

# Tratamiento de la Variación Sintáctica mediante un Modelo de Recuperación Basado en Localidad\*

Jesús Vilares y Miguel A. Alonso

Departamento de Computación, Universidade da Coruña

Campus de Elviña s/n, 15071 - A Coruña

{jvilar,alonso}@udc.es

**Resumen:** La aplicación de información sintáctica en el modelo de recuperación basado en documentos imperante en la actualidad ha sido probada sin excesivo éxito en numerosas ocasiones, debido mayormente a los problemas que supone la integración de este tipo de información en el modelo. En este artículo proponemos el empleo de un modelo basado en localidad aplicado a la reordenación de resultados, el cual aborda el problema de la variación lingüística sintáctica mediante medidas de similaridad basadas en la distancia entre palabras. Se estudian dos aproximaciones cuya efectividad ha sido evaluada sobre el corpus CLEF de documentos en español.

**Palabras clave:** Recuperación de Información, variación lingüística sintáctica, modelo basado en localidad, fusión de datos.

**Abstract:** To date, attempts for applying syntactic information in the document-based retrieval model dominant have led to little practical improvement, mainly due to the problems associated with the integration of this kind of information into the model. In this article we propose the use of a locality-based retrieval model for reranking, which deals with syntactic linguistic variation through similarity measures based on the distance between words. We study two approaches whose effectiveness has been evaluated on the CLEF corpus of Spanish documents.

**Keywords:** Information Retrieval, syntactic linguistic variation, locality-based model, data fusion.

## 1. Introducción

El procesamiento sintáctico ha sido empleado repetidamente en el ámbito de la Recuperación de Información (RI) para hacer frente a la *variación lingüística sintáctica* presente en los textos (Perez-Carballo y Strzalkowski, 2000; Hull et al., 1997), si bien su empleo en el caso del español ha sido poco estudiado hasta ahora (Alonso, Vilares, y Darriba, 2002; Vilares y Alonso, 2003). Estas técnicas precisan de algún tipo de analizador sintáctico, para lo cual es necesario contar con una gramática apropiada, por sencilla que sea. Sin embargo, aún cuando dicha información sintáctica pueda ser convenientemente extraída del texto, persiste todavía el problema de cómo incorporar dicha información al sistema. La aproximación más común, consistente en una combinación ponderada de

términos simples y términos multipalabra — formados por términos simples relacionados sintácticamente—, no logra siempre resolver adecuadamente los problemas derivados de la sobrevaloración que el sistema tiende a dar a los términos complejos en detrimento de los términos simples (Mitra et al., 1997).

En este contexto, el empleo de técnicas pseudo-sintácticas basadas en distancias entre términos se presenta como una alternativa práctica que evita dichos problemas, al no ser necesaria gramática o analizador alguno, y al integrar de modo consistente la información obtenida, tanto a nivel de la aparición de los términos en sí, como de su proximidad, frecuentemente ligada a la existencia de una relación sintáctica entre los mismos.

En este artículo proponemos la utilización de un *modelo basado en localidad*, sustentado sobre similitudes basadas en distancias, como complemento a las técnicas clásicas de Recuperación de Información basadas en la indexación de términos simples, con el fin de

\* Parcialmente financiado por el Ministerio de Educación y Ciencia y FEDER (TIN2004-07246-C03-02), y por la Xunta de Galicia (PGIDIT05PXIC30501PN, PGIDIT05PXIC10501PN, PGIDIT05SIN044E).

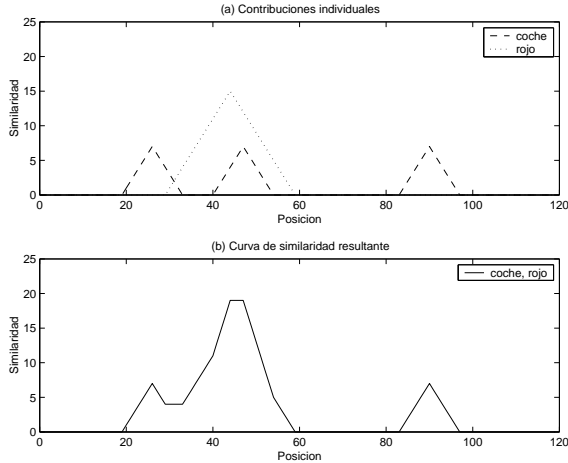


Figura 1: Modelo basado en localidad: (a) posiciones del texto con aparición de términos de la consulta y sus áreas de influencia; y (b) curva de similitud resultante

incrementar la precisión de los documentos devueltos por el sistema.

## 2. Recuperación de Información Basada en Localidad

### 2.1. Modelo de Recuperación

En el modelo de recuperación imperante en RI, denominado *basado en documentos*, el usuario solicita del sistema los documentos relevantes a su consulta o necesidad de información. Por otra parte, el modelo *basado en localidad* propuesto por de Kretser y Moffat (de Kretser y Moffat, 1999b; de Kretser y Moffat, 1999a) va un paso más allá y busca las *posiciones* concretas del texto que pueden resultar relevantes al usuario.

La *Recuperación de Pasajes* (Kaszkiel y Zobel, 2001) es una aproximación intermedia que persigue identificar aquellas secciones del documento —*pasajes*— relevantes para la consulta. En este modelo, una vez que el documento original ha sido dividido en pasajes, éstos son procesados y ordenados mediante técnicas tradicionales. Sin embargo restan por resolver problemas acerca de cómo definir el concepto de pasaje, su tamaño, grado de superposición, etc. (Llopis, 2003).

Por el contrario, el modelo basado en localidad considera la colección a indexar no como un conjunto de documentos, sino como una secuencia de palabras donde cada aparición de un término de la consulta ejerce una influencia sobre los términos circundantes. Dichas influencias son aditivas, de forma que la contribución de diferentes apariciones

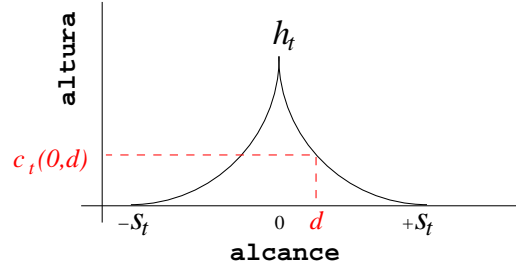


Figura 2: Función de contribución de similitud  $c_t$  de forma circular.

de términos de la consulta pueden sumarse, dando lugar a una medida de similitud, tal y como muestra la figura 1. Aquellas áreas del texto con una mayor densidad de términos de la consulta, o con términos de mayor peso, darán lugar a picos en la curva de influencia resultante, señalando posiciones del texto potencialmente relevantes. Todo ello sin necesidad de particionar artificialmente el documento como en el caso de la Recuperación de Pasajes.

### 2.2. Cálculo de Similaridades

A continuación, describiremos el modelo basado en localidad propuesto originalmente por de Kretser y Moffat (de Kretser y Moffat, 1999b; de Kretser y Moffat, 1999a). En este modelo la medida de similitud o relevancia es calculada únicamente sobre aquellas posiciones donde aparecen términos de la consulta, reduciendo de este modo el coste computacional asociado.

La contribución a dicha similitud por parte de un término de la consulta viene dada por una *función de contribución de similitud*  $c_t$  definida en base a los siguientes parámetros (de Kretser y Moffat, 1999a):

- La *forma* de la función, siendo la misma para todos los términos.
- La *altura máxima*  $h_t$  de la función, que se da en la posición del término que ejerce la influencia.
- El *alcance*  $s_t$  de la función, es decir, su radio de influencia.
- La distancia en palabras entre los dos términos considerados,  $d = |x - l|$ , donde  $l$  es la posición del término de la consulta que ejerce la influencia y  $x$  la posición sobre la que se desea calcular la medida de similitud.

Si bien en (de Kretser y Moffat, 1999a) se describen diversas formas de función, experimentos previos mostraron un mejor comportamiento de la función circular (*cir*) en el caso del español. Dicha función, cuya representación gráfica se muestra en la figura 2, se define mediante la ecuación:

$$c_t(x, l) = h_t \cdot \sqrt{1 - (d/s_t)^2} \quad (1)$$

con  $c_t(x, l) = 0$  para  $|x - l| > s_t$ , y equivalente a los cuadrantes de dos círculos normalizados con centros en  $(h_t, -s_t)$  y  $(h_t, s_t)$ .

Por su parte, la altura máxima  $h_t$  asociada a un término  $t$  se calcula como función inversa de su frecuencia en la colección:

$$h_t = f_{Q,t} \cdot \log_e \left( \frac{N}{f_t} \right) \quad (2)$$

donde  $N$  es el número total de términos en la colección,  $f_t$  el número de apariciones del término  $t$  en la colección y  $f_{Q,t}$  la frecuencia del término  $t$  en la consulta  $Q$ .

En lo que respecta al alcance  $s_t$  de la influencia de un término  $t$ , ésta viene dada también por el inverso de su frecuencia en la colección, pero normalizada en base a la frecuencia media:

$$s_t = \frac{n}{N} \cdot \frac{N}{f_t} = \frac{n}{f_t} \quad (3)$$

siendo  $n$  el número de términos diferentes en la colección, es decir, el tamaño del vocabulario.

De este modo, la medida de similaridad  $C_Q(x)$  asignada a la posición  $x$  del documento en la cual aparece un término de la consulta  $Q$  se calcula como:

$$C_Q(x) = \sum_{t \in Q} \sum_{\substack{l \in I_t \\ |l-x| \leq s_t \\ \text{term}(x) \neq \text{term}(l)}} c_t(x, l) \quad (4)$$

donde  $I_t$  es el conjunto de posiciones donde ocurre un término  $t$  de la consulta  $Q$ , y donde  $\text{term}(w)$  denota el término asociado a la posición  $w$ . En otras palabras, la medida de similaridad o relevancia asociada a una posición es la suma de las influencias ejercidas por los demás términos de la consulta presentes en el documento y dentro de cuyo alcance se encuentra, exceptuando otras apariciones del término existente en la posición considerada (de Kretser y Moffat, 1999b).

Finalmente, la medida de relevancia  $\text{sim}(D, Q)$  asignada a un documento  $D$  respecto a una consulta  $Q$  vendrá dada en fun-

ción de las similaridades asignadas a las apariciones de términos de la consulta que dicho documento contenga. Este punto se comenta en mayor detalle en el siguiente apartado.

### 2.3. Adaptaciones del Modelo.

Dado que el modelo basado en localidad permite trabajar a un nivel de detalle mayor que las técnicas clásicas de RI, al identificar no sólo los documentos relevantes sino también concretar las posiciones de interés dentro de los mismos, hemos optado por emplear este modelo en nuestros experimentos. Para ello ha sido necesario realizar ciertas adaptaciones de acuerdo con nuestras necesidades, las cuales nos diferencian del planteamiento original del modelo.

El planteamiento elegido a la hora de integrar la similaridad basada en distancias dentro de nuestro sistema de RI, ha sido el del postprocesado de los documentos previamente obtenidos mediante un sistema de recuperación clásico basado en documentos, con intención de incrementar la precisión de los primeros documentos devueltos. Este primer conjunto de documentos devuelto por el sistema es a continuación procesado empleando el modelo basado en localidad, tomando la ordenación final obtenida en base a distancias como aquella a devolver al usuario.

Otra de las principales diferencias respecto al modelo original es el del empleo de la lematización (Graña, Chappelier, y Vilares, 2001) en lugar del *stemming* a la hora de la normalización de consultas y documentos, dado su mejor comportamiento en el caso del español (Vilares et al., 2002).

Por otra parte, debemos señalar que los parámetros de altura máxima  $h_t$  y alcance  $s_t$  utilizados durante la reordenación se calculan en base a los parámetros globales de la colección, y no en base a los parámetros locales al subconjunto de documentos devueltos, para así evitar los problemas derivados de la correlación que esto conllevaría.<sup>1</sup>

Finalmente, a la hora de calcular la relevancia  $\text{sim}(D, Q)$  de un documento  $D$  respecto a una consulta  $Q$ , en lugar del algoritmo iterativo del modelo original (de Kretser y Moffat, 1999a), nuestra solución calcula dicha medida de relevancia como la suma de

<sup>1</sup>Por ejemplo, el parámetro  $f_t$  de número de apariciones de un término  $t$  es el número de apariciones de  $t$  en toda la colección, no el número de apariciones de  $t$  en el conjunto de documentos a reordenar.

las medidas de similaridad individuales de las apariciones de términos de la consulta en dicho documento:

$$\text{sim}(D, Q) = \sum_{\substack{x \in D \\ \text{term}(x) \in Q}} C_Q(x) \quad (5)$$

### 3. Resultados Experimentales con Distancias

Nuestra aproximación ha sido probada sobre el corpus monolingüe para español del CLEF<sup>2</sup>, conformado por las siguientes colecciones de documentos y *topics* asociados a partir de los cuales generar las consultas:

**CLEF 2001-02-A:** colección de entrenamiento y estimación de parámetros, formada por 215.738 teletipos de la agencia española de noticias EFE<sup>3</sup> correspondientes al año 1994, siendo sus *topics* asociados aquéllos de número impar empleados en las ediciones 2001 y 2002 del CLEF.

**CLEF 2001-02-B:** colección de evaluación similar a la anterior, si bien emplea los *topics* de número par.

**CLEF 2003:** colección de evaluación, formada por 454.045 teletipos de EFE correspondientes a los años 1994 y 1995, y que usa los *topics* del CLEF 2003.

Dichos *topics* están formados por tres campos: *título*, un breve título como su nombre indica; *descripción*, una somera frase de descripción; y *narrativa*, un pequeño texto especificando los criterios de relevancia.

Debemos precisar que aquellos *topics* con menos de 6 documentos relevantes fueron eliminados, ya que en dichos casos la modificación en la posición de uno o dos documentos devueltos puede acarrear cambios muy marcados en los resultados obtenidos para dicha consulta, distorsionando así los resultados globales (Hull et al., 1997).

Asimismo se emplearon dos tipos de consultas, las denominadas consultas *cortas*, generadas a partir de los campos *título* y *descripción*, y las denominadas consultas *largas*, que emplean la totalidad de los campos del *topic*. En el caso de las consultas largas, se ha doblado la relevancia asignada al campo

*título*, al concentrar éste la semántica básica de la consulta.

En lo que respecta a la indexación inicial de términos lematizados (*lem*), se empleó el conocido motor de indexación vectorial SMART (Buckley, 1985), empleando un esquema de pesos *atn-ntc* (Salton y Buckley, 1988). Por otra parte, con objeto de mejorar en lo posible el rendimiento final del sistema resultante, se partirá de un conjunto inicial de documentos tan bueno como sea posible. Para ello aplicaremos realimentación mediante expansión de consultas con el algoritmo de Rocchio (Rocchio, 1971):

$$Q_1 = \alpha Q_0 + \beta \sum_{k=1}^{n_1} \frac{R_k}{n_1} - \gamma \sum_{k=1}^{n_2} \frac{S_k}{n_2} \quad (6)$$

donde  $Q_1$  es vector de la consulta final,  $Q_0$  es el vector de la consulta inicial,  $R_k$  es el vector del  $k$ -ésimo documento relevante,  $S_k$  es el vector del  $k$ -ésimo documento no relevante,  $n_1$  es el número de documentos relevantes examinados,  $n_2$  es el número de documentos no relevantes examinados, y  $\alpha$ ,  $\beta$  y  $\gamma$  son, respectivamente, los parámetros que controlan las contribuciones relativas de la consulta original, los documentos relevantes, y los documentos no relevantes. En concreto, expandiremos automáticamente la consulta inicial con los  $t=10$  mejores términos de los cinco primeros documentos devueltos ( $n_1=5$ ), con unas contribuciones relativas  $\alpha=0.8$ ,  $\beta=0.1$ ,  $\gamma=0$  para consultas cortas y  $\alpha=1.2$ ,  $\beta=0.1$ ,  $\gamma=0$  para consultas largas.

Los resultados obtenidos se muestran en la tabla 1. El rendimiento del sistema se ha medido en base a los parámetros recogidos en cada fila: número de consultas empleadas, número de documentos devueltos, número de documentos relevantes esperados, número de documentos relevantes devueltos, precisión media no interpolada para todos los documentos relevantes, precisión-R, precisión en los primeros niveles estándar de cobertura, y precisión a los  $n$  documentos devueltos. La primera columna de cada grupo recoge los resultados de la línea base, la indexación de lemas con realimentación (*lem*), mientras que la segunda columna muestra los resultados obtenidos tras la ordenación de *lem* mediante distancias (*cir*). Para cada parámetro se han marcado en negrita los valores para los que se ha obtenido una mejora respecto a la línea base.

<sup>2</sup><http://www.clef-campaign.org>

<sup>3</sup><http://www.efe.es>

<i>corpus</i>	<i>CLEF 2001-02-A</i>				<i>CLEF 2001-02-B</i>				<i>CLEF 2003</i>			
<i>consulta</i>	<i>cortas</i>		<i>largas</i>		<i>cortas</i>		<i>largas</i>		<i>cortas</i>		<i>largas</i>	
<i>técnica</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>
#consultas	46	=	46	=	45	=	45	=	47	=	47	=
#docs. dev.	46k	=	46k	=	45k	=	45k	=	47k	=	47k	=
#rlvs. esp.	3007	=	3007	=	2513	=	2513	=	2335	=	2335	=
#rlvs. dev.	2767	=	2779	=	2376	=	2406	=	2240	=	2223	=
Pr. no int.	.5220	.4668	.5604	.4714	.4773	.4278	.5392	.4831	.5024	.3924	.5207	.4005
Pr.-R	.4990	.4651	.5366	.4652	.4599	.4205	.5104	.4592	.4912	.3921	.4871	.3911
Pr. a 0 %	.8221	<b>.8835</b>	.8895	<b>.8979</b>	.8210	<b>.8233</b>	.8710	.8678	.8145	<b>.8230</b>	.8301	<b>.8415</b>
Pr. a 10 %	.7490	<b>.7870</b>	.8028	<b>.8143</b>	.6861	<b>.7197</b>	.7619	<b>.8084</b>	.7369	.6626	.7421	.6518
Pr. a 20 %	.6866	<b>.6883</b>	.7352	.7017	.6319	<b>.6378</b>	.6929	.6808	.6632	.5663	.6758	.5737
Pr. a 30 %	.6573	.6148	.6996	.6066	.5688	.5464	.6497	.6000	.6019	.5098	.6304	.5030
Pr. a 40 %	.5997	.5267	.6541	.5372	.5289	.4827	.6202	.5438	.5638	.4439	.5975	.4400
Pr. a 50 %	.5456	.4656	.6005	.4728	.5017	.4322	.5733	.4987	.5410	.4077	.5479	.3956
Pr. a 5 docs.	.6609	<b>.6913</b>	.6957	<b>.7261</b>	.5956	<b>.6000</b>	.6844	.6533	.5872	.5532	.6213	.5574
Pr. a 10 docs.	.6457	.6391	.6848	.6522	.5600	.5444	.6178	.6089	.5596	.5064	.5872	.4979
Pr. a 15 docs.	.5884	<b>.5899</b>	.6435	.5971	.5274	.5111	.5822	.5556	.5305	.4624	.5504	.4652
Pr. a 20 docs.	.5630	.5446	.6043	.5674	.5011	.4822	.5533	.5189	.4883	.4181	.5266	.4277
Pr. a 30 docs.	.5225	.4848	.5580	.4971	.4444	.4215	.5081	.4733	.4433	.3702	.4667	.3780
Pr. a 100 docs.	.3507	.3052	.3598	.3048	.2940	.2780	.3191	.3022	.2770	.2400	.2853	.2404

Tabla 1: Resultados obtenidos mediante reordenación por distancias (*cir*) de la lematización con realimentación (*lem*)

Como muestran los resultados, la reordenación por distancias ha producido una disminución general del rendimiento del sistema, salvo para los primeros niveles de cobertura y primeros documentos devueltos, donde en algunos casos los resultados son similares o incluso mejores. Podemos concluir, pues, que esta primera aproximación no ha demostrado ser de demasiado interés práctico.

#### 4. Fusión de Datos mediante Intersección

##### 4.1. Justificación

Dado que el número de documentos relevantes devueltos es el mismo, la caída en el rendimiento del sistema en ésta primera aproximación sólo puede deberse a una mala ordenación de los resultados por el modelo basado en distancias. Por esta razón decidimos estudiar la variación en la distribución de documentos relevantes y no relevantes en los  $K$  primeros documentos devueltos. Comentaremos únicamente los resultados obtenidos empleando consultas cortas para el corpus CLEF 2001-02-A, mostrados en la tabla 2, ya que los resultados obtenidos en dicho estudio son muy similares para los demás corpus y tipos de consultas.

Cada fila muestra los resultados obtenidos al comparar los  $K$  primeros documentos devueltos por el sistema mediante inde-

xación de lemas con realimentación (*lem*) —conjunto de resultados  $L$ — con aquéllos devueltos tras su reordenación mediante distancias (*cir*) —conjunto de resultados  $D$ . Las columnas muestran los resultados obtenidos para cada uno de los parámetros considerados: número medio de nuevos relevantes obtenidos mediante distancias ( $D \setminus L$ ), número medio de relevantes perdidos con distancias ( $L \setminus D$ ), número medio de relevantes que se mantienen ( $L \cap D$ ), coeficiente de superposición de relevantes ( $R_{sup}$ ), precisión de *lem* a los  $K$  primeros documentos ( $Pr(L)$ ), precisión a los  $K$  documentos tras la reordenación por distancias ( $Pr(D)$ ), y precisión en los documentos comunes a ambas aproximaciones dentro de sus  $K$  primeros documentos ( $Pr(L \cap D)$ ). En la parte derecha de la tabla se muestran sus equivalentes para el caso de los documentos no relevantes: número medio de no relevantes añadidos, perdidos y comunes, y grado de superposición de no relevantes.

A partir de estos resultados se pueden extraer diversas conclusiones de importancia. En primer lugar, observamos que el número de documentos relevantes obtenidos por ambas aproximaciones dentro de sus  $K$  primeros documentos es muy similar —si bien algo menor para las distancias—, tal como se puede apreciar en las cifras absolutas de documentos relevantes entrantes y salientes y en

$K$	Docs. relevantes							Docs. no relevantes			
	$D \setminus L$	$L \setminus D$	$L \cap D$	$R_{sup}$	$Pr(L)$	$Pr(D)$	$Pr(L \cap D)$	$D \setminus L$	$L \setminus D$	$L \cap D$	$N_{sup}$
5	1.93	1.78	1.52	0.45	<b>0.66</b>	<b>0.69</b>	<b>0.80</b>	1.15	1.30	0.39	0.24
10	3.24	3.30	3.15	0.49	<b>0.65</b>	<b>0.64</b>	<b>0.76</b>	2.61	2.54	1.00	0.28
15	4.17	4.15	4.67	0.53	<b>0.59</b>	<b>0.59</b>	<b>0.72</b>	4.35	4.37	1.80	0.29
20	4.59	4.96	6.30	0.57	<b>0.56</b>	<b>0.54</b>	<b>0.72</b>	6.65	6.28	2.46	0.28
30	5.61	6.74	8.93	0.59	<b>0.52</b>	<b>0.48</b>	<b>0.68</b>	11.22	10.09	4.24	0.28
100	7.00	11.48	23.52	0.72	<b>0.35</b>	<b>0.31</b>	<b>0.49</b>	45.43	40.96	24.04	0.36
200	5.54	9.63	37.35	0.83	<b>0.23</b>	<b>0.21</b>	<b>0.36</b>	90.50	86.41	66.61	0.43
500	2.35	3.39	52.67	0.95	<b>0.11</b>	<b>0.11</b>	<b>0.16</b>	167.43	166.39	277.54	0.62

Tabla 2: Distribución de documentos relevantes y no relevantes tras la reordenación mediante distancias. Corpus CLEF 2001-02-A, consultas cortas

las precisiones a los  $K$  documentos de ambas aproximaciones. Esto nos permite confirmar que se trata de un problema de mala ordenación de los documentos.

En segundo lugar debemos referirnos a los coeficientes de superposición de documentos relevantes ( $R_{sup}$ ) y no relevantes ( $N_{sup}$ ). Estos coeficientes, definidos en (Lee, 1997), indican el grado de superposición entre el conjunto de documentos relevantes o no relevantes de dos conjuntos de documentos devueltos. Para dos ejecuciones  $run_1$  y  $run_2$ , dichos coeficientes se definen como:

$$R_{sup} = \frac{2 |Rel(run_1) \cap Rel(run_2)|}{|Rel(run_1)| + |Rel(run_2)|} \quad (7)$$

$$N_{sup} = \frac{2 |Nonrel(run_1) \cap Nonrel(run_2)|}{|Nonrel(run_1)| + |Nonrel(run_2)|} \quad (8)$$

donde  $Rel(X)$  y  $Nonrel(X)$  representan, respectivamente, el conjunto de documentos relevantes y no relevantes devueltos en la ejecución  $X$ .

Como podemos apreciar en la tabla 2, los factores de superposición de los documentos relevantes son considerablemente mayores que los de los no relevantes. De esta forma ambas aproximaciones devuelven un conjunto similar de documentos relevantes, pero un conjunto diferente de documentos irrelevantes. Se cumple, pues, la denominada *propiedad de la superposición desigual* (Lee, 1997), que dice que diferentes ejecuciones deben devolver conjuntos similares de documentos relevantes a la vez que devolver conjuntos disimilares de no relevantes como primer indicador de la efectividad que tendría la fusión de datos de ambas.

En tercer lugar, y en relación al punto anterior, puede verse que la precisión en los documentos comunes a ambas aproximaciones dentro de sus  $K$  primeros documentos

( $Pr(L \cap D)$ ) es mayor que las precisiones alcanzadas tanto por lemas ( $Pr(L)$ ) como por distancias ( $Pr(D)$ ); o lo que es lo mismo, la probabilidad de que un documento sea relevante es mayor cuando es devuelto por ambas aproximaciones.

Conforme a estas observaciones, se planteó una nueva aproximación para la reordenación, esta vez basada en la fusión de datos.

## 4.2. Descripción del Algoritmo

La *fusión de datos* es una técnica de combinación de evidencias consistente en la combinación de resultados devueltos empleando diferentes representaciones de consultas o documentos, o mediante múltiples técnicas de recuperación (Fox y Shaw, 1994; Lee, 1997).

En nuestro caso hemos optado por una aproximación basada no en la combinación de puntuaciones en base a similitudes (Fox y Shaw, 1994; Lee, 1997) o rango (Lee, 1997), sino en un criterio booleano para el cual, una vez fijado un valor  $K$ , los documentos son devueltos en el siguiente orden:

1. En primer lugar, los documentos pertenecientes a la intersección de los  $K$  primeros documentos de ambas aproximaciones:  $L_K \cap D_K$ . El objetivo perseguido es el de incrementar la precisión en los primeros documentos devueltos.
2. A continuación, los documentos pertenecientes a los  $K$  primeros documentos de ambas aproximaciones que no estén en la intersección:  $(L_K \cup D_K) \setminus (L_K \cap D_K)$ . El objetivo es añadir a los primeros documentos devueltos aquellos documentos relevantes devueltos únicamente mediante la aproximación basada en distancias, sin perjudicar la ordenación de aquéllos devueltos únicamente por la indexación de lemas.

- Finalmente, los restantes documentos devueltos por los lemas:  $L \setminus (L_K \cup D_K)$ .

donde  $L$  es el conjunto de resultados devuelto por  $lem$ ,  $L_K$  el conjunto de  $K$  primeros resultados devuelto con  $lem$ , y  $D_K$  el conjunto de  $K$  primeros resultados devuelto mediante la reordenación por distancias.

Con respecto a la ordenación interna de los resultados, se tomará como referencia, por sus mejores resultados, la ordenación obtenida mediante la indexación de lemas ( $lem$ ). De esta forma cuando se devuelva un subconjunto  $S$  de resultados, los documentos que lo conforman se devolverán en el mismo orden relativo que existía entre ellos cuando eran devueltos por  $lem$ .<sup>4</sup>

### 5. Resultados Experimentales con Fusión de Datos

Tras experimentos previos de puesta a punto, se optó finalmente por emplear un valor  $K = 30$  en el caso de consultas cortas y  $K = 50$  en el caso de las largas.

La tabla 3 recoge los resultados obtenidos para la nueva aproximación. Al igual que antes, la primera columna de cada grupo muestra los resultados para la línea base, la indexación de lemas con realimentación ( $lem$ ), mientras que la segunda columna muestra los resultados obtenidos tras su ordenación mediante fusión por intersección ( $cir$ ).

Podemos apreciar que las mejoras obtenidas con la reordenación mediante fusión son consistentes, especialmente en el caso de la precisión a los  $n$  documentos devueltos —tal como se pretendía—, si bien dichas mejoras también se extienden al resto de parámetros estudiados, siendo algo menores en el caso del corpus CLEF 2003.

### 6. Conclusiones y Trabajo Futuro

A lo largo de este artículo se ha planteado la utilización de un modelo de recuperación basado en distancias entre palabras, o basado en localidad, para tratar el problema de la variación lingüística de carácter sintáctico presente en los textos.

Se han considerado dos aproximaciones, ambas enfocadas a la reordenación de resultados, en este caso obtenidos mediante inde-

<sup>4</sup>Es decir, si la secuencia original en  $lem$  era  $d2-d3-d1$  y se toma un subconjunto  $\{d1, d2\}$  a devolver, los documentos se obtendrían en el mismo orden relativo original:  $d2-d1$ .

xación de lemas de palabras con contenido. La primera aproximación, que asumía el orden obtenido mediante la aplicación del modelo basado en localidad como el orden final a devolver, no obtuvo buenos resultados. Tras analizar el comportamiento del sistema se optó por emplear una aproximación basada en la fusión de datos, y que emplea la intersección de los conjuntos de documentos devueltos por ambos sistemas como guía para la reordenación. Esta segunda aproximación resultó fructífera, obteniendo una mejora consistente y general.

En lo que respecta al trabajo futuro, pretendemos dar capacidad al sistema para el tratamiento de variantes morfosintácticas de una expresión (Jacquemin y Tzoukermann, 1999), así como de relaciones de sinonimia ponderada (Fernández-Lanza, Graña, y Sobrino, 2003).

### Bibliografía

- Alonso, M. A., J. Vilares, y V. M. Darriba. 2002. On the usefulness of extracting syntactic dependencies for text indexing. En vol. 2464 de *Lecture Notes in Artificial Intelligence*. Springer-Verlag, pág. 3–11.
- Buckley, C. 1985. Implementation of the SMART information retrieval system. Technical Report TR85-686, Department of Computer Science, Cornell University.
- de Kretser, O. y A. Moffat. 1999a. Effective document presentation with a locality-based similarity heuristic. En *Proc. of SIGIR '99, Berkeley, USA*, pág. 113–120.
- de Kretser, O. y A. Moffat. 1999b. Locality-based information retrieval. En *Proc. of 10th Australasian Database Conference (ADC '99), Auckland, New Zealand*, pág. 177–188.
- Fernández-Lanza, S., J. Graña, y A. Sobrino. 2003. Introducing FDSA (Fuzzy Dictionary of Synonyms and Antonyms): Applications on Information Retrieval and Stand-Alone Use. *Mathware & Soft Computing*, 10(2-3):57–70.
- Fox, E. A. y J. A. Shaw. 1994. Combination of multiple searches. En *The 2nd Text REtrieval Conference (TREC-2)*, Gaithersburg, USA, pág. 243–252.
- Graña, J., J.-C. Chappelier, y M. Vilares. 2001. Integrating external dictionaries into stochastic part-of-speech taggers. En

<i>corpus</i>	<i>CLEF 2001-02-A</i>				<i>CLEF 2001-02-B</i>				<i>CLEF 2003</i>			
<i>consulta</i>	<i>cortas</i>		<i>largas</i>		<i>cortas</i>		<i>largas</i>		<i>cortas</i>		<i>largas</i>	
<i>técnica</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>	<i>lem</i>	<i>cir</i>
#consultas	46	=	46	=	45	=	45	=	47	=	47	=
#docs. dev.	46k	=	46k	=	45k	=	45k	=	47k	=	47k	=
#rlvs. esp.	3007	=	3007	=	2513	=	2513	=	2335	=	2335	=
#rlvs. dev.	2767	=	2779	=	2376	=	2406	=	2240	=	2223	=
Pr. no int.	.5220	<b>.5327</b>	.5604	.5589	.4773	.4768	.5392	<b>.5497</b>	.5024	.4977	.5207	.5167
Pr.-R	.4990	<b>.5126</b>	.5366	<b>.5433</b>	.4599	.4551	.5104	<b>.5188</b>	.4912	.4737	.4871	.4865
Pr. a 0%	.8221	<b>.8386</b>	.8895	<b>.9091</b>	.8210	<b>.8248</b>	.8710	<b>.8751</b>	.8145	<b>.8163</b>	.8301	.8257
Pr. a 10%	.7490	<b>.7758</b>	.8028	<b>.8256</b>	.6861	<b>.7191</b>	.7619	<b>.7740</b>	.7369	.7283	.7421	<b>.7540</b>
Pr. a 20%	.6866	<b>.7193</b>	.7352	<b>.7528</b>	.6319	<b>.6426</b>	.6929	<b>.7188</b>	.6632	<b>.6737</b>	.6758	<b>.6834</b>
Pr. a 30%	.6573	<b>.6844</b>	.6996	.6922	.5688	<b>.5818</b>	.6497	<b>.6784</b>	.6019	.6015	.6304	<b>.6391</b>
Pr. a 40%	.5997	<b>.6164</b>	.6541	<b>.6610</b>	.5289	<b>.5470</b>	.6202	<b>.6460</b>	.5638	<b>.5672</b>	.5975	.5876
Pr. a 50%	.5456	<b>.5644</b>	.6005	<b>.6026</b>	.5017	.4909	.5733	<b>.5996</b>	.5410	.5359	.5479	.5327
Pr. a 5 docs.	.6609	<b>.6739</b>	.6957	<b>.7217</b>	.5956	<b>.6178</b>	.6844	<b>.6933</b>	.5872	<b>.6298</b>	.6213	<b>.6553</b>
Pr. a 10 docs.	.6457	<b>.6761</b>	.6848	<b>.7065</b>	.5600	<b>.5756</b>	.6178	<b>.6400</b>	.5596	<b>.5745</b>	.5872	<b>.5979</b>
Pr. a 15 docs.	.5884	<b>.6188</b>	.6435	<b>.6449</b>	.5274	<b>.5393</b>	.5822	<b>.6000</b>	.5305	<b>.5390</b>	.5504	<b>.5560</b>
Pr. a 20 docs.	.5630	<b>.5826</b>	.6043	<b>.6185</b>	.5011	<b>.5089</b>	.5533	<b>.5722</b>	.4883	<b>.5074</b>	.5266	.5170
Pr. a 30 docs.	.5225	.5225	.5580	<b>.5652</b>	.4444	.4444	.5081	<b>.5148</b>	.4433	.4433	.4667	<b>.4716</b>
Pr. a 100 docs.	.3507	.3502	.3598	.3539	.2940	<b>.3011</b>	.3191	<b>.3304</b>	.2770	<b>.2789</b>	.2853	.2809

Tabla 3: Resultados obtenidos mediante reordenación por fusión con intersección (*cir*) de la lematización con realimentación (*lem*)

- Proc. of RANLP 2001, Tzigov Chark, Bulgaria*, pág. 122–128
- Hull, D. A., G. Grefenstette, B. M. Schulze, E. Gaussier, H. Schütze, y J. O. Pedersen. 1997. Xerox TREC-5 site report: Routing, filtering, NLP, and Spanish tracks. En *The 5th Text REtrieval Conference (TREC-5), Gaithersburg, USA*, pág. 167–180.
- Jacquemin, C. y E. Tzoukermann. 1999. NLP for term variant extraction: synergy between morphology, lexicon and syntax. En *Natural Language Information Retrieval*, vol. 7 de *Text, Speech and Language Technology*. Kluwer Academic Publishers, pág. 25–74.
- Kaszkiel, M. y J. Zobel. 2001. Effective ranking with arbitrary passages. *Journal of the American Society of Information Science*, 52(4):344–364.
- Lee, J. H. 1997. Analyses of multiple evidence combination. En *Proc. of SIGIR '97, Philadelphia, USA*, pág. 267–276. ACM Press.
- Llopis, F. 2003. *IR-n: Un sistema de Recuperación de Información basado en Pasajes*. PhD. Thesis, Universidad de Alicante.
- Mitra, M., C. Buckley, A. Singhal, y C. Cardie. 1997. An analysis of statistical and syntactic phrases. En *Proc. of RIAO-97, 5th International Conference "Recherche d'Information Assistee par Ordinateur", Montreal, Canada*, pág. 200–214.
- Perez-Carballo, J. y T. Strzalkowski. 2000. Natural language information retrieval: progress report. *Information Processing and Management*, 36(1):155–178.
- Rocchio, J.J., 1971. *The SMART Retrieval System—Experiments in Automatic Document Processing*, capítulo Relevance feedback in information retrieval, pág. 313–323. Prentice-Hall.
- Salton, G. y C. Buckley. 1988. Term-weighting approaches in automatic retrieval. *Information Processing & Management*, 24(5):513–523.
- Vilares, J., M. A. Alonso, F. J. Ribadas, y M. Vilares. 2002. COLE experiments at CLEF 2002 Spanish monolingual track. En *Working Notes for the CLEF 2002 Workshop, Rome, Italy*, pág. 153–160.
- Vilares, J. y M. A. Alonso. 2003. A grammatical approach to the extraction of index terms. En *Proc. of RANLP 2003, Borovest, Bulgaria*, pág. 500–504.