

*Lenguas
y computación*

Antonio Moreno Sandoval



EDITORIAL
SINTESIS

Índice

Prólogo	9
1. Introducción al tratamiento computacional de las lenguas	13
1.1. El modelo computacional del lenguaje.....	14
1.1.1. <i>Relaciones multidisciplinares</i>	17
1.2. Evolución histórica de la lingüística computacional.....	19
1.3. ¿Qué se puede hacer con los ordenadores?	22
1.3.1. <i>Procesamiento de voz</i>	23
1.3.2. <i>Otras aplicaciones de la lingüística computacional</i>	25
1.4. Arquitectura de un sistema de procesamiento del lenguaje natural.....	25
1.4.1. <i>Conocimiento lingüístico y procesamiento informático</i>	28
1.4.2. <i>Análisis y generación</i>	29
1.4.3. <i>Procesamiento basado en la sintaxis y procesamiento basado en la semántica</i>	30
Ideas fundamentales del capítulo.....	31
2. Modelos simbólicos	33
2.1. Introducción	33
2.1.1. <i>Perspectiva histórica</i>	33
2.1.2. <i>Algunas características de los modelos simbólicos</i>	36
2.2. Gramáticas formales.....	38
2.2.1. <i>Limitaciones de las gramáticas formales</i>	40
2.2.2. <i>La gramaticalidad</i>	41
2.2.3. <i>Requisitos de una gramática formal</i>	42
2.2.4. <i>Tipos de gramáticas formales</i>	43
2.3. Gramáticas, lenguas y autómatas	46
2.3.1. <i>Gramáticas generativas</i>	46
2.3.2. <i>La jerarquía de Chomsky</i>	48
2.3.3. <i>Gramáticas irrestrictas</i>	50

2.3.4. Gramáticas dependientes del contexto	50
2.3.5. Gramáticas independientes del contexto	51
2.3.6. Gramáticas regulares o de estados finitos	51
2.3.7. La adecuación formal de las gramáticas	52
Ideas fundamentales del capítulo.....	55
3. Ejemplos y limitaciones de los modelos simbólicos	57
3.1. Un ejemplo de gramática formal básica del español.....	57
3.1.1. Gramática 1: independiente del contexto	57
3.1.2. Gramática 2: unificación y rasgos	67
3.2. Limitaciones de las gramáticas formales en el procesamiento automático.....	77
Ideas fundamentales del capítulo.....	79
4. Modelos estadísticos	81
4.1. Conceptos básicos de probabilidad y estadística.....	83
4.1.1. Probabilidad y estadística.....	83
4.1.2. Relevancia de la estadística para el estudio del lenguaje.....	93
4.1.3. Teoría de información aplicada al lenguaje natural	97
4.2. Técnicas estadísticas: estimación y evaluación de la probabilidad.....	101
4.2.1. ¿Cómo se calcula la probabilidad de una unidad?	101
4.2.2. Probabilidad condicionada	105
4.2.3. Evaluación de la predicción.....	109
4.2.4. Cálculo de n-gramas	112
Ideas fundamentales del capítulo.....	115
5. Aplicaciones y limitaciones de los modelos estadísticos	117
5.1. Tres tipos de aplicaciones de los modelos estadísticos en la lingüística computacional.....	117
5.1.1. Estadística de corpus	118
5.1.2. Modelos generativos con probabilidades.....	121
5.1.3. Aprendizaje automático: clasificadores y clusters.....	129
5.1.4. Herramientas para desarrollar modelos predictivos	139
5.2. Limitaciones de los modelos estadísticos.....	140
Ideas fundamentales del capítulo.....	143
6. Análisis morfológico	145
6.1. Análisis morfológico y su descripción.....	146
6.1.1. Localidad, alomorfía y procesos concatenativos y no concatenativos.....	148

Índice

6.2. Métodos de procesamiento morfológico.....	155
6.2.1. <i>Analizadores morfológicos basados en estados finitos</i>	156
6.2.2. <i>Analizadores morfológicos basados en paradigmas</i>	160
6.3. La desambiguación en los análisis morfológicos.....	163
6.4. Algunos casos problemáticos.....	165
6.4.1. <i>Las unidades léxicas ambiguas</i>	166
6.4.2. <i>Las unidades léxicas pluriverbales</i> (multiword expressions).....	166
6.4.3. <i>Las palabras desconocidas</i>	168
Ideas fundamentales del capítulo.....	170
7. Análisis sintáctico	171
7.1. El problema de la asignación de estructuras sintácticas.....	171
7.1.1. <i>Robustez</i>	173
7.1.2. <i>Desambiguación</i>	174
7.1.3. <i>Acierto</i>	179
7.2. Dos modelos de representación de estructuras sintácticas.....	181
7.3. <i>Treebanks</i> y la inducción gramatical.....	188
7.4. Estructuras relevantes para el análisis sintáctico.....	192
7.4.1. <i>Estructuras de constituyentes básicas</i>	193
7.4.2. <i>Estructuras de dependencias básicas</i>	195
7.4.3. <i>Estructuras sintácticas características del español</i>	196
Ideas fundamentales del capítulo.....	199
8. Interpretación: semántica y pragmática	201
8.1. Conceptos básicos de representación del significado.....	202
8.1.1. <i>¿Qué es el significado y cómo se puede representar?</i>	202
8.1.2. <i>Problemas interpretativos: ambigüedad, metáfora</i> <i>e ironía</i>	204
8.2. Tratamiento computacional de la semántica oracional.....	206
8.3. Tratamiento computacional de la semántica léxica.....	208
8.3.1. <i>Modelos simbólicos de representación léxica</i>	208
8.3.2. <i>Modelos estadísticos de representación léxica</i>	215
8.3.3. <i>Comparación entre ambos modelos</i> <i>de representación léxica</i>	218
8.4. Tratamiento computacional del discurso	220
Ideas fundamentales del capítulo.....	226
9. Traducción automática	227
9.1. Evolución de los sistemas de traducción automática.....	230
9.1.1. <i>Los primeros tiempos de la traducción automática</i>	230

9.1.2. <i>Los primeros sistemas comerciales</i>	231
9.1.3. <i>La revolución de la traducción automática estadística</i>	233
9.1.4. <i>La traducción automática en el nuevo milenio</i>	233
9.2. <i>Distintas aproximaciones a la traducción automática</i>	234
9.2.1. <i>La traducción automática basada en reglas</i>	234
9.2.2. <i>La traducción automática estadística</i>	236
9.3. <i>Problemas y limitaciones actuales de la traducción automática</i>	238
9.3.1. <i>La ambigüedad léxica</i>	238
9.3.2. <i>La ambigüedad sintáctica</i>	239
9.3.3. <i>Los aspectos de estilo (pragmáticos) de la traducción</i>	239
9.4. <i>Evaluación de la calidad de las traducciones automáticas</i>	240
9.4.1. <i>Evaluación humana de la calidad de las traducciones</i>	240
9.4.2. <i>Evaluación automática de la calidad de las traducciones</i>	242
9.5. <i>Sistemas más populares en la actualidad</i>	242
9.5.1. <i>Systran</i>	243
9.5.2. <i>Apertium</i>	244
9.5.3. <i>Google Translate</i>	245
9.5.4. <i>Moses</i>	246
Ideas fundamentales del capítulo.....	248
10. Gestión de la información y analítica de textos	249
10.1. <i>Recuperación y extracción de información</i>	250
10.1.1. <i>Técnicas generales</i>	253
10.1.2. <i>Búsqueda de respuestas</i>	263
10.1.3. <i>Generación automática de resúmenes</i>	267
10.2. <i>Minería de datos y analítica de textos</i>	269
10.2.1. <i>Clasificación de textos</i>	270
10.2.2. <i>Análisis de opiniones y sentimientos</i>	271
10.3. <i>Perspectivas de la gestión de información: campos de aplicación</i>	273
Ideas fundamentales del capítulo.....	277
Conclusiones: pasado, presente y futuro	279
Lingüistas e informáticos.....	280
Relación entre LC y otras ramas de la lingüística	281
Bibliografía seleccionada	285

Prólogo

Si utilizásemos una metáfora culinaria, podríamos decir que este no es un libro de recetas de cocina, sino de consejos generales para escoger apropiadamente el menú en función de los comensales.

En 1998, la Editorial Síntesis me dio la oportunidad de publicar el libro *Lingüística computacional: una introducción a los modelos simbólicos, estadísticos y biológicos*, de cuyo prólogo están tomadas estas palabras. Dos décadas después, y dentro de una nueva colección dirigida por el profesor Juan Carlos Moreno Cabrera, se me brinda de nuevo la ocasión para hablar de lenguas y sistemas informáticos que procesan lenguas. Para mí es un tema apasionante, al que he dedicado ya unos cuantos años y me gustaría dedicar otros tantos. La disciplina cuenta con más de sesenta años de desarrollo y he tenido la oportunidad de conocer a alguno de sus grandes actores, cuyo magisterio es necesario reconocer y transmitir a las próximas generaciones. Este reconocimiento es justo y útil al mismo tiempo, pues nos permitirá entender mejor los éxitos y fracasos de las diferentes estrategias que han abordado el tratamiento automático de las lenguas naturales.

La lingüística computacional es ya una disciplina madura con un enfoque teórico (o de ciencia básica y fundamental) y otro aplicado (o de ingeniería). La rama fundamental está aportando cada vez más a la lingüística teórica, no solo en conceptos y herramientas de experimentación con teorías formales, sino también en métodos cuantitativos y estadísticos para estudiar los fenómenos de variación. La fabulosa potencia de análisis de datos que proporcionan las aplicaciones computacionales que se utilizan en la lingüística de corpus es un buen ejemplo de ello. La rama ingenieril, por su parte, está contribuyendo a la lingüística aplicada de diversas maneras: con programas que ayudan a la traducción, a los terminólogos y lexicógrafos, a los lingüistas forenses, a los sociolingüistas y a los psicolingüistas. La enseñanza de lenguas por ordenador es otro de los campos que se han beneficiado de los avances de la lingüística computacional.

El objetivo principal de este libro es proporcionar el contexto necesario para entender la relación entre todas las ramas de la lingüística. El futuro de las ciencias del lenguaje está, creo, en la apertura metodológica y conceptual a las aportaciones de diferentes especialidades. La historia de la lingüística computacional proporciona excelentes ejemplos de colaboración multidisciplinar entre lingüistas, informáticos, lógicos, psicólogos, traductores, ingenieros y matemáticos. La complementariedad permite superar los desconocimientos individuales. Y esto es necesario porque la complejidad de formalizar el lenguaje para tratarlo computacionalmente supera los conocimientos individuales, como tendremos ocasión de comprobar a lo largo del libro.

Mi libro de 1998 estaba dirigido a estudiantes de Lingüística y Computación que estuvieran interesados en una introducción al estado de la cuestión a finales de los noventa. El libro que tiene el lector en sus manos está dirigido fundamentalmente a lingüistas, traductores, psicólogos y en general cualquier humanista o científico social que quiera introducirse en los conceptos y métodos para poder interactuar con informáticos. Durante años, en mis exámenes finales de la asignatura Lingüística Computacional, formulaba la siguiente pregunta:

¿Qué aproximación le parece más acorde con sus intereses: formarse como lingüista y entenderse con los informáticos o formarse como lingüista computacional autónomo?

La respuesta más frecuente era la primera opción y coincido con esa apreciación. Formarse como lingüista computacional integral supone adquirir conocimientos sustanciales de diferentes especialidades: formalización lingüística en múltiples aspectos, desde el acústico hasta el interpretativo, programación informática eficiente, robusta y distribuida, conocimientos lógicos para la interpretación, así como computación estadística para modelar el aprendizaje automático. Sencillamente, el campo se ha ampliado tanto que es casi imposible abarcarlo de forma unitaria e individual. Es un signo claro de la madurez de la disciplina.

Si en los años ochenta la investigación estaba dirigida por los lingüistas gracias a las aportaciones de las gramáticas de unificación y rasgos, a partir de los noventa la disciplina ha estado dominada por los informáticos y, más en concreto, por los modelos estadísticos. Este cambio metodológico sustancial se puede comprobar en los contenidos de los manuales más conocidos: Grishman (1986) o Gazdar y Mellish (1989) frente a Jurafsky y Martin (2000; 2008). La revolución de los modelos estadísticos ha transformado la lingüística computacional y este libro persigue dotar a los lectores con conocimientos para entender sus conceptos, herramientas e implicaciones. El libro de 1998 incluía una parte dedicada a los modelos estadísticos, pero entonces estaban todavía en sus inicios y ahora presentamos una actualización que recoge la madurez de lo que se ha investigado en los últimos veinte años.

Por otra parte, lo que el lector tiene delante no es un libro de texto con ejercicios para aprender de manera pautada a elaborar gramáticas computacionales, lexicones, corpus anotados, etc. Un buen ejemplo de este tipo de manuales es el de Bird, Klein y Loper (2009). Tampoco es una enciclopedia con capítulos escritos por diferentes autores especialistas en sus dominios, para dar una visión lo más general y completa posible de todo lo que concierne a la lingüística computacional (ejemplos de esto último son los *handbook* de Jurafsky y Martin, 2008, o de Indurkha y Damerou, eds., 2010).

Este libro pretende ser una fuente para conocer los conceptos y métodos claves en el campo. Se complementa con referencias especializadas y selectas en cada caso. Una parte significativa del libro tratará la conexión entre la lingüística computacional y otras áreas de la lingüística, tanto teórica como aplicada.

La estructura del libro se divide en dos partes. En la primera se abordan los fundamentos, que incluyen los modelos simbólicos y modelos estadísticos (dividido cada uno en dos partes para guardar un equilibrio con la extensión de los otros capítulos), así como los diferentes niveles desde la morfología a la interpretación. La segunda parte se dedica a las aplicaciones: dos capítulos se dedican a la traducción automática y a las distintas aplicaciones de la gestión de información (recuperación, extracción y análisis de opinión). En el último capítulo, de conclusiones, se repasan los hallazgos y limitaciones actuales para hacer luego una previsión de los desarrollos futuros de la disciplina.

Me he tomado la libertad, por sugerencia de la editorial, de traducir al español todas las citas de otros autores, con el objetivo de facilitar la lectura.

Finalmente, unas palabras de agradecimiento. Si bien mis maestros y compañeros en los proyectos iniciales se merecen mi sincero reconocimiento (se aprende más de los pares que de los maestros, como nos decía en clase un catedrático de Teoría de la Literatura), mis últimos veinte años como investigador responsable de los proyectos del Laboratorio de Lingüística Informática de la UAM me han hecho cambiar de perspectiva: los discípulos son los verdaderos responsables de que este libro tenga un alcance multilingüístico y multinivel. Sus tesis doctorales me han permitido trabajar (y aprender) con ellos sobre temas que se habrían escapado a mis posibilidades o intereses. Corren malos tiempos para la ciencia y debo manifestar mi convicción de que la parte más sustancial de la investigación científica que se desarrolla en los departamentos y laboratorios universitarios la realizan los becarios y estudiantes de doctorado. El relevo generacional es imprescindible. Sin ellos, no habrá progreso en el conocimiento futuro.