



DH@Madrid Summer School 2016

Procesamiento del Lenguaje Natural

y Poesía

Núria Bel (UPF)

nuria.bel@upf.edu

Madrid, 27 junio al 1 de julio de 2016







PLN y poesía

Núria Bel



Procesamiento del Lenguaje Natural

Objetivo

Realizar tareas que asociamos a leer, entender y extraer información

¿Qué tareas?

Todas las que podamos resolver como una secuencia de acciones planteadas en términos de "si ... entonces ..."



"si ... entonces ..."???

Los programas informáticos son colecciones de instrucciones diseñadas por humanos que se basan en la declaración de condiciones que se han de cumplir para tomar decisiones y manipular el input.

Por ejemplo:

Corrector ortográfico: si una palabra no está en la lista de palabras de la lengua, entonces señala que es un error.

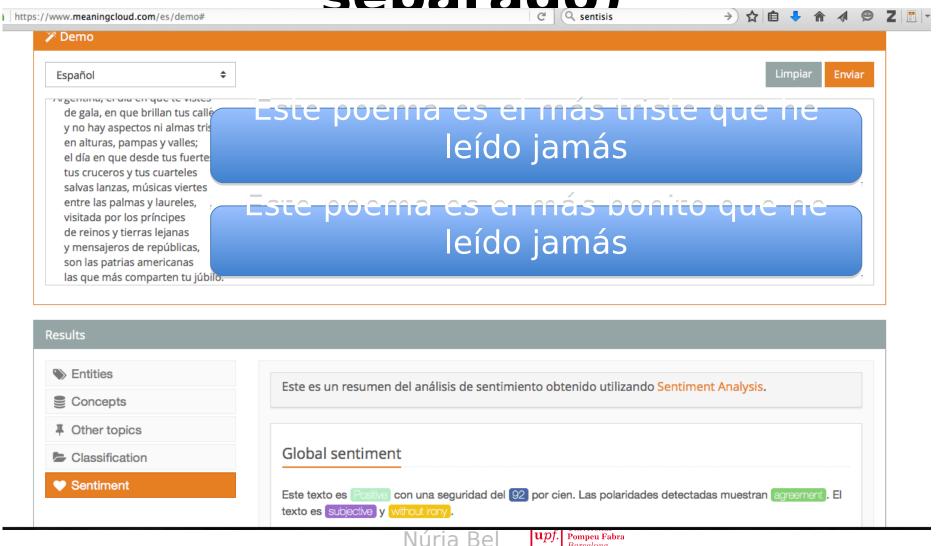


Ejercicio 1 (1/5)

- Mira este programa. https://www.meaningcloud.com/demo
- Busca "sentiment" en las opciones de análisis
- Nos dice si un texto puede ser interpretado como una opinión positiva o negativa.
- ¿cómo crees que puede decirlo?
 - Piensa en el corrector ortográfico y en plantearlo como "si ... entonces ..."
- Prueba estos dos textos. ¿Hace lo que esperas?
 - Este poema es el más bonito que he leído jamás
 - Este poema es el peor que he leído jamás



Ejercicio 1 (2/5) Prueba los textos ... (por separado)



Barcelona

Ejercicio 1 (3/5)

Respuesta: efectivamente, como en el caso del corrector ortográfico, tiene una lista de palabras que considera "positivas" y otras "negativas".

De este extracto del poema de Ruben Darío, Canto a la Argentina, ¿qué palabras pondrías en la lista de palabras positivas y cuáles en la de negativas?

Argentina, el día en que te vistes de gala, en que brillan tus calles y no hay aspectos ni almas tristes en alturas, pampas y valles; el día en que desde tus fuertes, tus cruceros y tus cuarteles salvas lanzas, músicas viertes entre las palmas y laureles, visitada por los príncipes de reinos y tierras lejanas y mensaieros de

repúblicas, son las

Ruben Darío, Canto a la Argentina (1914), Extraído de Proyecto Gutemberg:

https://www.gutenberg.org/files/51458/51458-0.txt

Ejercicio 1 (4/5) Se calcula la ratio positivas/negativas

Argentina, el día en que te vistes de gala, en que brillan tus calles y no hay aspectos ni almas tristes en alturas, pampas y valles; el día en que desde tus fuertes, tus cruceros y tus cuarteles salvas lanzas, músicas viertes entre las palmas y laureles, visitada por los príncipes de reinos y tierras lejanas y mensaieros de

Palabras positivas	Palabras negativas
gala	tristes
brillan	
laureles	
júbilo	

Ruben Darío, Canto a la Argentina (1914), Extraído de Proyecto Gutemberg:

https://www.gutenberg.org/files/51458/51458-0.txt

Ejercicio 1 (5/5) ¿Preguntas?

¿Y las otras palabras?

Si no están en las listas, asignadas a una información **explícita**, el programa no las puede reconocer ni tenerlas en cuenta.

¿Más preguntas?



Procesamiento del Lenguaje Natural

Procesar: **reconocer** palabras y asignarles una representación (información) que permita poder tomar decisiones.

Representación: **añadir** (substituir) a la cadena de caracteres que forma una palabra (o secuencia de palabras) información explícita de sus características para una tarea determinada

¿**Qué** información explícita? Depende de la tarea ..



PLN, tareas más conocidas

- Análisis de opinión
- Corrección gramatical
- Traducción automática
- Búsqueda y recuperación de información
- Extracción de información
- Resumen automático
- Respuesta a preguntas y asistentes virtuales
- Análisis lingüístico

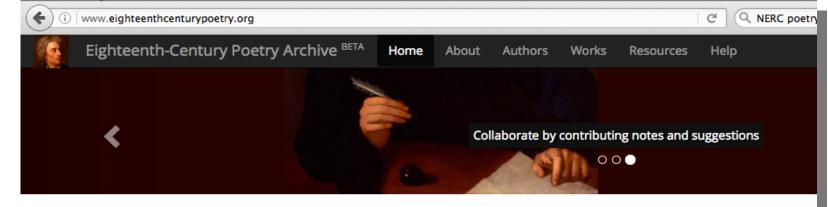


¿¿¿Y para poesía???

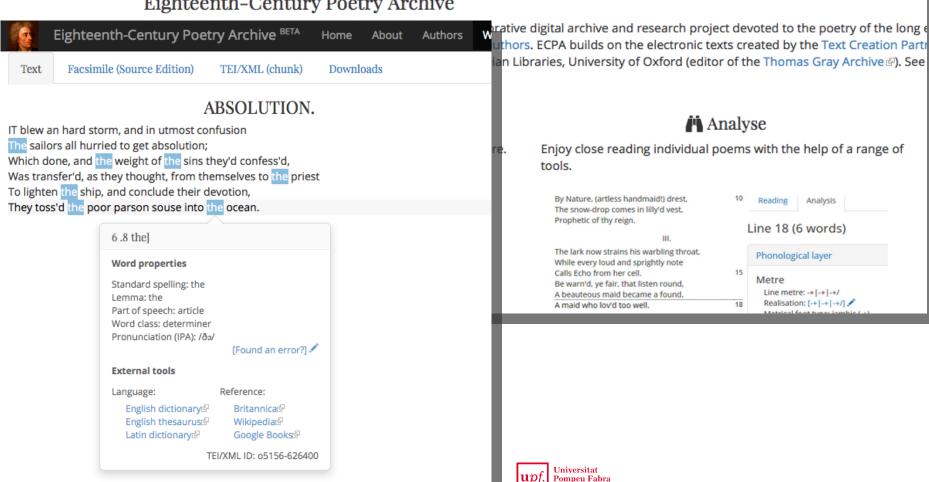
¿Extracción de información? ¿Análisis de opinión? ¿Análisis lingüístico?

pero hay que saber qué puede hacer y qué no ...





Eighteenth-Century Poetry Archive



Barcelona

Las palabras y la tecnología ...

- Empezamos por "palabras en una lista"
- Hay que hacer un programa para reconocer palabras en un texto
 - si QUÉ entonces "es una palabra"

- Hay que saber un poco más ...
- Haremos un pequeño paréntesis:
 Formato y codificación de textos



Formato y codificación de textos (1/11) Palabras son...

- Una palabra es una secuencia de caracteres
- Un caracter es un código
- Los espacios en blanco también son código!
- Una palabra será una secuencia de códigos de caracteres entre dos códigos de espacio en blanco.

textos (2/11) Código para todo ... para espacio, tabulador, final de línea ...

```
he aquí el paraíso terrestre, ¶
he aquí la ventura esperada, ¶
he aquí el Vellocino de Oro, ¶
he aquí Canaán la preñada, \P
la Atlántida resucitada; ¶
he aquí los campos del Toro¶
y del Becerro simbólicos; ¶
he aquí el existir que en sueños ¶
miraron los melancólicos, ¶
los clamorosos, los dolientes ¶
poetas visionarios ¶
que en sus olimpos o calvarios ¶
amaron a todas las gentes. ¶
```

Formato y codificación de textos (3/11) Un caracter también es código: una secuencia de biolio 0001 a

0100 0001 A

✓ Para buscar "triste" en el texto:

"La nena triste quiere un caramelo"

- Se comparará cada palabra del texto para decidir si es o no es lo que busca
- Se comparará cada caracter con el criterio de búsqueda
- ✓ Se comparará cada bit ...

```
"triste" - "La "
"c" - "L"
"01100011" - "01001100"
"0" AND "0" = 1
"1" AND "1" = 1
"1" AND "0" = 0 Núria Bel
```



Formato y codificación de textos (4/11) Codificación de caracteres

- ✓ "Codificación de caracteres" se refiere a las "correspondencias entre cada caracter ("a") y una secuencia única de código binario (01100001)"
- ✓ Hay diferentes listas de correspondencias, (todas estándares):

ASCII (para el alfabeto del inglés)

ISO-8859 (para el de las lenguas europeas)

UTF8 - UNICODE (para todas las del mundo)

http://www.unicode.org



Formato y codificación de textos (5/11)

Una aplicación o programa que no tiene la misma correspondencia que la de la codificación del texto que queremos analizar

... no lo leerá bien!

```
he aquí el para√≠so terrestre, ¶
 he aquí la ventura esperada, ¶
he aquí el Vellocino de Oro, ¶
he aquí Cana√on la pre√±ada, ¶
la Atl√ontida resucitada; ¶
he aquí los campos del Toro¶
y del Becerro simbólicos; ¶
he aquí el existir que en sue√±os¶
miraron los melancólicos, ¶
los clamorosos, los dolientes ¶
poetas visionarios ¶
 que en sus olimpos o calvarios ¶
  amaron a todas las gentes. T
```

Muchas
aplicaciones
para el
inglés no
pueden ni
leer los
textos en
otras
lenguas que
necesitan
otros
caracteres!

Formato y codificación de textos (6/11)

Lo mejor es guardar documento en "formato texto" ".txt" y "codificación de texto": UNICODE UTF8

Conversión de archivo - UNED-	-2016.txt
Advertencia: al guardar como archivo de texto se perderá todo el forn Codificación de texto:	nato, las imágenes y los objetos del archivo.
Mac OS (predeterminado) MS-DOS • Otra codificación:	Occidental (Windows latino 1) Turco (Mac OS)
Opciones: Insertar saltos de línea	Turco (Windows latino 5) Unicode 6.0 Unicode 6.0 (Little-Endian)
Terminar líneas con: Sólo retorno de carro \$	Unicode 6.0 UTF-8
Permitir la sustitución de caracteres	

Esto es lo que sale en mi ordenador al guardar un documento en formato ".txt". Puede salir de otra manera en un ordenador windows o linux.

Formato y codificación de textos (7/11) ¿Por qué "formato de texto"?

- ✓ El texto ha de ser **procesable** y no contener código (para la impresora o pantalla)
- √ "Texto sin formato" es sin ese código
- ✓ PDF no es un formato pr (está pensado para la in

Se pueden utilizar conversores, pero.... suelen dar problemas.

EK^C^D^T^@^@^H^@^@Yf<89>?^R2^L'^@^@^@'^@^@^H^@^@^@mimetypeapplication/vnd.oas

s.opendocument.textPK^C^D^T^@*H^H^H^H^@Yf<889>?^@^@^@^@^@^@^@^@^@^@^@^@^@^@^@^@
ent.xmlilirā6^R\$iS°<94>001;{<8a>4^>%i>a*w*nY'5ycO6vR[{IÁs(a8B8\^P<94>à-2ec-rŮ^GØ-Á^^
£)\R^ULI-9e>d^(å-8f5<89>-929-HÜ#<8b>*V-93>-8a>-AY@w^?<8d>1F<8b>f59\$ñ\m<8a>y@^
=n<98>IVCÂÔb6iāāR7W_è<83>RéÉ_>a<8e>C,<?<99>^Uz<98>

ŸĎTÄO^M_18~8a>F<8f>^[!\$#<86>^B^RC&C;\dop\n

Ne<9b>^HÏ]¿ÝåhB:æ¶<9d>K

iáf^?^x^Z<8a>(u\<82>Ÿd<97>¤′±:øÆÇ<9c>HM<90>+^]A÷^80^Z8^G6G^YÌÂ<9 '³<c°9b⊳Lr^C^[*<93><8c>^DG4<90><94>^PZäP¬,Ø=h*<t¹,Ö8<9a>^P3a>^PZäP¬,

0><85>u^[[npòI^TßòÜZt-ō:n\s85>µ.o\k^FÀA,K"<ah.^L^Z<9b>§yC æ*<85> 1>Ï<82><8f><97> £»^MmarÉ <8f>1^E^BØ¿<9C>v<88>.PØDX^På!<88>^Så^N%

1/06U, V-885+288-H-<81-34b6<99><8c-u½-786i`Òa-1Í-93>-:<91>ói ÁMf9V,µ

A, !66uiG\<97>^P\lÄfF2±,<8e>+f\q^Ra=82>\degx\0^1c3>\degx\0^1c3>\degx\0^1c3>\degx\0^1c3>\degx\0^1c3\degx\0^

Formato y codificación de textos (8/11)

¿Dónde encontrar textos en formato texto? ¡un problema!

Aquí http://www.gutenberg.org/ Online Book Catalog Download Bibrec Download This eBook Format (2) Read this book online: HTML RUBÉN DARÍO ILUSTRACIONES EPUB (with images) ENRIQUE OCHOA EPUB (no images) tas. Administración: Edi-Kindle (with images) Kindle (no images) Plain Text UTF-8 Q# 🕶 🚐 More Files...



textos (9/11) También el .html (página web) Se puede intentar ódigidar página en formato texto" desde el navegador



Formato y codificación de textos (10/11) También .html (página web) da más problemas:

caracteres que no son del texto, pero que se leen como si lo fueran

	Égloga II		
	[Poema: Texto co	mpleto.]	
		Garcilaso de la V	⁷ ega
Albanio			
		1	
En medio de	el invierno está templad	la	
el agua dulc	e desta clara fuente,	'	
y en el verar	no más que nieve helad	a. '	
•	ndas, cómo veo presen		
•	, la memoria d'aquel d	-	
de que el alr	na temblar y arder se s	iente!	
	laridad vi mi alegría	'	
	toda v enturbiarse:	i	

Formato y codificación de textos (11/11) Ejercicio 2

- Escojan su obra favorita en castellano y descárguenla
- Pueden buscar en https://www.gutenberg.org/browse/languages/es
- Por favor, ¡que sea larguito!
- Si no está en .txt, codificación utf8, abran el documento en un procesador (word, OpenOffice) y guarden una versión .txt, codificación utf8.



Las palabras y la tecnología

volvamos a lo nuestro ...

"si ... entonces ..."????

El programa añade información en forma de etiquetas

si entonces ETIQUETA

si una secuencia de caracteres está entre caracteres de espacio en blanco, o de puntuación (o...), entonces añadir etiqueta: "es un TOKEN-palabra".

si una secuencia de TOKEN-palabra está entre TOKEN-puntuación-inicio-de-frase y TOKEN-puntuación-final-de-frase entonces añadir etiqueta: es un TOKEN-oración.



Marcas XML, TEI, por ejemplo pero son etiquetas para

▼ XML output

```
<document>
  <wordcount>70</wordcount>
  <cputime>0.015291</cputime>
  <paragraph>
    <sentence id="1">
    <token form="he" id="t1.1">
```

Si reconocemos palabras ... entonces...

Cuando reconocemos palabras, es posible:

- √ contarlas
- √ saber con qué otras palabras aparecen
- ✓identificar qué características lingüísticas tienen:
 - · categoría gramatical: Nombre, Verbo, Adjetivo, ...
 - características morfosintácticas: Femenino, Pretérito, ..
 - otras características: nombre propio de persona,

Son programas de **análisis lingüístico**, la base del Procesamiento del Lenguaje Natural



Programas de análisis lingüístico (1/5)

Análisis lingüístico para saber:

- palabras más frecuentes
- combinaciones (bigramas) más frecuentes de palabras
- combinaciones significativas de palabras: colocaciones

¿cómo sabe el programa qué palabra es un adjetivo?

Si entonces ADJ Si está en una lista de palabras ADJ (diccionario)

Pero hay algunas complicaciones

Herramientas de análisis lingüístico (2/5) http://nlp.lsi.upc.edu/freeling/node/1

FreeLing Home Page

Hooked on a FreeLing

Programa de código abierto, descargable, instalable en windows, mac y linux, gratis para investigación.

Main menu

- Home
- Features
- Linguistic Data
- Contributions
- License
- Installing
- Documentation
- Contributing
- Download
- Source code
- References
- Web Links
- Forum & FAQs

References

To cite FreeLing in your academic works, please reference the following papers:

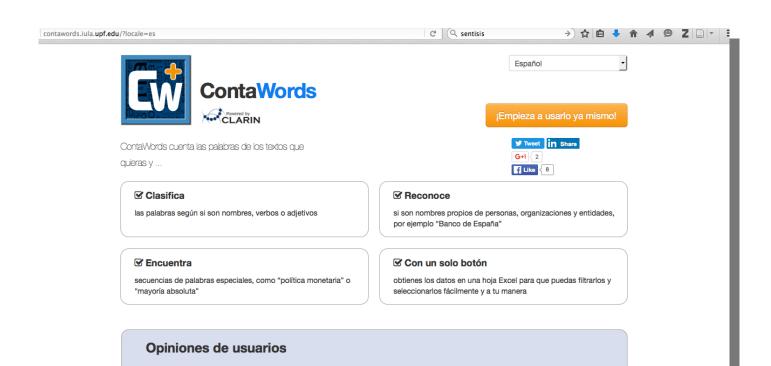
About FreeLing as a whole:

Lluís Padró and Evgeny Stanilovsky.
 FreeLing 3.0: Towards Wider Multilinguality
 Proceedings of the Language Resources and Evaluation
 Conference (LREC 2012) ELRA.
 Istanbul, Turkey. May, 2012.
 [pdf] [bibtex]

lingüístico (3/5) http://contawords.iula.upf. edu

- utiliza Freeling
- es una aplicación wel

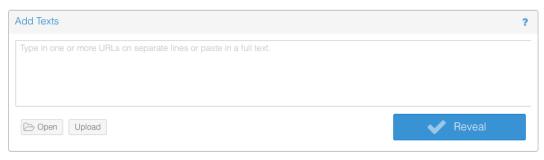
No necesita instalación, gratis para investigación



Herramientas de análisis lingüístico (4/5) http://voyant-tools.org/

- es una aplicación web
- tiene funciones interesantes





Herramientas de análisis lingüístico (5/5)

Para obtener:

- ✓ Información cuantitativa
- ✓ Reconocer nombres propios y clasificarlos (lugar, persona, etc.)
- ✓ Anotación morfosintáctica para saber formas de pasado, en primera persona

Se puede hacer a mano, pero depende del texto se puede hacer muy pesado, ¿o no?

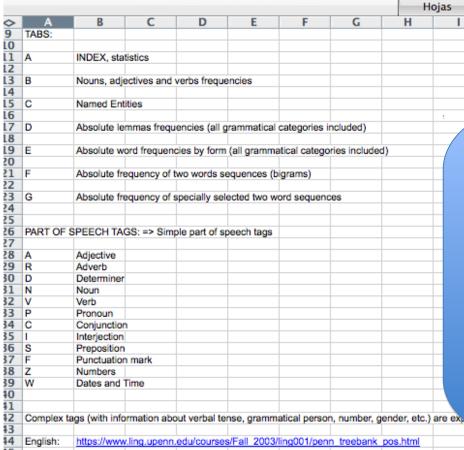


ContaWords

Gráficos

http://contawords.iula.upf.edu/more







Dado uno o más documentos:

Gráfico

- extrae listas de palabras según su categoría,
- lista las palabras por frecuencia,
- identifica y clasifica nombres propios y
- cuenta las combinaciones de dos palabras

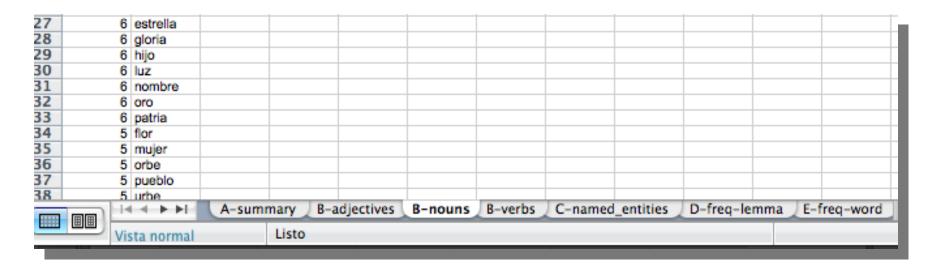
y lo pone todo en una hoja excel

http://nlp.lsi.upc.edu/freeling/doc/tagsets/tagset-ca.html



Distingue entre Nombres, Verbos, Adjetivos, y reconoce "Entidades con Nombre"





Entidades con Nombre



¡Argentina, región de la Ción

En inglés, Named Entity Recognition (NERC)

Usa información del contexto donde

aparece

región de la aurora!

¡Argentina! tu ser no abriga

El Plata, padre extraordinario,

10	Argental Mete enforces o jes ele	Getyrzphic :
	América	Geographic
6	Argentina	Organization
6	Centenario	Others
4	España	Geographic
	Italia	Geographic
	Buenos Aires	Geographic
	Dios	Geographic
	Sol	Person
	Dios	Others
	Exodos	Person
	Libertad	Others
2	Oid	Person
2	Plata	Person
	Roma	Geographic
2	Roma	Person
-		



Entidades con Nombre "Extracción de información"



Comete errores, sí, pero también el error está en la palabra

	Α	В	Ċ	D
	10	Argentina	Geographic	
		América	Geographic	
		Argentina	Organizatio	n
		Centenario	Others	
		España	Geographic	
		Italia	Geographic	
		Buenos Aires	Geographic	
¡Exodos! ¡Exodos!		Dios	Geographic	
Evodoci		Sol	Person	
ILXUUUS:		Dios	Others	
		Exodos	Person	
		Libertad	Others	
		Oid	Person	
0:4	2	Plata	Person	
Oid		Roma	Geographic	
	2	Roma	Person	

Análisis Lingüístic Nombres (comunes frecuencia de lema



Α	В	С	D
	argentina		
	tierra		
14	imagen		
13	fiesta		
	paz		
11	dios		
11	hombre		
	sol		
	himno		
9	libertad		
	hermano		
	sangre		
7	centenario		
	día		
7	grito		
7	mar		
7	nación		
7	pampa		
	sombra vida		
- 7	voz		
	VUZ		lúrio Dol

Atención!!
Son lemas, si queremos
formas ir a la pestaña Efreq-word

Lema: representante del paradigma flexivo



Análisis Lingüístico Verbos: frecuencia de



lem

		_	
Α	В	С	
40	ser		
	haber		
16			
12	cantar		
11	dar		
8	estar		
	pasar		
7	libertar		
7	oír		
	tener		
7	ver		
6	llegar		
	mirar		
	abrir		
5	amar		
5	crear		
5	decir		
_			

información asignada: lema
e información
morfosintáctica:
arrebata arrebatar V
arrebatas arrebatar V
Pero hay palabras que tienen
más de una etiqueta en la
lista:

paso pasar V paso paso N El programa elige la que es en contexto:

"El paso de" vs. "Yo

oaso"



Análisis lingüístico Etiquetado morfosintáctico

- Para cada palabra, el programa accede al diccionario y recupera etiquetas asociadas.
- Elige la correcta en contexto (v es el

Ejemplo de etiquetas asociadas a palabras. Ambigüedad categorial Sentences Sentence 1 Vellocino de Oro he aquí el paraíso terrestre he aquí ventura esperada aguí he paraíso haber haber terrestre ventura esperar haber vellocino de oro VAIP1S0 DA0MS0 NCMS000 AQ0CS00 VAIP1S0 DA0FS0 VAIP1S0 DA0MS0 NP00000 0.962485 0.962485 0.98926 0.638706 0.962485 haber haber venturo haber VMIP1S0 VMIP1S0 PP3FSA0 AQ0FS00 VMIP1S0 0.010734 0.0363423 0.0363423 0.361294 0.0363423 NCMS000 0.00117233 0.00117233 6.20105e-06 0.00117233

Análisis lingüístico El diccionario



El diccionario es una lista (finita) de palabras con información lingüística que se ha hecho manualmente Si una palabra no está en la lista, puede suponerse la categoría por el sufijo

in the dictionary.

- The Asturian dictionary contains more than 140,000 forms corresponding to some 40,000 combinations lemma-PoS.
- The Catalan dictionary contains more than 520,000 forms corresponding to 71,000 combinations lemma-PoS.
- The Welsh dictionary contains some 345,000 forms, corresponding to over 8,500 lemma-PoS combinations.
- The German dictionary contains near 395,000 forms, corresponding to almost 130,000 lemma-PoS combinations.
- The English dictionary was automatically extracted from WSJ and other corpuses, with accurate manual post-edition and completion.
 - It contains about 68,000 forms corresponding to some 37,000 different
 - combinations lemma-PoS.
- The Spanish dictionary contains over 555,000 forms corresponding to more than 76.000 lemma-PoS combinations.
- The French dictionary contains over 350,000 forms corresponding to more than 54.000 lemma-PoS

Diccionarios disponibles en FreeLing Si la aplicación no tiene el diccionario de una lengua, no puede analizar textos de esa lengua

Análisis lingüístico Las etiquetas

La información lingüística, las características morfosintácticas de las palabras, es la etiqueta asociada

```
abalanzara abalanzar VMIC1SO | abalanzar VMIC3SO bajo bajar VMIP1SO | bajo AQOMSO | bajo
```

NCMS000 | bajo SPS00

La etiqueta intenta ser comprensible' Hay estándares de cómo han de ser las etiquetas.



Análisis Lingüístico Las etiquetas (Tagset): para cada lengua, y en estilos diferentes

- Penn Treebank: 45 tags
- Lancaster BNC/Full (C5/C7): 61/146 tags
- German Stuttgart-Tübingen Tag Set: 50 tags
- Swedish SUC Tag Set: 25 tags
- EAGLES: 254 tags (las que utiliza





Análisis Lingüístico Las etiquetas: Penn Treebank

una convención .. pero cuestionable

	Tag	Description	Example	Tag	Description	Example
	CC	Coordin. Conjunction	and, but, or	SYM	Symbol	+,%, &
	CD	Cardinal number	one, two, three	TO	"to"	to
	DT	Determiner	a, the	UH	Interjection	ah, oops
	$\mathbf{E}\mathbf{X}$	Existential 'there'	there	VB	Verb, base form	eat
	FW	Foreign word	теа сиІра	VBD	Verb, past tense	ate
	IN	Preposition/sub-conj	of, in, by	VBG	Verb, gerund	eating
	JJ	Adjective	yellow	VBN	Verb, past participle	eaten
	JJR.	Adj., comparative	bigger	VBP	Verb, non-3sg pres	eat
	JJS	Adj., superlative	wildest	VBZ	Verb, 3sg pres	eats
	LS	List item marker	1, 2, One	WDT	Wh-determiner	which, that
	MD	Modal	can, should	WP	Wh-prencun	what, who
	NN	Noun, sing. or mass	llama	WP\$	Possessive wh-	whose
	NNS	Noun, plural	llamas	WRB	Wh-adverb	how, where
	NNP	Proper noun, singular	IBM	\$	Dollar sign	\$
	NNPS	Proper noun, plural	Carolinas	#	Pound sign	#
1	PDT	Predeterminer	all, both	DC.	Left quote	(" or ")
	POS	Possessive ending	'S	30	Right quote	(' or '')
	PP		I, you, he	(Left parenthesis	([, (, {, <)
	PP\$	Possessive pronoun	your, one's)	Right parenthesis	$(],),\},>)$
	RB	Adverb	quickly, never	,	Comma	,
	RBR	Adverb, comparative	-		Sentence-final punc	
	RBS		fastest	:	Mid-sentence punc	(: ;)
	RP	Particle	up, off			



Análisis Lingüístico Penn TreeBank Tag Set for

Snanish

ACRNM acronym (ISO, CEI)

ADJ Adjectives (mayores, mayor)

ADV Adverbs (muy, demasiado, cómo)

ALFP Plural letter of the alphabet (As/Aes,

bes)

ALFS Singular letter of the alphabet (A, b)

ART Articles (un, las, la, unas)

BACKSLASH backslash (\)

CARD Cardinals

CC Coordinating conjunction (y, o)

CCAD Adversative coordinating conjunction (pero)

CCNEG Negative coordinating conjunction (ni)

. . .

PAL Portmanteau word formed by a and el

PDEL Portmanteau word formed by de and el



Análisis lingüístico Las etiquetas: EAGLES

- Las etiquetas EAGLES codifican la información en una secuencia donde la posición se relaciona con el atributo del que se codifican los valores.
- La primera posición siempre codifica la categoría gramatical.
- Si una forma no tiene un valor para un determinado atributo, se codifica '0'.

 Dependiendo de la categoría gramatical, habrá un número determinado de posiciones/valor

 Está preparado para cubrir las necesida lenguas europeas
 category
 N:noun

0	eas	category	N:noun
	1	type	C:common; P:proper
	2	case	N:nominative; G:genitive; D:dative; F:accusative;
	3	gen	F:f; M:m; C:c
	4	num	S:s; P:p; N:n

Ncms

paso

Ejercicio 2: EAGLES Verbos

Codifica las etiquetas para las siguientes formas verbales:

cantaría abriríamos maté leyera correrás oídas

inicial = 0, y aquí decimos = 1.

No importa el número, ha de

Pos.	Atributo	Valor	Código
1	Categoría	Verbo	V
2	Tipo	Principal	M
		Auxiliar	Α
		Semiauxiliar	S
3	Modo	Indicativo	I
		Subjuntivo	S
		Imperativo	M
		Infinitivo	N
		Gerundio	G
		Participio	P
4	Tiempo Presente		P
		Imperfecto	I
		Futuro	F
		Pasado	S
		Condicional	С
		-	0
5	Persona	Primera	1
		Segunda	2
		Tercera	3
6	Número	Singular	S
		Plural	P
7	Género	Masculino	M
		Femenino	F

Análisis lingüístico Etiquetado morfosintáctico automático FreeLing

Se basa en información de contexto ...

A partir de datos correctos (humanos) aprende la probabilidad de que, dada una etiqueta X, la anterior sea una etiqueta Y de entre las posibles para esa palabra.

▼ :	Sentences															
S	Sentence 1															
	he	aquí	el	paraíso	terrestre	,	he	aquí	la	ventura	esperada	,	he	aquí	el	Vellocino_de_Oro
	haber VAIP1S0 0.962485	aquí RG 1	el DAOMSO	paraíso NCMS000	terrestre AQ0CS00	Fc 1	haber VAIP1S0 0.962485	aquí RG 1	el DA0FS0 0.98926	ventura NCFS000 0.638706	esperar VMP00SF	Fc 1	haber VAIP1S0 0.962485	aquí RG 1	el DAOMSO	vellocino_de_oro NP00000 1
	haber VMIP1S0 0.0363423					_	haber VMIP1S0 0.0363423		lo PP3FSA0 0.010734	venturo <i>AQ0FS00</i> 0.361294		_	haber VMIP1S0 0.0363423			
	he / 0.00117233						he / 0.00117233		la NCMS000 6.20105e-06				he / 0.00117233			

Análisis lingüístico FreeLing Etiquetado morfosintáctico automático Select output Morphological Analysis

http://nlp.lsi.upc.edu/freeling/demo/demo.php

Dependency Parsing Sentences Coreferences Semantic Graph Sentence 1 Animará la virgen tierra la sangre de los finos brutos que da la animar el virgen tierra sangre de el fino bruto aue dar el NCCS000 AQ0MP00 VMIF3S0 DA0FS0 NCFS000 DA0FS0 NCFS000 SPDAOMPO AQ0MP00 PROCNO0 VMIP3S0 DA0FS0 0.98926 0.638702 0.98926 0.999961 0.992728 0.97619 0.998555 0.973684 0.550139 0.98926 bruto lo virgen sangrar de aue dar PP3FSA0 AQ0CS00 PP3FSA0 VMM03S0 NCFS000 **РРЗМРАО** NCMP000 ĊS VMM02S0 PP3FSA0 0.010734 0.361298 0.010734 0.00877193 3.85246e-05 0.0072574 0.0238095 0.449861 0.00144509 0.010734 la la sangrar NCMS000 NCMS000 VMSP3S0 NCMP000 NCMS000 6.20105e-6.20105e-06 6.20105e-06 0.00877193 1.44858e-05 sangrar VMSP1S0 0.00877193

Sentence 1 Animará la virgen tierra la sangre de los finos brutos que da la pecuaria Inglaterra tierra el animar virgen sangre de el fino bruto que dar el pecuario inglaterra **DAOMPO** VMIF3S0 DA0FS0 AQ0CS00 NCFS000 DA0FS0 NCFS000 AQ0MP00 AQ0MP00 **PROCNOO** VMIP3S0 DA0FS0 AQ0FS00 NP00000

Morphological Analysis

PoS Tagging

Shallow Parsing Full Parsing

ContaWords

http://contawords.iula.upf.edu/more_info

Bigramas y su frecuencia:

Α	В	С	D	E
4	grito	N	sagrado	Α
	gran	Α	Dios	N
1	último	Α	vigor	N
1	íntimo	A	eslabona	V
1	ínclitos	A	nombra	V
1	ímpetu	N	exterior	A
1	áureas	Α	alegra	V
- 4	Andread	k.i	In a consideration	A

Podemos ver las combinaciones más frecuentes de dos palabras

I VINCUIO	IN	mumo	A
1 vuelo	N	inspirado	V
1 voto	N	cordial	A
1 voluntad	N	extirpe	V
	NUITA DOL	Barcelona	

ContaWords

http://contawords.iula.upf.edu/more_inf

0



Bigramas

Podemos ver la combinación para una determinada palabra

	tonos	N	distintos	A
1	tigres	N	marciales	Α
1	tierras	N	lejanas	Α
1	tierra	N	sagrado	Α
1	tierra	N	libre	Α
1	tierra	N	labrada	V
1	tierra	N	feraz	Α
1	tierra	N	austral	Α
1	tierra	N	abierta	٧
1	testa	N	imperiosa	Α
1	tesoro	N	año	N
1	terruño	N	convierte	V



ContaWords

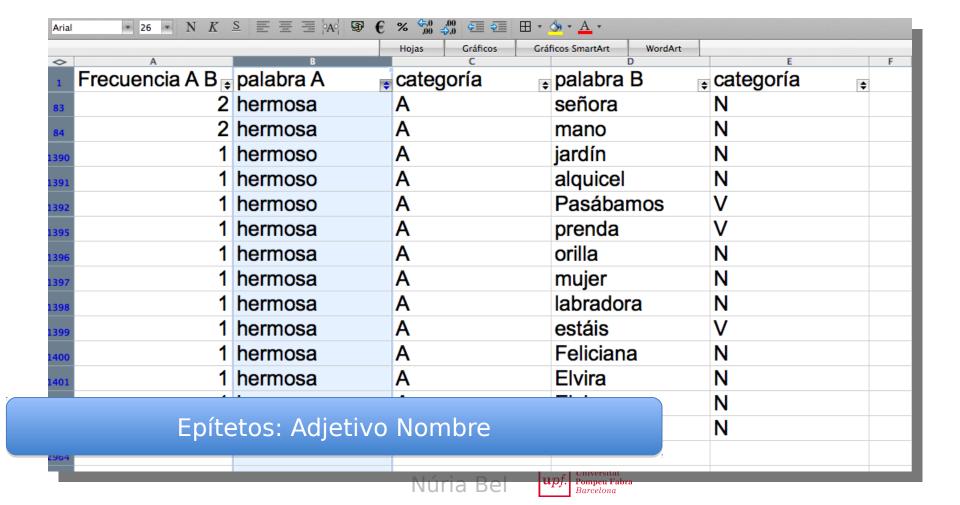
Con Excel (o similar) se pueden hacer "Filtros" para seleccionar vistas de los datos: Para ver las combinaciones de N-ADJ más frecuentes

- insertar fila (que sea la primera!)
- buscar la opción "Filtro" en el Menú Datos y marcar "Autofiltro"
- seleccionar: las combinaciones N-Adj o Adj-N o todas las palabras que combinan con ... hermosa?

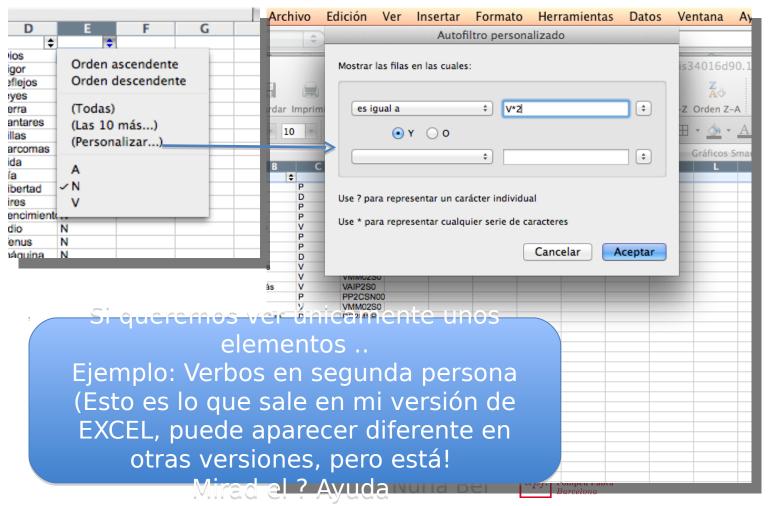




ContaWords Filtros hoja Excel ra inspeccionar datos



ContaWords Autofiltro para hacer búsquedas



Ejercicio 4 (1/2)

- Utiliza ContaWords en: http://contawords.iula.upf.edu
- Sigue las instrucciones para subir el archivo que has guardado antes como .txt, utf8 y seleccionar la lengua.
- Encuentra y acciona el comando "Ejecutar".
- Descarga el resultado cuando aparezca en pantalla. Puede tardar un poco.
- Cambia el nombre al documento Excel cuando lo guardes en tu ordenador.



Ejercicio 4 (2/2)

- Inspecciona los resultados del análisis de ContaWords y responde a estas preguntas:
 - En la pestaña de nombres, ¿por qué están todos en singular?
 - Encuentra cuáles son las 10 palabras más frecuentes del texto ¿Tienen alguna característica en común?
- Crea un Filtro en la pestaña "G-freq-bigram" siguiendo las instrucciones de las transparencias del curso.
- Con el filtro creado, selecciona las opciones para que se muestren únicamente las combinaciones de Adjetivo-Nombre, y las de Nombre-Adjetivo.



Herramientas de análisis http://voyant-tools.org/

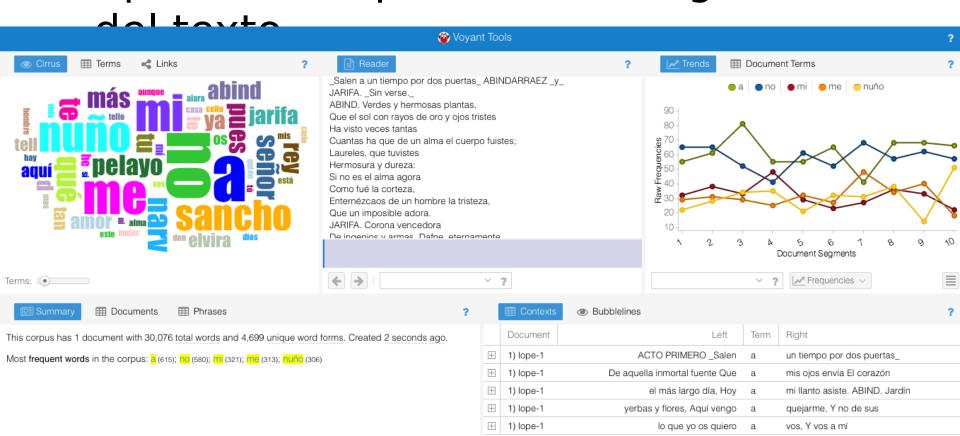
- es una aplicación web
- tiene funciones interesantes





Otra herramienta útil Voyant Tools

- Visualización palabras en nube
- Aparición de palabras en segmentos



voyant-tools.org



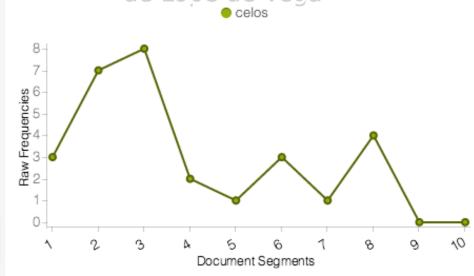
Funciones similares, pero algunas cosas nuevas: frecuencias de palabras y por segmentos del texto

Aparición de la palabra "celos" en El remedio en la desdicha, de Lope de Vega

Antes me ofende y resfría.

JARIFA. No es justo que en el amor,
Abindarráez, tan justo
De hermanos, halles disgusto,
Siendo el más limpio y mejor.
Amor que celos no sabe,
Amor que pena no tiene,
A mayor perfeción viene,
Y a ser más dulce y suave.
Quiéreme bien como hermano:
No te aflijas ni desueles,
Sique el camino que sueles.

mi X

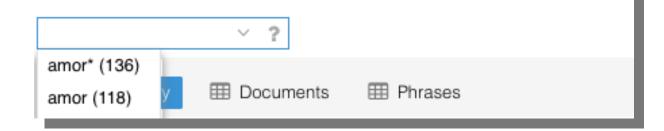


voyant-tools.org

Voyant Tools

Word Tree permite visualizar las secuencias de palabras (bigramas)





Atención!! palabras más frecuentes

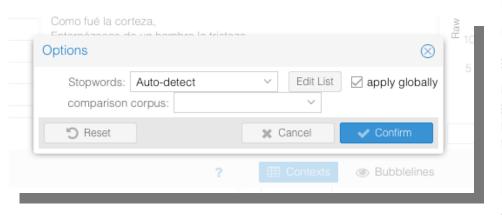


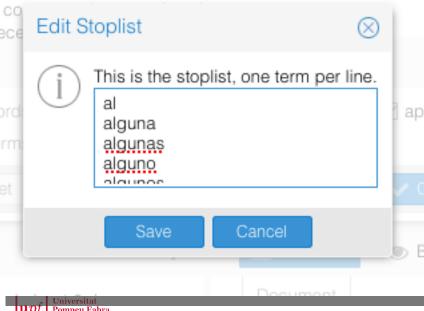
This corpus has 1 document with 30,076 total words and 4,699 unique word forms. Created 2 seconds ago.

Most frequent words in the corpus: a (615); no (580); mi (321); me (313); nuño (306)

Es posible eliminar palabras poco interesantes, como las preposiciones,

las negaciones, pronombres. Se conocen como "palabras vacías" o "storeca Edit Stoplist words".





Más herramientas!!

Descúbrelas! TAPoR3, CLARIN.EU, y

más en el sicuiente curso:

About Tour Contact Tools Lists Useful links Logir TAPoR 3 Discover research tools for studying texts. Borja Navarr Vovant 2.0 Terms Radio shows the high frequency terms scrolling by as if broadcast. You can select the words to watch rise and fall. You can control the speed and number of Home About

▼ Services

▼ Events News Contact CLARIN Call for Papers: CLARIN Annual Conference 2016 CALL FOR PAPERS Call for Papers for the 5th CLARIN Annual Read more > Search for Language Resources Deposit your resources Featured Resource Search in the Virtual Language Czech National Corpus Observatory for language resources in the CLARIN repositories:

Resumen y conclusiones (1/3)

- ✓ Procesamiento del Lenguaje Natural, una tecnología para el análisis de datos textuales, también para la filología (y la poesía).
- En principio, para grandes cantidades de datos, pero ... también para acelerar el análisis de cualquier texto.
- Cualquier texto, pero todavía hay que estar seguros de:
 - √ codificación de caracteres (utf8)
 - √ formato del documento (no pdf, mejor .txt)
- ✓ Actualmente, la digitalización de textos literarios no prioriza la creación de textos procesables.



Resumen y conclusiones (2/3)

- ✓ Anotación: añadir información con etiquetas específicas a una tarea. Para cada herramienta hay que buscar información sobre las etiquetas.
- ✓ Anotación: hay que verificar que el texto está escrito en la lengua que está prevista por la herramienta: palabras, ortotipografía, versos y mayúsculas, puede dar problemas ...
- ✓ El PLN para poesía podría dar herramientas más adecuadas: se necesitan recursos específicos (diccionarios, analizadores, etc.) para este género literario.

Resumen y conclusiones (3/3)

Dificultades

- ✓ Identificación y localización de herramientas, saber qué hay ...
- ✓ Instalación de herramientas, no apto para ...
- ✓ Comprender su "lenguaje": etiquetarios, etc.
- ✓ Comprender sus limitaciones: qué lenguas, cobertura histórica, etc.
- ✓ Función de infraestructuras de investigación



Infraestructura Lingüística



CATALOGO

HERRAMIENTAS

AREAS +

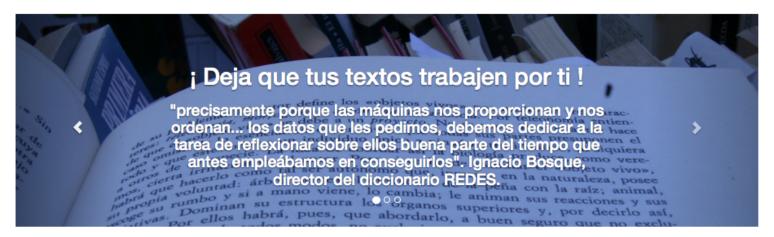
Centro-K CLARIN

Contacto +

AREAS+

Idioma +

El Centro de Competencias CLARIN del IULA-UPF, con el lema "Deja que tus textos trabajen por ti", tiene la misión de promover y asesorar el uso de tecnología y herramientas de análisis de textos en la investigación en Humanidades y Ciencias Sociales.













Centro de competencias UPF-CLARIN http://www.clarin-es-lab.org /

El análisis depende de la disponibilidad de "recursos lingüísticos", listas de palabras y textos con información explícita. Todos hemos de contribuir. ¿Tienes palabras? ¿Tienes textos? Escríbeme a nuria.bel@upf.edu

Gracias!