}_{\mathbf{e}' \sim \mathcal{U}_{\mathcal{D}}(\mathcal{E})} [\log \sigma(-\psi(\mathbf{e}_1, \mathbf{r}, \mathbf{e}'; \boldsymbol{\theta}))] \\ & - \sum_{j=1}^k \mathbb{E}_{\mathbf{e}' \sim \mathcal{U}_{\mathcal{D}}(\mathcal{E})} [\log \sigma(-\psi(\mathbf{e}', \mathbf{r}, \mathbf{e}_2; \boldsymbol{\theta}))] \end{aligned} \right]$$

1. Take a sample uniformly from the dataset.
2. Sample a relation  $\mathbf{r}$  from the output of the PCNN classifier.
3. Increase the energy of this fact.



## Degenerate distributions



## Desired distribution



## VAE Model Reminder (Marcheggiani)

$$\overbrace{P(e_{-i} | s, e_i)}^{\text{fill-in-the-blank}} = \sum_{r \in \mathcal{R}} \overbrace{P(r | s)}^{\text{classifier}} \overbrace{P(e_{-i} | r, e_i)}^{\text{entity predictor}}$$

$$\mathcal{L}_{\text{VAE REG}}(\phi) = D_{\text{KL}}(Q(r | \mathbf{e}; \phi) \| \mathcal{U}(\mathcal{R}))$$



## Degenerate distributions

$$P(r | s_1) = \text{[uniform distribution]}$$

$$P(r | s_2) = \text{[uniform distribution]}$$

$$P(r | s_3) = \text{[uniform distribution]}$$

$$P(r | s_4) = \text{[uniform distribution]}$$

$$P(r | s_1) = \text{[distribution with high peak at center]}$$

$$P(r | s_2) = \text{[distribution with high peak at center]}$$

$$P(r | s_3) = \text{[distribution with high peak at center]}$$

$$P(r | s_4) = \text{[distribution with high peak at center]}$$

Problem: **Marcheggiani's model cannot handle deep encoder.**

## Desired distribution

$$P(r | s_1) = \text{[distribution with peak at right end]}$$

$$P(r | s_2) = \text{[distribution with peak at left end]}$$

$$P(r | s_3) = \text{[distribution with peak at left end]}$$

$$P(r | s_4) = \text{[distribution with peak at right end]}$$

$$\vdots$$

## VAE Model Reminder (Marcheggiani)

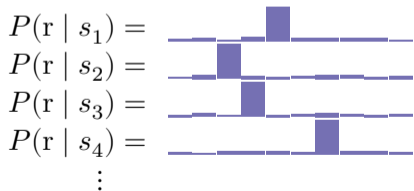
$$\overbrace{P(e_{-i} | s, e_i)}^{\text{fill-in-the-blank}} = \sum_{r \in \mathcal{R}} \overbrace{P(r | s)}^{\text{classifier}} \overbrace{P(e_{-i} | r, e_i)}^{\text{entity predictor}}$$

$$\mathcal{L}_{\text{VAE REG}}(\phi) = D_{\text{KL}}(Q(r | \mathbf{e}; \phi) \| \mathcal{U}(\mathcal{R}))$$

Degenerate distributions:



Desired distributions:

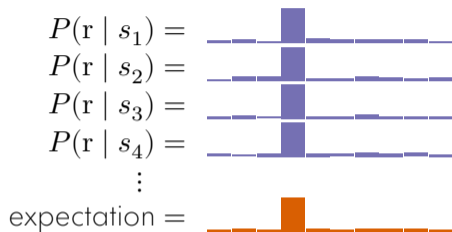


## Ensure Confidence

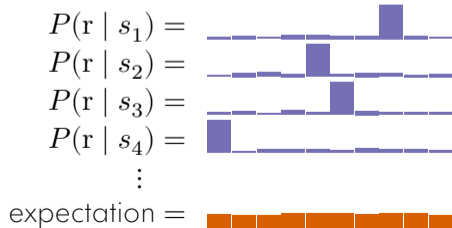
$$\mathcal{L}_S(\phi) = \mathbb{E}_{(s, \mathbf{e}) \sim \mathcal{U}(\mathcal{D})} [\mathbf{H}(\mathbf{R} | s, \mathbf{e}; \phi)]$$

The entropy of the relation distribution must be low for each sample.

Degenerate distributions:



Desired distributions:



### Ensure Diversity

$$\mathcal{L}_D(\phi) = D_{\text{KL}}(P(\mathbf{R} | \phi) \| \mathcal{U}(\mathcal{R}))$$

At the level of the dataset (or mini-batch) the distribution of relations must be uniform.

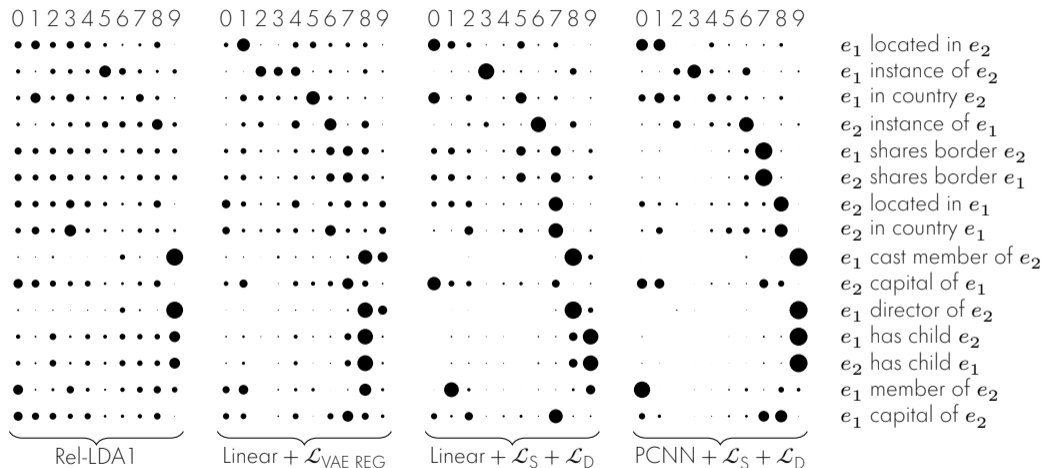
Model		$B^3$			V-measure			ARI
Classifier	Reg.	$F_1$	Prec.	Rec.	$F_1$	Hom.	Comp.	
rel-LDA		29.1	24.8	35.2	30.0	26.1	35.1	13.3
rel-LDA1		36.9	30.4	47.0	37.4	31.9	45.1	24.2
Linear	$\mathcal{L}_{VAE REG}$	35.2	23.8	67.1	27.0	18.6	49.6	18.7
PCNN	$\mathcal{L}_{VAE REG}$	27.6	24.3	31.9	24.7	21.2	29.6	15.7
Linear	$\mathcal{L}_S + \mathcal{L}_D$	37.5	31.1	47.4	<b>38.7</b>	32.6	47.8	27.6
PCNN	$\mathcal{L}_S + \mathcal{L}_D$	<b>39.4</b>	32.2	50.7	38.3	32.2	47.2	<b>33.8</b>
BERTcoder	$\mathcal{L}_S + \mathcal{L}_D$	41.5	34.6	51.8	39.9	33.9	48.5	35.1
BERTcoder	SelfORE	<b>49.1</b>	47.3	51.1	<b>46.6</b>	45.7	47.6	<b>40.3</b>

Model		$B^3$			V-measure			ARI
Classifier	Reg.	$F_1$	Prec.	Rec.	$F_1$	Hom.	Comp.	
rel-LDA		29.1	24.8	35.2	30.0	26.1	35.1	13.3
rel-LDA1		36.9	30.4	47.0	37.4	31.9	45.1	24.2
Linear	$\mathcal{L}_{VAE REG}$	35.2	23.8	67.1	27.0	18.6	49.6	18.7
PCNN	$\mathcal{L}_{VAE REG}$	27.6	24.3	31.9	24.7	21.2	29.6	15.7
Linear	$\mathcal{L}_S + \mathcal{L}_D$	37.5	31.1	47.4	<b>38.7</b>	32.6	47.8	27.6
PCNN	$\mathcal{L}_S + \mathcal{L}_D$	<b>39.4</b>	32.2	50.7	38.3	32.2	47.2	<b>33.8</b>
BERTcoder	$\mathcal{L}_S + \mathcal{L}_D$	41.5	34.6	51.8	39.9	33.9	48.5	35.1
BERTcoder	SelfORE	<b>49.1</b>	47.3	51.1	<b>46.6</b>	45.7	47.6	<b>40.3</b>

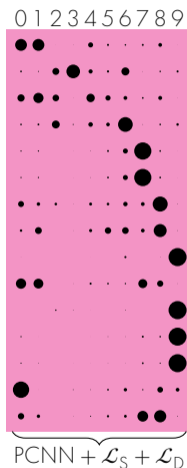
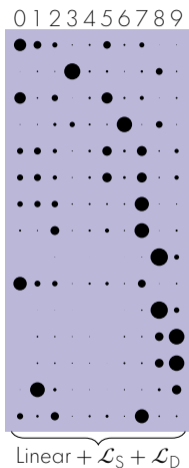
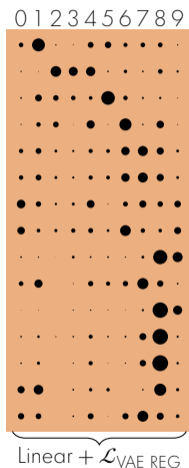
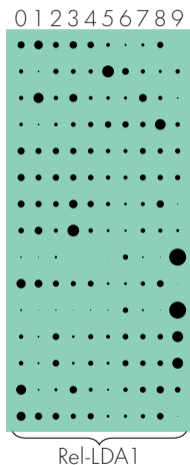
Model		$B^3$			V-measure			ARI
Classifier	Reg.	$F_1$	Prec.	Rec.	$F_1$	Hom.	Comp.	
rel-LDA		29.1	24.8	35.2	30.0	26.1	35.1	13.3
rel-LDA1		36.9	30.4	47.0	37.4	31.9	45.1	24.2
Linear	$\mathcal{L}_{VAE REG}$	35.2	23.8	67.1	27.0	18.6	49.6	18.7
PCNN	$\mathcal{L}_{VAE REG}$	27.6	24.3	31.9	24.7	21.2	29.6	15.7
Linear	$\mathcal{L}_S + \mathcal{L}_D$	37.5	31.1	47.4	<b>38.7</b>	32.6	47.8	27.6
PCNN	$\mathcal{L}_S + \mathcal{L}_D$	<b>39.4</b>	32.2	50.7	38.3	32.2	47.2	<b>33.8</b>
BERTcoder	$\mathcal{L}_S + \mathcal{L}_D$	41.5	34.6	51.8	39.9	33.9	48.5	35.1
BERTcoder	SelfORE	<b>49.1</b>	47.3	51.1	<b>46.6</b>	45.7	47.6	<b>40.3</b>

Model		$B^3$			V-measure			ARI
Classifier	Reg.	$F_1$	Prec.	Rec.	$F_1$	Hom.	Comp.	
rel-LDA		29.1	24.8	35.2	30.0	26.1	35.1	13.3
rel-LDA1		36.9	30.4	47.0	37.4	31.9	45.1	24.2
Linear	$\mathcal{L}_{VAE REG}$	35.2	23.8	67.1	27.0	18.6	49.6	18.7
PCNN	$\mathcal{L}_{VAE REG}$	27.6	24.3	31.9	24.7	21.2	29.6	15.7
Linear	$\mathcal{L}_S + \mathcal{L}_D$	37.5	31.1	47.4	<b>38.7</b>	32.6	47.8	27.6
PCNN	$\mathcal{L}_S + \mathcal{L}_D$	<b>39.4</b>	32.2	50.7	38.3	32.2	47.2	<b>33.8</b>
BERTcoder	$\mathcal{L}_S + \mathcal{L}_D$	41.5	34.6	51.8	39.9	33.9	48.5	35.1
BERTcoder	SelfORE	<b>49.1</b>	47.3	51.1	<b>46.6</b>	45.7	47.6	<b>40.3</b>

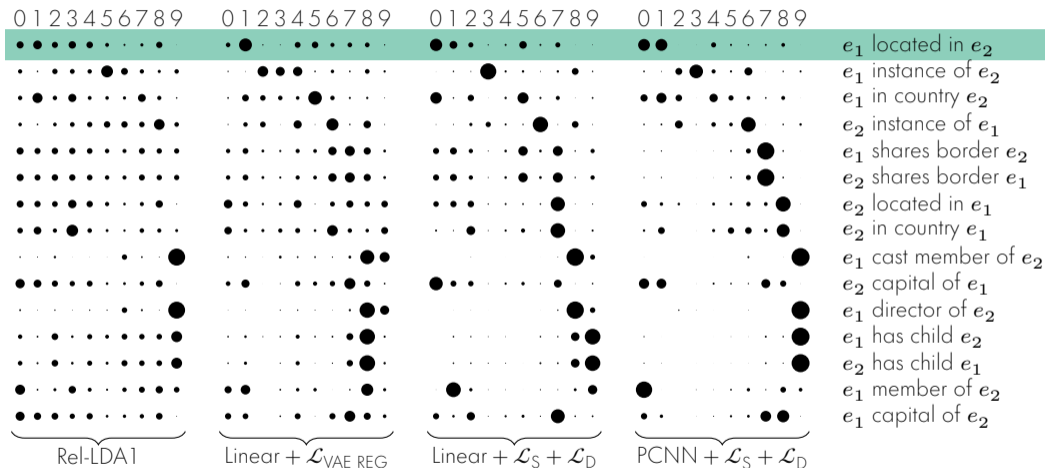
Hu et al. "SelfORE: Self-supervised Relational Feature Learning for Open Relation Extraction" EMNLP 2020

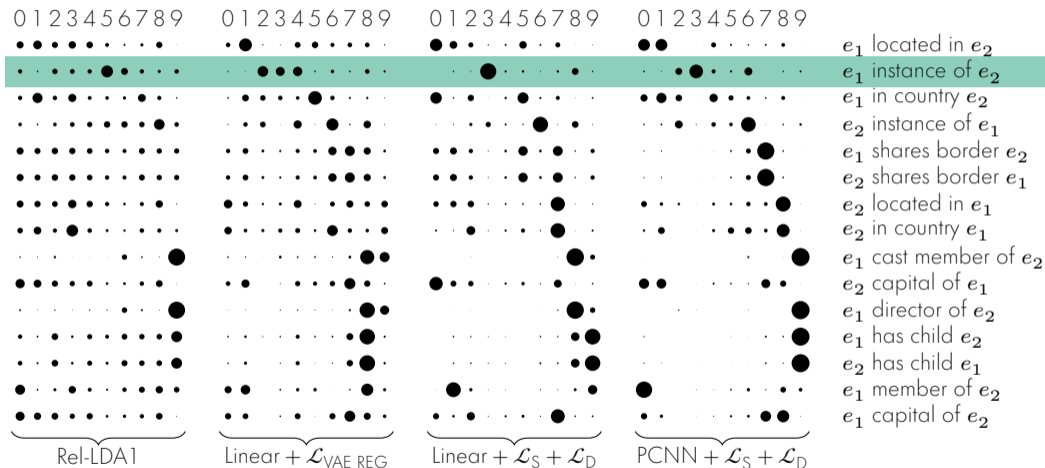


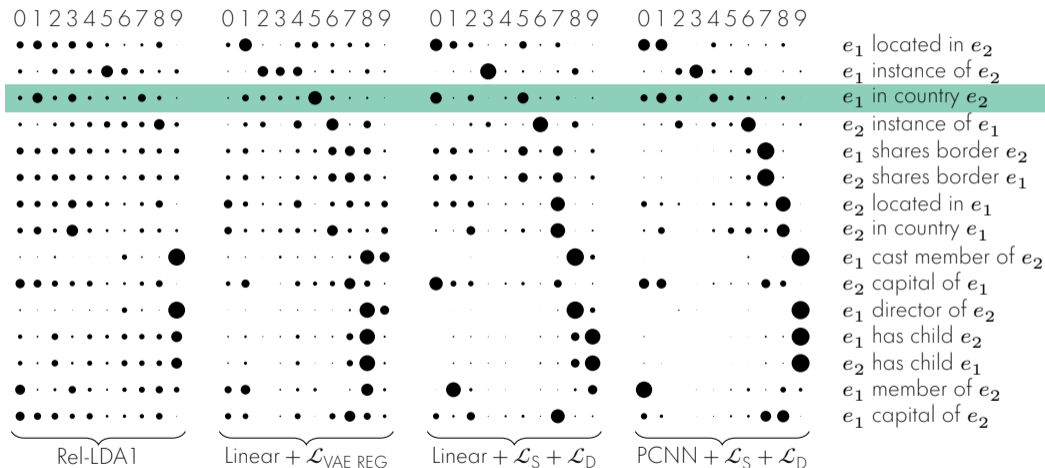


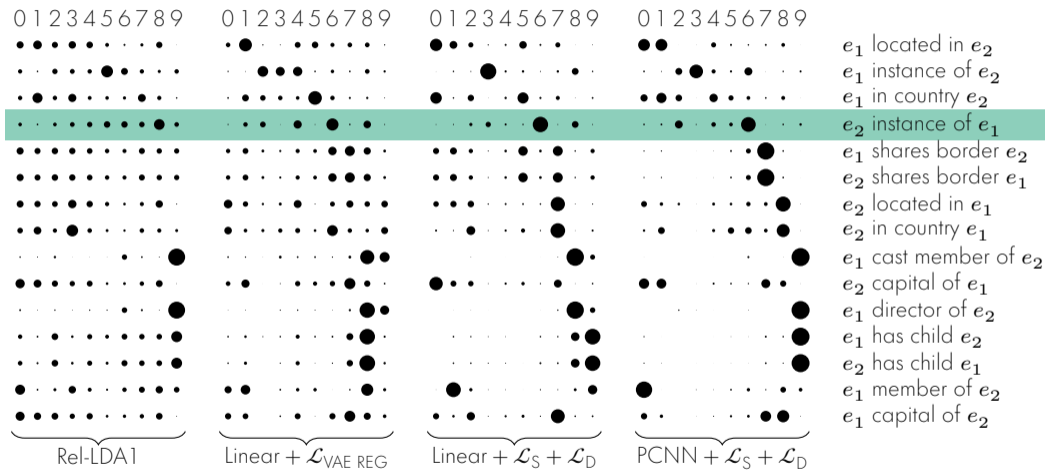


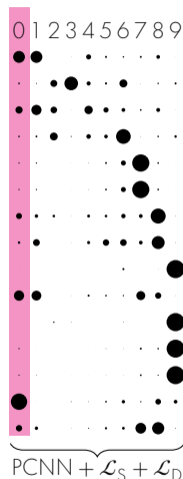
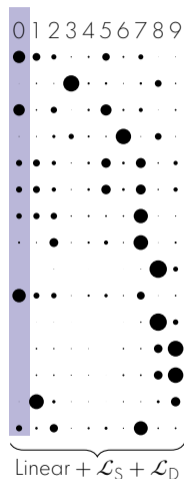
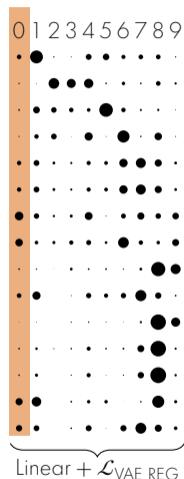
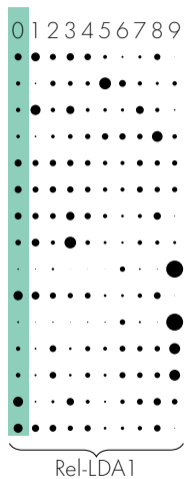
- $e_1$  located in  $e_2$
- $e_1$  instance of  $e_2$
- $e_1$  in country  $e_2$
- $e_2$  instance of  $e_1$
- $e_1$  shares border  $e_2$
- $e_2$  shares border  $e_1$
- $e_2$  located in  $e_1$
- $e_2$  in country  $e_1$
- $e_1$  cast member of  $e_2$
- $e_2$  capital of  $e_1$
- $e_1$  director of  $e_2$
- $e_1$  has child  $e_2$
- $e_2$  has child  $e_1$
- $e_1$  member of  $e_2$
- $e_1$  capital of  $e_2$



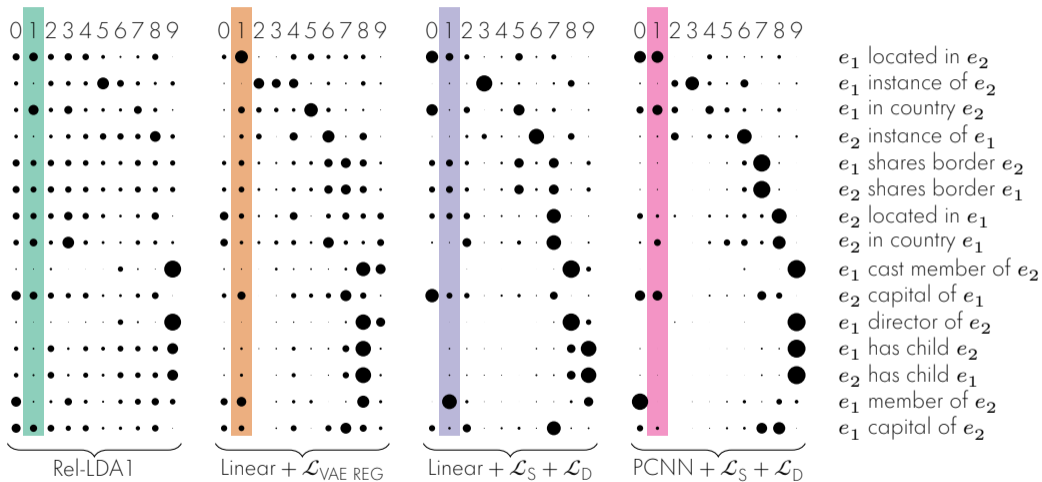


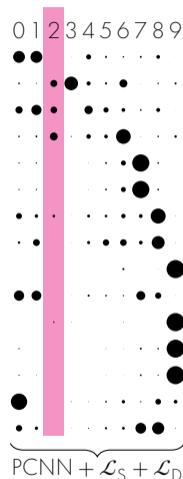
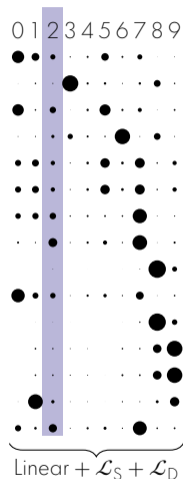
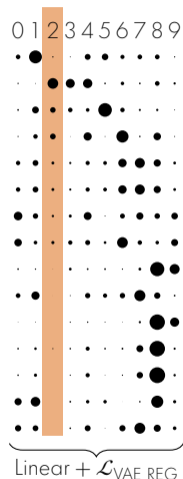
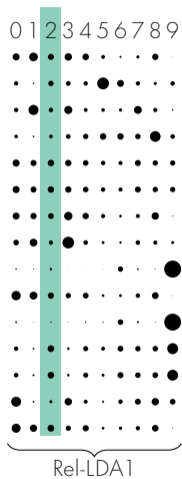






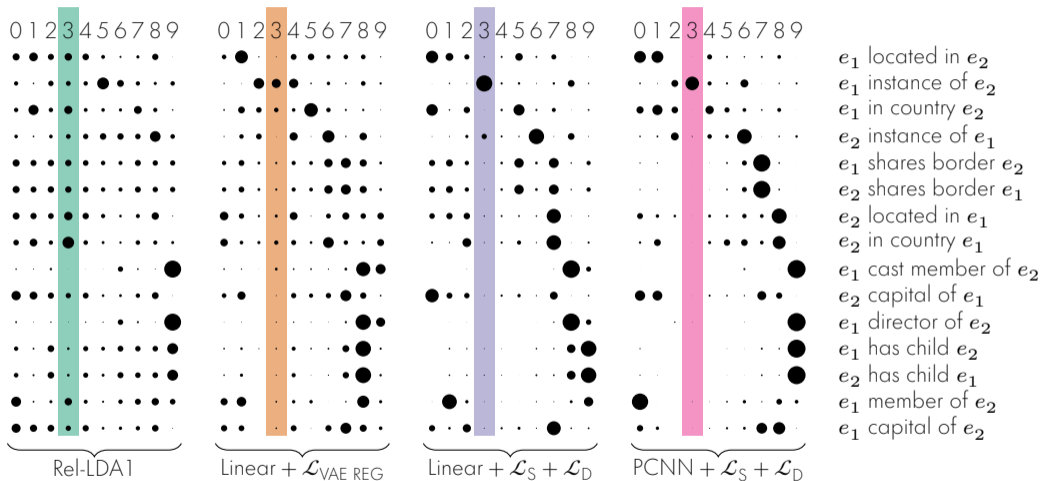
- $e_1$  located in  $e_2$
- $e_1$  instance of  $e_2$
- $e_1$  in country  $e_2$
- $e_2$  instance of  $e_1$
- $e_1$  shares border  $e_2$
- $e_2$  shares border  $e_1$
- $e_2$  located in  $e_1$
- $e_2$  in country  $e_1$
- $e_1$  cast member of  $e_2$
- $e_2$  capital of  $e_1$
- $e_1$  director of  $e_2$
- $e_1$  has child  $e_2$
- $e_2$  has child  $e_1$
- $e_1$  member of  $e_2$
- $e_1$  capital of  $e_2$



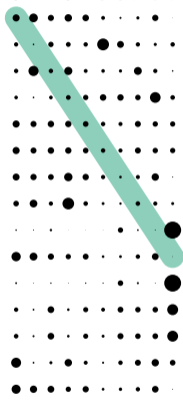


- $e_1$  located in  $e_2$
- $e_1$  instance of  $e_2$
- $e_1$  in country  $e_2$
- $e_2$  instance of  $e_1$
- $e_1$  shares border  $e_2$
- $e_2$  shares border  $e_1$
- $e_2$  located in  $e_1$
- $e_2$  in country  $e_1$
- $e_1$  cast member of  $e_2$
- $e_2$  capital of  $e_1$
- $e_1$  director of  $e_2$
- $e_1$  has child  $e_2$
- $e_2$  has child  $e_1$
- $e_1$  member of  $e_2$
- $e_1$  capital of  $e_2$



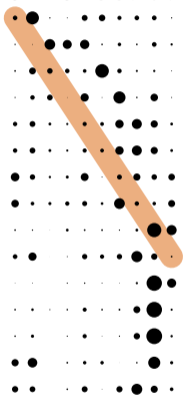


0 1 2 3 4 5 6 7 8 9

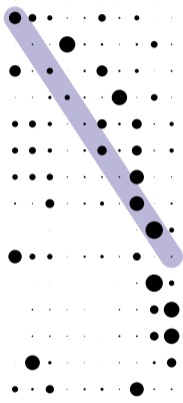


Rel-LDA1

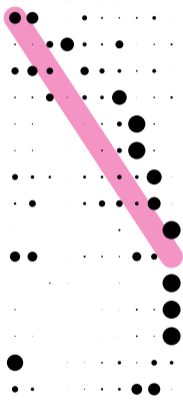
0 1 2 3 4 5 6 7 8 9

Linear +  $\mathcal{L}_{VAE REG}$ 

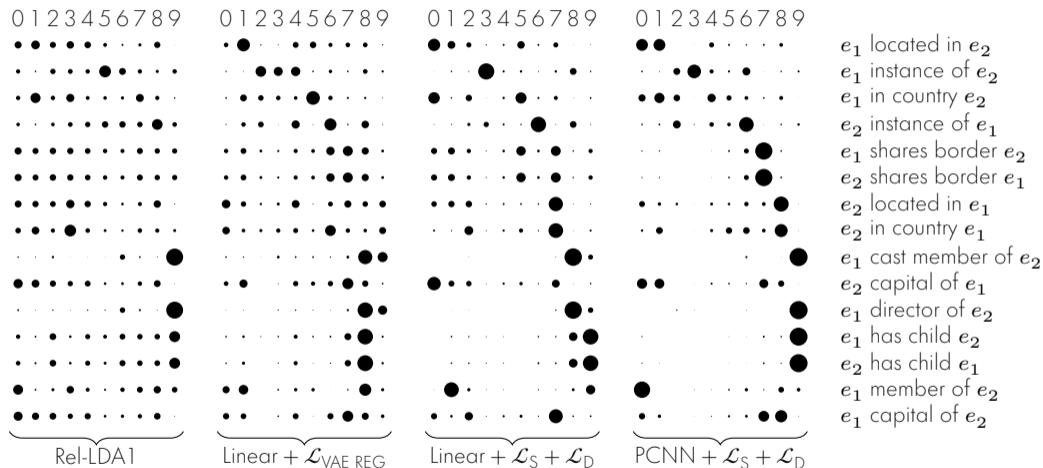
0 1 2 3 4 5 6 7 8 9

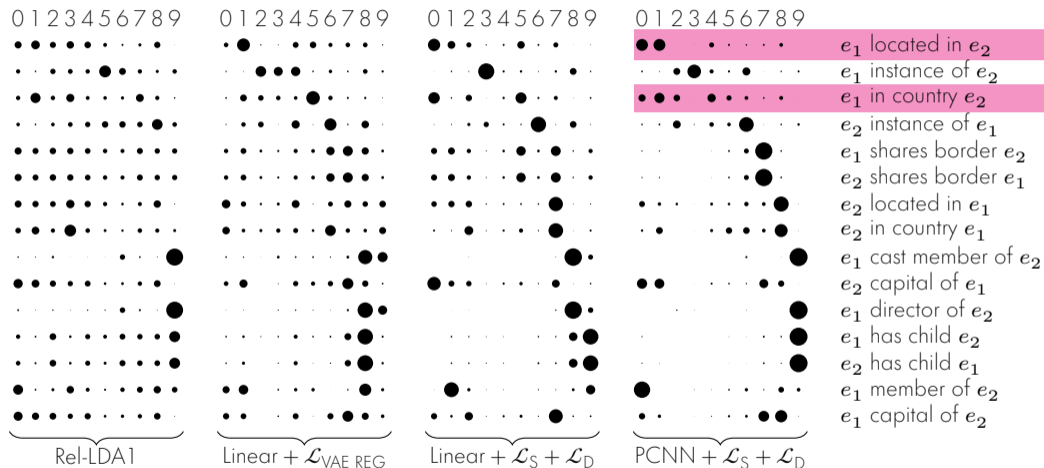
Linear +  $\mathcal{L}_S + \mathcal{L}_D$ 

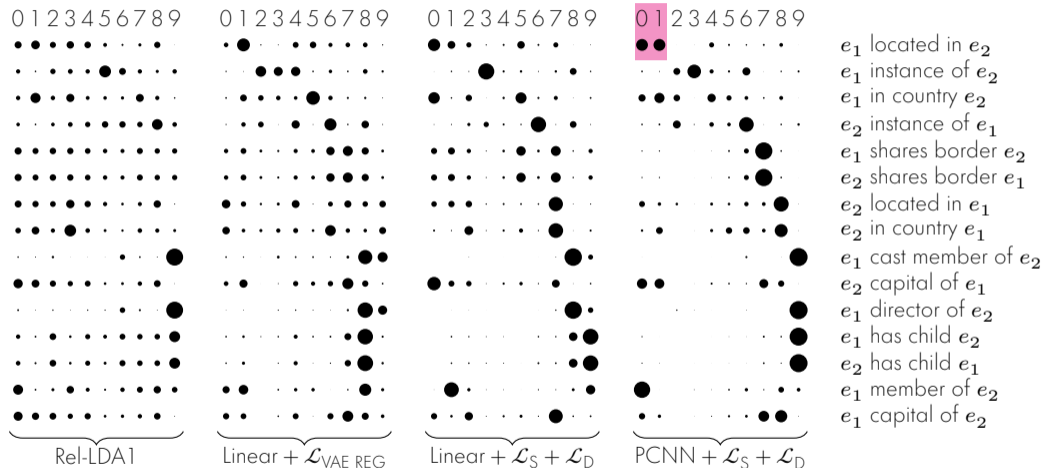
0 1 2 3 4 5 6 7 8 9

PCNN +  $\mathcal{L}_S + \mathcal{L}_D$ 

- $e_1$  located in  $e_2$
- $e_1$  instance of  $e_2$
- $e_1$  in country  $e_2$
- $e_2$  instance of  $e_1$
- $e_1$  shares border  $e_2$
- $e_2$  shares border  $e_1$
- $e_2$  located in  $e_1$
- $e_2$  in country  $e_1$
- $e_1$  cast member of  $e_2$
- $e_2$  capital of  $e_1$
- $e_1$  director of  $e_2$
- $e_1$  has child  $e_2$
- $e_2$  has child  $e_1$
- $e_1$  member of  $e_2$
- $e_1$  capital of  $e_2$







## Take-home Message

Selecting good regularizations to enforce modeling hypotheses enables us to train a deep classifier.

## Contributions

- Train a PCNN without supervision
- Designed two regularization losses (Skewness, Distribution distance)
- Introduced new datasets (T-RExes)
- Evaluated using additional metrics (V-measure, ARI)

## Graph-based Aggregate Extraction

---

Megrez <sub>$e_1$</sub> <sup>Q850779</sup> is a star in the northern circum-polar constellation of Ursa Major <sub>$e_2$</sub> <sup>Q10460</sup>. }  $x_1$

Posidonius <sub>$e_1$</sub> <sup>Q185770</sup> was a Greek philosopher, astronomer, historian, mathematician, and teacher native to Apamea, Syria <sub>$e_2$</sub> <sup>Q617550</sup>. }  $x_2$

Hipparchus <sub>$e_1$</sub> <sup>Q159905</sup> was born in Nicaea, Bithynia <sub>$e_2$</sub> <sup>Q739037</sup>, and probably died on the island of Rhodes, Greece. }  $x_3$

Learn a similarity function  
 $\text{sim} : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$

$\text{sim}(x_1, x_2) < \text{sim}(x_2, x_3)$

$\text{sim}(x_1, x_3) < \text{sim}(x_2, x_3)$

5 way 1 shot: given 1 query and 5 candidates, which of the candidates is most similar to the query?

Evaluated using accuracy.



**Sentential approaches:** extract sentences' relation independently ( $\mathcal{S} \times \mathcal{E}^2 \rightarrow \mathcal{R}$ )

**Aggregate approaches:** maps a set of sentences to a set of facts ( $2^{\mathcal{S} \times \mathcal{E}^2} \rightarrow 2^{\mathcal{E}^2 \times \mathcal{R}}$ )

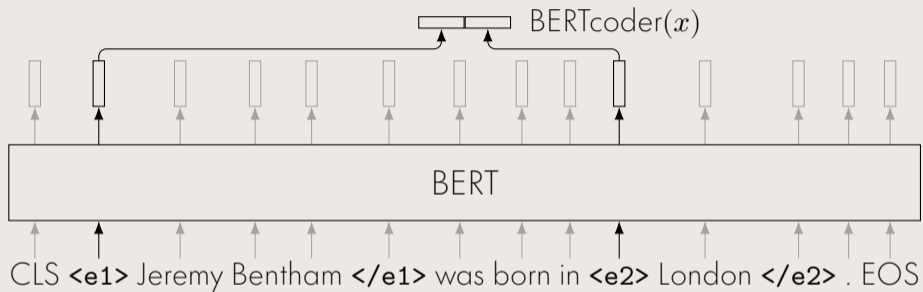
## Goal

Exploit dataset-level regularities to leverage additional information

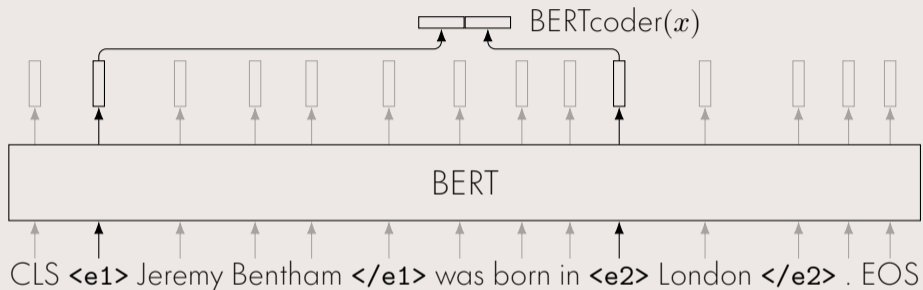
## Plan

1. Model datasets as graphs
2. Related relation extraction work only uses **linguistic** similarities
3. Proof that **topological** information can be used
4. How topological features are usually extracted (GCN)
5. How to extract them differently (WL isomorphism test)
6. Experimental results
7. Perspective

## BERTcoder (linguistic)



## BERTcoder (linguistic)

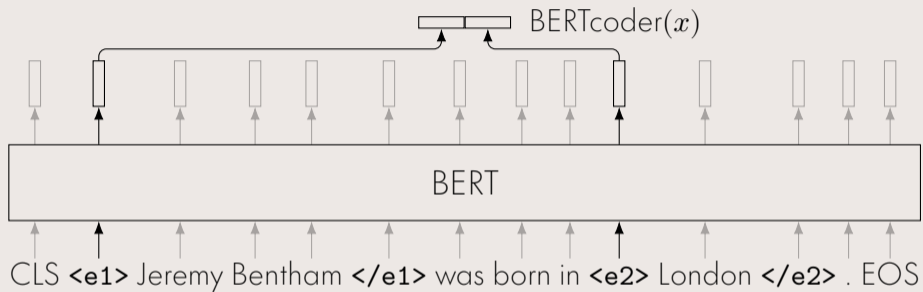


## Prediction

Compare samples using:

$$\text{sim}(x, x') = \text{sigmoid}(\text{BERTcoder}(x)^T \text{BERTcoder}(x'))$$

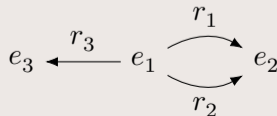
## BERTcoder (linguistic)



## Prediction

Compare samples using:  
 $\text{sim}(x, x') = \text{sigmoid}(\text{BERTcoder}(x)^T \text{BERTcoder}(x'))$

## Hypotheses



MTB assumes:

$$r_1 = r_2 \ (\mathcal{H}_{1\text{-ADJACENCY}})$$

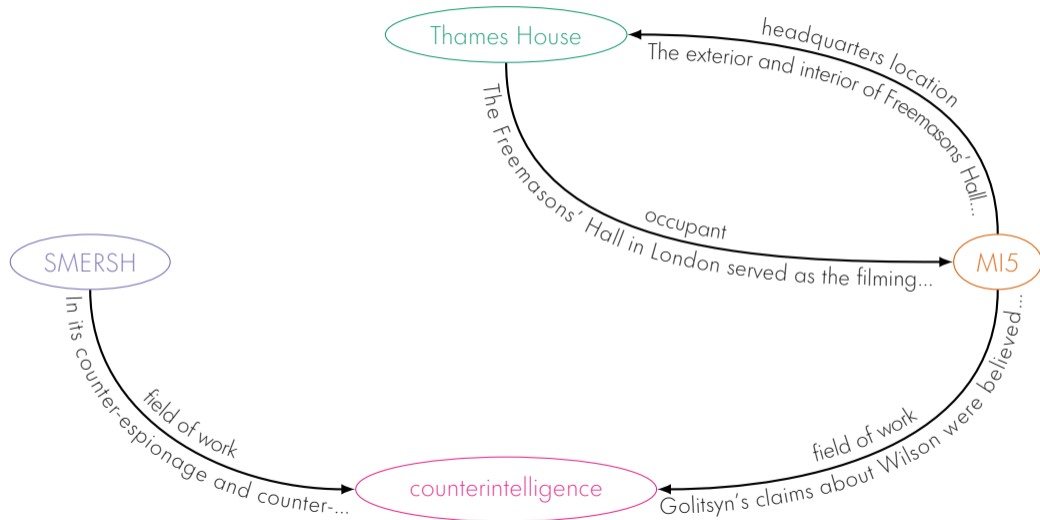
$$r_3 \neq r_1 \wedge r_3 \neq r_2 \ (\mathcal{H}_{1 \rightarrow 1})$$

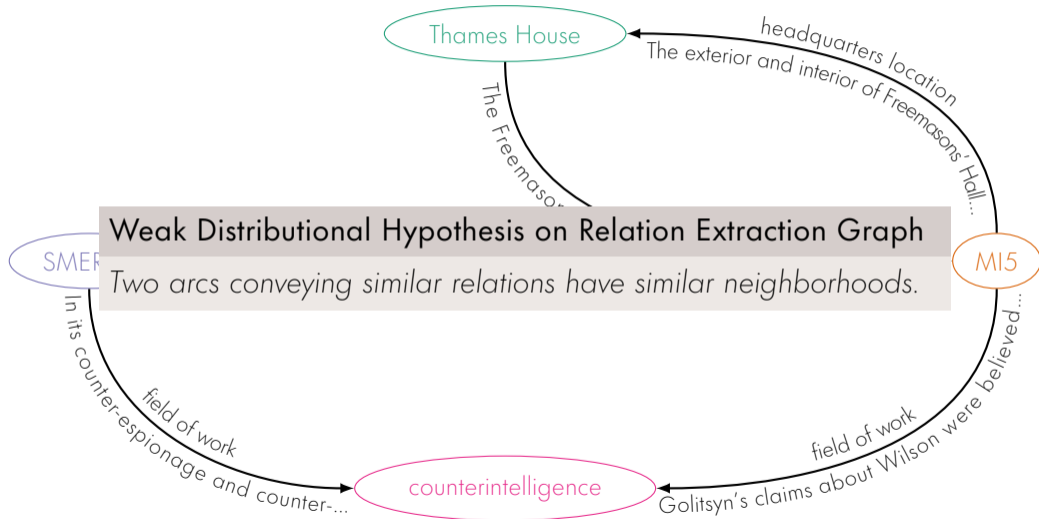
The exterior and interior of Freemasons' Hall continued to be a stand-in for Thames House<sub>e<sub>2</sub></sub>, the headquarters of MI5<sub>e<sub>1</sub></sub>.

Golitsyn's claims about Wilson were believed in particular by the senior MI5<sub>e<sub>1</sub></sub> counterintelligence<sub>e<sub>2</sub></sub> officer Peter Wright.

In its counter-espionage<sub>e<sub>2</sub></sub> and counter-intelligence roles, SMERSH<sub>e<sub>1</sub></sub> appears to have been extremely successful throughout World War II.

The Freemasons' Hall in London served as the filming location for Thames House<sub>e<sub>1</sub></sub>, the headquarters for MI5<sub>e<sub>2</sub></sub>.





### Proposition

Given the **path**  $e_1 \xrightarrow{r_1} e_2 \xrightarrow{r_2} e_3 \xrightarrow{r_3} e_4$ , we expect  $r_1 \not\perp r_2 \not\perp r_3$ .

### Goal

Compute the mutual information  $I(r_2; r_1, r_3)$



## Proposition

Given the **path**  $e_1 \xrightarrow{r_1} e_2 \xrightarrow{r_2} e_3 \xrightarrow{r_3} e_4$ , we expect  $r_1 \not\perp r_2 \not\perp r_3$ .

## Goal

Compute the mutual information  $I(r_2; r_1, r_3)$

## Path Counting Algorithm

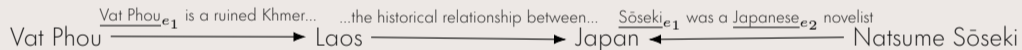
We can (slowly) sample **walks** using power of the adjacency matrix.

1. Sample a walk by chaining neighbors
2. Reject non-path
3. Count the accepted paths weighted by importance

## Path Frequency

Frequency	Relation Surface forms	Relation Identifiers
31.696%	<i>country • diplomatic relation • citizen of</i>	<b>P17 • P530 • P27</b>

Example of path:



## Path Frequency

Frequency	Relation Surface forms	Relation Identifiers
31.696‰	<i>country • diplomatic relation • citizen of</i>	<b>P17 • P530 • P27</b>

Example of path:



## Summary Statistics

$$I(r_2; r_1, r_3) = \mathbb{E}_{r_1, r_3} [\mathbb{H}_{P(r_2)}(r_2 | r_1, r_3)] - \mathbb{H}(r_2 | r_1, r_3)$$

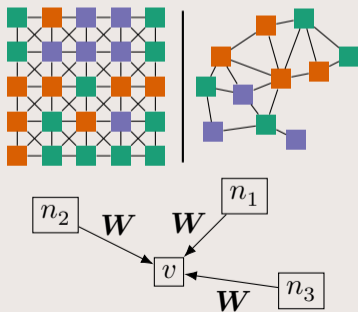
$$\begin{array}{ccc} \approx & \approx & \approx \\ 6.95 \text{ bits} & 8.01 \text{ bits} & 1.06 \text{ bits} \end{array}$$

## Modeling Hypothesis

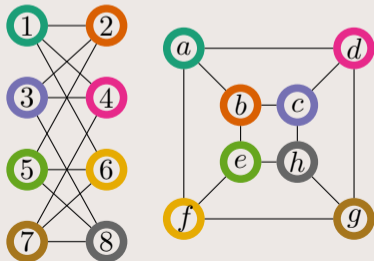
$\mathcal{H}_{1\text{-NEIGHBORHOOD}}$ : Two samples with the same neighborhood in the relation extraction graph convey the same relation.

$$\forall a, a' \in \mathcal{A}: \mathcal{N}(a) = \mathcal{N}(a') \implies \rho(a) = \rho(a')$$

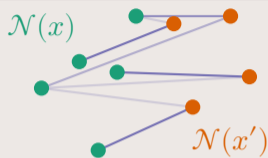
## Graph Convolutional Network



## Graph Isomorphism



## Earth Mover Distance



## Compare Topological Features

Skip recoloring, directly compare neighborhoods in  $\mathbb{R}^d$ :

$S(x, k) =$  samples at distance  $k$  of  $x$

$\mathfrak{S}(x, k) =$   
 $\{ \text{BERTcoder}(y) \in \mathbb{R}^d \mid y \in S(x, k) \}$

$$W_1(\mathfrak{S}(x, 1), \mathfrak{S}(x', 1))$$

**algorithm** WEISFEILER-LEMAN

*Inputs:*  $G = (V, E)$  graph

$k$  dimensionality

*Output:*  $\chi_\infty$  coloring of  $k$ -tuples

$\chi_0(\mathbf{x}) \leftarrow \text{iso}(\mathbf{x}) \quad \forall \mathbf{x} \in V^k$

**for**  $\ell = 1, 2, \dots$  **do**

$\mathfrak{I}_\ell \leftarrow$  new color index

**for all**  $\mathbf{x} \in V^k$  **do**

$c_\ell(\mathbf{x}) \leftarrow$

$\{ \chi_{\ell-1}(\mathbf{y}) \mid \mathbf{y} \in N^k(\mathbf{x}) \}$

$\chi_\ell(\mathbf{x}) \leftarrow$

$(\chi_{\ell-1}(\mathbf{x}), c_\ell(\mathbf{x}))$  in  $\mathfrak{I}_\ell$

**until**  $\chi_\ell = \chi_{\ell-1}$

**output**  $\chi_\ell$

## Redefining similarity

We keep the **linguistic** similarity from MTB:

$$\text{sim}_{\text{ling}}(x, x') = \text{sigmoid}(\text{BERTcoder}(x)^{\top} \text{BERTcoder}(x'))$$

But also define a **topological** similarity:

Either using GCN:

$$\text{sim}_{\text{topo}}^{\text{GCN}}(x, x') = \text{sigmoid}(\text{GCN}(G)_x^{\top} \text{GCN}(G)_{x'})$$

Or 1-Wasserstein:

$$\text{sim}_{\text{topo}}^{W_1}(x, x') = -W_1(\mathfrak{S}(x, 1), \mathfrak{S}(x', 1))$$

Define the **topolinguistic** similarity as:

$$\text{sim}_{\text{topoling}}(x, x') = \text{sim}_{\text{ling}}(x, x') + \lambda \text{sim}_{\text{topo}}(x, x')$$

Model	Accuracy
<b>Pre-trained</b>	
Linguistic (BERT)	69.46
Topological ( $W_1$ )	65.75
Topolinguistic	72.18
<b>Fine-tuned</b>	
MTB	78.83
MTB GCN-Chebyshev	76.10

### Few-Shot Evaluation

1 query

5 candidates

Which candidate conveys the same relation as the query?

Random model score 20% accuracy.

Model	Accuracy
<b>Pre-trained</b>	
Linguistic (BERT)	69.46
Topological ( $W_1$ )	65.75
Topolinguistic	72.18
<b>Fine-tuned</b>	
MTB	78.83
MTB GCN-Chebyshev	76.10

### Few-Shot Evaluation

1 query

5 candidates

Which candidate conveys the same relation as the query?

Random model score 20% accuracy.

Soares et al. "Matching the Blanks: Distributional Similarity for Relation Learning"  
ACL 2019



Model	Accuracy
<b>Pre-trained</b>	
Linguistic (BERT)	69.46
Topological ( $W_1$ )	65.75
Topolinguistic	72.18
<b>Fine-tuned</b>	
MTB	78.83
MTB GCN-Chebyshev	76.10

### Few-Shot Evaluation

1 query

5 candidates

Which candidate conveys the same relation as the query?

Random model score 20% accuracy.

Model	Accuracy
<b>Pre-trained</b>	
Linguistic (BERT)	69.46
Topological ( $W_1$ )	65.75
Topolinguistic	72.18
<b>Fine-tuned</b>	
MTB	78.83
MTB GCN-Chebyshev	76.10

### Few-Shot Evaluation

1 query

5 candidates

Which candidate conveys the same relation as the query?

Random model score 20% accuracy.

Model	Accuracy
<b>Pre-trained</b>	
Linguistic (BERT)	69.46
Topological ( $W_1$ )	65.75
Topolinguistic	72.18
<b>Fine-tuned</b>	
MTB	78.83
MTB GCN-Chebyshev	76.10

### Few-Shot Evaluation

1 query

5 candidates

Which candidate conveys the same relation as the query?

Random model score 20% accuracy.

Soares et al. "Matching the Blanks: Distributional Similarity for Relation Learning"  
ACL 2019

Model	Accuracy
<b>Pre-trained</b>	
Linguistic (BERT)	69.46
Topological ( $W_1$ )	65.75
Topolinguistic	72.18
<b>Fine-tuned</b>	
MTB	78.83
MTB GCN-Chebyshev	76.10

### Few-Shot Evaluation

1 query

5 candidates

Which candidate conveys the same relation as the query?

Random model score 20% accuracy.

## Take-home Message

Topological information can be leverage for unsupervised relation extraction.

## Contributions

- Explicitly modeled the aggregate setup for the unsupervised problem.
- Provided proof on the quality of topological information.
- Proposed an approach to exploit the mutual information between topological and linguistic features.

Several directions still need to be explored.

Use the topological features to identify the relational information in the linguistic features.

$$\mathcal{L}_{\text{LT}}(x_1, x_2, x_3) = \max \left( \begin{array}{l} 0, \zeta + 2(\text{sim}_{\text{ling}}(x_1, x_2) - \text{sim}_{\text{topo}}(x_1, x_2))^2 \\ \quad - (\text{sim}_{\text{ling}}(x_1, x_2) - \text{sim}_{\text{topo}}(x_1, x_3))^2 \\ \quad - (\text{sim}_{\text{ling}}(x_1, x_3) - \text{sim}_{\text{topo}}(x_1, x_2))^2 \end{array} \right)$$

Use the topological features to identify the relational information in the linguistic features.

$$\mathcal{L}_{\text{LT}}(x_1, x_2, x_3) = \max \left( \begin{array}{l} 0, \zeta + 2 \left( \text{sim}_{\text{ling}}(x_1, x_2) - \text{sim}_{\text{topo}}(x_1, x_2) \right)^2 \\ \quad - \left( \text{sim}_{\text{ling}}(x_1, x_2) - \text{sim}_{\text{topo}}(x_1, x_3) \right)^2 \\ \quad - \left( \text{sim}_{\text{ling}}(x_1, x_3) - \text{sim}_{\text{topo}}(x_1, x_2) \right)^2 \end{array} \right)$$

- Ideally we want to align the two similarities. ←

Use the topological features to identify the relational information in the linguistic features.

$$\mathcal{L}_{\text{LT}}(x_1, x_2, x_3) = \max \left( \begin{array}{l} 0, \zeta + 2 \left( \text{sim}_{\text{ling}}(x_1, x_2) - \text{sim}_{\text{topo}}(x_1, x_2) \right)^2 \\ - \left( \text{sim}_{\text{ling}}(x_1, x_2) - \text{sim}_{\text{topo}}(x_1, x_3) \right)^2 \\ - \left( \text{sim}_{\text{ling}}(x_1, x_3) - \text{sim}_{\text{topo}}(x_1, x_2) \right)^2 \end{array} \right)$$

- Ideally we want to align the two similarities. ←
- However to stabilize the loss we need to use negative samples. ←



Use the topological features to identify the relational information in the linguistic features.

$$\mathcal{L}_{\text{LT}}(x_1, x_2, x_3) = \max \left( 0, \zeta + 2 \left( \text{sim}_{\text{ling}}(x_1, x_2) - \text{sim}_{\text{topo}}(x_1, x_2) \right)^2 - \left( \text{sim}_{\text{ling}}(x_1, x_2) - \text{sim}_{\text{topo}}(x_1, x_3) \right)^2 - \left( \text{sim}_{\text{ling}}(x_1, x_3) - \text{sim}_{\text{topo}}(x_1, x_2) \right)^2 \right)$$

- Ideally we want to align the two similarities.
- However to stabilize the loss we need to use negative samples.
- Up to a margin  $\zeta$ .

Questions?

---

## Supplementary Material

---

$\mathcal{H}_{\text{DISTANT}}$ 

A sentence conveys all the possible relations between all the entities it contains.

$$\mathcal{D}_{\mathcal{R}} = \mathcal{D} \bowtie \mathcal{D}_{\text{KB}}$$

where  $\bowtie$  denotes the natural join operator:

$$\mathcal{D} \bowtie \mathcal{D}_{\text{KB}} = \{ (s, e_1, e_2, r) \mid (s, e_1, e_2) \in \mathcal{D} \wedge (e_1, e_2, r) \in \mathcal{D}_{\text{KB}} \}.$$

1. the bag of words of the infix;
2. the surface form of the entities;
3. the lemma words on the dependency path;
4. the POS of the infix words;
5. the type of the entity pair (e.g. person–location);
6. the type of the head entity (e.g. person);
7. the type of the tail entity (e.g. location);
8. the words on the dependency path between the two entities.

$$B^3 \text{ precision}(g, c) = \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \mathcal{U}(\mathcal{D}_{\mathcal{X}})} P(g(\mathbf{X}) = g(\mathbf{Y}) \mid c(\mathbf{X}) = c(\mathbf{Y}))$$

$$B^3 \text{ recall}(g, c) = \mathbb{E}_{\mathbf{X}, \mathbf{Y} \sim \mathcal{U}(\mathcal{D}_{\mathcal{X}})} P(c(\mathbf{X}) = c(\mathbf{Y}) \mid g(\mathbf{X}) = g(\mathbf{Y}))$$

$$B^3 F_1(g, c) = \frac{2}{B^3 \text{ precision}(g, c)^{-1} + B^3 \text{ recall}(g, c)^{-1}}$$

$$\text{homogeneity}(g, c) = 1 - \frac{H(c(\mathbf{X}) | g(\mathbf{X}))}{H(c(\mathbf{X}))}$$

$$\text{completeness}(g, c) = 1 - \frac{H(g(\mathbf{X}) | c(\mathbf{X}))}{H(g(\mathbf{X}))}$$

$$\text{V-measure}(g, c) = \frac{2}{\text{homogeneity}(g, c)^{-1} + \text{completeness}(g, c)^{-1}}$$

$$\text{RI}(g, c) = \mathbb{E}_{\mathbf{X}, \mathbf{Y}} [P(c(\mathbf{X}) = c(\mathbf{Y}) \Leftrightarrow g(\mathbf{X}) = g(\mathbf{Y}))]$$

$$\text{ARI}(g, c) = \frac{\text{RI}(g, c) - \mathbb{E}_{c \sim \mathcal{U}(\mathcal{R}^{\mathcal{D}})} [\text{RI}(g, c)]}{\max_{c \in \mathcal{R}^{\mathcal{D}}} \text{RI}(g, c) - \mathbb{E}_{c \sim \mathcal{U}(\mathcal{R}^{\mathcal{D}})} [\text{RI}(g, c)]}$$



$$\pi_r = \frac{(\exp(y_r) + G_r) / \tau}{\sum_{r' \in \mathcal{R}} (\exp(y_{r'}) + G_{r'}) / \tau}$$

Confidence	B <sup>3</sup>			V-measure			ARI
	F <sub>1</sub>	Prec.	Rec.	F <sub>1</sub>	Hom.	Comp.	
$\mathcal{L}_S$ regularization	39.4	32.2	50.7	38.3	32.2	47.2	33.8
Gumbel-Softmax	35.0	29.9	42.2	33.2	28.3	40.2	25.1

$$P(r = r \mid s, \mathbf{e}; \boldsymbol{\theta}, \boldsymbol{\phi}) = P(r_s = r \mid s; \boldsymbol{\phi})P(r_e = r \mid \mathbf{e}; \boldsymbol{\theta})$$

$$\mathcal{L}_{\text{ALIGN}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = -\log \sum_{r \in \mathcal{R}} P(r \mid s, \mathbf{e}; \boldsymbol{\theta}, \boldsymbol{\phi}) + \mathcal{L}_{\text{D}}(\boldsymbol{\theta}) + \mathcal{L}_{\text{D}}(\boldsymbol{\phi}).$$

Model	B <sup>3</sup>			V-measure			ARI
	F <sub>1</sub>	Prec.	Rec.	F <sub>1</sub>	Hom.	Comp.	
$\mathcal{L}_{\text{EP}} + \mathcal{L}_{\text{S}} + \mathcal{L}_{\text{D}}$	39.4	32.2	50.7	38.3	32.2	47.2	33.8
$\mathcal{L}_{\text{ALIGN}}$ average	37.6	30.3	49.7	39.4	33.1	48.8	20.3
$\mathcal{L}_{\text{ALIGN}}$ maximum	41.2	33.6	53.4	43.5	36.9	53.1	29.5
$\mathcal{L}_{\text{ALIGN}}$ minimum	34.5	26.5	49.3	35.9	29.6	45.7	15.3

## Spectral (convolution is multiplication in Fourier space)

	Graph	Euclidean
Laplacian	$L = D - M$	$\nabla^2$
$\hookrightarrow$ Eigenfunctions	$U$ s.t. $L = U\Lambda U^{-1}$	$\xi \mapsto e^{2\pi i \xi x}$
Fourier transform	$U^T \mathbf{f}$	$\mathcal{F}(f) = \int_{-\infty}^{\infty} f(x) e^{2\pi i \xi x} dx$
Convolution	$U(U^T \mathbf{w} U^T \mathbf{f})$	$\mathcal{F}^{-1}(\mathcal{F}(w) \mathcal{F}(f))$

## Spatial

$$\text{GCN}(\mathbf{X}; \mathbf{W})_v = \text{ReLU} \left( \frac{1}{|N(v)|} \sum_{n_i \in N(v)} \mathbf{W} \mathbf{X}_{n_i} \right)$$

