

Departamento de Lenguajes y Computación

Universidad de Almería

Práctica 1 Informática Documental

Profesor: Antonio Becerra Terón

Aula de Prácticas: Aula 4, CITE III

Octubre, 2007

PRÁCTICA 1. Implementación de un sistema de recuperación espacio-vectorial I. Ponderación de términos e indización de documentos

9 horas

Objetivo general de la práctica

OBJETIVO

El objetivo básico de esta práctica es abordar la primera fase en la implementación de un sistema de recuperación de información sencillo, basado en el modelo espacio-vectorial. Este objetivo abarca la selección de los términos de indización, ponderación automática, generación de estructuras de indización y, por último, almacenamiento de estas estructuras.

Descripción de la práctica

DESCRIPCIÓN

Esta práctica pretende que el alumno implemente el proceso de selección de términos de indización, ponderación automática e indización, para un sistema de recuperación textual.

La indexación se refiere a la caracterización de documentos para la tarea de la recuperación de información. Qué sería lo ideal? Que un documento indexado sea una representación eficiente de los contenidos semánticos del documento original. Normalmente, la indexación consiste en obtener una lista de términos significativos (también denominados conceptos) con información asociada (frecuencia en el documento, frecuencia en la base de datos documental, etc). Qué necesitamos ahora? Una técnica que nos permita obtener estas listas de términos.

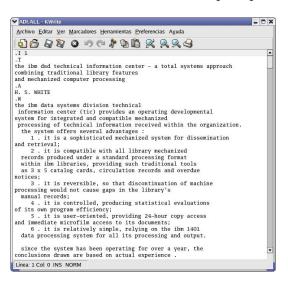
La técnica clásica de recuperación de información supone considerar cada documento como un conjunto de cadenas concatenadas, intentando eliminar las ambiguedades del lenguaje natural que puedan aparecer en los documentos. Las actividades que se realizan en este enfoque de la recuperación de información son:

- Eliminar las palabras que no expresan ningún contenido (stopwords o palabras vacías);
- Reducir las palabras a algún formato base como puede ser la raíz de cada una de ellas (stemming o lematización);
- Obtener un peso para los términos de indexación de acuerdo a su relevancia (frecuencia en el documento, frecuencia en la base de datos);

Para el desarrollo de la práctica, en primer lugar el alumno deber diponer de una colección de documentos en formato texto a almacenar en el sistema de recuperación que queremos implementar. Estos documentos proceden de la colección

del grupo de investigación *Information Retrieval* de la Universidad de Glasgow, liderado por el profesor *C. J. van RIJSBERGEN*. La colección la podemos encontrar en la pgina Web de la asignatura, concretamente en la URL http://www.ual.es/~abecerra/ID/. Esta colección contiene un total de 82 documentos estructurados de la siguiente forma

Figura 1: Documentos en la colección para práctica 1



Aquí, podemos comprobar que I representa el identificador de documento, T el título, A el autor, y W el texto a indexar de forma automática utilizando las técnicas de búsqueda en texto completo.

Una vez entregado el fondo documental, al alumno debe desarrollar un programa que permita realizar una selección automatizada de los términos de indización, y generación de las estructuras de indización. La figura 2 muestra, de forma gráfica, el proceso de esta práctica.

El proceso de indexación de un fondo documental conlleva realizar el proceso de filtrado y el proceso de extracción de los términos relevantes o keywords. El proceso de filtrado supone eliminar todas las cadenas de texto no relevantes y reducir las palabras a un formato base estándar. Este proceso de filtrado está involucrado por los siguientes subprocesos:

- (a) Listar las palabras de los documentos
- (b) Poner todo el texto en minúscula, y eliminar los símbolos especiales;
- (c) Eliminar de los documentos el conjunto de palabras vacías o stopwords;
- (d) Realizar la reducción de términos por lematización o stemming;

A continuación, desglosamos las actividades más importantes.

(c) Para la eliminación de las palabras vacías o *stopwords*, el alumno podrá disponer de la lista de palabras vacías propuesta por C. J. van Rijsbergen.

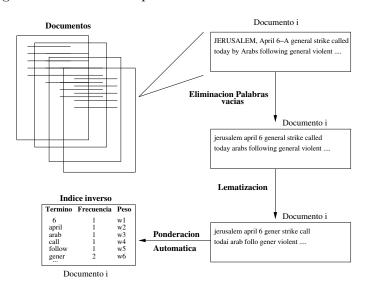


Figura 2: Proceso de la práctica 1 de Informática Documental

Esta lista contiene un conjunto bastante completo de palabras que no tienen contenido semántico para el documento. Entre estas palabras, podemos incluir artcíulos, preposiciones, conjunciones, adverbios, adjetivos, palabras muy cortas, etc.

El alumno debe implementar un programa que lea el texto del conjunto de documentos incluidos en el sistema de recuperación documental, y elimine todas las palabras vacías que aparezcan en dichos documentos. Para la implementación de este programa, se proporciona al alumno un archivo tipo texto, que contiene la lista de palabras vacías que debe comparar. Esta lista de palabras vacías se encuentra en el web de la asignatura, dirección URL http://www.ual.es/~abecerra/ID/. La implementación del programa deber utilizar una estructura de indexación adecuada para la comparación de las palabras leídas en los documentos con las palabras que aparecen en la lista.

La salida que debe generar el programa, es el conjunto de documentos de la colección, donde se han eliminado todas las palabras vacías, y, por tanto, las palabras que aparecen son susceptibles de considerarse como palabras relevantes o *keywords* para la indización de dicho documento.

(d) El proceso de lematización consiste en reducir un término hasta quedarnos con su raíz, que es, precisamente, la que aporta mayor contenido semántico. Debemos tener en cuenta que el número de palabras derivadas a partir de una misma raíz puede ser muy elevado. Esto supone que indizar documentos por la raíz de los términos proporciona una elevada disminución en el número de términos a considerar.

Para este proceso de lematización o *stemming*, el alumno deber implementar el algoritmo de lematización más utilizado en este proceso dentro del

contexto de la recuperación de información. Se proporciona al alumno este algoritmo en pseudo-cdigo, y el alumno debe aplicarlo a las palabras posiblemente relevantes de los documentos resultantes del proceso de filtrado anterior.

A continuación, mostramos, el pseudo-cdigo del algoritmo que facilitamos al alumno para su implementación:

```
Consonante: es una letra distinta de A, E, I, O, U, y distinta de Y precedida por una consonante;
Vocal: letra que no es consonante;
Una consonante se denota por c, una vocal por v. Una lista ccc. . . de longitud superior a 0 se denota
por C, y una lista vvv...por V; tenemos, por tanto, que cualquier palabra o parte de una palabra, tiene
uno de los siguientes cuatro formatos:
\mathsf{CVCV} \dots \mathsf{C}
CVCV ... V
VCVC ...C
VCVC ... V
Esto se puede representar de una única forma como sigue
[C]VCVC . . . [V]
donde [ ] los corchetes indican presencia arbitraria de contenido. Ahora, utilizando (VC){m} para denotar
VC repetida m veces, tenemos que esto lo podemos escribir como:
[C](VC)\{m\}[V].
m se denomina la medida de la palabra. El caso m=0 representa palabras nulas. Algunos ejemplos seran:
m=0
                TR, EE, TREE, Y, BY
m=1
                TROUBLE, OATS, TREES, IVY,
                TROUBLES, PRIVATE, OATEN, ORRERY.
m=2
Ahora, las reglas para eliminar un sufijo serán de la forma:
(condición) S_1 \rightarrow S_2
Esto significa que si una palabra finaliza con el sufijo \mathsf{S}_1, y el lema anterior a \mathsf{S}_1 satisface la
condición, entonces S_1 es reemplazado por S_2. La condición, normalmente, se proporciona en términos
de m; por ejemplo (m > 1) EMENT \rightarrow. Permitiría reemplazar REPLACEMENT por REPLAC,
dado que REPLAC es una parte de palabra con m=2.
La condición de la regla puede tener los siguientes formatos:
*S - El lema finaliza con S (igual para el resto de letras)
v^* - El lema contiene una vocal.
*d - El lema finaliza con una consonante doble (i.e. -TT, -SS).
*o - El lema finaliza con cvc, donde la segunda c no es W, X, o Y (i.e. -WIL, -HOP).
Adems, la condición puede incluir expresiones con operadores and, or y not. De forma que,
(\mathsf{m} > 1 \; \mathsf{and} \; (*\mathsf{S} \; \mathsf{or} \; *\mathsf{T}))
comprueba para un lema con m > 1, si finaliza en S o T, mientras que
(*d and not (*L or *S or *Z))
comprueba si un lema finaliza con una doble consonante distinta de L, S o Z.
Podemos encontrar regla de reemplazamiento con condiciones nulas; por ejemplo, con
\mathsf{SSES} \to \mathsf{SS}
\mathsf{IES} \to \mathsf{I}
SS \rightarrow SS
CARESSES se corresponde con CARESS, Igualmente CARESS se corresponde con CARESS, y,
por último, CARES con CARE. A continuación, comenzamos a presentar el algoritmo.
```

```
\begin{array}{cccc} \underline{\mathsf{Paso}\ 1a} \\ \mathsf{SSES} \to \mathsf{SS} & \mathsf{caresses} \to \mathsf{caress} \\ \mathsf{IES} \to \mathsf{I} & \mathsf{poines} \to \mathsf{poni} \\ & \mathsf{ties} \to \mathsf{ti} \\ \mathsf{SS} \to \mathsf{SS} & \mathsf{caress} \to \mathsf{caress} \\ \mathsf{S} \to & \mathsf{cats} \to \mathsf{cat} \end{array}
```

Si la segunda o tercera de las reglas en el Paso 1b tiene éxito, entonces realizar lo siguiente:

```
\begin{array}{lll} \mathsf{AT} \to \mathsf{ATE} & \mathsf{conflat(ed)} \to \mathsf{conflate} \\ \mathsf{BL} \to \mathsf{BLE} & \mathsf{troubl(ed)} \to \mathsf{trouble} \\ \mathsf{IZ} \to \mathsf{IZE} & \mathsf{siz(ed)} \to \mathsf{size} \\ (\mathsf{*d} \ \mathsf{and} \ \mathsf{not} \ (\mathsf{*L} \ \mathsf{or} \ \mathsf{*S} \ \mathsf{or} \ \mathsf{*Z})) \to \mathsf{una} \ \mathsf{letra} & \mathsf{hopp(ing)} \to \mathsf{hop} \\ & \mathsf{tann(ed)} \to \mathsf{tan} \\ & \mathsf{fall(ing)} \to \mathsf{fall} \\ & \mathsf{hiss(ing)} \to \mathsf{hiss} \\ & \mathsf{fizz(ed)} \to \mathsf{fizz} \\ (\mathsf{m} = 1 \ \mathsf{and} \ \mathsf{*o}) \to \mathsf{E} & \mathsf{fail(ing)} \to \mathsf{fail} \\ & \mathsf{fil(ing)} \to \mathsf{file} \\ \end{array}
```

Esta regla de hacer corresponder una única letra provoca la eliminación de uno de los pares de letras dobles.

$$\begin{array}{c} \underline{\mathsf{Paso}\ 1c} \\ (^*\mathsf{v}^*)\ \mathsf{Y} \to \mathsf{I} & \mathsf{happy} \to \mathsf{happi} \\ \mathsf{sky} \to \mathsf{sky} \end{array}$$

Esta regla trata con los plurales y los participios del pasado. Los siguientes pasos del algoritmo son mucho más directos.

```
(m > 0) ATIONAL \rightarrow ATE
                                                \mathsf{relational} \to \mathsf{relate}
(m>0) \ \mathsf{TIONAL} \to \mathsf{TION}
                                                \mathsf{conditional} \to \mathsf{condition}
                                                rational \rightarrow rational
(m > 0) ENCI \rightarrow ENCE
                                                \mathsf{valenci} \to \mathsf{valence}
(m > 0) ANCI \rightarrow ANCE
                                                he sitanci \rightarrow he sitance
(m > 0) IZER \rightarrow IZE
                                                \mathsf{digitizer} \to \mathsf{digitize}
(m > 0) ABLI \rightarrow ABLE
                                                conformabli → conformable
(m>0) \; \mathsf{ALLI} \to \mathsf{AL}
                                                \mathsf{radicalli} \to \mathsf{radical}
(m > 0) ENTLI \rightarrow ENT
                                                \mathsf{differently} \to \mathsf{different}
(m>0)\ ELI \to E
                                                \mathsf{vileli} \to \mathsf{vile}
(m > 0) OUSLI \rightarrow OUS
                                                analogousli \rightarrow analogous
                                                {\sf vietnamization} \to {\sf vietnamize}
(m>0) \; IZATION \rightarrow IZE
(m > 0) ATION \rightarrow ATE
                                                predication → predicate
(m > 0) ATOR \rightarrow ATE
                                                operator \rightarrow operate
(m > 0) ALISM \rightarrow AL
                                                \mathsf{feudalism} \, \to \, \mathsf{feudal}
(m>0) \ \mathsf{IVENESS} \to \mathsf{IVE}
                                                \mathsf{decisiveness} \to \mathsf{decisive}
(m>0) \; \mathsf{FULNESS} \to \mathsf{FUL}
                                                \mathsf{hopefulness} \to \mathsf{hopeful}
(m>0) \; \mathsf{OUSNESS} \to \mathsf{OUS}
                                               \mathsf{callousness} \to \mathsf{callous}
                                                \mathsf{formaliti} \to \mathsf{formal}
(m>0) \ ALITI \to AL
(m > 0) \text{ IVITI} \rightarrow \text{IVE}
                                                \mathsf{sensitiviti} \to \mathsf{sentitive}
(m>0) \; \mathsf{BILITI} \to \mathsf{BLE}
                                                sensibiliti \rightarrow sensible
```

```
\begin{array}{lll} & \underline{\mathsf{Paso}\ 3} \\ & (m>0)\ \mathsf{ICATE} \to \mathsf{IC} \\ & (m>0)\ \mathsf{ATIVE} \to \\ & (m>0)\ \mathsf{ALIZE} \to \mathsf{AL} \\ & (m>0)\ \mathsf{ALIZE} \to \mathsf{AL} \\ & (m>0)\ \mathsf{ICITI} \to \mathsf{IC} \\ & (m>0)\ \mathsf{ICAL} \to \mathsf{IC} \\ & (m>0)\ \mathsf{ICAL} \to \mathsf{IC} \\ & (m>0)\ \mathsf{FUL} \to \\ & (m>0)\ \mathsf{NESS} \to \\ \end{array}
```

```
Paso 4
(m > 1) AL \rightarrow
                                                  revival → reviv
(m > 1) ANCE \rightarrow
                                                  \mathsf{allowance} \to \mathsf{allow}
(m > 1) ENCE \rightarrow
                                                  inference → infer
(m > 1) ER \rightarrow
                                                  \mathsf{airliner} \to \mathsf{airlin}
(\mathsf{m}>1)~\mathsf{IC}\to
                                                  \mathsf{gyroscopic} \to \mathsf{gyroscop}
(m > 1) ABLE \rightarrow
                                                  adjustable \rightarrow adjust
(m > 1) IBLE \rightarrow
                                                  defensible → defens
(m > 1) ANT \rightarrow
                                                  \mathsf{irritant} \to \mathsf{irrit}
(m > 1) EMENT \rightarrow
                                                  replacement \rightarrow replac
(m > 1) MENT \rightarrow
                                                  \mathsf{adjustment} \to \mathsf{adjust}
(\mathsf{m}>1)~\mathsf{ENT}\to
                                                  \mathsf{depedent} \to \mathsf{depend}
(m > 1 and (*S or *T)) ION \rightarrow
                                                  \mathsf{adoption} \to \mathsf{adopt}
(m>1)~\text{OU} \rightarrow
                                                  homologou 	o homolog
(m > 1) ISM \rightarrow
                                                  communism → commun
(\mathsf{m}>1)\;\mathsf{ATE}\to
                                                  activate → activ
(m > 1) ITI \rightarrow
                                                  angulariti → angular
(m > 1) OUS \rightarrow
                                                  homologous \rightarrow homolog
(m > 1) IVE \rightarrow
                                                  effective \rightarrow effect
(m > 1) IZE \rightarrow
                                                  bowdlerize -- bowdler
```

```
\begin{array}{c} \underline{\mathsf{Paso}\;5a} \\ (\mathsf{m}>1)\;\mathsf{E}\to \\ \\ (\mathsf{m}=1\;\mathsf{and}\;\mathsf{not}\;\mathsf{*o})\;\mathsf{E}\to \\ \\ \underline{\mathsf{Paso}\;5b} \\ (\mathsf{m}>1)\;\mathsf{and}\;\mathsf{*d}\;\mathsf{and}\;\mathsf{*L})\to\mathsf{letra}\;\mathsf{nica} \\ \\ \mathsf{roll}\to\mathsf{roll}\to\mathsf{roll} \end{array}
```

Una vez implementado el algoritmo, utilizamos el lematizador sobre los documentos para generar el conjunto de palabras relevantes o *keywords* dentro de cada documento. Estas palabras constituyen una representación del texto, tal que los términos se pueden comparar con los términos de otros documentos y los términos que aparecen en las consultas.

La siguiente parte de la práctica se centra en el proceso de extracción de las keywords. En este caso, tenemos que los términos proporcionados por el lematizador serán utilizados como términos de indización para los documentos. A cada término de un documento se le asigna un peso de acuerdo a su frecuencia en el documento y en la colección. Antes de introducir lo que debe realizar el alumno en esta última parte de la práctica, definimos los conceptos de frecuencia de un término un documento y frecuencia de documentos de un término.

- Frecuencia de un término i en un documento j: número de ocurrencias del término i en el documento j. Cuantos más aparezca en un documento un término determinado, tenemos que este documento tratará sobre el concepto representado por el término;
- Frecuencia de documentos para un término i: número de documentos indizados por el término i, es decir número de documentos en la colección que contienen el término. Si tenemos términos que aparecen en muchos documentos, entonces estos términos no serán discriminantes para estos documentos.

A continuación, especificamos, de forma más detallada, el proceso de extracción de las palabras relevantes y su ponderación. Este proceso de ponderación debe iniciarse con el cómputo automático de la relevancia de los términos de indización presentes en cada uno de los documentos. La idea es poder generar vectores de documentos

$$d_i \rightarrow \overrightarrow{d_i} = (w(t_1, d_i), \dots, w(t_k, d_i))$$

donde d_i es el documento i, t_1, \ldots, t_k son los k términos de indización obtenidos después de aplicar el lematizador, y $w(t_j, d_i)$ expresa la relevancia del término de indización t_i en el documento d_i .

Por tanto, el alumno debe implementar un programa que le permita calcular, de forma automática, los vectores de ponderación o relevancia de los documentos originales con respecto a los términos de indización. Este programa calculará la ponderación o relevancia de los términos con la fórmula de ponderación que se utiliza habitualmente en el contexto de la recuperación de información.

$$w(t_i, d_j) = w_{i,j} = \frac{F_{i,j} \times Idf_i}{|\overrightarrow{d_j}|} = \frac{F_{i,j} \times \log \frac{N}{n_i}}{\sqrt{\sum_{s=1...k} (F_{s,j} \times \log \frac{N}{n_s})^2}}$$

donde $F_{i,j}$ representa la frecuencia del término i en el documento j, Idf_i la inversa de la frecuencia de documentos del término i, $|\overrightarrow{d_j}|$ representa el módulo del vector $\overrightarrow{d_j}$, N el número de documentos en la colección, y n_i es la cantidad de documentos donde aparece t_i .

Este mecanismo de cómputo de la ponderación implica recorrer las palabras obtenidas del lematizador para cada documento, y calcular, en función de la frecuencia de aparición de un término lematizado, su importancia o relevancia dentro de un documento. Además, si un término aparece mucho en un documento, se supone importante, pero si aparece en muchos documentos entonces no se considera útil y su medida de relevancia será mínima.

Una vez generados los vectores de documentos, al alumno organizará esta información en una matriz de asociación término-documento, utilizando las estructuras de datos necesarias para generar el conjunto completo de índices inversos de los documentos. Esta matriz organiza los documentos por filas y los términos por columnas, de forma que la posición M[i,j] de la matriz representa la relevancia del término de indexación j en el documento i.

El resultado de esta práctica será un programa que permite eliminar palabras vacías, lematizar y ponderar los términos en función de la frecuencia. Además, generará un índice inverso (estructuras de indización adecuadas) para cada documento, incluyendo la lista completa de términos de indización con la relevancia en el documento.

Recursos:

RECURSOS

Recursos software.

- $\bullet\,$ Compilador de Borland C++
- Java Development Kit