

Recursos lingüísticos

Miguel A. Alonso Jorge Graña Jesús Vilares

Departamento de Computación, Facultad de Informática, Universidade da Coruña

- 1 Tipos de recursos lingüísticos
- 2 Ejemplos de corpus etiquetados
 - El corpus ITU o CRATER
 - CORGA
- 3 Ejemplos de Treebanks
 - Susanne

Tipos de principales recursos lingüísticos

- Corpus etiquetados
- Diccionarios (lexicones)
- Bancos de árboles (Treebanks)

Corpus etiquetado

- Un **corpus** es una colección de documentos
- Un **corpus etiquetado** es aquél en el que todas y cada una de las palabras aparecen acompañadas al menos de su **etiqueta correcta**
- También suele aparecer el **lema**
- Sirven como referencia de un estilo de uso de una lengua durante el proceso de ajuste o puesta a punto de las herramientas de etiquetación

Diccionario

- En su forma más simple, una lista de palabras
- Generalmente acompañada de sus posibles etiquetas o papeles candidatos que puede desempeñar dentro de la frase
- También puede contener acepciones semánticas
- También puede contener relaciones entre palabras
 - derivativas (poco común)
 - relaciones semánticas (WordNet, EuroWordNet, ...)
 - ...
- También puede contener información de un dominio específico (SentiWordNet, ...)

Treebank

- Cada frase aparece completamente analizada sintácticamente
- Dos grandes paradigmas (no excluyentes):
 - Treebanks de constituyentes
 - Treebanks de dependencias
- Permiten extraer:
 - las reglas de la gramática a la que obedece un determinado idioma, o al menos un subconjunto relevante de él, incluso con una probabilidad de uso asignada a cada una de esas reglas.
 - Las relaciones de dependencias entre palabras, así como sus probabilidades relativas

Ejemplo: El corpus ITU o CRATER

- Corpus etiquetado **paralelo** inglés-francés-español (contiene los mismos documentos en varios idiomas)
- Dominio de los textos: telecomunicaciones (ITU es la International Telecommunications Union)
- Tamaño:
 - Texto original: 5 millones de palabras
 - Texto anotado automáticamente: 5 millones de palabras
 - Textos anotado **revisado**: medio millón de palabras

En/PREP/en esta/DMPXFS/este colaboración/NCFS/colaboración debe/VMPI3S/deber
reconocerse/VCLI/reconocer el/ARTDMS/el carácter/NCMS/carácter consultivo/ADJGMS/consultivo
de/PREP/de las/ARTDFP/el organizaciones/NCFP/organización que/CQUE/que participan/VLPI3P/participar
en/PREP/en los/ARTDMP/el trabajos/NCMP/trabajo del/PDEL/del CCITT/NPTOS/CCITT ,/CM/,
en.particular/ADVN/en.particular a/PREP/a la/ARTDFS/el ISO/ACRNM/ISO ,/CM/, desde/PREP/desde el/ARTDMS/el
punto/NCMS/punto de/PREP/de vista/NCFS/vista de/PREP/de su/PPOSPS/suyo labor/NCFS/labor
con.respecto.a/PREP/con.respecto.a los/ARTDMP/el sistemas/NCMP/sistema de/PREP/de datos/NCMP/dato y/CC/y
a/PREP/a las/ARTDFP/el comunicaciones/NCFP/comunicación ./FS/. < > < > < > < > < > < > < >

Estadísticas

- 486.073 palabras (7.564.781 bytes)
- 14.919 frases (media de 33 palabras por frase)
- cada línea es una frase:
palabra/etiqueta/lema palabra/etiqueta/lema ...
- 12.462 lemas, 17.138 formas con 18.917 etiquetas posibles:
 - 15.745 formas con 1 etiqueta
 - 1.097 formas con 2 etiquetas
 - 223 formas con 3 etiquetas
 - 61 formas con 4 etiquetas
 - 8 formas con 5 etiquetas
 - 3 formas con 6 etiquetas
 - 1 forma con 7 etiquetas
- 8, 13% de formas ambiguas
- 1, 10 etiquetas por forma

Estadísticas

- Estadísticas sobre palabras:
 - 263.592 palabras con 1 etiqueta
 - 107.746 palabras con 2 etiquetas
 - 79.751 palabras con 3 etiquetas
 - 7.013 palabras con 4 etiquetas
 - 18.949 palabras con 5 etiquetas
 - 8.528 palabras con 6 etiquetas
 - 494 palabras con 7 etiquetas
- 486.073 palabras con 895.760 etiquetas posibles.
- 45,77% de palabras ambiguas
- 1,84 de etiquetas por palabra

Tagset

ACRNM	Sigla (ADN)
ADJCP	Adjetivo comparativo general plural (mayores, menores)
ADJCS	Adjetivo comparativo general singular (mayor, menor)
ADJGFP	Adjetivo positivo general femenino plural
ADJGFS	Adjetivo positivo general femenino singular
ADJGMP	Adjetivo positivo general masculino plural
ADJGMS	Adjetivo positivo general masculino singular
ADJSFP	Adjetivo superlativo general femenino plural (máximas, mínimas)
ADJSFS	Adjetivo superlativo general femenino singular (máxima, mínima)
ADJSMP	Adjetivo superlativo general masculino plural (máximos, mínimos)
ADJSMS	Adjetivo superlativo general masculino singular (máximo, mínimo, grandísimo)
ADVGR	Adverbio con grado positivo (muy, demasiado, mucho)
ADVGRC	Adverbio con grado comparativo (más, menos)
ADVGRS	Adverbio con grado superlativo (abundantísimamente)
ADVINT	Adverbio interrogativo (cómo)
ADVL	Adverbio locativo direccional (abajo)
ADVLD	Adverbio locativo dinámico (adelante)
ADVLE	Adverbio locativo estático (dentro)
ADVLLD	Adverbio locativo interrogativo dinámico (adónde)
ADVLLN	Adverbio locativo interrogativo (dónde)
ADVLP	Adverbio locativo con deixis próxima (aquí)
ADVLR	Adverbio locativo con deixis remota (allí)
ADVLRD	Adverbio locativo relativo dinámico (adonde)
ADVLRN	Adverbio locativo relativo direccional (donde)
ADVLMR	Adverbio relativo modal (como)
ADVNL	Adverbio general (salvajemente, bien, probablemente)
ADVNEG	Adverbio general negativo (tampoco)
ADVT	Adverbio temporal (ahora, ayer)
ADVTIN	Adverbio temporal interrogativo (cuándo)
ADVTNE	Adverbio temporal negativo (nunca)
ADVTRE	Adverbio temporal relativo (cuando)
ALFP	Letra del alfabeto en plural (As/Aes, bes)

Ejemplo: El CORGA

- COrpus de Referencia do Galego Actual
- <http://corpus.cirp.es/corga/>
- Documentos en XML
- Dominio general: 457 periódicos, 103 revistas, 415 libros (novela, ensayo, relato corto y teatro) desde 1975 a la actualidad.
- Tamaño
 - 25 millones de palabras
 - 59.193 palabras etiquetadas y revisadas manualmente(767 noticias)

Tagset

			m= masculino f= femenino a= masculino ou feminino n= neutro 0= non aplica	s= singular p= plural a= singular ou plural 0= non aplica	c= comparativo s= superlativo 0= non aplica	1= primeira 2= segunda 3= terceira a= primeira ou terceira 0= non aplica	n= nominativo a= acusativo d= dativo p= preposicional f= fuxante léxico 0= outros casos	p= presente i= copretérito 1= antepretérito s= futuro e= pretérito c= posesitivo 0= non aplica	i= indicativo s= subjuntivo m= imperativo f= infinitivo X= xerando p= participio	s= singular p= plural a= singular ou plural	d= determinante n= non determinante a= determinante ou non determinante
Categoría	Tipo	Subtipo	Sexo	Número	Grado	Persona	Caso	Tempo Verbal	Modo	Posición	Valor
Substantivo S	c= común p= propio		m/f/a/n	s/p							
Adectivo A			m/f/a/n	s/p	c/s						
Verbo V			m/f/a/n	s/p		1/2/3/a		p/f/c/e	i/s/m/f/X/p		
Preposición P											
Conxunción C	c= coordinante a= subordinante										
Adverbio W	m= nuclear n= modificador a= nuclear ou modificador r= relativo q= interrogativo/exclamativo										
Artigo D	d= determinado i= indeterminado		m/f	s/p							
Demostrativo E			m/f/n	s/p							d/s
Relativo T			m/f/a/n	s/p							a/s
Posestivo M			m/f	s/p		1/2				s/p	d/s
Indefinido I			m/f/a/n	s/p							d/s
Numeral N	c= cardinal s= ordinal		m/f	s/p							d/s
Pronome R	t= tónico a= áctico		m/f/n	s/p		1/2					
Interrogativo/exclamativo G			m/f	s/p							d/s
Locución L	a= adverbial p= preposicional c= conxectiva	c= coordinante s= subordinante 0= non aplica									
Intersección Y											

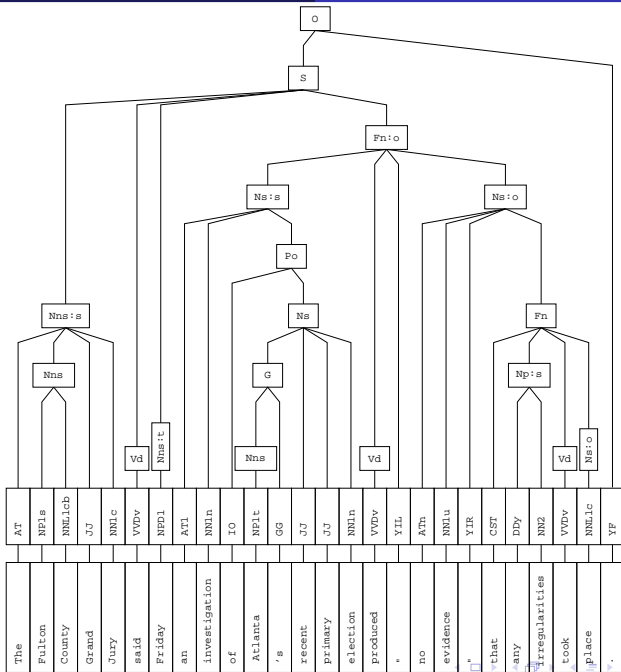
Ejemplo: El corpus Susanne

- Treebank de constituyentes del inglés
- Tamaño: 150.000 palabras (10.500 lemas diferentes) del corpus BROWN para el inglés americano

```

...
A01:0010b - AT      The          the          [O[S[Nns:s.
A01:0010c - NP1s   Fulton      Fulton      [Nns.
A01:0010d - NNL1cb County     county      .Nns]
A01:0010e - JJ      Grand       grand       .
A01:0010f - NN1c   Jury        jury        .Nns:s]
A01:0010g - VVDv   said        say         [Vd.Vd]
A01:0010h - NPD1   Friday      Friday      [Nns:t.Nns:t]
A01:0010i - AT1    an          an          [Fn:o[Ns:s.
A01:0010j - NN1n   investigation investigation .
A01:0020a - IO     of          of          [Po.
A01:0020b - NP1t   Atlanta     Atlanta     [Ns[G[Nns.Nns]
A01:0020c - GG     +<apos>s   -          .G]
A01:0020d - JJ      recent      recent      .
A01:0020e - JJ      primary     primary     .
A01:0020f - NN1n   election    election    .Ns]Po]Ns:s]
A01:0020g - VVDv   produced    produce     [Vd.Vd]
A01:0020h - YIL    <ldquo>     -          .
A01:0020i - ATn    +no        no          [Ns:o.
A01:0020j - NN1u   evidence    evidence    .
A01:0020k - YIR    +<rdquo>   -          .
A01:0020m - CST    that        that        [Fn.
A01:0030a - DDy    any         any         [Np:s.
A01:0030b - NN2    irregularities irregularity .Np:s]
A01:0030c - VVDv   took        take        [Vd.Vd]
A01:0030d - NNL1c place       place       [Ns:o.Ns:o]Fn]Ns:o]Fn:o]S]

```



Tagset

APPGf	her como posesivo ≠ PPHO1f
APPGh1	its
APPGh2	their
APPGi1	my como posesivo
APPGi2	our
APPGm	his excepto como pronombre ≠ PPGm
APPGy	your
AT	the como determinante
AT1	Artículo indefinido a, an
AT1e	every
ATn	no como determinante o calificador ≠ UH
BTO21	in en la secuencia in order + infinitivo
BTO22	order en la secuencia in order + infinitivo
CC	Conjunción coordinada
CC31	Conjunción coordinada de 3 palabras (1a. palabra), e.g. <u>as</u> well as
CC32	Conjunción coordinada de 3 palabras (2a. palabra), e.g. as <u>well</u> as
CC33	Conjunción coordinada de 3 palabras (3a. palabra), e.g. as well <u>as</u>
CCB	but como conjunción coordinada ≠ ICSx RR
CCn	nor
CCr	or
CS	Conjunción subordinada
CS21	Conjunción subordinada de 2 palabras (1a. palabra), e.g. <u>even</u> though
CS22	Conjunción subordinada de 2 palabras (2a. palabra), e.g. even <u>though</u>
CS31	Conjunción subordinada de 3 palabras (1a. palabra), e.g. <u>as</u> long as
CS32	Conjunción subordinada de 3 palabras (2a. palabra), e.g. as <u>long</u> as
CS33	Conjunción subordinada de 3 palabras (3a. palabra), e.g. as long <u>as</u>
CSA	as como conjunción subordinada o como preposición en sentido comparativo ≠ Ila RGA
CSN	than en todos los usos
CST	that como conjunción subordinada ≠ DD1a
CST21	<u>as</u> how
CST22	as <u>how</u>
CSW	whether en todos los usos

Extracción de reglas gramaticales

- El símbolo [indica el comienzo de una regla: se toma el siguiente símbolo del árbol parentizado, se añade a la parte derecha de la regla que está en la cima de la pila, y se crea un nuevo elemento en la pila para una nueva regla cuya parte izquierda es dicho símbolo
- El símbolo] indica el final de una regla: la regla que está en la cima de la pila en está ya completa, por lo que se extrae y se elimina de la pila
- Cualquier otro símbolo se añade a la parte derecha de la regla que está en la cima de la pila

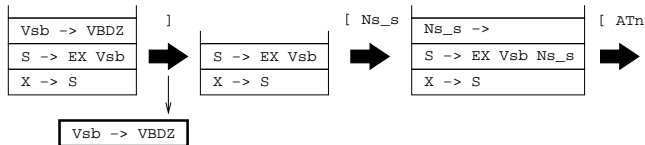
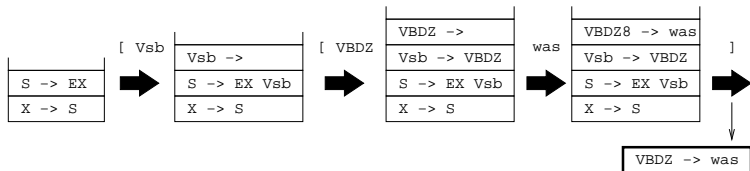
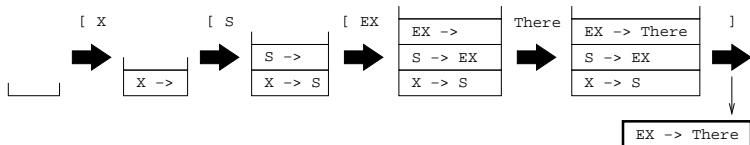
Algoritmo de extracción de reglas gramaticales

```

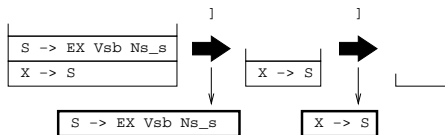
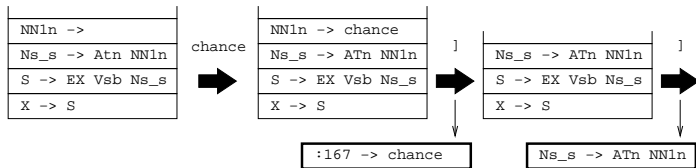
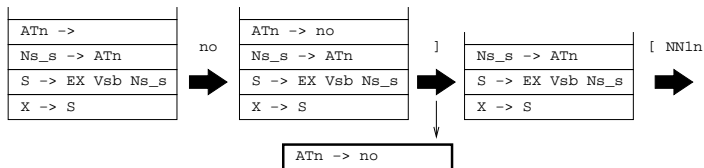
P ← Pila_Vacía;
R ← Conjunto_Vacío;
while (queden símbolos del Árbol_Parentizado por procesar) do begin
  A ← siguiente símbolo del Árbol_Parentizado;
  case A of
    [ : A ← siguiente símbolo del Árbol_Parentizado;
      Regla ← pop (P);
      Añadir A en la parte derecha de Regla;
      push (P, Regla);
      push (P, A → );
    ] : Regla ← pop (P);
      Añadir Regla a R;
    cualquier otro símbolo :
      Regla ← pop (P);
      Añadir A en la parte derecha de Regla;
      push (P, Regla)
  endcase
endwhile;
return R

```

Ejemplo de extracción de reglas gramaticales



Ejemplo de extracción de reglas gramaticales



Treebank SUSANNE

Corpus sin trazas	Corpus con trazas
4.292 frases	2.188 frases
77.275 palabras	60.759 palabras
18 palabras por frase	28 palabras por frase
44,93% de las palabras son ambiguas	43,50% de las palabras son ambiguas
2,30 etiquetas por palabra	2,39 etiquetas por palabra

