

PTD Algebra

Alberto Simões <ambs@ilch.uminho.pt>
José João Almeida <jj@di.uminho.pt>

November 23rd, 2011



- PTD stands for Probabilistic Translation Dictionary;
- Can be seen as a common translation dictionary, but . . .
 - . . . is extracted automatically from parallel corpora;
 - . . . for each word suggests a set of possible translations;
 - . . . together with each translation presents a certainty level;
 - . . . does not present just translations, but also “related words;”
- Can be extracted by some different tools:
 - NATools (?);
 - GIZA++ (?);



europe =>	stupid =>
count => 42853,	count => 180,
trans =>	trans =>
europa => 94.7 %	estúpido => 47.6 %
europeus => 3.4 %	estúpida => 11.0 %
europeu => 0.8 %	estúpidos => 7.4 %
europeia => 0.1 %	avisada => 5.6 %
	direita => 5.6 %
	impasse => 4.5 %
	ocupado => 3.8 %



Extraction process

$$\text{align} : TU^* \longrightarrow \text{Dic}(\mathcal{L}_\alpha, \mathcal{L}_\beta) \times \text{Dic}(\mathcal{L}_\beta, \mathcal{L}_\alpha)$$

$$TU \quad \cong \text{sentence}_{\mathcal{L}_\alpha} \times \text{sentence}_{\mathcal{L}_\beta}$$

$$\text{Dic}(A, B) \quad \cong \text{word}_A \rightharpoonup \text{wordInfo}(B)$$

$$\text{wordInfo}(B) \quad \cong \text{occur} \times \text{word}_B \rightharpoonup \text{probabilidade}$$

Dictionary structure

$$\text{word}_{\mathcal{L}_\alpha} \mapsto (\text{occurrences} \times \text{word}_{\mathcal{L}_\beta} \mapsto \mathcal{P}(\mathcal{T}(\text{word}_{\mathcal{L}_\alpha}) = \text{word}_{\mathcal{L}_\beta}))$$



Per-Fide

PTDC/CLE-LLI/108948/2008

Select Type

Select language

Select corpora

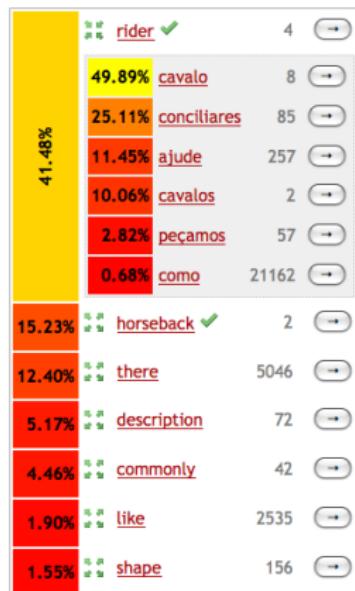
- DGT-TM ([info](#))
- EuroParl ([Info](#))
- JRC-Acquis ([Info](#))
- Vatican v1 ([Info](#))
- EurLex v1 ([Info](#))
- ECB v1 ([Info](#))
- EMEA 0.3 ([Info](#))
- Comboni ([Info](#))
- zenity ([Info](#))

PT query

EN query

PTD for Vatican v1: PT → EN

cavalo (8 occurrences)





There is a list of operations one can perform with PTDs:

Description	Notation
union	$d1 \cup d2$
interception	$d1 \cap d2$
domain restrict	$d1 d2$
domain subtract	$d1 \setminus d2$
PTD composition	$d1 \circ d2$
PTD distance	$distance(d1, d2)$
PTD totalize	$totalize(d1)$
PTD filtering	$filter(d, entry \rightarrow bool)$



- PTD Union **sums** two or more dictionaries;
- The base of NATools algorithm for scalability (?);
- This operation:
 - is applied to dictionaries with the same source and target languages;
 - gives different weights to probabilities accordingly to corpora sizes and word occurrence counts;
 - is applied to multi-sets. This means that

$$D \cup D = 2 \times D,$$

where the multiplication is the duplication of occurrences counts (translations probabilities are kept).

$$\underbrace{\left\{ \begin{array}{ll} \text{patient} & 80.4 \% \\ \text{patients} & 4.5 \% \\ \text{card} & 0.6 \% \\ \text{ill} & 0.6 \% \\ \text{you} & 0.4 \% \\ \text{she} & 0.3 \% \end{array} \right.}_{\text{doente} (10964)} \cup \underbrace{\left\{ \begin{array}{ll} \text{patient} & 24.2 \% \\ \text{ill} & 18.7 \% \\ \text{sick} & 17.7 \% \\ \text{patients} & 10.1 \% \\ \text{illness} & 2.0 \% \\ \text{well} & 1.7 \% \end{array} \right.}_{\text{doente} (404)} =$$

$$\underbrace{\left\{ \begin{array}{ll} \text{patient} & 79.8 \% \\ \text{patients} & 4.6 \% \\ \text{ill} & 0.8 \% \\ \text{card} & 0.6 \% \\ \text{you} & 0.4 \% \\ \text{shee} & 0.3 \% \end{array} \right.}_{\text{doente} (11368)}$$



Intersection is a dictionary enhancing mechanism:

- intersects the domain of both dictionaries (removing any word that does not occur in one of the dictionaries);
- uses, for each maintained word, the minimum number of occurrences in the two dictionaries;
- intersects the probable translation sets, maintaining words that are probable translations on both dictionaries;
- associates for each translation a probability that is the minimum in the two dictionaries.

This process:

- makes a stronger dictionary;
- reduces probability values;
- resulting values lose significance as a probability;



Intersection can be useful for:

- the intersection of dictionaries obtained from different domain corpora, to compute the shared or **base** lexicon;
- the intersection with a specific small hand-controlled dictionary can be used to tune PTD extraction algorithms;



$$\underbrace{\left\{ \begin{array}{ll} \text{patient} & 80.4 \% \\ \text{patients} & 4.5 \% \\ \text{card} & 0.6 \% \\ \text{ill} & 0.6 \% \\ \text{you} & 0.4 \% \\ \text{she} & 0.3 \% \end{array} \right.}_{\text{doente} (10964)} \cap \underbrace{\left\{ \begin{array}{ll} \text{patient} & 24.2 \% \\ \text{ill} & 18.7 \% \\ \text{sick} & 17.7 \% \\ \text{patients} & 10.1 \% \\ \text{illness} & 2.0 \% \\ \text{well} & 1.7 \% \end{array} \right.}_{\text{doente} (404)} =$$

$$\underbrace{\left\{ \begin{array}{ll} \text{patient} & 80.4 \% \\ \text{patients} & 4.5 \% \\ \text{ill} & 0.6 \% \end{array} \right.}_{\text{doente} (404)}$$



Domain restriction,

- restricts the PTD domain to a specific set of words;
- can be used to restrict a dictionary to some specific words;

Domain subtraction,

- removes from the PTD domain a specific set of words;
- can be used to remove from a dictionary some common words;
- ... and that can be used to detect terminology;

Recalculate probabilities:

- Most of the operations described lose translations, or change translation probabilities;
- For example, when removing some words from the list of possible translations, the sum of translation probabilities is no longer 100 %;
- This operation recalculate all probabilities in a way they sum up 100 % again.

$$\text{totalize } \underbrace{\left(\begin{array}{l} \overbrace{\text{doente}}^{(404)} \\ \text{patients} \\ \text{ill} \end{array} \right) \left\{ \begin{array}{ll} \text{patient} & 80.4 \% \\ \text{patients} & 4.5 \% \\ \text{ill} & 0.6 \% \end{array} \right.} = \underbrace{\left(\begin{array}{l} \overbrace{\text{doente}}^{(404)} \\ \text{patients} \\ \text{ill} \end{array} \right) \left\{ \begin{array}{ll} \text{patient} & 94.0 \% \\ \text{patients} & 5.3 \% \\ \text{ill} & 0.7 \% \end{array} \right.}$$



PTD composition is the most interesting operation:

- take two dictionaries, one $A \rightarrow B$ and another $B \rightarrow C$;
- compose the two dictionaries, obtaining a $A \rightarrow C$ dictionary;
- can be useful to calculate translation dictionaries for one pair of languages that does not have a parallel corpora available (?);

It can be used in other situations, like:

- take one dictionary, $A \rightarrow B$;
- compute $A \rightarrow A$ (compose with inverse dictionary);
- can be useful to calculate related-word sets.



PTD Composition: Triangulation Example

afluencia	influx	18.6 %	<table><tr><td>afflusso</td><td>48.9 %</td></tr><tr><td>flusso</td><td>12.7 %</td></tr><tr><td>flussi</td><td>4.7 %</td></tr></table>	afflusso	48.9 %	flusso	12.7 %	flussi	4.7 %
afflusso	48.9 %								
flusso	12.7 %								
flussi	4.7 %								
flow	12.9 %	<table><tr><td>flusso</td><td>46.9 %</td></tr><tr><td>flussi</td><td>9.9 %</td></tr><tr><td>gravi</td><td>1.7 %</td></tr></table>	flusso	46.9 %	flussi	9.9 %	gravi	1.7 %	
flusso	46.9 %								
flussi	9.9 %								
gravi	1.7 %								
inflow	6.1 %	<table><tr><td>sfogo</td><td>24.2 %</td></tr><tr><td>afflusso</td><td>16.8 %</td></tr><tr><td>ascritto</td><td>14.7 %</td></tr></table>	sfogo	24.2 %	afflusso	16.8 %	ascritto	14.7 %	
sfogo	24.2 %								
afflusso	16.8 %								
ascritto	14.7 %								
flood	5.9 %	<table><tr><td>inondazioni</td><td>5.6 %</td></tr><tr><td>flusso</td><td>4.4 %</td></tr><tr><td>alluvione</td><td>2.8 %</td></tr></table>	inondazioni	5.6 %	flusso	4.4 %	alluvione	2.8 %	
inondazioni	5.6 %								
flusso	4.4 %								
alluvione	2.8 %								
flows	4.7 %	<table><tr><td>flussi</td><td>72.3 %</td></tr><tr><td>flusso</td><td>1.6 %</td></tr><tr><td>onde</td><td>1.5 %</td></tr></table>	flussi	72.3 %	flusso	1.6 %	onde	1.5 %	
flussi	72.3 %								
flusso	1.6 %								
onde	1.5 %								



afluencia	afflusso	10.08 %
	flusso	8.73 %
	flussi	5.51 %
	sfogo	1.46 %
	ascritto	0.89 %
	inondazioni	0.33 %
	gravi	0.22 %
	alluvione	0.16 %
	onde	0.07 %



Composing with the Inverse Dictionary

casa	house	48 %	{ casa house casiña vivenda }	76 % 5 % 2 % 2 %	= 36.5 % = 2.4 % = 1.0 % = 1.0 %
	home	35 %	{ casa fogar país domicilio terra }	52 % 11 % 2 % 2 % 1 %	= 18.2 % = 3.9 % = 0.7 % = 0.7 % = 0.4 %
	cottage	2 %	{ casa cabana }	79 % 9 %	= 1.60 % = 0.20 %
	houses	1 %	{ casa casopa moradas }	80 % 5 % 2 %	= 0.80 % = 0.05 % = 0.02 %
	homes	1 %	{ casa asilos domicilios }	19 % 5 % 3 %	= 0.20 % = 0.05 % = 0.03 %





Result of inverse composition

casa	57.3 %
fogar	3.9 %
house	2.4 %
casiña	1.0 %
vivenda	1.0 %
país	0.7 %
domicilio	0.7 %
terra	0.4 %
cabana	0.2 %
casopa	0.05 %
asilos	0.05 %
domicilios	0.03 %
moradas	0.02 %





- big corpora PTDs tend to be big;
- PTDs tend to be used for word translations:
 - so we can remove non-words, like numbers and punctuation;
- PTDs have a probability value associated with translations:
 - so we can use it to reduce the PTD size, removing non probable translations;
- PTDs have a number of occurrences:
 - so we can use it to remove infrequent words...

$$\text{filter} \left(\underbrace{\begin{array}{ll} \text{doente} & \left\{ \begin{array}{ll} \text{patient} & 80.4 \% \\ \text{patients} & 4.5 \% \\ \text{card} & 0.6 \% \\ \text{ill} & 0.6 \% \\ \text{you} & 0.4 \% \\ \text{she} & 0.3 \% \end{array} \right. \\ (10964) \end{array}}_{(10964)} \right) = \underbrace{\begin{array}{ll} \text{doente} & \left\{ \begin{array}{ll} \text{patient} & 80.4 \% \\ \text{patients} & 4.5 \% \end{array} \right. \\ (10964) \end{array}}_{(10964)}$$

- PTDs can be used for terminology extraction (?; ?);
- The process can be quite good if used together with a morphological analyzer;
- It uses translation patterns:

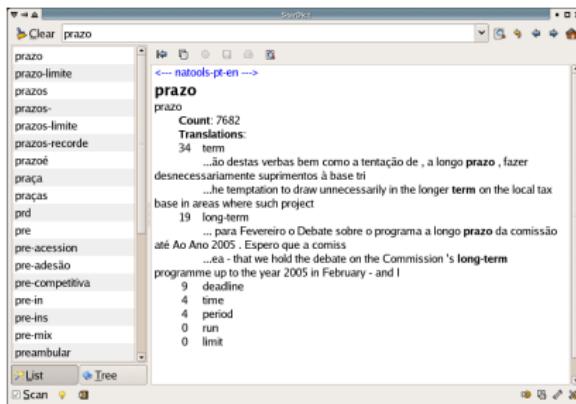
$$\mathcal{T}(A \cdot \text{"de"} \cdot B) = \mathcal{T}(B) \cdot \mathcal{T}(A)$$

	Human	Rights
Direitos		X
do		
Homem	X	



PTD can also be used directly by end-users.

- We've shown how to query them using Per-Fide corpora interface;
- We can also use a command-line tool: **nat-ptd query**;
- We can create dictionaries to be queried offline:





Bibliography

