

Reseña de / Review of / Ressenya de

## **Rojo, Guillermo (2021). *Introducción a la lingüística de corpus en español*. Routledge**

ISBN: 9780367635848 (tapa blanda); 9780367635855 (tapa dura); 9781003119760 (libro electrónico);  
DOI: <http://doi.org/10.4324/9781003119760>

Rogelio NAZAR  
Pontificia Universidad Católica de Valparaíso (Chile)  
[rogelio.nazar@pucv.cl](mailto:rogelio.nazar@pucv.cl)  
<https://orcid.org/0000-0002-8853-1353>

**Palabras clave:** metodología; lingüística de corpus en castellano; lexicología y gramática basadas en corpus; investigación lingüística

*English:*

**Keywords:** *methodology; corpus linguistics in Spanish; corpus-based lexicology and grammar; linguistic research*

*Català:*

**Paraules clau:** *metodologia; lingüística de corpus en castellà; lexicologia i gramàtica basada en corpus; recerca lingüística*

El libro de Guillermo Rojo que aquí se reseña se ubica en una tradición de volúmenes que introducen la metodología de trabajo con corpus, en la que podemos encontrar obras tales como las de Sinclair (1991), Stubbs (1996), McEnery y Hardie (2012) y Brezina (2018), entre otras, pero este tiene la particularidad de estar destinado a un lector hispanohablante. Viene a cubrir una necesidad, ya que no existía un manual para utilizar en cursos de grado o posgrado sobre estos temas, adaptado a las características de la lengua castellana. Será de ayuda en materias relacionadas con la lingüística de corpus pero también con la metodología de investigación lingüística. Sin duda, ejercerá influencia en investigadores de áreas diversas, debido a que las destrezas técnicas que expone son transversales, por ejemplo, a estudios vinculados con el léxico, la terminología y la gramática, ya sea que trabajen con corpus monolingüe o bilingüe, en primera (L1) o segundas lenguas (L2). Se trata de una obra de calidad, redactada con maestría y autoridad por quien posee la capacidad de procesar datos de primera mano. Al mismo tiempo, el libro es también testimonio de cómo ha cambiado la lingüística española en menos de 20 años ya que, por entonces, hablar de algunos de los temas que aquí se discuten podía atraer miradas de extrañeza. Con este libro, muchos conceptos nuevos y su terminología entran al *mainstream* de la lingüística en español.

En los primeros tres capítulos, Rojo presenta los aspectos básicos del trabajo con corpus. En el primero, titulado “La explotación básica de los corpus”, cubre tanto temas técnicos como conceptos lingüísticos fundamentales.

Es aqu3 donde delimita el 3rea de trabajo del libro y prepara al lector para los cap3tulos siguientes. Explica qu3 es la lingüística de corpus, qu3 es un corpus y los distintos tipos de técnicas y tecnologías disponibles en la actualidad para la investigaci3n en este 3mbito. Con el cap3tulo 2, “La lingüística de corpus y la metodolog3a de la investigaci3n lingüística”, la introducci3n a la lingüística de corpus se articula con el tema m3s general de la metodolog3a de la investigaci3n cient3fica. En este lugar se discute sobre cuestiones epistemol3gicas y se dialoga con otras aproximaciones metodol3gicas ya establecidas. Cabe destacar aqu3, adem3s, el enfoque hist3rico en el desarrollo de la metodolog3a lingüística. En el cap3tulo 3, “Diseño, construcci3n y explotaci3n de corpus”, incluye aspectos ya un poco m3s operativos como la anotaci3n, codificaci3n y an3lisis de datos. El cap3tulo incluye tambi3n aspectos legales y 3ticos que, tal como ya hab3an seÑalado McEnery y Hardie (2012), conviene considerar. Se declara como un manual orientado a la docencia, y por ello en los cap3tulos iniciales incluye algunos temas ya transitados, como las dos culturas de Snow (1959), la enciclopedia china de Borges (1952) y la clasificaci3n de las ciencias de Bunge (1972).

En los cap3tulos siguientes, el libro comienza a profundizar en los temas centrales de la investigaci3n basada en corpus. En el cap3tulo 4, con el t3tulo “Recuperaci3n de informaci3n contenida en corpus textuales: el l3xico”, se trabaja primero el concepto de palabra y se destaca la importancia del estudio de la frecuencia l3xica, un aspecto de la lexicolog3a que ha atra3do cada vez m3s atenci3n desde que la lingüística de corpus se consolid3 como metodolog3a de trabajo. Un campo que se desarrolla ampliamente en este cap3tulo es el estudio de la variaci3n del vocabulario, desde los puntos de vista diat3pico, diacr3nico, diastr3tico y diaf3sico. Posiblemente, algunos lectores encontrar3n que un aspecto menos desarrollado es el an3lisis sem3ntico basado en corpus. Menciona el tema por ejemplo al presentar el estudio de la coaparici3n en el an3lisis de la polisemia. Explica el caso de *bombilla*, que aparecer3 junto a distintas palabras dependiendo de que se use con el sentido de foco o como la bombilla del mate, etc. Sin embargo, se trata de un tema complejo que podr3a haberse desarrollado de manera m3s extensiva.

El cap3tulo 5, “Recuperaci3n de informaci3n contenida en corpus textuales: fen3menos gramaticales”, es probablemente el m3s atractivo del libro. Resultar3 3til como gu3a para investigadores interesados en aplicar la lingüística de corpus al estudio de los fen3menos gramaticales, pero sin duda ser3 iluminador tambi3n para los gram3ticos que a3n basan su metodolog3a de trabajo en la introspecci3n. El inter3s radica en que, en lingüística de corpus, el an3lisis gramatical es generalmente m3s complejo que el l3xico, ya que muchos de los fen3menos gramaticales no son directamente observables o presentan mayor dificultad para su medici3n con medios mec3nicos. En este cap3tulo se estudia, por ejemplo, el comportamiento verbal, los adverbios en *-mente*, los fen3menos de concordancia, la adaptaci3n de pr3stamos, las relaciones sint3cticas y algunos aspectos de la variaci3n diacr3nica de la gram3tica espaÑola, entre otra diversidad de temas que pueden ser muy estimulantes particularmente para investigadores que se encuentren en las primeras etapas de su carrera y est3n a3n definiendo su l3nea de trabajo o tema de investigaci3n.

El cap3tulo tambi3n dedica atenci3n al trabajo con el espaÑol como L2. Describe técnicas que pueden ser de utilidad para la docencia, ya que pueden ser fuente de inspiraci3n para proponer ejercicios a los estudiantes y de esa

forma llevar el corpus al aula. Interesará particularmente a los lectores de la revista *TEISEL* porque formula preguntas y expone una amplia variedad de posibilidades de investigación en L2. Presenta los usos de corpus formados por producciones orales o escritas de aprendientes de español como L2 y explica, por ejemplo, ya sea la relación entre las variables nivel de conocimiento de la L2 y la riqueza de vocabulario o la frecuencia de aparición de secuencias gramaticalmente incorrectas, como artículo + demostrativo + sustantivo, etc., o bien la relación entre las variables L1 y L2 en la interferencia lingüística. En este último caso, se observan regularidades en el desempeño de los aprendientes de castellano en función de su L1, en cuanto a características y tipos de error, tal como la selección incorrecta de preposiciones en combinación con verbos.

En el capítulo 6, “Otras cuestiones centrales en lingüística de corpus”, el autor aborda una serie de reflexiones generales sobre la lingüística de corpus, entre las que cabe destacar, nuevamente, los antecedentes históricos de la disciplina. Es también el lugar en el que presenta la discusión de algunos aspectos de lingüística cuantitativa y la estructura estadística de la lengua. Aquí establece, además, la relación con conceptos generales de la estadística como el muestreo, ya que explica los problemas de la determinación del tamaño de la muestra, su representatividad y equilibrio. Presenta también una reflexión en profundidad sobre el presente y el futuro de la lingüística de corpus, así como su diálogo con otras ramas lingüísticas y otras disciplinas.

El capítulo 7, por último, se titula “Herramientas de recuperación de datos: resumen y ampliación” y está dedicado a cuestiones más avanzadas de informática, como el uso de expresiones regulares y algunos comandos de Linux. El capítulo no incluye un desarrollo en profundidad de los aspectos más técnicos como el tratamiento estadístico de datos, medidas de evaluación, aplicación de etiquetadores morfosintácticos ni conceptos básicos de programación. Esto se debe seguramente al tipo de público al que va dirigida la obra, ya que la idea es que pueda comprenderse sin necesidad de ser experto en lingüística computacional. Sin embargo, tal vez hubiese sido deseable aprovechar la oportunidad de presentar un capítulo algo más desafiante para este tipo de lector.

Posiblemente, un mayor desarrollo de conceptos de programación, más allá de lo que explica sobre AWK (Aho, Kernighan y Weinberger, 1988), hubiese sido beneficioso, ya que esta habilidad representa una forma de alfabetización y de acceso al pensamiento lógico. Más allá de la cuestión meramente operativa, resulta conveniente tener la capacidad de pensar con símbolos y lenguajes formales. En este ámbito, además, la cuestión del sistema operativo no es solo un aspecto técnico: un sistema operativo es prácticamente una ideología, y como tal puede inhibir la libertad y la creatividad. Cierta tipo de herramientas, más allá de sus interfaces gráficas amigables, no ofrecen al investigador ni la autonomía ni la potencia de procesamiento de lenguajes como Perl o Python que, a la vez, son simples en comparación con otros. Rojo no ignora este hecho, y lo explicó en una charla en Barcelona allá por 2010, en la que intentó convencer a la audiencia asegurando que “nadie se ha muerto por aprender un poquito de Perl”. La gente se reía pero la cuestión es bastante seria. Las habilidades de programación ubican al investigador en una situación distinta a la

de quien utiliza interfaces que solo permiten realizar las operaciones decididas por quien diseñó el programa y que, además, tienden a orientarse hacia el trabajo cualitativo.

Otra cuestión que podría echarse a faltar en el libro ya en general es un mayor énfasis en el contraste entre métodos cuantitativos y cualitativos. En lingüística de corpus es muy frecuente el enfoque cualitativo, y el corpus suele utilizarse más como medio de comprobación de teorías que como medio de descubrimiento, aunque el tema no se discute mayormente. Por ejemplo, cuando en el libro se estudia la alternancia *cocodrilo/crocodilo*, resulta necesario saber, de manera previa, que tal alternancia existe, para buscar cada una de las formas y observar su frecuencia de aparición en el corpus. Desde el punto de vista cuantitativo, en cambio, cabe otro tipo de preguntas y métodos. Es posible, por ejemplo, encontrar estas alternancias aplicando medidas de similitud ortográfica o semántica. Otro ejemplo similar se da en la descripción que ofrece de la interfaz de búsqueda de enigramas de Google Books. La interfaz permite buscar palabras para obtener curvas de distribución de frecuencias en la serie temporal. Sin embargo, lo verdaderamente interesante del recurso es descargar la totalidad del corpus, ya que entonces las posibilidades se multiplican. Con métodos cuantitativos y la utilización de un lenguaje de *scripting* es posible, por ejemplo, proceder al revés, y así introducir curvas para obtener las palabras, claramente un resultado que no puede obtenerse con el método cualitativo.

Es evidente que un libro no puede desarrollar todos los temas con la misma profundidad y siempre habrá algunos que sacrificar. Así, debe valorarse por lo que es y no por lo que podría haber sido, y la valoración es muy positiva. El lector modelo es un individuo con formación en humanidades, sin un nivel avanzado de alfabetización tecnológica. Pero también podría ser útil para otra audiencia constituida por el grupo diametralmente opuesto, es decir, investigadores con formación en informática que tengan interés en dedicarse al procesamiento del lenguaje natural y necesiten fortalecer sus conocimientos de lingüística. Cualquiera sea el caso, se trata de un manual útil para la lingüística de corpus en español y también para la docencia en lingüística general. Dada su selección y organización de temas, puede adaptarse a distintas necesidades, pero podría ser especialmente útil para vertebrar el programa de un curso de lingüística de corpus de un semestre. Proporciona la base que permitirá luego a cada cual elegir un recorrido de lecturas. Para el desarrollo de destrezas técnicas, se requerirá inevitablemente de un periodo de ejercitación práctica. Para la formación avanzada en lingüística cuantitativa, el libro puede complementarse con otras lecturas como el manual de Manning y Schütze (1999) o el de Herdan (1964). Para finalizar, cabe señalar que esta obra es un recordatorio del privilegio de los lingüistas de la época actual, que disponen de tantas herramientas y posibilidades.

## Referencias

- Aho, Alfred; Kernighan, Brian; Weinberger, Peter (1988). *The AWK programming language*. Addison-Wesley Publishing Company.
- Borges, Jorge Luis (1952). El idioma analítico de John Wilkins. En *Otras inquisiciones* (pp. 121-125). Sur.

- Brezina, Vaclav (2018). *Statistics in corpus linguistics: A practical guide*. Cambridge University Press. <https://doi.org/10.1017/9781316410899>
- Bunge, Mario (1972). *La ciencia: su método y su filosofía*. Siglo veinte.
- Herdan, Gustav (1964). *Quantitative linguistics*. Butterworths.
- Manning, Christopher; Schütze, Hinrich (1999). *Foundations of statistical natural language processing*. MIT Press.
- McEnery, Tony; Hardie, Andrew (2012). *Corpus linguistics: Method, theory and practice*. Cambridge University Press.
- Sinclair, John (1991). *Corpus, concordance, collocation*. Oxford University Press.
- Snow, Charles Percy (1959). *The two cultures*. Cambridge University Press.
- Stubbs, Michael (1996). *Text and corpus analysis*. Blackwell.

