IJCAI 2015

# Iterative Learning of Parallel Lexicons and Phrases from Non-Parallel Corpora

Meiping Dong, Yang Liu, Huanbo Luan, Maosong Sun, Tatsuya Izuha, and Dakun Zhang

# Learning Lexicons from Parallel Corpora

- Parallel corpora are important for learning lexicons

Garcia y asociados

Garcia and associates

sus asociados no son fuertes

his associates are not strong

Garcia y sus asociados no son enemigos

Garcia and his associates are not enemies

parallel corpora

(Brown et al., 1993)

# Learning Lexicons from Parallel Corpora

- Parallel corpora are important for learning lexicons

Garcia y asociados
Garcia and associates

sus asociados no son fuertes
his associates are not strong

Garcia y sus asociados no son enemigos
Garcia and his associates are not enemies

parallel corpora

(Brown et al., 1993)

# Learning Lexicons from Parallel Corpora

- Parallel corpora are important for learning lexicons

Garcia y asociados
Garcia and associates

sus asociados no son fuertes
his associates are not strong

Garcia y sus asociados no son enemigos
Garcia and his associates are not enemies

parallel corpora

| Spanish | English |
|---|---|
| Garcia | Garcia |
| y | and |
| asociados | associates |
| sus | his |
| no | not |
| son | are |
| fuertes | strong |
| enemigos | enemies |

parallel lexicon

(Brown et al., 1993)

# Limitation of Parallel Corpora

- Parallel corpora are limited in both quantity and coverage

# Limitation of Parallel Corpora

- Parallel corpora are limited in both quantity and coverage



**Can we learn lexicons from non-parallel corpora?**

# Previous Work

- Joint parallel sentence and lexicon extraction (Fung and Cheung, 2004)

- Parallel sub-sentential fragment extraction (Munteanu and Marcu, 2006; Quirk et al., 2007; Cettolo et al., 2010)

- Extracting parallel phrases using lexicons (Zhang and Zong, 2013)

- Translation as decipherments (Ravi and Knight, 2011; Nuhn et al., 2012; Dou and Knight, 2012)

- Transductive learning on monolingual corpora (Ueffing et al., 2007; Bertoldi and Federico, 2009)

# Our Work

- Learning lexicons and phrases from non-parallel corpora

Garcia y asociados

sus asociados no son fuertes

no son enemigos

Buenos Ares

his associates are not strong

Messi scored a goal

Garcia and associates

# Our Work

- Learning lexicons and phrases from non-parallel corpora

Garcia y asociados

sus asociados no son fuertes

no son enemigos

Buenos Ares

his associates are not strong

Messi scored a goal

Garcia and associates

# Our Work

- Learning lexicons and phrases from non-parallel corpora

Garcia y asociados

Buenos Ares

sus asociados no son fuertes

his associates are not strong

Messi scored a goal

no son enemigos

Garcia and associates

# Our Work

- Learning lexicons and phrases from non-parallel corpora

Garcia y asociados

sus asociados no son fuertes

no son enemigos

Buenos Ares

his associates are not strong

Messi scored a goal

Garcia and associates

# Our Work

- Learning lexicons and phrases from non-parallel corpora

Garcia y asociados

sus asociados no son fuertes

no son enemigos

Buenos Ares

his associates are not strong

Messi scored a goal

Garcia and associates

# Idea

- "**find one, get more**" (Fung and Cheung, 2004) **iteratively**

Garcia y asociados

Buenos Ares

`sus` asociados `no` `son` fuertes

`his` associates `are` `not` strong

Messi scored a goal

no son enemigos

Garcia and associates

**iteration 1**

6

# Idea

- "**find one, get more**" (Fung and Cheung, 2004) **iteratively**

Garcia y asociados

sus asociados no son fuertes

no son enemigos

Buenos Ares

his associates are not strong

Messi scored a goal

Garcia and associates

**iteration 2**

7

# Idea

- "**find one, get more**" (Fung and Cheung, 2004) **iteratively**

Garcia y asociados

sus asociados no son fuertes

no son enemigos

Buenos Ares

his associates are not strong

Messi scored a goal

Garcia and associates

**iteration 3**

8

# Idea

- "**find one, get more**" (Fung and Cheung, 2004) **iteratively**

Garcia y asociados

sus asociados no son fuertes

no son enemigos

Buenos Ares

his associates are not strong

Messi scored a goal

Garcia and associates

**iteration 4**

9

# Model

- Extending IBM models in a non-parallel scenario

**phrase matching**

$$P(F|E; \boldsymbol{\theta}) = \sum_{\mathbf{m}} \frac{p(T|S)}{(S+1)^T} \prod_{t=1}^{T} P(\mathbf{f}^{(t)}|\mathbf{e}^{(\mathbf{m}_t)}; \boldsymbol{\theta})$$

**word alignment**

$$P(\mathbf{f}^{(t)}|\mathbf{e}^{(\mathbf{m}_t)}; \boldsymbol{\theta})$$

$$= \quad \delta(\mathbf{m}_t, 0)\epsilon +$$

$$(1 - \delta(\mathbf{m}_t, 0)) \frac{p(J^{(t)}|I^{(\mathbf{m}_t)})}{(I^{(\mathbf{m}_t)}+1)^{J^{(t)}}} \prod_{j=1}^{J^{(t)}} \sum_{i=0}^{I^{(\mathbf{m}_t)}} p(\mathbf{f}_j^{(t)}|\mathbf{e}_i^{(\mathbf{m}_t)})$$

# Training Objective

- Objective: maximizing the likelihood of training data

- We use a seed dictionary to inject prior knowledge

$$
\begin{aligned}
J(\boldsymbol{\theta}) \;=\; & P(F|E;\boldsymbol{\theta}) - \\
& \sum_f \sum_e \sigma(f,e,\mathbf{d}) \log \frac{\sigma(f,e,\mathbf{d})}{p(f|e)} - \quad \text{seed dictionary} \\
& - \sum_I \lambda_I \left( \sum_J p(J|I) - 1 \right) - \\
& - \sum_e \gamma_e \left( \sum_f p(f|e) - 1 \right)
\end{aligned}
$$

# Training Algorithm

- It is intractable to calculate expected counts exactly

- We resort to a Viterbi EM algorithm

1: **procedure** $\text{VITERBIEM}(F, E, \mathbf{d})$
2:     Initialize $\boldsymbol{\theta}^{(0)}$
3:     **for all** $k = 1, \ldots, K$ **do**
4:         $\hat{\mathbf{m}}^{(k)} \leftarrow \text{ALIGN}(F, E, \boldsymbol{\theta}^{(k-1)})$
5:         $\boldsymbol{\theta}^{(k)} \leftarrow \text{UPDATE}(F, E, \mathbf{d}, \hat{\mathbf{m}}^{(k)})$
6:     **end for**
7:     **return** $\hat{\mathbf{m}}^{(K)}, \boldsymbol{\theta}^{(K)}$
8: **end procedure**

# Calculating Viterbi Matching

- Viterbi matching can be calculated independently

- We use information retrieval techniques for speed-up

$$
\hat{\mathbf{m}}_t = \begin{cases} \tilde{\mathbf{m}}_t & \text{if } P(\mathbf{f}^{(t)}|\mathbf{e}^{(\tilde{\mathbf{m}}_t)};\boldsymbol{\theta}) > \epsilon \\ 0 & \text{otherwise} \end{cases}
$$

$$
\tilde{\mathbf{m}}_t = \operatorname*{argmax}_{s \in \{1,\dots,S\}} \left\{ \frac{p(J^{(t)}|I^{(s)})}{(I^{(s)}+1)^{J^{(t)}}} \prod_{j=1}^{J^{(t)}} \sum_{i=0}^{I^{(s)}} p(\mathbf{f}_j^{(t)}|\mathbf{e}_i^{(s)}) \right\}
$$

# Updating Model Parameters

- Model parameters are updated based on the Viterbi matching

| parameter | expected count |
|---|---|
| $p(J\|I)$ | $\displaystyle\sum_{t=1}^{T}(1-\delta(\hat{\mathbf{m}}_t,0))\delta(J^{(t)},J)\delta(I^{(\hat{\mathbf{m}}_t)},I)$ |
| $p(f\|e)$ | $\displaystyle\sum_{t=1}^{T}(1-\delta(\hat{\mathbf{m}}_t,0))\frac{p(f\|e)}{\sum_{i=0}^{I^{(\hat{\mathbf{m}}_t)}}p(f\|\mathbf{e}_i^{(\hat{\mathbf{m}}_t)})}$ $\displaystyle\times\sum_{j=1}^{J^{(t)}}\delta(f,\mathbf{f}_j^{(t)})\sum_{i=0}^{I^{(\hat{\mathbf{m}}_t)}}\delta(e,\mathbf{e}_i^{(\hat{\mathbf{m}}_t)})$ $+\sigma(f,e,\mathbf{d})$ |

# Matching Evaluation

- Non-parallel corpora: 20K Chinese + 40K English phrases
- seed dictionary: 1K entries

# Translation Evaluation

- Using learned parallel phrases to train translation models

| iteration | corpus | lexicon | BLEU |
|-----------|--------|---------|------|
| 0 | 7.4K | 4.1K | 7.68 |
| 1 | 43.4K | 6.6K | 11.21 |
| 2 | 102.4K | 9.3K | 12.65 |
| 3 | 135.3K | 10.5K | 12.77 |
| 4 | 148.2K | 11.0K | 13.23 |
| 5 | 153.7K | 11.3K | 13.40 |

# Examples

| Chinese | 其 主要 产品 是 食物 |
|---------|---------------------|
| English | its main products are food |
| Chinese | 美国 既定 的 战略 计划 |
| English | the set US strategic plan |
| Chinese | 全球化 的 趋势 |
| English | the trends of globalization |
| Chinese | 强化 联合 作战 功能 |
| English | strengthening joint combat functions |
| Chinese | 伊斯兰教 ， 天主教 和 基督教 |
| English | Islam , Catholicism , and Christianity |

# Conclusion and Future Work

- **Conclusion**

  - We have presented a framework for learning parallel lexicons and phrases from non-parallel corpora

  - Useful for machine translation and bilingual lexicography

- **Future Work**

  - Extend to more sophisticated alignment models

  - Scale to large-scale data

# Thanks