

PROPUESTA PARA UN ANALISIS AUTOMATICO DEL CONTENIDO EN DOCUMENTACION DE EST

Félix de Moya (Universidad de Granada)
José Manuel Muñoz (Universidad de Córdoba)
Pedro Hípola (Universidad de Granada)

Moya Anegón, F.; Muñoz Muñoz, J. M.; Hípola, P. (1988). [Propuesta para un análisis automático del contenido en documentación de EST](#). In M. Laurén, Christer; Nordman (Ed.), *Special Language: From Humans Thinking To Thinking Machines* (pp. 448-452). Clevedon: Multilingual Matters Ltd. Clevedon.

<http://books.google.es/books?id=oKRVbYGv3fEC&lpg=PA3&ots=r82HshWzX4&lr&hl=es&pg=PA448#v=onepage&q&f=false>

Las características del sistema que presentamos en estas páginas han determinado la participación en él de un grupo de trabajo con carácter marcadamente interdisciplinar. Se ha tratado de hacer confluír la tarea de estudio sobre materiales de EST para fines docentes con las técnicas usuales en indexación automática. Los componentes del grupo coinciden en unos intereses de trabajo comunes, encaminados a obtener el aislamiento de términos especialmente relevantes dentro de textos de Ciencia y Tecnología escritos en inglés, que sirvan en primer lugar para el establecimiento y puesta al día de vocabularios utilizables en la enseñanza del inglés con fines específicos, los cuales a su vez puedan ser utilizados dentro de una aplicación que se encargue de realizar la indexación documental de forma automática sobre palabras clave extraídas de los textos escritos que se presenten al sistema.

La metodología empleada para cubrir los objetivos definidos se basa en procedimientos usuales dentro del ámbito de la estadística lingüística, especialmente los más extendidos en el tratamiento de unidades léxicas, completados con la utilización del análisis morfográfico necesario. En la base del funcionamiento del sistema se encuentra el recurso a la comparación de las frecuencias del léxico empleado en textos ingleses de Ciencia y Tecnología con las frecuencias del vocabulario de la lengua standard, para efectuar así el aislamiento de los términos que se desea localizar. Establecidos de manera semiautomática repertorios léxicos pertenecientes a dominios temáticos específicos, los vocabularios obtenidos pueden ser utilizados como objeto de estudio lingüístico, como auxiliar en la enseñanza de "English for Specific Purposes" y además como término de comparación para el siguiente objetivo, el de la indexación automática.

Por este tipo de procedimiento se facilita en primer lugar que la selección del vocabulario que se va a usar en la enseñanza de ESP se realiza no sólo a través de un criterio de estimación subjetiva por parte del profesor, sino con la ayuda de un

estudio objetivo, que le orienta en la determinación de qué tipo de material se debe emplear y además le ofrece datos para establecer prioridades en el léxico que se enseña. Por otra parte, desde el punto de vista documental el análisis que se produce podrá ser clasificado, dentro del esquema que ha señalado Van Slype², como "indización automática selectiva en lenguaje natural".

Se ha recurrido a algoritmos y criterios de ordenación del material basados en la ley de Zipf y en la fórmula del test de Pearson, además de diversos cálculos de probabilidades. Las herramientas informáticas empleadas están constituidas por un sistema de gestión de base de datos relacional. Para preparar las rutinas de tratamiento del material lingüístico se ha utilizado el lenguaje de comandos incorporado en el mismo gestor.

Presentamos en esta comunicación el proceso de preparación de uno de los vocabularios de tema específico, el correspondiente a textos sobre temas informáticos, realizado a partir de un corpus de más de doscientos cincuenta mil tokens, del que, desechadas fórmulas, nombres de variables, etc., quedaron aproximadamente unas ciento cincuenta mil formas léxicas. El interés de contar con estudios serios de vocabulario es especialmente patente en el área temática de la computación, por el alto nivel de especificidad léxica que caracteriza a los textos de este ámbito, su considerable ritmo en la creación de nuevos vocablos y por la gran cantidad de préstamos que en este dominio científico suelen tomar del inglés y otras lenguas, entre ellas el español.

Tras llevar a cabo la recopilación de los textos que componen el corpus, se realizó la carga de los mismos y, previa depuración formal de diversos caracteres, se efectuó el recuento de frecuencias de cada una de las formas léxicas³. Independientemente de esto, se ha procedido a la carga del diccionario de frecuencias elaborado por Kucera y Francis⁴. Los datos así extraídos fueron incorporados al sistema completo de base de datos. Por otra parte, de forma manual, se ha introducido en cada registro un campo que contiene un puntero que remite, si procede, a la forma canónica correspondiente, o, en caso contrario, una marca que indica que se trata de un término lexical canónico. El sistema de punteros no está sólo concebido para agrupar flexiones de un mismo lema, sino también con vistas a compilar grupos de palabras que, generalmente por estar construidas en torno a la misma lexema, constituyen un subconjunto léxicosemántico uniforme cuya reducción simplifica el proceso de indización documental⁵. Por otra parte, ese mismo campo se ha utilizado para anotar diversas informaciones de tipo gramatical y lexicográfico tales como marcas para identificar acrónimos, palabras "vacías", etc.

En lo que se refiere al establecimiento del antídiconario de palabras vacías, se han seguido diversos criterios. El planteamiento de la cuestión partió de considerar qué tipo de sistema se deseaba implementar. En primer lugar fueron marcadas como vacías las formas habitualmente consideradas como de valor primordialmente gramatical o al menos de escasa relevancia significativa: preposiciones, conjunciones, determinantes, verbos auxiliares y modales, copulas, etc. Asimismo se ha intentado señalar el mayor número posible de palabras "atemáticas", de sentido impreciso o excesivamente general, utilizando también como punto de referencia

critérios estadísticos. Resulta delicado estabelecer o grupo de termos "omnibus" que devem formar parte do antídicionário, pois é difícil determinar o grau de generalidade das palavras. Em nosso caso, quando se ha planteado uma dúvida sobre si considerar ou não vacuo algum termo, na maioria das ocasiões se ha preferido optar por incluí-lo no antídicionário, pois o tipo de indexação que se pretende alcançar por agora é de tipo "indicativo", mais que "informativo". Isto é, se trata de obter um conjunto de palavras chave lo mais reducido posible que consigne os aspectos mais essenciais do conteúdo do documento que se processe, não de fazer uma avaliação exaustiva. Assim, se han considerado "atemáticos" termos como _get_, _go_, _take_, _give_, _make_, etc., e incluso outros de carácter polissémico, como _even_, uma de cujas acepções (_even number_) puede resultar relevante no análisis de determinados documentos.

A la hora de evaluar o sistema se han observado determinadas limitaciones, especialmente en lo que se refiere a la gestión de homógrafas -y por tanto a la red de reenvíos entre flexiones y lexicales canónicas-, que son origen de diversos ruidos no el proceso. Otro tipo de problemas surgen do dicionário de frecuencias que se ha seleccionado como término de comparación, pois por las especiales características do corpus procesado no resulta un reflejo exacto do inglés standard. Sin embargo, nos parece que los resultados alcanzados a través de la metodología que hemos empleado se pueden considerar satisfactorios, pois se obtienen unas indexaciones con un grau de ruido no muy elevado y de suficiente exhaustividad, todo ello por procedimientos automáticos, lo que supone una solución a la falta de coherencia -y a la carga económica- que suele caracterizar el análisis de los indexadores humanos. Dadas las características do sistema, su actualización resulta posible sin excesiva intervención humana. Y esto aparece tanto más necesario cuanto más se considera cuál es la rapidez con la que se modifica el subsistema lingüístico de la Ciencia y la Tecnología. La creciente disponibilidad de material escrito en soporte legible por ordenador aumenta las facilidades para hacerlo. La opinión de los autores es que la investigación actual, en lo que se refiere a los proyectos de la inteligencia artificial, no debería marginar, a la hora de preparar bases de conocimientos para sistemas intérpretes de lenguaje natural, los avances clásicos de la estadística lingüística.

NOTAS

(1) La utilización de este tipo de técnicas ha sido muy experimentada durante las últimas cuatro décadas con vistas a desarrollar una "estilística computacional". Puede verse el reciente artículo de YANG HUIZHONG, "A New Technique for Identifying Scientific/Technical Terms and Describing Science Texts", en *Literary and Linguistic Computing*, vol. 1, n. 2, Oxford University Press, 1986.

(2) Cfr. VAN SLYPE, G., *Les langages d'indexation*, Paris, 1987, pág. 168 y ss. La indexación selectiva, bien en lenguaje natural, bien en lenguaje controlado, se caracteriza porque "toma en cuenta sólo ciertos términos, considerados por el algoritmo del sistema como los más representativos del contenido del documento" (pág. 169).

(3) Puede encontrarse una información más detallada en: Moya Anegón, F.; Muñoz Muñoz, J. M.; Hípola, P. "Análisis automático de documentación técnica informática", V Congreso de la AESLA, Pamplona, 1987.

(4) KUCERA, H. y FRANCIS, W. N., *Computational Analysis of Present-day American English*, Rhode Island, 1967.

(5) Deweze utiliza el término "forma semántica" para referirse a este tipo de agrupamientos. Cfr. A. DEWEZE, *Informatique documentaire*, Paris, 1985, pág. 129 y ss.