

Mark Davies

Uso del *Corpus del Español* y los corpus relacionados para la lexicografía histórica española

Resumen: El *Corpus del Español* original de 100 millones de palabras (siglos XIII–XX) permite a los investigadores llevar a cabo investigaciones avanzadas que incluyen colocativos, n-gramas, sinónimos, así como comparaciones a través de períodos históricos, que no son posibles con otros corpus como *CORDE*. Además, los datos del siglo XX permiten comprender la variación en el léxico basada en el género textual, lo que no es posible con otros corpus como *CREA*. En 2016, el *Corpus del Español* se amplió considerablemente para incluir dos mil millones de palabras de 20 países de habla hispana, lo que permite que la investigación sobre la variación regional del léxico sea mucho más profunda que con otros corpus como *CORPES*. Finalmente, el próximo corpus *NOW-Spanish* (que se lanzará a finales de 2018), con seis mil millones de palabras en español, es completamente original en la forma en que los investigadores pueden rastrear los cambios léxicos en curso (incluso semana por semana).

Palabras clave: Corpus Léxico, Histórico, Géneros, Dialecto

Abstract: The original 100 million word *Corpus del Español* (1200s–1900s) allows researchers to carry out advanced research involving collocates, n-grams, synonyms—as well as comparisons across historical periods—in ways that are not possible with other corpora such as *CORDE*. In addition, the data from the 1900s allows insight into genre-based variation in lexis, which is not possible with other corpora like *CREA*. In 2016 the *Corpus del Español* was greatly enlarged to include two billion words of data from 20 Spanish-speaking countries, and it allows research on regional variation in lexis that is much more powerful than with other corpora like *CORPES*. Finally, the upcoming *NOW-Spanish* corpus (to be released in late 2018) is completely unique in the way that researchers can track ongoing changes in lexis (even week-by-week) in a six billion-word corpus of Spanish.

Keywords: Corpus, Lexis, Historical, Genres, Dialect

1 Introducción

Hasta hace unos veinte años, la lexicografía histórica solía llevarse a cabo mediante minuciosos y laboriosos análisis de textos impresos. Sin embargo, a partir de la llegada de grandes corpus en línea, la tarea de quienes estudian el

cambio y la variación léxica se ha convertido en algo mucho más fácil. Los investigadores pueden lograr ahora en solo unos minutos lo que antes podría haber necesitado días o incluso semanas.

En este trabajo, consideraremos con cierto detalle cómo el *Corpus del Español* se puede utilizar para examinar la variación y el cambio léxico¹. Nos centraremos en las dos partes del corpus: el corpus de 100 millones de palabras «histórico/basado en el género» que se lanzó en 2002, y el corpus de dos mil millones de palabras «basado en web/dialectal» que se lanzó en 2016. También trataremos brevemente un corpus de cinco mil millones de palabras que se lanzó en 2018, y que permitirá a los investigadores examinar los neologismos y el cambio léxico actual con increíble detalle. A medida que analicemos estos diferentes corpus, también trataremos brevemente los corpus similares de la Real Academia Española: *CORDE* (textos históricos), *CREA* (principalmente textos de finales del siglo XX) y *CORPES* (textos de principios del siglo XXI).

2 Los corpus textuales

El *Corpus del Español* se compone de dos partes diferentes, ambas diseñadas para dar respuesta a diferentes preguntas sobre el cambio y la variación del lenguaje. La primera parte es el corpus «histórico/género». Se completó en 2002 y fue revisado a fines de 2007. Está compuesto por aproximadamente 100 millones de palabras, desde el español antiguo hasta finales del siglo XX, tiene aproximadamente 18 millones de palabras del siglo XIII al XV, 42 millones de palabras del siglo XVI al XVIII, y alrededor de 40 millones de palabras de los siglos XIX y XX. Está compuesto de textos de una amplia gama de géneros, incluidos más de cinco millones de palabras de transcripciones de conversaciones habladas de finales del siglo XX. Para el siglo XX, están divididos en partes iguales entre los siguientes géneros: hablado, ficción, prensa y académico. Los detalles completos sobre cada uno de los casi 14 000 textos se pueden encontrar a través del enlace «Textos» en el sitio web del corpus, y los usuarios pueden descargar un archivo de Excel que enumera todos esos textos.

La segunda parte del *Corpus del Español* fue lanzada en 2016. Está orientada a estudios sobre la variación dialectal en español o que necesitan un corpus mucho más grande que el corpus original de 100 millones de palabras. El corpus contiene más de dos mil millones de palabras, lo que lo hace 100 veces más grande que la porción del siglo XX del corpus original. Además, los usuarios pueden

realizar comparaciones entre los 21 países del corpus (España, México, Colombia, Argentina, etc.) para ver la frecuencia por dialecto y buscar directamente todas las palabras que son más frecuentes en un país que en otro.

En este artículo, comparemos la parte histórica (y basada en el género) del *Corpus del Español* con el *CORDE* (*Corpus diacrónico del español*) y el *CREA* (*Corpus de referencia del español actual*) de la Real Academia Española. A continuación, en las Secciones 3–11, comparemos la funcionalidad de la parte histórica del *Corpus del Español* con el *CORDE*, por lo que daremos una visión general del *CORDE*.

El *CORDE* se creó a finales de los años 90 del siglo XX y fue el primer gran corpus histórico del español. Está compuesto por aproximadamente 250 millones de palabras de texto, con buena representación a lo largo de los diferentes periodos históricos, y un buen equilibrio entre géneros, que incluyen poesía, escritos históricos, literatura, materiales didácticos, etc. Con 250 millones de palabras, es de dos a tres veces más grande que el *Corpus del Español*. Sin embargo, como veremos, el tamaño no es todo. Sin el tipo correcto de arquitectura e interfaz de corpus, los datos textuales están, en esencia, «atrapados» dentro del corpus y no están disponibles para los usuarios finales.

3 Descripción general del uso de corpus para la lexicografía histórica

Los corpus históricos, para ser útiles, deberían permitir a los usuarios realizar investigaciones sobre varios aspectos diferentes del cambio léxico². Estos incluyen lo siguiente:

- Encuentra resultados representativos. En el nivel más básico, los usuarios pueden buscar una palabra o frase, encontrar la primera aparición de la palabra o frase y ver todas las ocurrencias en contexto.
- Frecuencia (más avanzada). Los usuarios pueden ver fácilmente la frecuencia de una palabra o frase a lo largo del tiempo, con frecuencias normalizadas³.
- Frecuencia (aún más avanzada). En lugar de tener que decirle al corpus qué palabras o frases específicas buscar, el corpus puede generar una lista de palabras cuya frecuencia coincide con ciertos criterios, como los nombres que

¹ Para estudios similares sobre el inglés, vid. Davies en prensa a y b.

² Para estudios similares con corpus históricos del inglés, vid. Davies 2012a y 2012b.

³ En otras palabras, frecuencia por millar o por millón de palabras de texto, para tener en cuenta el diferente tamaño del corpus en diferentes periodos históricos.

se incorporaron al léxico en el siglo XVII o todas las palabras que se usan al menos cinco veces más en el siglo XIII que en el XIV

- Forma de palabra (morfológica). Los usuarios deben poder buscar por prefijos, sufijos y raíces, y ver la frecuencia de cada forma coincidente en los diferentes periodos históricos, así como la frecuencia general de todas las formas en cada período histórico.
- Significado de la palabra (semántica): simple. Los usuarios pueden encontrar las colocaciones más frecuentes (palabras cercanas) de una palabra o frase determinada, lo que obviamente proporciona una muy buena comprensión del significado de la palabra. Prácticamente cualquier arquitectura e interfaz de corpus permite a los usuarios ver las palabras cercanas caso por caso, pero los corpus realmente útiles resumen toda esta información sobre las colocaciones para todas las ocurrencias de una palabra o frase determinada.
- Semántico (más avanzado). Suponiendo que el corpus pueda encontrar colocaciones, debería ser posible comparlas a través de períodos históricos o entre diferentes géneros. Los cambios en las colocaciones a lo largo de períodos sirven a menudo como marcadores de cambio semántico.
- Semántica (aún más avanzada). En lugar de simplemente buscar palabras y frases, los usuarios pueden buscar por campo semántico. P. ej., si un repertorio está integrado en el corpus, o si los usuarios pueden crear listas personalizadas de palabras, entonces podrían crear una búsqueda donde cualquier palabra en un campo semántico es parte de la consulta. Un ejemplo de esto podría ser [miembro de la familia] seguido de [sinónimo de *pedir*] seguido de [sinónimo de *limpiar*], o [hora del día] cerca de [sinónimo de *lánguete*]. Del mismo modo, se podría comparar la frecuencia de todas las palabras o frases en un campo semántico completo, y comparar la frecuencia y la distribución de cada miembro a lo largo del tiempo.

En los siguientes apartados, proporcionaré ejemplos concretos de cómo estos dos grandes corpus históricos, el *Corpus del Español* y el *CORDE*, pueden usarse (o no) para investigar la amplia gama de fenómenos enumerados anteriormente. (Para obtener una visión general previa de la funcionalidad del *Corpus del Español* con arquitectura e interfaz más antiguas, v. Davies 2002, 2005a, 2005b, 2008, 2010). Algunos lectores pueden, así, comenzar a obtener una perspectiva completamente nueva de lo que se puede hacer con los corpus históricos. Si han utilizado corpus con arquitecturas e interfaces limitadas, quizá los hayan usado solo para encontrar las ocurrencias de una palabra o frase específica. Sin embargo, una vez que una persona ha utilizado un corpus de los que permiten una amplia gama de consultas, se da cuenta de que hay innumerables temas de lingüística histórica que podrían estudiarse con un corpus completo.

4 Encontrar y mostrar resultados representativos

En el nivel más básico, un corpus debería permitir al investigador buscar una palabra o frase, encontrar la primera ocurrencia de la palabra o frase o ver todos los resultados en contexto (y quizás limitar las ocurrencias a un período histórico determinado). Los programas para realizar tales búsquedas son abundantes y bastante rápidos: 1–2 segundos (en el caso de la búsqueda simple) incluso para un corpus de 100 millones de palabras.

El *CORDE* puede hacer estas búsquedas básicas bastante bien. P. ej., supongamos que el usuario quiere encontrar todas las apariciones de la palabra *braveza*. Después de enviar la búsqueda, el usuario ve que hay 273 ocurrencias en 86 documentos. Al hacer clic en «Obtención de ejemplos», el usuario ve entradas de «Palabra clave en contexto» (KWIC, *Keyword in Context*) como las siguientes:

CONCORDANCIA
 a que de auto dicha es mourendo con saña: o con bravaça o con mal querencia como quiera que pena q1 1491
 lxo el rey malnon ayal en la yta del rey como la bravaça del. Jorpe que q2 1491
 que que las duena por al salasso. Lienamente e syn bravaça ninguna duena. Jor cobrea q3 1491
 que que que q4 1491
 tca esta rey q5 1491
 que q6 1491
 verde e apouena e asin tuerto q7 1491
 tezar al rey deue lo fazer en la bravaça por cobdiçia e por cobdiçia. Mas democantiga q8 1491
 q9 1491
 q10 1491
 q11 1491
 q12 1491
 q13 1491
 q14 1491
 q15 1491
 q16 1491
 q17 1491
 q18 1491
 q19 1491
 q20 1491
 q21 1491
 q22 1491
 q23 1491
 q24 1491
 q25 1491
 q26 1491
 q27 1491
 q28 1491
 q29 1491
 q30 1491
 q31 1491
 q32 1491
 q33 1491
 q34 1491
 q35 1491
 q36 1491
 q37 1491
 q38 1491
 q39 1491
 q40 1491
 q41 1491
 q42 1491
 q43 1491
 q44 1491
 q45 1491
 q46 1491
 q47 1491
 q48 1491
 q49 1491
 q50 1491
 q51 1491
 q52 1491
 q53 1491
 q54 1491
 q55 1491
 q56 1491
 q57 1491
 q58 1491
 q59 1491
 q60 1491
 q61 1491
 q62 1491
 q63 1491
 q64 1491
 q65 1491
 q66 1491
 q67 1491
 q68 1491
 q69 1491
 q70 1491
 q71 1491
 q72 1491
 q73 1491
 q74 1491
 q75 1491
 q76 1491
 q77 1491
 q78 1491
 q79 1491
 q80 1491
 q81 1491
 q82 1491
 q83 1491
 q84 1491
 q85 1491
 q86 1491
 q87 1491
 q88 1491
 q89 1491
 q90 1491
 q91 1491
 q92 1491
 q93 1491
 q94 1491
 q95 1491
 q96 1491
 q97 1491
 q98 1491
 q99 1491
 q100 1491
 q101 1491
 q102 1491
 q103 1491
 q104 1491
 q105 1491
 q106 1491
 q107 1491
 q108 1491
 q109 1491
 q110 1491
 q111 1491
 q112 1491
 q113 1491
 q114 1491
 q115 1491
 q116 1491
 q117 1491
 q118 1491
 q119 1491
 q120 1491
 q121 1491
 q122 1491
 q123 1491
 q124 1491
 q125 1491
 q126 1491
 q127 1491
 q128 1491
 q129 1491
 q130 1491
 q131 1491
 q132 1491
 q133 1491
 q134 1491
 q135 1491
 q136 1491
 q137 1491
 q138 1491
 q139 1491
 q140 1491
 q141 1491
 q142 1491
 q143 1491
 q144 1491
 q145 1491
 q146 1491
 q147 1491
 q148 1491
 q149 1491
 q150 1491
 q151 1491
 q152 1491
 q153 1491
 q154 1491
 q155 1491
 q156 1491
 q157 1491
 q158 1491
 q159 1491
 q160 1491
 q161 1491
 q162 1491
 q163 1491
 q164 1491
 q165 1491
 q166 1491
 q167 1491
 q168 1491
 q169 1491
 q170 1491
 q171 1491
 q172 1491
 q173 1491
 q174 1491
 q175 1491
 q176 1491
 q177 1491
 q178 1491
 q179 1491
 q180 1491
 q181 1491
 q182 1491
 q183 1491
 q184 1491
 q185 1491
 q186 1491
 q187 1491
 q188 1491
 q189 1491
 q190 1491
 q191 1491
 q192 1491
 q193 1491
 q194 1491
 q195 1491
 q196 1491
 q197 1491
 q198 1491
 q199 1491
 q200 1491
 q201 1491
 q202 1491
 q203 1491
 q204 1491
 q205 1491
 q206 1491
 q207 1491
 q208 1491
 q209 1491
 q210 1491
 q211 1491
 q212 1491
 q213 1491
 q214 1491
 q215 1491
 q216 1491
 q217 1491
 q218 1491
 q219 1491
 q220 1491
 q221 1491
 q222 1491
 q223 1491
 q224 1491
 q225 1491
 q226 1491
 q227 1491
 q228 1491
 q229 1491
 q230 1491
 q231 1491
 q232 1491
 q233 1491
 q234 1491
 q235 1491
 q236 1491
 q237 1491
 q238 1491
 q239 1491
 q240 1491
 q241 1491
 q242 1491
 q243 1491
 q244 1491
 q245 1491
 q246 1491
 q247 1491
 q248 1491
 q249 1491
 q250 1491
 q251 1491
 q252 1491
 q253 1491
 q254 1491
 q255 1491
 q256 1491
 q257 1491
 q258 1491
 q259 1491
 q260 1491
 q261 1491
 q262 1491
 q263 1491
 q264 1491
 q265 1491
 q266 1491
 q267 1491
 q268 1491
 q269 1491
 q270 1491
 q271 1491
 q272 1491
 q273 1491
 q274 1491
 q275 1491
 q276 1491
 q277 1491
 q278 1491
 q279 1491
 q280 1491
 q281 1491
 q282 1491
 q283 1491
 q284 1491
 q285 1491
 q286 1491
 q287 1491
 q288 1491
 q289 1491
 q290 1491
 q291 1491
 q292 1491
 q293 1491
 q294 1491
 q295 1491
 q296 1491
 q297 1491
 q298 1491
 q299 1491
 q300 1491
 q301 1491
 q302 1491
 q303 1491
 q304 1491
 q305 1491
 q306 1491
 q307 1491
 q308 1491
 q309 1491
 q310 1491
 q311 1491
 q312 1491
 q313 1491
 q314 1491
 q315 1491
 q316 1491
 q317 1491
 q318 1491
 q319 1491
 q320 1491
 q321 1491
 q322 1491
 q323 1491
 q324 1491
 q325 1491
 q326 1491
 q327 1491
 q328 1491
 q329 1491
 q330 1491
 q331 1491
 q332 1491
 q333 1491
 q334 1491
 q335 1491
 q336 1491
 q337 1491
 q338 1491
 q339 1491
 q340 1491
 q341 1491
 q342 1491
 q343 1491
 q344 1491
 q345 1491
 q346 1491
 q347 1491
 q348 1491
 q349 1491
 q350 1491
 q351 1491
 q352 1491
 q353 1491
 q354 1491
 q355 1491
 q356 1491
 q357 1491
 q358 1491
 q359 1491
 q360 1491
 q361 1491
 q362 1491
 q363 1491
 q364 1491
 q365 1491
 q366 1491
 q367 1491
 q368 1491
 q369 1491
 q370 1491
 q371 1491
 q372 1491
 q373 1491
 q374 1491
 q375 1491
 q376 1491
 q377 1491
 q378 1491
 q379 1491
 q380 1491
 q381 1491
 q382 1491
 q383 1491
 q384 1491
 q385 1491
 q386 1491
 q387 1491
 q388 1491
 q389 1491
 q390 1491
 q391 1491
 q392 1491
 q393 1491
 q394 1491
 q395 1491
 q396 1491
 q397 1491
 q398 1491
 q399 1491
 q400 1491
 q401 1491
 q402 1491
 q403 1491
 q404 1491
 q405 1491
 q406 1491
 q407 1491
 q408 1491
 q409 1491
 q410 1491
 q411 1491
 q412 1491
 q413 1491
 q414 1491
 q415 1491
 q416 1491
 q417 1491
 q418 1491
 q419 1491
 q420 1491
 q421 1491
 q422 1491
 q423 1491
 q424 1491
 q425 1491
 q426 1491
 q427 1491
 q428 1491
 q429 1491
 q430 1491
 q431 1491
 q432 1491
 q433 1491
 q434 1491
 q435 1491
 q436 1491
 q437 1491
 q438 1491
 q439 1491
 q440 1491
 q441 1491
 q442 1491
 q443 1491
 q444 1491
 q445 1491
 q446 1491
 q447 1491
 q448 1491
 q449 1491
 q450 1491
 q451 1491
 q452 1491
 q453 1491
 q454 1491
 q455 1491
 q456 1491
 q457 1491
 q458 1491
 q459 1491
 q460 1491
 q461 1491
 q462 1491
 q463 1491
 q464 1491
 q465 1491
 q466 1491
 q467 1491
 q468 1491
 q469 1491
 q470 1491
 q471 1491
 q472 1491
 q473 1491
 q474 1491
 q475 1491
 q476 1491
 q477 1491
 q478 1491
 q479 1491
 q480 1491
 q481 1491
 q482 1491
 q483 1491
 q484 1491
 q485 1491
 q486 1491
 q487 1491
 q488 1491
 q489 1491
 q490 1491
 q491 1491
 q492 1491
 q493 1491
 q494 1491
 q495 1491
 q496 1491
 q497 1491
 q498 1491
 q499 1491
 q500 1491
 q501 1491
 q502 1491
 q503 1491
 q504 1491
 q505 1491
 q506 1491
 q507 1491
 q508 1491
 q509 1491
 q510 1491
 q511 1491
 q512 1491
 q513 1491
 q514 1491
 q515 1491
 q516 1491
 q517 1491
 q518 1491
 q519 1491
 q520 1491
 q521 1491
 q522 1491
 q523 1491
 q524 1491
 q525 1491
 q526 1491
 q527 1491
 q528 1491
 q529 1491
 q530 1491
 q531 1491
 q532 1491
 q533 1491
 q534 1491
 q535 1491
 q536 1491
 q537 1491
 q538 1491
 q539 1491
 q540 1491
 q541 1491
 q542 1491
 q543 1491
 q544 1491
 q545 1491
 q546 1491
 q547 1491
 q548 1491
 q549 1491
 q550 1491
 q551 1491
 q552 1491
 q553 1491
 q554 1491
 q555 1491
 q556 1491
 q557 1491
 q558 1491
 q559 1491
 q560 1491
 q561 1491
 q562 1491
 q563 1491
 q564 1491
 q565 1491
 q566 1491
 q567 1491
 q568 1491
 q569 1491
 q570 1491
 q571 1491
 q572 1491
 q573 1491
 q574 1491
 q575 1491
 q576 1491
 q577 1491
 q578 1491
 q579 1491
 q580 1491
 q581 1491
 q582 1491
 q583 1491
 q584 1491
 q585 1491
 q586 1491
 q587 1491
 q588 1491
 q589 1491
 q590 1491
 q591 1491
 q592 1491
 q593 1491
 q594 1491
 q595 1491
 q596 1491
 q597 1491
 q598 1491
 q599 1491
 q600 1491
 q601 1491
 q602 1491
 q603 1491
 q604 1491
 q605 1491
 q606 1491
 q607 1491
 q608 1491
 q609 1491
 q610 1491
 q611 1491
 q612 1491
 q613 1491
 q614 1491
 q615 1491
 q616 1491
 q617 1491
 q618 1491
 q619 1491
 q620 1491
 q621 1491
 q622 1491
 q623 1491
 q624 1491
 q625 1491
 q626 1491
 q627 1491
 q628 1491
 q629 1491
 q630 1491
 q631 1491
 q632 1491
 q633 1491
 q634 1491
 q635 1491
 q636 1491
 q637 1491
 q638 1491
 q639 1491
 q640 1491
 q641 1491
 q642 1491
 q643 1491
 q644 1491
 q645 1491
 q646 1491
 q647 1491
 q648 1491
 q649 1491
 q650 1491
 q651 1491
 q652 1491
 q653 1491
 q654 1491
 q655 1491
 q656 1491
 q657 1491
 q658 1491
 q659 1491
 q660 1491
 q661 1491
 q662 1491
 q663 1491
 q664 1491
 q665 1491
 q666 1491
 q667 1491
 q668 1491
 q669 1491
 q670 1491
 q671 1491
 q672 1491
 q673 1491
 q674 1491
 q675 1491
 q676 1491
 q677 1491
 q678 1491
 q679 1491
 q680 1491
 q681 1491
 q682 1491
 q683 1491
 q684 1491
 q685 1491
 q686 1491
 q687 1491
 q688 1491
 q689 1491
 q690 1491
 q691 1491
 q692 1491
 q693 1491
 q694 1491
 q695 1491
 q696 1491
 q697 1491
 q698 1491
 q699 1491
 q700 1491
 q701 1491
 q702 1491
 q703 1491
 q704 1491
 q705 1491
 q706 1491
 q707 1491
 q708 1491
 q709 1491
 q710 1491
 q711 1491
 q712 1491
 q713 1491
 q714 1491
 q715 1491
 q716 1491
 q717 1491
 q718 1491
 q719 1491
 q720 1491
 q721 1491
 q722 1491
 q723 1491
 q724 1491
 q725 1491
 q726 1491
 q727 1491
 q728 1491
 q729 1491
 q730 1491
 q731 1491
 q732 1491
 q733 1491
 q734 1491
 q735 1491
 q736 1491
 q737 1491
 q738 1491
 q739 1491
 q740 1491
 q741 1491
 q742 1491
 q743 1491
 q744 1491
 q745 1491
 q746 1491
 q747 1491
 q748 1491
 q749 1491
 q750 1491
 q751 1491
 q752 1491
 q753 1491
 q754 1491
 q755 1491
 q756 1491
 q757 1491
 q758 1491
 q759 1491
 q760 1491
 q761 1491
 q762 1491
 q763 1491
 q764 1491
 q765 1491
 q766 1491
 q767 1491
 q768 1491
 q769 1491
 q770 1491
 q771 1491
 q772 1491
 q773 1491
 q774 1491
 q775 1491
 q776 1491
 q777 1491
 q778 1491
 q779 1491
 q780 1491
 q781 1491
 q782 1491
 q783 1491
 q784 1491
 q785 1491
 q786 1491
 q787 1491
 q788 1491
 q789 1491
 q790 1491
 q791 1491
 q792 1491
 q793 1491
 q794 1491
 q795 1491
 q796 1491
 q797 1491
 q798 1491
 q799 1491
 q800 1491
 q801 1491
 q802 1491
 q803 1491
 q804 1491
 q805 1491
 q806 1491
 q807 1491
 q808 1491
 q809 1491
 q810 1491
 q811 1491
 q812 1491
 q813 1491
 q814 1491
 q815 1491
 q816 1491
 q817 1491
 q818 1491
 q819 1491
 q820 1491
 q821 1491
 q822 1491
 q823 1491
 q824 1491
 q825 1491
 q826 1491
 q827 1491
 q828 1491
 q829 1491
 q830 1491
 q831 1491
 q832 1491
 q833 1491
 q834 1491
 q835 1491
 q836 1491
 q837 1491
 q838 1491
 q839 1491
 q840 1491
 q841 1491
 q842 1491
 q843 1491
 q844 1491
 q845 1491
 q846 1491
 q847 1491
 q848 1491
 q849 1491
 q850 1491
 q851 1491
 q852 1491
 q853 149

1. 12. Capítulos e informaciones de Sinto y V. A. B. C. ¿ambos se encuentran en promedio anónimo medio menor como a dar una impresión con palabras? ¿una impresión o el estado que tiempo y presencia que vez?
2. 11. Capítulos e informaciones de Sinto y V. A. B. C. ¿ambos se encuentran en promedio anónimo medio menor como a dar una impresión con palabras? ¿una impresión o el estado que tiempo y presencia que vez?
3. 11. Capítulos e informaciones de Sinto y V. A. B. C. ¿ambos se encuentran en promedio anónimo medio menor como a dar una impresión con palabras? ¿una impresión o el estado que tiempo y presencia que vez?
4. 12. Capítulos e informaciones de Sinto y V. A. B. C. ¿ambos se encuentran en promedio anónimo medio menor como a dar una impresión con palabras? ¿una impresión o el estado que tiempo y presencia que vez?
5. 12. Capítulos e informaciones de Sinto y V. A. B. C. ¿ambos se encuentran en promedio anónimo medio menor como a dar una impresión con palabras? ¿una impresión o el estado que tiempo y presencia que vez?
6. 12. Capítulos e informaciones de Sinto y V. A. B. C. ¿ambos se encuentran en promedio anónimo medio menor como a dar una impresión con palabras? ¿una impresión o el estado que tiempo y presencia que vez?
7. 12. Capítulos e informaciones de Sinto y V. A. B. C. ¿ambos se encuentran en promedio anónimo medio menor como a dar una impresión con palabras? ¿una impresión o el estado que tiempo y presencia que vez?
8. 12. Capítulos e informaciones de Sinto y V. A. B. C. ¿ambos se encuentran en promedio anónimo medio menor como a dar una impresión con palabras? ¿una impresión o el estado que tiempo y presencia que vez?
9. 12. Capítulos e informaciones de Sinto y V. A. B. C. ¿ambos se encuentran en promedio anónimo medio menor como a dar una impresión con palabras? ¿una impresión o el estado que tiempo y presencia que vez?

Gráfico 3: Pantalla de palabras clave en contexto con el *Corpus del Español*

un pico de alrededor de 12 ocurrencias por millón de palabras en el siglo XIV). Además, uno puede ver la palabra clave en contexto para cualquier palabra, no solo aquellas con baja frecuencia (como es el caso del *CORDE*).

En el *Corpus del Español*, al hacer clic en los números en cualquier columna se mostrará la palabra clave en el contexto de ese siglo, o se pueden ver todas las entradas a la vez haciendo clic en TOTAL. Luego se puede hacer clic en otras entradas para ver el contexto extendido (aproximadamente 200 palabras en total).

La visualización de palabra clave en contexto en el *Corpus del Español* permite algunas funcionalidades importantes que no son posibles con el *CORDE*. Primero, los usuarios pueden guardar ocurrencias (en tres grupos diferentes de ocurrencias) haciendo clic en A, B o C. Posteriormente, pueden organizar estas listas de ocurrencias (lo que incluye moverlas entre diferentes listas), lo que facilita en gran medida el almacenamiento y análisis de datos.

Hasta este punto, entonces, las búsquedas en los dos corpus son bastante similares. El *CORDE* tiene la ventaja de ser el corpus más grande, mientras que el *Corpus del Español* tiene la ventaja de mostrar la frecuencia en cada siglo y de mostrar la palabra clave en contexto para todas las palabras, independientemente de la frecuencia.

5 Frecuencia de palabra: datos básicos

Sin embargo, además de obtener simplemente todas las apariciones de una palabra o frase determinada, los usuarios a menudo quieren saber cuán frecuente eran en diferentes siglos o en determinados periodos históricos. Es en este punto donde el *CORDE* comienza a mostrar algunas debilidades serias. P. ej., después de buscar *braveza* y luego seleccionar «Ver estadística», el usuario visualiza:

Estadísticas			
Año	%	Casos	País
1627	20.95	35	ESPAÑA
1547	20.35	34	MÉXICO
1610	17.96	30	PERÚ
1632	8.98	15	BOLIVIA
1566	4.79	8	COLOMBIA

Tema	%	Casos
22. - Verso narrativo	26.37	72
12. - Prosa narrativa	20.51	56
19. - Prosa histórica	17.58	48
21. - Verso lírico	12.45	34
16. - Prosa de sociedad	7.69	21

Gráfico 4: Frecuencia (por año) con el *CORDE*

Esta tabla nos dice los años específicos en los cuales la palabra o frase es más común, pero es imposible ver la frecuencia por década o por siglo. No sirve de nada mostrar que la palabra fue la más frecuente en 1627 si, de hecho, es mucho menos común en el siglo XVII que en el XIII o el XIV. El otro problema importante es que las cifras no están relativizadas. En otras palabras, vemos la frecuencia absoluta por año, pero una palabra o frase puede ser más común en ese año simplemente porque hay más palabras para ese año en el corpus. Cualquier comparación sería de frecuencia requiere que los resultados se relativicen a lo largo de periodos históricos para que podamos tener en cuenta los diferentes tamaños del corpus en diferentes periodos históricos, y ver la frecuencia de la palabra o frase por millón de palabras.

El *Corpus del Español* permite este tipo de búsqueda con bastante facilidad. P. ej., con *braveza* en el *Corpus del Español*, podemos ver la «visualización de tabla» (como en la tabla 2 anterior), o una visualización en gráfico:

SECCIÓN	S13	S14	S15	S16	S17	S18	S19	S20
OCURENCIAS	48	32	27	32	5	0	0	0
POR MILL	7.15	11.99	3.31	1.88	0.40	0.00	0.00	0.00

Gráfico 5: *Corpus del Español*: Frecuencia de palabra y frase por siglo

Se nos muestra aquí la frecuencia absoluta (p. ej., 32 ocurrencias en siglo XIV), así como la importante frecuencia relativa («por millón»), que tiene en cuenta el tamaño de la sección en millones de palabras. P. ej., hay 32 ocurrencias en los 3 millones de palabras del siglo XIV, o 10,8 ocurrencias por millón de palabras. Un gráfico como este es la única manera de ver realmente los cambios en la frecuencia de una palabra, frase o construcción, y solo es posible con el *Corpus del Español*.

6 Frecuencia de palabra: comparación de períodos históricos

En lugar de tener que decirle al corpus qué palabras o frases específicas buscar, un corpus con arquitectura e interfaz bien diseñadas generaría una lista de palabras cuya frecuencia coincida con ciertos criterios. P. ej., podría encontrar todos los sustantivos que entraron en la lengua del siglo XVII o las palabras que se usaron al menos cinco veces más en el siglo XIII que en el XIV. Tal consulta es completamente imposible con el *CORDE*. Todo lo que se puede hacer es buscar palabras y frases específicas. Si el *CORDE* reconoce la frecuencia de todas las palabras y frases en todos los períodos históricos, ciertamente no permite a los investigadores usar esa información mediante una consulta.

Con el *Corpus of Spanish*, por otro lado, tales consultas son bastante sencillas. P. ej., se puede buscar simplemente [nn *] (sustantivos) y seleccionar [siglo XIII] (1200–1299) «SECCION 1» para comparar con [siglo XIV] (1300–1399) «SECCION 2». En uno o dos segundos, el usuario ve la siguiente lista⁴:

Tabla 1: *Corpus del Español*. Comparación de la frecuencia de palabras por siglo (todas las palabras a la vez)

Siglo XIII	Siglo XIV
capitolo, ascendente, ladeza, saturnus, orizon, morauedis, roque, xaque, acendent, armella, murcia, baldouin, ascendent, segonda, gudufre, significador, ferrando, coronan, afffil, zonte, dond, iudga, tiempro, boy monte, catamieto, infortunas, unie, caput, canyaron, sacrificio, declinacion, sacrificios, ygnador, juppiet, algauru, himas, delllos, hienusalem, decina.	armadas, osso, ome, avia, paris, abe, collado, elena, falcon, yuierno, verano, goncales, encarnacion, pase, bien, venja, avras, falcones, ynfante, facer, puero, ynfanta, ynfanta, ynfanto, yvno, venjdo, sembrar, fojas, enqina, ynfante, mjel, menhalo, syempre, dolencia, ssiete, avedes, castilla, muria, aujdo, peca, arroyo, vuas, qienca, termjino, tenjan.

4 Téngase en cuenta que en la versión web hay frecuencias –en bruto y normalizadas– para cada palabra, así como enlaces para ver la palabra en contexto, como se muestra en el Gráfico 2. En la Tabla 1, hemos simplificado la visualización.

SEC1 (119005): 22,892,256 PALABRAS				SEC2 (119005): 19,297,249 PALABRAS								
PALABRA/FRASE	OCCURR 1	OCCURR 2	PROPORCION	PALABRA/FRASE	OCCURR 1	OCCURR 2	PROPORCION					
1 SECTOR	2540	0	111.3	0.0	11,729.5	1	VENTURA	1237	25	64.1	1.1	58.5
2 TELEVISION	2119	0	92.8	0.0	9,284.8	2	NUMI	1025	21	55.1	0.9	57.7
3 SECTORES	1960	0	59.6	0.0	5,959.1	3	RODRICO	697	19	36.1	0.8	43.4
4 FUTUOL	1208	0	52.9	0.0	5,293.1	4	VULCO	662	20	34.3	0.9	35.1
5 LIDER	1119	0	49.0	0.0	4,903.1	5	PAE	611	21	31.7	0.9	34.4
6 PROTELINAS	753	0	33.0	0.0	3,299.4	6	MATLIDE	609	24	34.7	1.1	33.0
7 IMPACTO	746	0	32.7	0.0	3,268.7	7	AOSENTO	1174	44	60.8	1.9	31.6
8 AEROPUERTO	721	0	31.6	0.0	3,159.2	8	HONRA	1652	63	85.6	2.8	31.0
9 NARCOTRAFICO	705	0	30.9	0.0	3,098.1	9	MENSTER	1545	60	80.1	2.6	30.5
10 INFACCION	704	0	30.8	0.0	3,084.7	10	MARGARITA	554	23	28.7	1.0	28.5
11 LIBERES	666	0	29.2	0.0	2,918.2	11	EMIRIANO	553	23	28.7	1.0	28.4
12 CAMPENATO	657	0	28.8	0.0	2,878.8	12	CONDESA	1829	78	94.8	3.4	27.7

Gráfico 6: *Corpus del Español*. Comparación de la frecuencia de las palabras (sustantivos en siglo XIX/siglo XX)

Obviamente, solo algunas de las palabras de esta lista son significativas. Muchas palabras son simplemente variantes ortográficas, otras son sustantivos propios que pueden aparecer en un puñado de textos de un siglo, pero no de otro.

Por supuesto, tales búsquedas no están limitadas solo al español antiguo, también se pueden llevar a cabo para períodos históricos más recientes. La siguiente tabla muestra (a la izquierda) nombres que son comunes en el siglo XX, pero no en el XIX y (a la derecha) aquellos que son comunes en el siglo XIX pero no en el siglo XX. Esta tabla muestra todas las palabras que son más comunes en un período que en otro, incluso si no aparecen en el segundo período. P. ej., no hay resultados de televisión o aeropuerto en el siglo XIX (lo que probablemente no sea sorprendente).

Para «comparar manzanas con manzanas», podría ser útil indicar que una palabra debe darse con al menos una frecuencia determinada en cada uno de los dos períodos. P. ej., la tabla 7 muestra los nombres que aparecen al menos diez veces en los siglos XIX y XX, pero en los que hay un aumento significativo a lo largo del tiempo⁵.

En resumen, debido a la arquitectura del *Corpus del Español*, donde el corpus «reconoce» la frecuencia de cada palabra y frase en cada período histórico, tales comparaciones son bastante simples. En cambio, con el *CORDE* el corpus aparentemente no reconoce la frecuencia de palabras y frases en cada sección

5 Téngase en cuenta que, debido a que el etiquetado no es siempre perfecto, hay algunas entradas extrañas, como *lic* en el siglo XX o el nombre propio de *Matilde* en el siglo XIX, pero la mayoría de las palabras son relevantes.

SEC 1119004: 22,822,256 PALABRAS				SEC 2118004: 19,297,249 PALABRAS			
PALABRA/FRASE	OCCURREN1	P/M 1	PROPORCIÓN	PALABRA/FRASE	OCCURREN1	P/M 1	PROPORCIÓN
1 PROYOTONISMO	175	0	7.7	0.0	766.8	1	0.0
2 MEBADOLISMO	148	0	6.5	0.0	648.5	2	0.0
3 URBANISMO	127	0	5.6	0.0	556.5	3	0.0
4 SUPERBOLISMO	118	0	5.2	0.0	517.0	4	0.0
5 ATLETISMO	111	0	4.9	0.0	486.4	5	0.0
6 MARRISMO	108	0	4.7	0.0	473.2	6	0.0
7 FASCISMO	108	0	4.7	0.0	473.2	7	0.0
8 CUBISMO	105	0	4.6	0.0	460.1	8	0.0
9 HINDUISMO	104	0	4.6	0.0	455.7	9	0.0
10 NEROSISMO	97	0	4.3	0.0	425.0	10	0.0
11 EPRENSIONISMO	79	0	3.5	0.0	346.2	11	0.0
1 PAUPERISMO	71	0	3.7	0.0	367.9	1	0.0
2 CACIQUISMO	123	14	6.4	0.6	10.4	2	0.0
3 DESPOTISMO	239	38	15.2	1.7	9.1	3	0.0
4 ABISMO	1030	186	53.4	6.1	6.5	4	0.0
5 PAGANISMO	88	16	4.6	0.7	6.5	5	0.0
6 FANATISMO	281	65	14.6	2.8	5.1	6	0.0
7 PATRIOTISMO	330	78	17.1	3.6	5.0	7	0.0
8 FATALISMO	50	19	2.6	0.8	3.1	8	0.0
9 SENTIMENTALISMO	57	22	3.0	1.0	3.1	9	0.0
10 CATORCISMO	304	138	15.8	6.0	2.6	10	0.0
11 CATECISMO	147	68	7.6	3.0	2.6	11	0.0

Gráfico 9: *Corpus del Español*: Comparación de formas (*ismo en siglos XIX/XX)

Esto muestra, p. ej., que *despotismo* se encuentra 293 veces en el siglo XIX, pero solo 38 veces en el XX, y que *fascismo* se encuentra 108 veces el siglo XX, pero ninguna en el XIX. La capacidad de comparar formas de palabras en diferentes siglos es una característica de gran alcance del *Corpus del Español*, pero no es posible con el *CORDE*.

8 Significado de la palabra (semántica): colocaciones básicas

Como les gusta señalar a los lingüistas de corpus, «se puede decir mucho sobre una palabra con las otras palabras con las que se junta» (Firth 1957). A veces, las colocaciones (palabras cercanas) simplemente confirman lo que ya sabemos. P. ej., las colocaciones nominales más comunes (palabras cercanas) para *selva* son *árboles*, *vegetación*, *sierra*, *bosque*, etc. Para una palabra menos concreta, a menudo es necesaria más perspicacia. P. ej., los sustantivos más comunes que aparecen con formas de *lúgubre* son *acento*, *voz*, *silencio*, *noche*, *eco*, *gemido*, etc.

La clave del significado, entonces, se suele encontrar en las colocaciones o en las palabras cercanas. Prácticamente cualquier arquitectura e interfaz de corpus permite a los usuarios buscar una palabra y luego ver esa palabra en contexto. El usuario del corpus siempre puede recorrer los ejemplos uno por uno, tomar notas sobre palabras cercanas comunes y luego tratar de usar estas colocaciones para discernir el significado. Sin embargo, esto puede consumir mucho tiempo para palabras comunes. Un enfoque más productivo consistiría en hacer que el corpus encuentre todas las colocaciones por sí mismo y luego presentarlas al usuario en orden de frecuencia.

En términos de cambio histórico, lo deseable sería poder encontrar las colocaciones de una palabra dada en diferentes períodos históricos. Al observar los

cambios en las colocaciones a lo largo del tiempo, podemos obtener información sobre los cambios en el significado y el uso de la palabra.

Consideremos brevemente cómo el *CORDE* y el *Corpus del Español* permiten a los usuarios encontrar y procesar las colocaciones para obtener una idea del significado de las palabras. En el caso del *CORDE*, supongamos que queremos examinar las casi 38 000 colocaciones de todas las formas de *duro* (*dura*, *duros*, etc.). Suponiendo que a un usuario le toma unos 20 segundos encontrar cada ocurrencia en contexto y anotar (lo que supone que son) las palabras cercanas relevantes, dedicaría alrededor de 26 días (a ocho horas diarias) a repasar todos los ejemplos relevantes; y esto, suponiendo que el usuario no decidiera cambiar el ancho de la «ventana de colocaciones» o buscarse un tipo diferente de colocación, en cuyo caso tendría que emplear, más o menos, otro mes.

El *CORDE* permite a los usuarios ver «agrupaciones» para una palabra determinada, como las de la forma simple *duro* en el siglo XIII:

De 2 palabras	%	Casos	De 3 palabras	%	Casos
<i>duro la</i>	6.97	30	<i>duro en el</i>	2.55	11
<i>duro el</i>	6.27	27	<i>duro la batalla</i>	1.86	8
<i>duro en</i>	5.34	23	<i>duro punto nado</i>	1.16	5
<i>duro fasta</i>	3.95	17	<i>duro esta batalla</i>	1.16	5
<i>duro esta</i>	3.48	15	<i>duro fasta el</i>	1.16	5

Gráfico 10: *CORDE*: *duro* + colocaciones

Pero tal listado es de poco valor, porque como el *CORDE* no tiene ninguna forma de discriminar qué palabras son relevantes, nos da colocaciones como *dura la*, *dura en*, etc. (ya que *la*, *en*, etc., concurren con frecuencia con casi cualquier término), que proporcionan poca o ninguna comprensión del significado de la palabra. En cualquier caso, solo lista las nueve colocaciones más frecuentes, lo que no es suficiente para obtener información completa sobre el significado. Estas búsquedas son mucho más fáciles con el *Corpus del Español*. Los usuarios simplemente introducen la «palabra del nodo» (p. ej., *duro*, o *lúgubre*, o *selva*) y pueden, opcionalmente, seleccionar la categoría gramatical de las colocaciones:

en aproximadamente dos o tres segundos tienen todas las colocaciones en orden. P. ej., suponemos que un usuario quiere encontrar colocaciones relacionadas con el concepto '*duro*' en español antiguo. Después de introducir [= *duro*] (*duro*, *duros*, etc.), y esperar unos dos segundos, aparece una lista como la siguiente:

	CONTEXTO	TODOS	s13	s14	s15	s16	s17
1	PIEDRA	188	33	7	52	56	40
2	QUEBRANTAR	91	89			2	
3	MIENTRAS	84			1	48	35
4	PIEDRAS	81	10	3	23	28	17
5	DURA	77			15	30	19
6	GOLPES	75		12	1	53	9
7	FUERT	71	71				
8	HIERRO	53			3	36	14
9	PEÑA	48			2	21	25
10	MÁRMOL	47				26	21
11	HADO	42				40	2
12	PEÑAS	42		1	1	24	16

Gráfico 11: *Corpus del Español*: colocaciones de *duro* en siglos XIII–XVIII

Esta tabla muestra la frecuencia de cada colocación en cada siglo (aquí solo se muestran los siglos XIII al XVIII). P. ej., *pietra* ocurre cerca de [*duro*] 33 veces en el siglo XIII y 7 veces más en el XV. Hay 188 apariciones totales de *pietra* entre los siglos XIII y XV, por lo que los 40 casos cerca de [*duro*] son aproximadamente el 1,8 % de todos los resultados. Esto se traduce en una puntuación de información mutua de 3.53, que muestra que la relación entre las dos palabras es significativa. Por lo tanto, con el *Corpus del Español*, podemos hacer en 2–3 segundos lo mismo que llevaría hacer con el CORDE un mes o más.

9 Cambio semántico: comparación de colocaciones en diferentes períodos históricos

Si tenemos una arquitectura e interfaz de corpus que nos permite encontrar fácilmente colocaciones (lo que es posible, como se ha visto, con el *Corpus del Español*), podemos utilizar esta información de manera ingeniosa para examinar el cambio semántico. La idea básica es que si las palabras «cercanas» a una

determinada palabra cambian con el tiempo, puede deberse a que la misma palabra ha cambiado de significado (o al menos se está utilizando de una manera diferente). P. ej., la siguiente tabla muestra (a la izquierda) los sustantivos que aparecen con [*duro*] en el siglo XX, pero que no son muy comunes en el XIX, y (a la derecha) en el XIX, pero no en el XX.

SEC 1 (1900): 22,822,256 PALABRAS					SEC 2 (1800): 19,297,249 PALABRAS					
PALABRA/FRESE	OCCURREN	P/M1	P/M2	PROPORCIÓN	PALABRA/FRESE	OCCURREN	P/M1	P/M2	PROPORCIÓN	
1. CRIEGAS	27	0	1.2	0.0	1	MILES	30	0	1.6	0.0
2. MADEBAS	19	0	0.8	0.0	2	RENTA	24	0	1.2	0.0
3. LINEA	18	0	0.8	0.0	3	MILLONES	20	0	1.0	0.0
4. REPRESION	16	0	0.7	0.0	4	TRANCE	18	0	0.9	0.0
5. DISCO	13	0	0.6	0.0	5	ART	18	0	0.9	0.0
6. REGALO	10	0	0.4	0.0	6	CARACTER	17	0	0.9	0.0
7. LUCHA	10	1	0.4	0.1	7	ENTRAÑAS	14	0	0.7	0.0
8. COMPETENCIA	15	2	0.7	0.1	8	DUROS	13	0	0.7	0.0
9. BATALLA	51	12	2.2	0.6	9	REALES	13	0	0.7	0.0
10. GOLPE	12	3	0.5	0.2	10	PESOS	12	0	0.6	0.0
11. MIES	12	3	0.5	0.2	11	PUMADO	10	0	0.5	0.0
12. FORMA	12	3	0.5	0.2	12	CANTIDAD	10	0	0.5	0.0
13. GOBIERNO	13	4	0.6	0.2	13	ALMA	10	0	0.5	0.0

Gráfico 12: *Corpus del Español*: comparación de las colocaciones de *duro*, siglos XIX/XX

P. ej., *disco* se produce cerca de *duro* 13 veces en el siglo XX, pero no hay resultados (lo que no sorprende) en el XIX. *Entrañas*, por otro lado, aparece 14 veces con [*duro*] en el XIX, pero no hay resultados en el siglo XX. Si *disco* está en el corpus en el siglo XIX y *entrañas* en el siglo XX, ¿por qué su frecuencia como colocación con *duro* cambia tanto de un siglo a otro? ¿Es porque el uso de *duro* puede haber cambiado en algo? Mostremos otro ejemplo: la siguiente es una lista parcial de las colocaciones adjetivales de *mujer* (y *mujeres*) en los siglos XX (izquierda) y XIX (derecha):

SEC 1 (1900): 22,822,256 PALABRAS					SEC 2 (1800): 19,297,249 PALABRAS					
PALABRA/FRESE	OCCURREN	P/M1	P/M2	PROPORCIÓN	PALABRA/FRESE	OCCURREN	P/M1	P/M2	PROPORCIÓN	
1. INTERNAZIONALE	17	0	0.7	0.0	1	HONRADA	59	0	3.0	0.0
2. SEXUAL	15	0	0.7	0.0	2	HONRADAS	35	0	1.8	0.0
3. CUBANAS	15	0	0.7	0.0	3	DIGNA	24	0	1.2	0.0
4. EMBARAZADAS	14	0	0.6	0.0	4	LIVANA	21	0	1.1	0.0
5. IMPORTANTE	13	0	0.6	0.0	5	INFAME	19	0	1.0	0.0
6. PROGRESIVALES	13	0	0.6	0.0	6	DESVENTURADA	19	0	1.0	0.0
7. LABORAL	12	0	0.5	0.0	7	PERIODAS	18	0	0.9	0.0
8. EX	12	0	0.5	0.0	8	VANO	18	0	0.9	0.0
9. ARABE	10	0	0.4	0.0	9	MALDITA	16	0	0.8	0.0
10. DERECHOS	10	0	0.4	0.0	10	SEMIVANTE	15	0	0.8	0.0
11. DIFERENTES	15	1	0.7	0.1	11	DESIGNADA	14	0	0.7	0.0
12. MADURA	15	1	0.7	0.1	12	LIVANAS	13	0	0.7	0.0
13. SEXUALES	12	1	0.5	0.1	13	INELICHS	13	0	0.7	0.0

Gráfico 13: *Corpus del Español*: comparación de colocaciones de *mujer*, siglos XIX/XX

Obsérvese cómo los adjetivos del siglo XIX (derecha) se refieren a las «virtudes morales» de las mujeres, que están casi por completo ausentes en el siglo XX (izquierda); en este siglo, por otro lado, son mucho más prosaicos y se refieren a clasificaciones sobre la nacionalidad, el empleo, etc. En este caso, los datos del corpus proporcionan información interesante del cambio en la forma de ver a las mujeres en estos dos siglos; podemos obtener esta útil información con una simple búsqueda de 1–2 segundos en el corpus.

Aplicado al español antiguo o en los textos de los siglos XVI al XVIII, se podría adoptar un enfoque similar. Usando la interfaz para el *Corpus del Español*, simplemente es necesario indicar qué palabras o conceptos son de interés, especificar el tipo de colocación (sustantivo, verbo, etc., si corresponde), y luego hacer clic una o dos veces más para mostrar qué dos periodos históricos deben ser comparados. En dos o tres segundos, se recopilan y resumen todos los datos relevantes. Usando el *CORDE*, en cambio, las búsquedas como esta serían muy difíciles o imposibles, ya que la arquitectura del *CORDE* no está diseñada para encontrar colocaciones.

10 Cambios léxicos en un campo semántico: sinónimos

Con la arquitectura de corpus adecuada, los usuarios podrían buscar por campos semánticos, en lugar de simplemente buscar palabras y frases. Así, en el *Corpus del Español*, los índices onomasiológicos se integran en la arquitectura del corpus; esto nos permite, en el nivel más básico, encontrar la frecuencia histórica de todas las palabras relacionadas con un concepto en particular. P. ej., los usuarios pueden introducir = *mujer* y ver la frecuencia de todos los sinónimos de *mujer* a lo largo del tiempo (y, p. ej., en diferentes géneros en el siglo XX).

Esta lista parcial de resultados muestra, p. ej., que *doncella* y *moza* han disminuido desde el siglo XVI (por cada millón de palabras), mientras que *chica* y *muchacha* han aumentado desde el siglo XVIII hasta el XIX/XX. En cuanto a los periodos medievales, lo que obviamente se necesitaría es algún tipo de «índice histórico» que por ahora no está disponible. Pero, en la medida en que lo estuviera, la arquitectura del corpus podría incorporarlo fácilmente.

Además de buscar la frecuencia de palabras sueltas, la información semántica de los repertorios o las listas de palabras personalizadas y definidas por el usuario se pueden integrar directamente en la sintaxis de la consulta. P. ej., en el *Corpus del Español*, es posible que los usuarios creen (a través de la interfaz web) listas personalizadas de palabras de un campo de interés semántico particular, como términos navales, palabras relacionadas con las emociones, una lista de términos relacionados con la estructura familiar o una lista de palabras relacionadas con

CONTEXTO	TODOS	913	914	915	916	917	918	919	920	ACAD	PER	FIG	ORAL
1 <input type="checkbox"/> MUJER (S)	35395	19	6	697	5707	8233	1962	11836	6945	478	803	3740	1924
2 <input type="checkbox"/> SEÑORA (S)	32812	161	955	1584	9385	6855	1704	7455	2223	43	314	1619	1924
3 <input type="checkbox"/> PERSONA (S)	24445	430	215	2637	5376	3430	2044	5586	4745	787	854	786	2318
4 <input type="checkbox"/> JOVEN (S)	11639					4	238	465	607	7193	3132	362	729
5 <input type="checkbox"/> ESPAÑOLA (S)	9400	121	36	219	1712	2899	365	2594	1444	196	400	572	196
6 <input type="checkbox"/> DAMA (S)	8287					1412	3879	332	2020	506	46	134	266
7 <input type="checkbox"/> MUCHACHA (S)	2266					1115	135	132	994	890	13	49	663
8 <input type="checkbox"/> DONCELLA (S)	3025					1003	752	224	987	65	8	11	42
9 <input type="checkbox"/> CHICA (S)	1725	29	43	95	97	59	96	390	916	19	98	406	333
10 <input type="checkbox"/> SEÑORITA (S)	1534					1	1	7	134	962	429	9	28
11 <input type="checkbox"/> HEMBRA (S)	1397	1	3	261	311	144	97	241	359	207	22	89	21
12 <input type="checkbox"/> MOZA (S)	1592					328	584	116	538	24	1	20	2

Gráfico 14: *Corpus del Español*: comparación de sinónimos de *mujer*

un concepto teológico particular. Esta lista personalizada de palabras se puede usar como parte de la sintaxis de la consulta. P. ej., si un usuario [andrés gómez] crea una lista de 100 palabras relacionadas con ‘emociones’ en español antiguo y otra lista de 70 palabras relacionadas con ‘relaciones familiares’ (*padre, hermano, nuera, etc.*), podría encontrar cada vez que aparezca una palabra en la Lista 1 cerca de la Lista 2. De esta forma, se pueden llevar a cabo búsquedas semánticas de gran alcance en el corpus.

El *Corpus del Español* puede incorporar este tipo de consultas orientadas semánticamente debido a la arquitectura subyacente del corpus, que se basa en bases de datos relacionales. Con este tipo de bases de datos relacionales es posible agregar cualquier cantidad de conjuntos de datos nuevos (repertorios, listas de palabras definidas por el usuario, etc.) y luego integrarlos sin problemas en la sintaxis de la consulta. La arquitectura para el *CORDE*, en cambio, no es «abierta» y no admite la incorporación de otros conjuntos de datos. Solo se puede buscar palabras o frases individuales, pero nada que se aproxime a un campo semántico completo o algo similar.

II Devió sincrónico # 1: género

Hasta este punto, hemos discutido los cambios históricos en el léxico y el significado, pero apenas hemos hecho alusión al género y a cómo este afecta a la frecuencia y al significado de las palabras. Sin embargo, obviamente, hay que considerar la importancia del género. P. ej., analicé los siguientes tres cuadros, que muestran la frecuencia de tres palabras del español moderno que son más comunes en el género de ficción (*borrachos*), académico (*proporcionar*) y el marcador de discurso bueno en español hablado (como en «lo haremos. Bueno, pero no es tan fácil»):

Tabla 2: Frecuencia de palabras de *-idad* en el *Corpus del Español* original (2002)

> 100	NONE
50–100	sonoridad, inmovilidad, expresividad, obscuridad
20–49	morbilidad, unicidad, posmodernidad, adaptabilidad, cotidianidad, discrecionalidad, selectividad, salubridad, elegibilidad, insensibilidad, disconformidad, mensuralidad, afectividad, deformidad
10–19	perpetuidad, receptividad, morosidad, anormalidad, hipersensibilidad, disparidad, extremidad, interioridad, permeabilidad, corresponsabilidad, perversidad, susceptibilidad, viscosidad, emotividad, natividad, plasticidad, asiduidad, promiscuidad, salinidad, virilidad, homogeneidad, inutilidad, frivolidad, perplejidad, voracidad
< 10	empleabilidad, trazabilidad, centralidad, consanguinidad, invisibilidad, inhabilidad, estadidad, habitabilidad, alteridad, sociabilidad, probidad, banalidad, materialidad

12 El nuevo *Corpus del Español* (Web/Dialectos)

Si bien es muy útil el componente original «Histórico/Género» de 100 millones de palabras, el *Corpus del Español* presenta algunas limitaciones. La primera es el tamaño: solo hay 20 millones de palabras del siglo XX y esta cifra resulta cuantitativamente escasa para una investigación en profundidad del léxico. La segunda es que el *Corpus del Español* original termina en 1999, no contiene textos del siglo XXI. Y la tercera es que el corpus original no permitía a los usuarios comparar el léxico en diferentes países.

Para abordar estas tres limitaciones creamos una gran extensión del *Corpus del Español* en 2014–2015. El nuevo corpus alcanza un tamaño de dos mil millones de palabras, lo que significa que es 100 veces el tamaño de la parte del español del siglo XX en el *Corpus del Español* original (2002). Todos los textos basados en la web se recopilieron en 2014–2015, por lo que el corpus representa mejor el español reciente. Finalmente, permite a los usuarios comparar frecuencias léxicas en 21 países diferentes de habla hispana.

En términos de tamaño, un corpus que es 100 veces más grande proporciona datos mucho más ricos y mayor información sobre la variación léxica. P. ej., consíderse la Tabla 2, que muestra aquellas palabras que terminan en *-idad*, que tienen una frecuencia de entre 2000 y 3000 ocurrencias en el nuevo corpus de dos mil millones de palabras. La tabla muestra cuántas veces aparecen en el corpus de más de 100 millones de palabras (donde 20 millones de palabras son del siglo XX).

Solo unas pocas palabras tienen una frecuencia de al menos 50 ocurrencias en el corpus anterior, la mayoría ocurre menos de 20 veces. Con un número tan limitado de apariciones, es muy poco lo que se puede hacer para investigar el significado y uso de las palabras o investigar su frecuencia en todos los géneros. Además, en algún punto hay tan pocos resultados que la misma palabra *frecuencia* se vuelve problemática. P. ej., una voz que aparece solo 15 veces puede aparecer en solo 3–4 textos diferentes, y la palabra podría no encontrarse en absoluto en un corpus que tuviera una composición textual ligeramente diferente. En algún punto, los términos con baja frecuencia son simplemente «ruidos».

Además de la frecuencia de voces aisladas, el nuevo (2016) corpus de dos mil millones de palabras también proporciona datos mucho más ricos para palabras en contexto con otras palabras. P. ej., en el *Corpus del Español* original, solo hay 67 adjetivos diferentes que aparecen cinco veces en el contexto gente + ADJ (gente *joven*, *pobre*, *rica*... *ponderosa*, *popular*), mientras que hay aproximadamente 1445 adjetivos diferentes en el corpus nuevo de dos mil millones de palabras (p. ej., *visionaria*, *desmotivada*, *sumisa*, *enfadada*, *gritona*, *carrutiva*). Del mismo modo, solo hay 24 nombres diferentes que aparecen cinco veces o más antes de una forma de *rico* en el corpus anterior (*países*, *comerciantes*, *alimentos*), mientras que hay alrededor de 650 nombres diferentes en el corpus de dos mil millones de palabras (*vocabulario*, *parientes*, *subsuelo*, *guion*, *piedras*). Como ejemplo final, solo hay unos diez sustantivos que aparecen al menos cinco veces en cuatro palabras después de una forma de *destronar* en el corpus anterior (p. ej., *corazón*, *alma*, *guerra*, *átomos*, *puerta*), pero hay más de 570 nombres en el nuevo corpus (p. ej., *canCIÓN*, *confianza*, *películas*, *cademas*, *pedacitos*, *hacha*). En la mayoría de los casos, hay 50–60 veces más colocaciones en el nuevo corpus de dos mil millones de palabras, lo que, por supuesto, proporciona a los investigadores datos mucho más detallados y útiles sobre el significado y el uso de estas palabras.

13 Búsqueda y estudio de neologismos

Desde el punto de vista del cambio léxico, el corpus nuevo también proporciona datos mucho más ricos para los neologismos. Esto es consecuencia del incremento del corpus (100 veces más grande que la parte del siglo XX del *Corpus del Español* original), pero también se debe a que el nuevo corpus contiene textos muy recientes. Los dos millones de textos fueron extraídos de la Web en 2014–2015, una década y media después de haber sido introducidos los últimos textos en el *Corpus del Español* original.

Hay muchos miles de palabras que aparecen al menos 500 veces en el corpus más reciente, pero menos de diez veces (en muchos casos nada) en el corpus

más antiguo y pequeño. Por supuesto, muchas de estas son palabras relacionadas con la tecnología (p. ej., *blog*, *web*, *internet*, *celular*, *navegador*, *correo electrónico*, *click*, *tweet*), pero otras son palabras relacionadas con otros campos (p. ej., *documental*, *implementación*, *fiscalía*, *inversionista*, *globalización*, *operativo*, *migrante*, *sostenibilidad*, *biodiversidad*).

Como una prueba más de la extensión de estos neologismos, los sustantivos que terminan en *-idad* y que se anotan a continuación aparecen al menos 700 veces en el nuevo corpus, pero menos de cinco (generalmente 0 veces) en el corpus antiguo: *interoperabilidad*, *catolicidad*, *virilidad*, *masividad*, *proactividad*, *escalabilidad*, *colonialidad*, *transversalidad*, *transsexualidad*, *digestibilidad*, *sincronicidad*, *ciberseguridad*, *emocionalidad*, *accidentalidad*, *inexequibilidad*, *tipicidad*, *ruralidad*. Las palabras que terminan en *-ismo* incluyen: *kirchnerismo*, *chavismo*, *fujimorismo*, *extractivismo*, *multiculturalismo*, *uribismo*, *massismo*, *madridismo*, *emprendedurismo*, *agnosticismo*, *biomagnetismo*, *ateísmo*, *vagivismo*, *veganismo*, *ventajismo*, *maerismo*, *bruxismo*, *cateterismo*.

Las palabras que terminan en *-ción* y aparecen al menos 700 veces en el corpus nuevo y más amplio, pero que prácticamente desaparecen en el corpus más antiguo y reducido son: *autenticación*, *fidelización*, *virtualización*, *mercantilización*, *victimización*, *geolocalización*, *encriptación*, *dispensación*, *precarización*, *suplementación*, *judicialización*, *visibilización*, *abducción*, *desafección*, *compartición*, *cosificación*, *procrastinación*, *remediación*, *iteración*, *resignificación*, *disrupción*, *deslocalización*, *feminización*, *extranjerización*, *redirección*, *evicción*, *invisibilización*, *bancarización*, *demonización*, *gamificación*, *previsualización*, *propiciación*.

Esto no quiere decir que todas estas palabras constituyan neologismos, solo que rara vez aparecen en el corpus más antiguo y pequeño y que, por el contrario, resultan frecuentes en el corpus más moderno y de mayor dimensión. Los lexicógrafos experimentados, por supuesto, querían investigar tales palabras con más detalle. Pero el punto principal es que el nuevo corpus de dos mil millones de palabras proporciona, además, los datos para la investigación de voces donde ha habido un relativo aumento de frecuencia en los últimos 15–20 años.

14 Desvío sincrónico # 1: dialectos

Además de identificar neologismos, una de las ventajas del nuevo corpus es que permite a los investigadores obtener la frecuencia de cada palabra o frase en 21 países diferentes de habla hispana⁶. Para algunos países, como España, el nuevo

corpus contiene hasta 459 millones de palabras (con 261 para México, 183 para Argentina y 180 para Colombia). Incluso entre los países con menor representación cuantitativa en el corpus, cada uno tiene al menos 30 millones de palabras (p. ej., 39 para El Salvador, 39 para Honduras, 37 para la República Dominicana, 36 para Puerto Rico, 35 para Nicaragua, 33 para Paraguay y 32 para Costa Rica). Las siguientes son solo algunas de las palabras que se registran con una frecuencia más alta en un país que en los otros⁷:

- Caribe
 - Puerto Rico *ay bendito*, *chavos*, *chirringa*, *mahones*, *habichuela* (+DR), *zafacón* (+DR); Cuba *guajiro*, *jimaguas*, *babalao*, *bitongo*, *pedir botella*; Rep. Dominicana *facú*, *tutumpote*, *manguilina*, *mofongo* (+PR).
- México y América Central
 - México *ándale*, *híjole*, *órale*, *güero*, (*muy*) *padre*, *chamaco* (Cam/Car), *piriche*, *popote*, *charola* Guatemala *huppil*, *canche*, *muchá*, *patojo*, *chafa* (+HN), *chirmol*, *canche*; El Salvador *cipote*, *chero*, *pupusa*, *cuillo*, *bayunco*, *piscuchá*; Honduras *catracho*, *papada*; Nicaragua *chavalo*, *majé* (+Cam), *pinol*, *pinolillo*, *chigüín*, *vigorón*, *gallo pinto* (+CR), *idlay* (+CR) Panamá *fulo*, *chombo*, *guan-dul*; Costa Rica *chinear*, *guila*, *chunche*.
- América del Sur
 - Colombia *cachaco*, *cachifo*, *verraquera*, *estar marnado*, *guandoca*, *biche*; Venezuela *bojote*, *coroto*, *catre*, *gafó*, *macundales*, *arepa*, *cachapa*, *canbur*, *carraotas*, *jojoto*; Ecuador *chumar*, *chulla*, *montuvio*, *omoto*; Perú *anticucho*, *jebe*, *chupe*, *pisco*, *jora*, *chompa* (+CL/EC), *choclo* (+CL/EC); Bolivia *opa*, *colla*, *chuno*, *lagar*; Chile *pololo**, *pololear*, *achuntar*, *bencina*, *bacón*, *jome*, *huaso*; Paraguay *ñembo*, *ñanduti*, *karai*, *yopará*, *mitái*; Uruguay *tropero*, *hacer** *sota*, *con fritas*; Argentina *pibe*, *fiaca*, *morfar*, *falopa*, *sobre el pucho*, *falluta*, *cafishi*.
- España *ordenador*, *aparcar*, *enfadar*, *gafas*, *zumo*, *chulo*, *guay*, *coger*, *boligráfico*, *patata*, *melocotón*, *echar de menos*, *vale*.

Además de poder analizar la frecuencia de una palabra o frase específica en los 21 países, también es posible hacer que el corpus produzca una lista de todas las palabras que son más comunes en un país (o conjunto de países) que en otro. P. ej., la siguiente tabla muestra palabras de *-ismo* que son más comunes en Venezuela (izquierda) que en Colombia, México, Argentina o España (derecha):

⁷ Téngase en cuenta que los países y la región entre paréntesis indican que la palabra también tiene una alta frecuencia relativa en estas otras zonas. La lista de palabras se corresponde con los datos léxicos aportados por Lipski 1996.

⁶ Para datos similares de un corpus de inglés análogo, vid. Davies/Fuchs 2015.

SECC 1 (Venezuela): 98,170,248 WORDS				SECC 2 (Argentina, Colombia, España...): 1,008,426,240 WORDS							
WORD/PHRASE	TOKENS	PM	PM/2	PM/1	RATIO	WORD/PHRASE	TOKENS	TOKENS/1	PM/2	PM/1	RATIO
1. JAABOLISMO	44	1	0.4	0.0	452.0	1. MASSISMO	1249	1	1.2	0.0	121.6
2. ESCLAUDISMO	21	1	0.2	0.0	215.7	2. PRISISMO	362	0	0.4	0.0	35.9
3. MADURISMO	141	10	1.4	0.0	144.8	3. EMPRENDIDORISMO	251	0	0.2	0.0	24.9
4. GOMECISMO	41	3	0.4	0.0	140.4	4. GARANTISMO	254	1	0.3	0.0	24.7
5. VENTAJISMO	812	63	8.3	0.1	132.4	5. PARKINSONISMO	219	0	0.2	0.0	21.7
6. PUNTOFISISMO	165	14	1.7	0.0	121.1	6. NEOPLATONISMO	218	1	0.2	0.0	21.2
7. RIQUEZISMO	22	3	0.2	0.0	75.3	7. CHARRISMO	204	1	0.2	0.0	19.9
8. HIRISMO	295	41	3.0	0.0	73.9	8. SOCORRISMO	196	0	0.2	0.0	19.4
9. CASTRO-COMUNISMO	55	8	0.6	0.0	70.6	9. LIBERTARISMO	184	0	0.2	0.0	18.2
10. ECO-SOCIALISMO	27	4	0.3	0.0	69.3	10. CRISTIANISMO	527	3	0.6	0.0	18.1
11. VENEZOLANISMO	27	4	0.3	0.0	69.3	11. COPENICANISMO	173	0	0.2	0.0	17.2
12. CASTROCOMUNISMO	58	9	0.6	0.0	66.2	12. MOTOTAXISMO	171	0	0.2	0.0	17.0

Gráfico 19: Comparación por dialecto en el corpus Web/Dialecto

15 Comparación con el CORPES

Al analizar los datos históricos del *Corpus del Español* original, los comparamos con los datos del corpus *CORDE* y, cuando discutimos los datos sincrónicos del *Corpus del Español* original, lo comparamos con el corpus *CREA*. Existe un tercer corpus de la Real Academia Española que es relevante, en términos de una comparación con la parte más moderna de dos mil millones de palabras «Web/Dialectos» del *Corpus del Español*, que hemos discutido en las Secciones 12–14; este tercer corpus es el *CORPES* (*Corpus del Español del Siglo XXI*).

En cuanto a posibilidades de estudios sobre léxico, existen algunas diferencias importantes entre el *CORPES* y la nueva extensión de dos mil millones de palabras del *Corpus del Español* (que llamaremos *CDE-2* en esta sección). La primera diferencia importante es el tamaño: el *CORPES* tiene aproximadamente 175 millones de palabras, que es menos del 10 % del tamaño del *CDE-2*. Esto es, en caso de que una palabra tenga 200 ocurrencias en *CDE-2* (probablemente un número suficiente para muchos tipos de investigación), en el *CORPES* probablemente tendría menos de 20 ocurrencias, lo que resulta mucho más problemático. En segundo lugar, los materiales que integran el *CDE-2* son más recientes, lo que resulta importante para detectar neologismos. Solo alrededor del 17 % del *CORPES* es posterior a 2010, mientras que el 100 % de los datos del *CDE-2* corresponde a ese período (los textos se recopilieron en 2014–2015).

También hay diferencias importantes en cuanto a la arquitectura e interfaz del corpus. Aunque el *CORPES* puede generar gráficos de frecuencia útiles para palabras individuales o para una frase exacta, no es posible encontrar la frecuencia de cadenas coincidentes en una búsqueda como «*menos * que*» (p. ej., *menos valor/intenso/arriesgado/sano que*) o «*adjetivo + ojos*» (p. ej., *ojos bellos/mublados/*

crystalinos/carinosos). En relación con esto, con el *CORPES* no es posible obtener la frecuencia de colocaciones significativas (es decir, es típico del *CORPES* mostrar artículos o preposiciones como las colocaciones más importantes de una palabra, como es también el caso del *CREA* y el *CORDE*, v. el Gráfico 10). El *CORPES* tampoco permite la comparación de colocaciones entre palabras diferentes (p. ej., *potente y poderoso*, o *iluminar y alumbrar*), aunque puedan ser búsquedas útiles para poner de manifiesto las diferencias de significado. Y, por último, el *CORPES* no permite comparaciones entre todas las palabras en diferentes países para encontrar las que son más comunes en un país que en otro (v. Gráfico 19). Por tanto, aunque el *CORPES* tiene una arquitectura e interfaz más avanzada que el *CORDE* o el *CREA*, todavía resulta bastante limitada en los tipos de búsquedas que permite.

16 Datos muy recientes y datos continuamente actualizados

En 2016 lanzamos un corpus de textos en inglés llamado *NOW* («Noticias en la web»). Este corpus crece automáticamente en tamaño en alrededor de 5–6 millones de palabras por día y está basado en aproximadamente 10 000 nuevas URL diarias extraídas de Google News. Con este corpus, pueden rastrearse los cambios en la frecuencia y el uso de palabras en el transcurso de meses, semanas e incluso días, lo que es obviamente muy útil para observar cambios extremadamente recientes en el lenguaje. P. ej., los investigadores pueden rastrear la frecuencia en el tiempo de las palabras «nuevas» en el idioma (que han surgido a partir de 2007), como *Brexit*, *manspreading*, *makerspace*, *gig economy*, *dadbod*, *monger*, *swatting*, *walkscore*, *trigger warning*, *mommy porn*, *normcore*, *listicle*, *suffertest*, *catfishing*, *sapiosexual*, *nomophobia*, *omnishambles*, *humblebrag*, *FOMO*, *precarriat*, *filter bubble*, *range anxiety*, *collaborative consumption*, *churrullism*, *birther*, *truthter*, *staycation*, *glamping*, *locavore*, *voluntourism*, *freggan*.

Hemos recopilado los textos para un corpus similar en español (que llamaremos aquí *NOW-Español*), que se lanzará en mayo de 2018. Cuando se publique contendrá aproximadamente 4800 millones de palabras desde enero de 2012 hasta mayo de 2018, y luego crecerá en tamaño alrededor de 140 millones de palabras cada mes (alrededor de 1600 millones de palabras cada año). De la misma manera que ya es posible para el inglés, los usuarios de *NOW-Español* podrán seguir la frecuencia de cualquier palabra, frase o construcción sintáctica a lo largo del tiempo (incluso a nivel de meses y semanas), y también tendrán las listas de neologismos en los textos que generará el corpus automáticamente.

Aunque *NOW-Español* aún no está disponible, podemos utilizar algunos datos preliminares para mostrar cómo *NOW-Español* se compara con la extensión de

dos mil millones de palabras en el *Corpus del Español* que, como hemos indicado, contiene textos recopilados de 2014–2015. P. ej., las siguientes palabras que terminan en *-idad* aparecen al menos 10 veces en NOW-Español, pero no figuran en los datos del *Corpus del Español* de 2014–2015: *poliautoinmunidad*, *banabilidad*, *promovilidad*, *superlatividad*, *sucrosidad*, *sucrenidad*, *inexportabilidad*, *teleguridad*. Las palabras terminadas en *-ismo* incluyen *narcouribismo*, *uribestialismo*, *cañavismo*, *cuñadismo*, *carnismo*, *larretismo*, *raspacapismo*, *chaveztidismo*, *urbanadismo*, *figurinismo*. Y las palabras terminadas en *-ción* incluyen *narcorevolución*, *uberización*, *multihabitación*, *posdesmovilización*, *incoración*, *cibervictimización*, *kilometración*, *deafmentación*, *electrorreducción*, *multilateración*, *desdibolización*, *termonebulización*, *reternalización*, *hipsterización*. Una vez más, no hay garantía de que todas estas palabras sean neologismos (es decir, que sean novedades de los últimos 3–4 años), solo de que aparecen en el español actual y no figuran en el *Corpus del Español* de dos mil millones de palabras. Los lexicógrafos tendrán acceso a esta rica información para observar los cambios más recientes en el lenguaje.

17 Conclusión

Los grandes corpus accesibles mediante Internet han ayudado a revolucionar el campo de la lexicografía histórica. Con solo unos pocos clics del ratón, los investigadores pueden buscar corpus que contengan cientos de millones (y ahora miles de millones) de palabras de miles (y ahora millones) de textos.

Cada corpus tiene sus propias fortalezas y debilidades. Los corpus *CORDE* y *CREA* de la Real Academia Española son bastante sólidos en términos de «corpus textual», pero tienen arquitecturas e interfaces muy limitadas y obsoletas que limitan el acceso de los investigadores a estos datos. Por otro lado, la parte «histórica/de género», original del *Corpus del Español* (que fue lanzado en 2002) constituye un corpus de dimensiones más reducidas que el *CORDE* o el *CREA*, pero permite varios tipos de investigación que no pueden llevarse a cabo con estos dos corpus, lo que incluye recuperar la frecuencia de palabras, frases y subcadenas (p. ej., prefijos o sufijos) por siglo o género, seguir los patrones de cambio de colocaciones para observar las modificaciones en el significado y uso, y utilizar el índice integrado y las listas de palabras personalizadas para ver cómo las voces compiten por el «espacio semántico» a lo largo del tiempo (o en diferentes géneros).

Uno de los usos más interesantes de los grandes corpus en línea es poner de relieve los cambios muy recientes en el lenguaje, como puede hacerse con la nueva extensión del *Corpus del Español*, que contiene dos mil millones de palabras de

21 países diferentes de habla hispana de 2014–2015. A partir de mayo de 2018, será posible utilizar NOW-Español para analizar los cambios que se producen virtualmente en «tiempo real», como ya es posible para NOW-English. En algún momento, incluso, sería posible crear corpus que monitoricen continuamente cómo las palabras y frases se propagan a lo largo del tiempo a través de distintas comunidades de habla y dialectos.⁸

En definitiva, los grandes corpus en línea nos permiten realizar muchos tipos de investigaciones que apenas eran imaginables hace 15 o 20 años y, con las mejoras y avances en la tecnología, podemos darnos cuenta de que solo estamos en el comienzo de lo que se puede llegar a hacer.

Referencias bibliográficas

- Davies, Mark (2002): «Un corpus anotado de 100.000.000 de palabras del español histórico y moderno», en *Actas del XVIII Congreso de la Sociedad Española para el Procesamiento del Lenguaje Natural*. Valladolid: SEPLN, 21–27.
- Davies, Mark (2005a): «Advanced research on syntactic and semantic change with the Corpus del Español», en Claus Fusch *et al.* (eds.), *Romance Corpus Linguistics II: Corpora and Diachronic Linguistics*. Tübingen: Gunter Naar, 203–214.
- Davies, Mark (2005b): «The advantage of using relational databases for large corpora: speed, advanced queries, and unlimited annotation», *International Journal of Corpus Linguistics* 10, 301–328.
- Davies, Mark (2008): «Spanish and Portuguese corpus linguistics», *Studies in Hispanic and Lusophone Linguistics* 1, 149–186.
- Davies, Mark (2010): «Creating useful historical corpora: A comparison of *CORDE*, the Corpus del Español, and the Corpus do Português», en Andrés Enrique-Arias (ed.), *Diachronía de las lenguas iberoromances: nuevas perspectivas desde la lingüística de corpus*. Fráncfort/Madrid: Vervuert/Iberoamericana, 137–166.
- Davies, Mark (2012a): «Expanding horizons in historical linguistics with the 400 million word Corpus of Historical American English», *Corpora* 7, 121–157.

8 Véase la Sección 7 de Davies 2015.

- Davies, Mark (2012b): «Examining Recent Changes in English: Some Methodological Issues», en Terttu Nevalainen y Elizabeth Closs Traugott (eds.), *The Oxford Handbook of the History of English*. Oxford: Oxford Univ. Press, 263–287.
- Davies, Mark (2015): «Corpora: An introduction», en Douglas Biber y Randi Reppen (eds.), *Cambridge Handbook of English Corpus Linguistics*. Cambridge: Cambridge University Press, 11–31.
- Davies, Mark (en prensa a): «Using large online corpora to examine lexical, semantic, and cultural variation in different dialects and time periods», en Eric Friginal et al. (eds.), *Corpus-Based Sociolinguistics*. Londres: Routledge.
- Davies, Mark (en prensa b): «Corpus-based studies of lexical and semantic variation: The importance of both corpus size and corpus design», en Carla Subir, Terttu Nevalainen e Irma Taavitsainen (eds.), *From data to evidence in English language research* (Digital Linguistics). Leiden: Brill.
- Davies, Mark/Robert Fuchs (2015): «Expanding Horizons in the Study of World Englishes with the 1.9 Billion Word Global Web-Based English Corpus (GloWbE)», *English World-Wide* 36, 1–28.
- Firth, J. R. (1957): *Papers in Linguistics 1934–1951*. Londres: Oxford University Press.
- Lipski, John (1996): *El español de América*. Madrid: Cátedra.

Virginia Bertolotti y Concepción Company Company

El corpus para América: *CORDIAM*

Resumen: El *Corpus Diacrónico y Diatópico del Español de América (CORDIAM)* fue creado para poder historiar el español de América y para hacer lingüística histórica general con datos de muchas variedades del español. Estos datos, exclusivamente americanos, se procesan a través de una interfaz amigable y eficiente, producto del trabajo conjunto de informáticos e investigadores en lingüística histórica e historia de la lengua. Este trabajo muestra las características lingüísticas y textuales del *CORDIAM*, sus características informáticas y da cuenta, además, del proceso de toma de decisiones para la creación de cada uno de los tres subcorpus que conforman el *CORDIAM: CORDIAM-Documentos, CORDIAM-Prensa* y *CORDIAM-Literatura*. Muestra también cómo los textos incluidos representan facetas diversas de la cultura en América.

Palabras clave: Lingüística de corpus, Corpus escrito, Lingüística histórica, Español en América

Abstract: The *Corpus Diacrónico y Diatópico del Español de América (CORDIAM)* is presented in this paper. It was created to make the history of Spanish in America and historical linguistics with data from many dialects of Spanish. These exclusively American data are processed by a friendly and efficient interface, which is the result of the joint work of computer engineers and investigators both in historical linguistics and in language history. This work shows the linguistic and textual characteristics of *CORDIAM*, its computing characteristics and, in addition, explains the decisions underlying process for the creation of each of the three subcorpus of *CORDIAM*, which are *CORDIAM-Documents*, *CORDIAM-Press*, and *CORDIAM-Literature*. It also shows how the included texts exhibit different aspects of the American culture.

Keywords: Corpus linguistics, Written corpus, Historical linguistics, Spanish in America

1 Presentación y objetivos

Los *corpus lingüísticos*, escritos y orales, diacrónicos y sincrónicos, las *ediciones críticas de textos*, antiguos y modernos, así como los *atlas lingüísticos*, son, como es sabido, los tres tipos de productos fundamentales que constituyen infraestructura para la investigación lingüística, y son, por ello, soportes esenciales para abrir nuevos horizontes de investigación, hallar nuevas evidencias, descripciones