

Parallel Corpora-based Dictionaries and Their Applications

Summer School of Linguistics 2013
Braga, 10th September 2013

Nuno Ramos Carvalho, Rui Brito, Luís Miguel Braga,
Alberto Simões and José João Almeida

<http://natura.di.uminho.pt>



Part 0: Background Review

- 1 Corpus and Parallel Corpora

Part I: Introduction to Probabilistic Translation Dictionaries

- 1 Definition
- 2 Query
- 3 Examples
- 4 Extraction Algorithms

Part II: Resources Created from PTD

- 1 Bi-Words (BWs)
- 2 Unambiguous-Concept Translation Sets (UCTS)
- 3 PTD Algebra

Part III: Other Applications

- 1 Triangulation and Transitivity
- 2 Probabilistic Syn Sets (PSS)
- 3 Creating Other Dictionaries (OmegaT resources)

Part 0

Background Review



- collection of structured texts
- usually large and domain restricted
- stored in easy to process formats

- collection of texts in different languages, where each of them is a translation of each other
- used in NLP and Computational Linguistics
- aligned at the sentence level
- different formats
 - Text Encoding Initiative (TEI)
 - Translation Memory eXchange (TMX)
 - XML Corpus Encoding Standard (XCES)
- tools for aligning
 - Sentence Alignment: Hunalign, Vanilla Aligner, WinAlign,...
 - Word Alignment: GIZA++, NATools, ...

Estes resultados constituem a base do Programa Europeu de defesa do Mar de Barents e, por esse motivo, peço-lhe que analise um projecto de carta que lhe expõe os factos mais importantes, e que, de acordo com as decisões do Parlamento, torne clara esta posição na Rússia.

No entanto, somos também da opinião de que deveria haver um debate sobre esta estratégia da comissão que seguisse um procedimento ordenado, e não só com base numa declaração oral pronunciada aqui no Parlamento Europeu, mas também com base num documento que seja decidido na comissão e que apresente uma descrição deste programa para um período de cinco anos.

These findings form the basis of the European Programmes to protect the Barents Sea, and that is why I would ask you to examine a draft letter setting out the most important facts and to make Parliament's position, as expressed in the resolutions which it has adopted, clear as far as Russia is concerned.

We believe, however, that the commission's strategic plan needs to be debated within a proper procedural framework, not only on the basis of an oral statement here in the European Parliament, but also on the basis of a document which is adopted in the commission and which describes this programme over the five-year period.

```
<?xml version="1.0"?>
<!DOCTYPE tmx SYSTEM "tmx11.dtd">
<tmx version="version 1.0">
<header creationtool="cwb-utils" creationtoolversion="1.0"
  segtype="sentence" adminlang="EN-US" srclang="pt">
</header>

<tu><!--1:1-->
  <tuv lang='pt'><seg>Constituição da República Portuguesa</seg></tuv>
  <tuv lang='es'><seg>CONSTITUCIÓN DE LA REPÚBLICA PORTUGUESA</seg></tuv>
</tu>

<tu><!--1:1-->
  <tuv lang='pt'><seg>IV REVISÃO CONSTITUCIONAL</seg></tuv>
  <tuv lang='es'><seg>Cuarta Revisión 1997</seg></tuv>
</tu>

<tu><!--1:1-->
  <tuv lang='pt'><seg>PREÂMBULO</seg></tuv>
  <tuv lang='es'><seg>PREÂMBULO</seg></tuv>
</tu>

<tu><!--1:1-->
  <tuv lang='pt'><seg>A 25 de Abril de 1974 , o Movimento das Forças Armadas ,
  <tuv lang='es'><seg>El 25 de Abril de 1974 , el Movimiento de las Fuerzas Ar
```

Part I

Introduction to Probabilistic Translation Dictionaries



- translation dictionaries
- probabilistic translations
- automatically extracted
- domain constrained
 - from parallel corpora, sentence aligned
- not a conventional dictionary
 - ... but can be transformed
- actually, a pair of dictionaries
 - example: $TMX_{pt \rightarrow en} \longrightarrow PTD_{pt \rightarrow en} \times PTD_{en \rightarrow pt}$
 - PTD from parallel corpora source to target language
 - PTD^{-1} (*inverse PTD*) from parallel corpora target to source language



```

europe => {      ocorr => 42853,
                trans => {      europa      => 0.9471,
                                europeus   => 0.0339,
                                europeu    => 0.0081,
                                europeia   => 0.0011,
                                },
                },
stupid => {      ocorr => 180,
                trans => {      estúpido  => 0.1755,
                                estúpida   => 0.1099,
                                estúpidos  => 0.0741,
                                avisada    => 0.0565,
                                direita    => 0.0558,
                                impasse    => 0.0448,
                                },
                },

```

CQuery

per-fide.di.uminho.pt/query/ptd/Vatican/PT/EN/mesa

Per-Fide

PTDC/CLE-LLI/108948/2008

Select Type
bilingual

Select language
PT-EN

Select corpora

- Comboni (info)
- Culinaria (info)
- DGT-Acquis (info)
- DGT-TM (info)
- ECB (info)
- ects (info)
- EMEA (info)
- EurLex-v1 (info)
- EuroParlV5 (info)
- JRC-Acquis-V3 (info)
- PressEU (info)
- Shakespeare (info)
- SoftwarePO-2 (info)
- thesis-abstract (info)
- Vatican (info)

PT query
mesa

EN query

Search

Entries per page: 20

PTD for : PT → EN

mesa (192 occurrences)

74.43%	table ✓	175
5.46%	banquet ✓	117
1.26%	altar ✓	361
1.09%	return	553
1.00%	both	2491
0.66%	start ✓	150
0.63%	a	58102
0.42%	presidency ✓	29
0.38%	liturgy	1745
0.17%	inviting ✓	86

QUERY> europa

Occurrences: 39917

Translations:

88.50%	europa
5.73%	european
2.37%	europa
1.16%	(none)
0.57%	eu
0.23%	unece
0.17%	the
0.16%	auto

QUERY> we

Occurrences: 300431

Translations:

17.81% (none)

8.25% que

6.02% temos

QUERY> read

Occurrences: 2435

Translations:

29.32% ler
 13.75% li
 8.36% read
 5.96% lido
 3.54% lemos
 1.60% leio
 1.46% estar
 1.45% leu

QUERY> represent

Occurrences: 2538

Translations:

17.87% representam
 11.57% representar
 8.93% represento
 7.54% representamos
 4.93% constituem
 3.63% representa
 3.37% (none)
 2.35% representante

QUERY> aceitável

Occurrences: 1713

Translations:

71.48% acceptable

8.56% unacceptable

$createPTD :: TU^* \longrightarrow PTD$

- 1 pre-processing stage
- 2 statistical algorithm stage
 - 1 build co-occurrence count table
 - 2 search highest values in matrix, to define *correct* relations
 - 3 when no highest values are found, use all others attenuated

- a flor cresce / a casa é grande / a casa azul tem flores
- the flower grows / the house is big / the blue house has flowers

Co-Occurrence Table

	a	flor	cresce	casa	é	grande	azul	tem	flores
the	3	1	1	2	1	1	1	1	1
flower	1	1	1	0	0	0	0	0	0
grows	1	1	1	0	0	0	0	0	0
house	2	0	0	2	1	1	1	1	1
is	1	0	0	1	1	1	0	0	0
big	1	0	0	1	1	1	0	0	0
blue	1	0	0	1	0	0	1	1	1
has	1	0	0	1	0	0	1	1	1
flowers	1	0	0	1	0	0	1	1	1

Currently working on:

- produce PTD from lemmas
- POS tagging

Part II

Resources Created from PTD



Bi-Words (BWs)

- strongly related word pair lists
- words may be translations
- example:
 - {informed}_{en} = {informados}_{pt}
 - {approved}_{en} = {aprovado}_{pt}
 - {modern}_{en} = {modernos}_{pt}

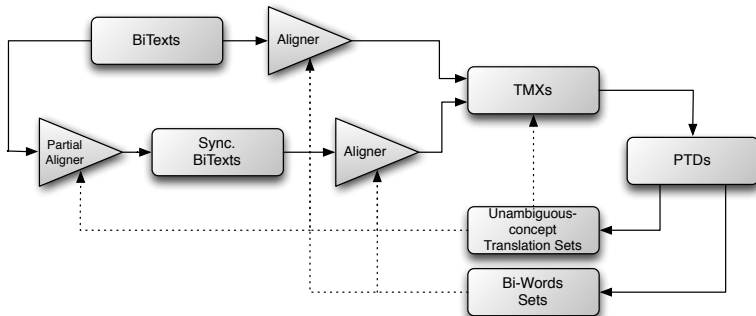


Unambiguous-Concept Translation Sets (UCTS)

- sets of equivalent words for two languages
- words that have always (or almost) the same translation
- stronger relations than BWs (smaller sets)
- examples:
 - $\{\text{London}\}_{en} = \{\text{Londres}\}_{pt}$
 - $\{\text{wolfram, tungsten}\}_{en} = \{\text{volfrâmio, tungsténio}\}_{pt}$
- useful for bootstrapping resources



- extracted automatically from PTD
- more important when few resources are available
- provide connection clues
- used for:
 - partial synchronization/alignment
 - alignment assessment
 - translation quality assurance



- PTD can be automatically extracted
- PTD formal definition (not discussed here)
- other resources can be created
- set of operations that described how these resources are created
- way to clearly share workflows



Creating a *PTD* composed only of verbs (PTD^V) can be defined using the filter function as:

$$PTD^V = filter(PTD, verb)$$

where, *verb* is a function defined as:

$$verb(entry) = \begin{cases} True & \text{if } word \text{ in } entry \text{ is a verb} \\ False & \text{otherwise} \end{cases}$$

Function	Syntax	Domain	Range
domain	$\text{dom}(-)$	$\text{ptd}_{A \rightarrow B}$	w_A^*
range	$\text{ran}(-)$	$\text{ptd}_{A \rightarrow B}$	w_B^*
union	$- \cup -$	$\text{ptd}_{A \rightarrow B} \times \text{ptd}_{A \rightarrow B}$	$\text{ptd}_{A \rightarrow B}$
intersection	$- \cap -$	$\text{ptd}_{A \rightarrow B} \times \text{ptd}_{A \rightarrow B}$	$\text{ptd}_{A \rightarrow B}$
composition	$- \circ -$	$\text{ptd}_{A \rightarrow B} \times \text{ptd}_{B \rightarrow C}$	$\text{ptd}_{A \rightarrow C}$
domain restrict	$- / -$	$\text{ptd}_{A \rightarrow B} \times w_A^*$	$\text{ptd}_{A \rightarrow B}$
domain subtract	$- \setminus -$	$\text{ptd}_{A \rightarrow B} \times w_A^*$	$\text{ptd}_{A \rightarrow B}$
totalization	$\text{tot}(-)$	$\text{ptd}_{A \rightarrow B}$	$\text{ptd}_{A \rightarrow B}$
filter	$\text{filter}(-, -)$	$\text{ptd}_{A \rightarrow B} \times (\text{entry}_{A \rightarrow B} \rightarrow \text{Bool})$	$\text{ptd}_{A \rightarrow B}$
map	$\text{map}(-, -)$	$\text{ptd}_{A \rightarrow B} \times (\text{entry}_{A \rightarrow B} \rightarrow \text{entry}_{A \rightarrow B})$	$\text{ptd}_{A \rightarrow B}$

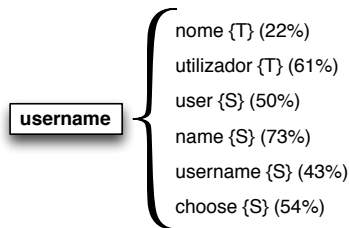
Part III

Other Applications

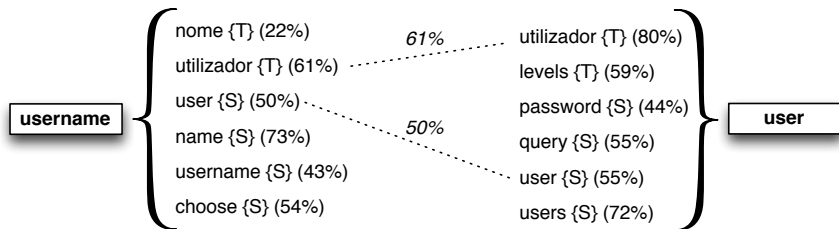
- some languages have a limited amount of available resource
- less available man power to create them
- come up with strategies to create them more efficiently

- start with: $PTD_{A \rightarrow B}$ and $PTD_{B \rightarrow C}$
- apply a composition function so that:
$$PTD_{A \rightarrow C} = PTD_{B \rightarrow C} \circ PTD_{A \rightarrow B}$$
- addressing number of occurrences, which words appear, and how probabilities are mapped
- result is: $PTD_{A \rightarrow C}$

- Syn Sets are required by many algorithms and techniques
- to discover words with the same meaning (or close)
- compute Probabilistic Syn Sets (PSS) from PTD
 - multi-lang Syn Sets, with weights, and domain constrained
- to compute $PSS(term)$ include:
 - all the possible translations and probabilities for $PTD(term)$
 - and, for each of these possible translations the inverse PTD translations and corresponding probabilities
- example:



- synonym search
- word similarity
- concepts similarity





Export *PTD* to use in OmegaT, as glossary

```
$ nat-ptd toTSV my_ptd.ptd glossary.tsv
```

Export *PTD* to use in OmegaT, as dictionary

```
$ nat-ptd toStarDict my_ptd.ptd star_dict
```




Software

- **Lingua::NATools**
`http://search.cpan.org/dist/Lingua-NATools/`
- **Lingua::PTD**
`http://search.cpan.org/dist/Lingua-PTD/`

Online

- `http://per-fide.di.uminho.pt/query`
- `http://ptd.natura.di.uminho.pt/`



- Structural Alignment of Plain Text Books
in LREC2012
- Defining a Probabilistic Translation Dictionaries Algebra
in EPIA2013 (*forthcoming*)
- Translation Dictionaries Triangulation
in Iberian SLTech Workshop 2012
- XML Schemas for Parallel Corpora
in XATA2011
- Extracção de Recursos de Tradução com base em Dicionários
Probabilísticos de Tradução
PhD thesis

Thank You!

Nuno Carvalho <narcarvalho@di.uminho.pt>