

Termonet: Construcción de terminologías a partir de WordNet y corpus especializados

Termonet: Terminology construction from WordNet and technical corpora

Miguel Anxo Solla Portela
Universidade de Vigo
Grupo TALG
miguelsolla@uvigo.es

Xavier Gómez Guinovart
Universidade de Vigo
Grupo TALG
xgg@uvigo.es

Resumen: En esta presentación, mostraremos la metodología y los recursos utilizados en el desarrollo de Termonet, una herramienta para la consulta y verificación en corpus de los léxicos de especialidad incluidos en WordNet. Termonet realiza una identificación en WordNet de los synsets pertenecientes a un ámbito terminológico a partir de las relaciones léxico-semánticas establecidas entre los synsets, y valida los términos identificándolos en un corpus especializado desambiguado semánticamente. La construcción de esta herramienta forma parte de las tareas del proyecto de investigación SKATeR-UVigo, orientado al desarrollo y aplicación de recursos para el procesamiento lingüístico del gallego.

Palabras clave: WordNet, lexicografía computacional, terminología computacional

Abstract: In this presentation, we review the methodology and the resources used in the development of Termonet, a tool for checking and verifying in a corpus the specialty lexicons embedded in WordNet. This tool performs an identification of the synsets in WordNet belonging to a terminological domain from the lexical-semantic relations established among synsets, and validates the terms identifying them by means of a semantically disambiguated specialized corpus. The construction of this tool is part of the tasks of the SKATeR-UVigo research project, aimed at the development and application of resources for Galician language processing.

Keywords: WordNet, computational lexicography, computational terminology

1 Introducción

En este artículo¹ se describen la metodología y los recursos utilizados en el desarrollo de Termonet², una herramienta para la consulta de los léxicos de especialidad incluidos en WordNet³ y para su verificación en corpus. La construcción de esta herramienta forma parte de los objetivos del proyecto de investigación SKATeR-UVigo, orientado al desarrollo y aplicación de recursos para el procesamiento lingüístico del gallego.

¹Esta investigación se realiza en el marco del proyecto *Adquisición de escenarios de conocimiento a través de la lectura de textos: Desarrollo y aplicación de recursos para el procesamiento lingüístico del gallego (SKATeR-UVigo)* financiado por el Ministerio de Economía y Competitividad, TIN2012-38584-C06-04.

²<http://sli.uvigo.es/termonet/termonet.php>

³<http://wordnet.princeton.edu>

Termonet se centra en la explotación de WordNet para la construcción de terminologías, mediante la exploración de las relaciones semánticas codificadas entre los nodos conceptuales (o synsets) de la ontología léxica. Como se explica con detalle más adelante, el funcionamiento de la aplicación se basa en que los términos propios de un ámbito terminológico incluidos en WordNet se localizan en synsets relacionados con un nodo raíz mediante ciertas configuraciones de relaciones semánticas y a determinadas distancias máximas de este nodo.

Termonet ofrece la posibilidad de explorar los distintos conjuntos de synsets asociados a un synset de origen en función de las configuraciones definidas por el usuario para la selección de relaciones exploradas y para el nivel máximo de exploración de cada re-

lación. La misma aplicación permite verificar los resultados de la exploración en un corpus de textos especializados.

2 Recursos

Las funcionalidades de Termonet se fundamentan en dos recursos básicos: un léxico WordNet y un corpus textual lematizado y desambiguado con respecto a los sentidos de WordNet. En la implementación actual de Termonet, diseñada para su aplicación en tareas terminológicas relacionadas con la ampliación del WordNet del gallego en el ámbito de la medicina, estos dos recursos son el léxico Galnet y el *Corpus Técnico do Galego*.

Galnet, la versión gallega de WordNet, se distribuye como parte del MCR (González Agirre, Laparra, y Rigau, 2012). Esta versión de Galnet, de 2012, incluye los Basic Level Concepts⁴, los ficheros lexicográficos de partes del cuerpo y de substancias, y la traducción parcial de los adjetivos. Además, contiene una primera ampliación realizada con el WN-Toolkit⁵ a partir de la Wikipedia⁶ y el *Diccionario CLUVI inglés-galego*⁷.

A partir de esta versión inicial, se ha seguido ampliando Galnet con el WN-Toolkit a partir de los diccionarios de Apertium⁸, Babelnet⁹ 2.0, Wiktionary¹⁰, Wikipedia, Geonames¹¹, Wikispecies¹² y los corpus SemCor inglés-gallego y CLUVI (Gómez Guinovart y O., 2014). También se ha realizado una expansión a partir del *Diccionario de sinónimos do galego*¹³ (Gómez Guinovart y Solla Portela, 2014). Finalmente, se han efectuado ampliaciones en el ámbito de la fraseología (locuciones verbales) y de la terminología (medicina y economía). Todas estas expansiones se pueden consultar en la interfaz web de Galnet¹⁴ utilizando la versión de desarrollo del recurso.

La implementación actual de Termonet usa la versión de desarrollo de Galnet 3.0.10 (2015), cuya extensión en número de synsets

(Syns) y variantes léxicas (Vars) se recoge en la Tabla 1 en comparación con la de la versión de 2002 distribuida con el MCR.

	MCR		3.0.10	
	Vars	Syns	Vars	Syns
N	18949	14285	27825	20621
V	1416	612	4199	1564
Adj	6773	4415	8086	5104
Adv	0	0	471	370
Total	27138	19312	40581	27659

Tabla 1: Extensión léxica de Galnet

Por su parte, el *Corpus Técnico do Galego* (CTG)¹⁵ es un corpus de orientación terminológica de 15 millones de palabras, formado por textos especializados del gallego contemporáneo en los ámbitos del derecho, informática, economía, ciencias ambientales, ciencias sociales y medicina. La sección del corpus de medicina del CTG (el subcorpus *Medigal*) utilizada en la implementación actual de Termonet totaliza 3.823.232 palabras. Para esta aplicación, se ha utilizado una versión del Medigal etiquetada mediante FreeLing¹⁶ y UKB (Agirre y Soroa, 2009), empleando Galnet 3.0.10 como léxico para la desambiguación semántica del corpus.

3 Funcionalidades

3.1 Construcción de terminologías

La función principal de Termonet consiste en facilitar la extracción de variantes de WordNet relacionadas con un ámbito de especialidad. Con este fin, Termonet ofrece un formulario de consulta que permite elegir un synset de la ontología léxica y, a partir de él, realizar una extracción de los términos relacionados en función de la configuración de relaciones semánticas que se seleccionen. Aunque Termonet permite realizar la extracción desde cualquier synset de la ontología, dada su orientación terminológica, la aplicación trata de sugerir siempre las variantes nominales más próximas cuando se propone un synset no nominal.

Como se ilustra en la parte superior de la Figura 1, Termonet permite indicar el synset de origen que definirá el ámbito de la extracción terminológica, y seleccionar el conjunto de relaciones semánticas que se utilizarán para la identificación de los términos de ese

⁴<http://adimen.si.ehu.es/web/BLC/>

⁵<http://sourceforge.net/projects/wn-toolkit/>

⁶<http://www.wikipedia.org>

⁷<http://sli.uvigo.es/diccionario/>

⁸<http://www.apertium.org>

⁹<http://www.babelnet.org>

¹⁰<http://www.wiktionary.org>

¹¹<http://www.geonames.org>

¹²<http://species.wikimedia.org>

¹³<http://sli.uvigo.es/sinonimos/>

¹⁴<http://sli.uvigo.es/galnet/>

¹⁵<http://sli.uvigo.es/CTG/>

¹⁶<http://nlp.lsi.upc.edu/freeling/>

ILL: ili-30-06045562-n indicar repeticiones

Filtro por distancia (nivel máximo de exploración de cada relación):

Synonyms 4 · Antonyms 4 · Hyperonyms Hyponyms

Holonyms Meronyms Related Verbs Domain

Glosses

has_hyperonym 1 · has_xpos_hyperonym 1 · has_hyponym 4 ·
has_xpos_hyponym 4 · has_holo_madeof 1 · has_holo_member 1 ·
has_holo_part 1 · has_mero_madeof 4 · has_mero_member 4 ·
has_mero_part 4 · has_derived 3 · has_pertainym 3 ·
is_derived_from 3 · pertains_to 3 · related_to 3 · see_also_wn15 3 ·
causes 3 · has_subevent 3 · is_caused_by 3 · is_subevent_of 3 ·
verb_group 3 · category 1 · category_term 4 · region 1 ·
region_term 1 · usage 1 · usage_term 1 · gloss 0 · rgloss 0

Filtro por relaciones (impide a exploración derivada das relacións seleccionadas):

Synonyms Antonyms Hyperonyms Hyponyms

Holonyms Meronyms Related Verbs Domain

Glosses

has_hyperonym has_xpos_hyperonym has_holo_madeof
has_holo_member has_holo_part has_derived has_pertainym
is_derived_from pertains_to related_to see_also_wn15
causes has_subevent is_caused_by is_subevent_of
verb_group category category_term region region_term
usage usage_term gloss rgloss

Vai -->

Figura 1: Consulta en Termonet.

ámbito, así como la distancia o nivel de profundidad hasta donde se desea desplegar cada tipo de relación. El concepto de distancia se refiere aquí al número de relaciones léxico-semánticas que unen dos synsets entre sí. De este modo, Termonet desplegará el árbol de relaciones desde el synset de origen a través de esa relación hasta alcanzar el nivel de profundidad determinado. Véase en la Figura 2, por ejemplo, la relación de hiponimia desplegada hasta el nivel 4 de profundidad en la terminología del ámbito de la medicina, construida a partir del synset *medical science* con los parámetros ilustrados en la Figura 1.

La aplicación cuenta también con un subformulario (parte inferior de la Figura 1) que permite restringir la extracción terminológica impidiendo la exploración derivada de las relaciones semánticas seleccionadas. Mediante este filtro, se trata de limitar la *toxicidad* de ciertas relaciones semánticas para la selección de los términos de un ámbito de especialidad, es decir, de reducir el impacto de las relaciones que introducen synsets que se desvían del campo conceptual. Según este criterio, la hiponimia, por ejemplo, se suele considerar una relación *tóxica*, ya que amplía la cobertura semántica inicial y tiende a introducir términos de campos conceptuales más amplios que los de partida.

Aunque la herramienta de extracción ter-

```
[0] 06045562-n medical_science | ***** { [2] biologist }
[+1] 1 06045562-n Hyperonyms (has_hyperonym) 06037298-n bioscience,
life_science | ***** { [1] biologist }
[+1] 2 06045562-n Hyponyms (has_hyponym) 06043075-n
medical_specialty, medicine | especialidade_médica (bootstrap), medicina
(bootstrap) { [0] medical_specialty }
[+2] 1 06043075-n Hyponyms (has_hyponym) 06046245-n allergology |
***** { [1] medical_specialty }
[+2] 2 06043075-n Hyponyms (has_hyponym) 06046383-n anesthesiology |
***** { [1] medical_specialty }
[+3] 1 06046383-n Related (related_to) 09793495-n anaesthetist,
anesthesiologist, anesthetist | anestesiista (wn6dic_02) { [1] medical_specialist }
[+2] 3 06043075-n Hyponyms (has_hyponym) 06046528-n angiology |
***** { [1] medical_specialty }
[+3] 1 06046528-n Related (related_to) 09793830-n angiologist |
***** { [1] doc }
[+2] 4 06043075-n Hyponyms (has_hyponym) 06046692-n bacteriology |
***** { [1] medical_specialty }
[+3] 1 06046692-n Related (has_pertainym) 02914740-a bacteriologic,
bacteriological | ***** { [2] medical_specialty }
[+3] 2 06046692-n Related (related_to) 02914740-a bacteriologic,
bacteriological | ***** { [2] medical_specialty }
[+3] 3 06046692-n Related (related_to) 09831411-n bacteriologist |
***** { [1] biologist }
[+3] 4 06046692-n Domain (category_term) 14899328-n
culture_medium, medium | medio (bootstrap), medio_do_cultivo
(bootstrap) { [2] substance [2] medical_specialty }
[+4] 1 14899328-n Hyponyms (has_hyponym) 14900184-n agar,
nutrient_agar | ágar-ágar (bootstrap), ágar_nutritivo (bootstrap),
placa_de_ágar-ágar (bootstrap) { [3] substance [3] medical_specialty }
[+4] 2 14899328-n Hyponyms (has_hyponym) 80000645-n
nutrient_broth | ***** { [3] substance [3] medical_specialty }
[+2] 5 06043075-n Hyponyms (has_hyponym) 06046898-n biomedicine |
***** { [1] medical_specialty }
```

Figura 2: Extracción de terminología.

minológica se encuentra aún en fase de desarrollo, en los experimentos se obtuvieron, con configuraciones muy simples de los parámetros de extracción, conjuntos de resultados con una congruencia mayor y cuantitativamente más significativos que la selección de variantes ligadas a un dominio de WordNet Domains¹⁷. Además, la extracción puede partir de cualquier synset y no está limitada a un dominio preestablecido, de modo que el procedimiento es idéntico para ámbitos conceptuales amplios, como la biología, y para campos más concisos, como la microbiología.

3.2 Verificación en corpus

Como ya se ha mencionado anteriormente, Termonet permite verificar los resultados de la extracción en un corpus textual lematizado y desambiguado con respecto a los sentidos de WordNet. En su implementación actual, permite contrastar los términos gallegos identificados en el corpus de medicina Medigal etiquetado con FreeLing y UKB.

El corpus desambiguado facilita el desarrollo de estrategias de verificación con base semántica para las variantes monoléxicas procedentes de Galnet, pero no para las pluriléxicas, que no cuentan con etiquetación semánti-

¹⁷<http://wndomains.fbk.eu>

Empregouse o corpus MEDIGAL
Termos monoléxicos (441 de 594, 74.24 %):
- Variantes galegas que coinciden cun lema con etiquetación semántica [ili_p]: 354 de 441 (80.27 %)
- Variantes galegas que coinciden cun lema coa etiquetación semántica con maior probabilidade [sense_p]: 338 de 441 (76.64 %)
- Promedio da frecuencia de variantes no corpus (valor máximo 1 para as variantes que se repiten 100 ou máis veces) [ili_f]: 0.4607 (46.07 %)
- Proporción das veces nas que o offset dunha variante está etiquetado como o de maior probabilidade polo UKB [sense_f]. Promedio de todos os valores sense_f: 0.8509 (85.09 %)
Termos pluriléxicos (153 de 594, 25.76 %):
- Variantes galegas pluriléxicas que coinciden con lemas sucesivos do corpus: 43 de 153 (28.1 %)
Ver as variantes pormenorizadamente

Figura 3: Verificación en corpus.

ca debido a las características de la lematización del corpus con FreeLing. Con el fin de comprobar de algún modo su presencia en el corpus, Termonet identifica las palabras léxicas de la variante en lemas sucesivos del corpus y calcula su frecuencia.

Termonet evalúa la presencia de cada término monoléxico en el corpus en base a cuatro criterios cuantificados de 0 a 1, y finalmente combina los resultados obtenidos por todos ellos en un índice general para cada criterio. Los criterios aplicados son:

1. La variante está presente (1) o no (0) como lema del corpus y con la etiqueta semántica del synset correspondiente.
2. La variante está presente como lema del corpus y con la etiqueta semántica más probable (1) o no (0) según UKB.
3. Frecuencia absoluta de la variante en el corpus, ponderando el valor máximo (1) para las variantes etiquetadas semánticamente que se repiten 100 veces o más, y el valor mínimo (0) para las variantes que no están presentes en el corpus.
4. Frecuencia con la que UKB le atribuye la mayor probabilidad a la etiqueta del synset de la variante, asignando el valor máximo (1) para la totalidad de las veces y el mínimo (0) para ninguna.

En la Figura 3 se muestran los índices globales obtenidos por la terminología construida a partir del synset *medical science* con los parámetros ilustrados en la Figura 1. A partir del análisis pormenorizado de las variantes (Figura 4), Termonet ofrece la posibilidad de comprobar sus contextos de uso en el corpus especializado (Figura 5), permitiendo así adquirir información terminológica muy valiosa sobre el uso real de los términos.

vasculite 14258176-n 32 <i>inflammation of a blood vessel or lymph duct</i>
<ul style="list-style-type: none"> • ili_p: 1 • sense_p: 1 • ili_f: 0.32 • sense_f: 1
apendicite 14258512-n 57 <i>inflammation of the vermiform appendix</i>
<ul style="list-style-type: none"> • ili_p: 1 • sense_p: 1 • ili_f: 0.57 • sense_f: 1
arterite 14258609-n 17 <i>inflammation of an artery</i>
<ul style="list-style-type: none"> • ili_p: 1 • sense_p: 1 • ili_f: 0.17 • sense_f: 1

Figura 4: Evaluación de los términos.

0.997333 - vasculite vasculite NCF5000 0.250038 14258176-n:0.0103699 ou ou CC 1 - a o DA0F50 0.696141 - miopattis miopatia NCFP000 1 14209201-n:0.0103244 conxéntas conxénito AQ0F50 1 01315844-a:0.0110103 . . Fp 1 -
3 [CTG 052/2148] - A o DA0F50 0.696141 - enfermidade enfermidade NCF5000 1 14070360-n:0.0125276/14061805-n:0.0103609/14055408-n:0.00737465/13923440-n:0.00637642 de de SPS00 0.997333 - Kawasaki kawasaki NP00000 1 - é ser VSIP30 1 00339934-v:0.00787303/02604760-v:0.0046272/02445925-v:0.00417368/02620587-v:0.00413664/02749904-v:0.00361052/01029368-v:0.00358414/02616386-v:0.00357256 unha un DIOF50 0.969159 - vasculite vasculite NCF5000 0.250038 14258176-n:0.0259059 sistémica sistémico AQ0F50 0.916667 - aguda agudo AQ0F50 1 00803038-a:0.00613675/01213197-a:0.00601892/00661885-a:0.00539746/00044760-a:0.00538203

Figura 5: Término en contexto.

4 Conclusiones

La verificación de los términos en un corpus desambiguado permite adquirir información muy valiosa sobre su uso real y constituye una fuente de conocimiento muy relevante en la expansión de Galnet guiada por campos conceptuales. Los resultados obtenidos en la extracción, avalados por su evaluación en corpus, nos animan a continuar investigando en esta dirección y a seguir completando el WordNet del gallego desde esta perspectiva.

Bibliografía

- Agirre, E. y A. Soroa. 2009. Personalizing PageRank for Word Sense Disambiguation. En *Proceedings of the 12th Conference of the European Chapter of the ACL*, págs. 33–41.
- Gómez Guinovart, X. y Antoni O. 2014. Methodology and evaluation of the Galician WordNet expansion with the WN-Toolkit. *Procesamiento del Lenguaje Natural*, 53:43–50.
- Gómez Guinovart, X. y M. A. Solla Portela. 2014. O dicionario de sinónimos como recurso para a expansión de WordNet. *Linguamática*, 6(2):69–74.
- González Agirre, A., E. Laparra, y G. Rigau. 2012. Multilingual Central Repository version 3.0. En *6th Global WordNet Conference*.