

*IX Congreso Internacional de la Asociación Asiática de Hispanistas*

*Universidad de Chulalongkorn Bangkok, Tailandia, del 22 al 24 de enero de 2016*

## **Tratamiento lingüístico y matemático de textos digitales españoles**

### **Presentación del Programa LEXIS-web**

**Hiroto Ueda**

**Universidad de Tokio**

# 1. Introducción

LETRAS-web (Tokio): <http://lecture.ecc.u-tokyo.ac.jp/~cueda/letras/>

LETRAS-web (Madrid): <http://shimoda.llf.uam.es/letras/>

NUMEROS-web (Tokio): <http://lecture.ecc.u-tokyo.ac.jp/~cueda/numeros/>

NUMEROS-web (Madrid): <http://shimoda.llf.uam.es/numeros/>

LEXIS-web (Tokio): <http://lecture.ecc.u-tokyo.ac.jp/~cueda/lexis/>

LEXIS-web (Madrid): <http://shimoda.llf.uam.es/lexis/>

Google: Hiroto Ueda

Google: Hiroto Ueda Letras

## 2. Interfaces de input y output

The screenshot displays the LEXIS-web interface, titled "LEXIS-web: Programa para análisis léxico de español". The interface is divided into several sections:

- Left Sidebar:** Contains settings for "Idioma" (set to "Español"), "Página de output" (set to "Esta página"), and links for "Cómo utilizar [PDF]" and "Cómo citar [PDF]". It also lists "LETRAS-web" and "NUMEROS-web", a "-Reiniciar-" button, and "COLABORACIÓN: LLI-UAM (Universidad Autónoma de Madrid)" with a small image of a building.
- [1] Input: Textbox-1:** A section for inputting text. It shows a list of textboxes (Textbox-1 to Textbox-7) and a "Total" dropdown menu with options: "Palabra", "C.S.(palabra)", "Palabra+C.S.", "Lema", "C.S.(lema)", and "Lema+C.S.". Below this is a label "(a) Archivo:" and a "Seleccione:" label.
- [2] Casilla de texto: Textbox-1:** A large text area containing a paragraph of text: "Alicia estaba empezando a cansarse de estar sentada en la orilla del río y no tener nada que hacer:había echado una ojeada con disimulo dos o tres veces el libro que estaba leyendo su hermana,pero no tenía ilustraciones o diálogos, y ¿de qué sirve un libro sin ilustraciones ni diálogos? Así que estaba considerando (hasta donde podía porque el calor la adormecía y atontaba),si valdría la pena levantarse y coger margaritas para hacerse una guirnalda, cuando de repente,un conejo blanco de ojos rosados pasó corriendo a su lado."
- [3] Output:** A section for output options, featuring a checkbox labeled "Archivo de texto (TSV)".
- Bottom:** A "EJECUTAR" button and a footer: "Hirotto Ueda, Universidad de Tokio (ver. 2016.1.8)".

**Fig. 1. Input de LEXIS-web**

## LEXIS-web: Program for lexical analysis of Spanish

Input: Textbox-1; Execution time: 0.407 s. /

Output lines: 107 /

Op	Palabra	C.S.(palabra)	Lema	C.S. (lema)	Homógrafos	Total	Prob.	Ip
1	Alicia	Xant	Alicia	Xant	1:Xant	-	-	1
2	estaba	Estar:IndImp13	estar	Estar	1:Estar	-	-	1
3	empezando	Ger	empezar	Inf	1:Ger	-	-	1
4	a	Prep	a	Prep	1:Prep	-	-	1
5	cansar	Inf	cansar	Inf	1:Inf	-	-	1
6	+se	Clit:3	se	Clit	1:Clit	-	-	1
7	de	Prep	de	Prep	1:Prep	-	-	1
8	estar	Estar	estar	Estar	1:Estar	-	-	1
9	*sentada	PP:fs	sentar	Inf	2:Sus/PP	100	0.556	1
10	en	Prep	en	Prep	1:Prep	-	-	1
11	*la	L:fs	el	L	2:L/Clit	180	0.643	1
12	*orilla	Sus:fs	orilla	Sus	2:Sus/V	140	0.933	1

**Fig. 2. Output de LEXIS-web**

<b>Op</b>	Número secuencial de output
<b>Palabra</b>	Palabra separada del texto, por ejemplo <i>del</i> se separa en <i>de +el</i> .
<b>C.S.(palabra)</b>	Categoría sintáctica de la palabra con informaciones gramaticales, por ejemplo, en el caso de <i>empezando</i> , «Ger» es abreviación de gerundio.
<b>Lema</b>	Forma representante, por ejemplo el lema de <i>empezando</i> es <i>empezar</i> , es decir, la forma canónica de una entrada de un diccionario.
<b>C.S.(lema)</b>	Categoría sintáctica del lema <i>empezar</i> «Inf», que es infinitivo
<b>Homógrafos</b>	Número y categorías posibles, por ejemplo <i>la</i> : «L» (forma femenina singular del artículo <i>el</i> ) «Clit» (forma femenina singular del clítico <i>lo</i> ).
<b>Total</b>	Frecuencia correspondiente.
<b>Prob(abilidad)</b>	Probabilidad dentro de las combinaciones calculadas en el mismo contexto.
<b>Ip.</b>	Número secuencial de línea de input.

## LEXIS-web: Programa para análisis léxico de español

Input: Textbox-1; Tiempo de ejecución: 0.476 s. /

Líneas de output: 77 /

Op ↕	Palabra ↕	F. A. ↕	P. M. U. ↕
1	Alicia	1	9.346
2	estaba	3	28.037
3	empezando	1	9.346
4	a	2	18.692
5	cansar	1	9.346
6	+se	3	28.037
7	de	5	46.729
8	estar	1	9.346

**Fig. 3. Output de LEXIS-web (2)**

### 3. Identificación léxica

Diccionario:

<i>caja</i>	Sus:fs
<i>cajamarca</i>	Xtop
<i>cajero</i>	Sus:ms
(...)	(...)

Abreviación (Abr.):

<b>Abrev.</b>	<b>Explicación</b>	<b>Ejemplo</b>
Adj	Adjetivo	<i>alto, interesante</i>
Adv	Adverbio	<i>abajo</i>
Clit	Clítico	<i>me, te, se, ..., lo, le, ...</i>
Comp	Comparativo	<i>más, menos</i>
Conj	Conj	<i>aunque, como, ...</i>

Det.dem	Determinante demostrativo	<i>este, ese, aquel</i>
Det.ind	Det.indefinido	<i>algún</i>
Det.pos	Det.posesivo	<i>mi, tu, su, ...</i>
Estar	Verbo <i>estar</i>	<i>estar</i>
Ger	Gerundio	<i>estando</i>
Haber	Verbo <i>haber</i>	<i>haber</i>
Inf	Verbo en infinitivo	<i>dar</i>
Int	Interjección	<i>hola, adiós, ...</i>
L	Artículo definido EL	<i>el (los, la, las, lo)</i>
Num	Numeral	<i>0, 1, 2, ..., uno, dos, ..., i, ii, ...</i>
Paren	Paréntesis	<i>( ) &lt; &gt; { } [ ] « »</i>
PP	Participio pasado	<i>estado</i>
Prep	Preposición	<i>a</i>
Pro.dem	Pronombre demostrativo	<i>aquel</i>



Pro.ind	Pronombre indefinido	<i>algo</i>
Pro.pers	Pronombre personal	<i>él</i>
Pro.prep	Pronombre prepositivo	<i>mí, ti, sí</i>
Punt	Puntuación	<i>. , : ; - ¿ ? ¡ !</i>
Q.adj	Interrogativo adjetival	<i>cuál</i>
Q.adv	Interrogativo adverbial	<i>cómo</i>
Q.pro	Interrogativo pronominal	<i>cuál</i>
Rel.adj	Relativo adjetival	<i>cuanto</i>
Rel.adv	Relativo adverbial	<i>cuando</i>
Rel.pro	Relativo pronominal	<i>cual</i>
S n	Sí o no	<i>sí, no</i>
Ser	Verbo <i>ser</i>	<i>ser, soy, eres, es, ...</i>
Signo	Signo	<i>#, \$, %, &amp;, +, -, =, *, /, ...</i>
Sus	Sus	<i>hombre, mujer, animal</i>

U	Artículo definido UN	<i>un (una, unos, unas)</i>
Xant	Xant	<i>abraham</i>
Xtop	Xtop	<i>cajamarca</i>
Y O	Y O	<i>y (e), o (ó, u)</i>

## 4. Separación:

KEY	ITM
al	a +el/a/Prep
ándola	ando +la/ar/.
ándolas	ando +las/ar/.
ándole	ando +le/ar/.
ádoles	ando +les/ar/.
ándolo	ando +lo/ar/.
(...)	(...)
rte	r +te/r/Inf

## 5. Lematización y asignación gramatical

<b>Abrev.</b>	<b>Explicación</b>	<b>Ejemplo</b>
«ms»	masculino singular	<i>libro</i>
«mp»	masculino plural	<i>ambos</i>
«cs»	común singular	<i>estudiante</i>
«V»	Verbo conjugado	<i>voy, comeremos</i>
«PP»	Participio pasado	<i>ido, comido</i>
«Ger»	Gerundio	<i>yendo, comiendo</i>
«Ind»	Indicativo	<i>sé, sabes</i>
«Sub»	Subjuntivo	<i>sepa, sepas</i>
«Fut»	Futuro	<i>sabré</i>
«Cond»	Condicional	<i>sabría</i>

«Pres»	Presente	<i>sé, sepa</i>
«Imp»	Imperfect	<i>sabía</i>
«Pas»	Pasado	<i>supe, supiera</i>
«1»	Primera persona	<i>yo, sé</i>
«2»	Segunda persona	<i>tú, sabes</i>
«3»	Tercera persona	<i>usted, él, ella, sabe</i>
«4»	Cuarta persona	<i>nosotros, nosotras, sabemos</i>
«5»	Quinta persona	<i>vosotros, vosotras, sabéis</i>
«6»	Sexta persona	<i>ustedes, ellos, ellas, saben</i>

KEY	ITM
a	/((Adj).*/\$1:fs#ó(Adj Det.pos Pro.ind Rel Q).*/\$1:fs#ar/(Inf:v Inf:r)/V:IndPres3#er/(Inf:v Inf:r)/V:SubPres13#ir/(Inf:v Inf:r)/V:SubPres13
á	/Inf/V:Fut3
aba	ar/(Inf ESTAR)/V:IndImp13
abais	ar/(Inf ESTAR)/V:IndImp5
ábamos	ar/(Inf ESTAR)/V:IndImp4
(...)	(...)

"-a" final de palabra:

(1) /(Adj).\*/\$1:fs

(2) o/(Adj|Det.pos|Pro.ind|Rel|Q).\*/\$1:fs

(3) ar/(Inf:v|Inf:r)/V:IndPres3

(4) er/(Inf:v|Inf:r)/V:SubPres13

(5) ir/(Inf:v|Inf:r)/V:SubPres13

## 6. Desambiguación

*sentada, la, orilla, etc:*

Secuencia	Frec.
Adj-Adj	60
Adj-Adv	30
Adj-Clit	60
Adj-Conj	60
(...)	(...)
Y o-Sus	80
Y o-U	50
Y o-V	60
Y o-Xant	70
Y o-Xtop	70



... *en la orilla*... Prep - {L/Clit} - Sus/V

(1) Secuencia anterior: *en {la}*

«Prep - {L}» (90)

«Prep - {Clit}» (0)

(2) Secuencia posterior: *{la} orilla*

«{L} - Sus» (90)

«{L} - V» (0)

«{Clit} - Sus» (10)

«{Clit} - V» (90)

=> {L}: 90 + 90 + 0 = 180 / {Clit}: 0 + 10 + 90 = 100

{L} Total: 180; Probabilidad 180 / (180 + 100) = .643

## 7. Final

Precisión, cantidad, rapidez, comodidad (sencillez)

98% de precisión, 10.000 palabras, menos de 5 segundos, ?

## Referencias:

- Almela, Ramón / Cantos, Pascual / Sánchez, Aquilino / Sarmiento, Ramón / Almela, Moisés. (2005). *Frecuencias del español. Diccionario y estudios léxicos y morfológicos*. Madrid : Universitas.
- Ávila Muñoz, Antonio Manuel. (199). *Léxico de frecuencia del español hablado en la ciudad de Málaga*. Málaga: Universidad de Málaga.
- Bull, William E. (1947). "Modern Spanish verb-form frequencies", *Hispania*, 451-466.
- Bybee, Joan. (2003). "Mechanisms of change in grammaticalization: the role of frequency", en Brian D. Joseph and Richard D. (eds.), *The Handbook of historical linguistics*. Oxford: Blackwell, 602-623.
- Company Company, Concepción. (2004). “¿Gramaticalización o desgramaticalización? Reanálisis y subjetivización de verbos como marcadores discursivos en la historia del español”, *Revista de Filología Española*, 84, 29-66.

- Davies, Mark. (2006). *A frequency dictionary of Spanish. Core vocabulary for learners*. New York: Routledge.
- García Hoz, Víctor. (1953). *Vocabulario usual, vocabulario común y vocabulario fundamental*. Madrid: Consejo Superior de Investigaciones Científicas.
- Gómez Díaz, Raquel. (2005). *La lematización en español. Una aplicación para la recuperación de información*. Gijón: Trea .
- Hopper, Paul J. / Traugott, Elizabeth Closs. (2003). *Grammaticalization*, 2nd ed. Cambridge: Cambridge University Press.
- Jiménez Juliá, Tomás. (2006). *El paradigma determinante en español. Origen nominativo, formación y características*. Verba, anexo 56, Santiago de Compostela: Universidade de Santiago de Compostela .
- Juilland, Alphonse / Chang-Rodríguez, Eugenio. (1964). *Frequency dictionary of Spanish words*. The Hague: Mouton .

- Lieberman, Erez / Michel, Jean-Baptiste / Jackson, Joe; Tang, T. / Nowak Martin A. (2007). "Quantifying the evolutionary dynamics of language", *Nature*, vol. 449, 713-716.
- Moreno Sandoval, Antonio (2014). "Desafíos de y para la lingüística de corpus", *Estudios Lingüísticos Hispánicos*, (Círculo de Estudios Lingüísticos Hispánicos de Tokio) 29, 69-85.
- Moreno Sandoval, Antonio / Guirao Miras, José María. (2008). "Frecuencia y distintividad en el uso lingüístico: casos tomados de la lematización verbal de corpus de distintos registros", *Actas del I Congreso Internacional de Lingüística de Corpus (CILC-09)*, Murcia: Universidad de Murcia. 195-210.
- Pagel, Mark. / Atkinson, Quentin D. / Meade Andrew. (2007). "Frequency of word-use predicts rates of lexical evolution throughout Indo-European history", *Nature*, 449, 717-720.

Ueda, Hiroto. (2015). "Frecuencia contrastiva, frecuencia ponderada y método de concentración. Aplicación al estudio de las dos formas prepositivas del español medieval «pora» y «para»", *Actas del IX Congreso Internacional de Historia de la Lengua Española (Cádiz, 2012)*, Madrid: Iberoamericana, 1139-1155.

Ueda (en prensa). "Analizador lingüístico común con reglas gramaticales y diccionario, preparados por el usuario: Una aplicación para el análisis tipológico del léxico español".

Ueda, Hiroto / Perea Maria Pilar. (2010). "Método general de lematización con una gramática mínima y un diccionario óptimo. Aplicación a un corpus dialectal escrito", en Moskowich-Spiegel Fandiño, I; Crespo García, B.; Lareo Martín, I.; Lojo, P. (eds.) *Visualización del lenguaje a través de corpus*. A Coruña: Universidade da Coruña , 919-932, .

Ueda, Hiroto / Rubio, Carlos. (2006). *Puerta al español. Nuevo diccionario español-japonés*. Tokio: Kenkyusha.

[Fin]