



RECUPERACI N DE INFORMACI N: UN  REA DE INVESTIGACI N EN CRECIMIENTO

Information retrieval: A Growing Area of research

Fernando R. A. Bordignon
Universidad Nacional de Luj n, Argentina

Gabriel H. Tolosa
Universidad Nacional de Luj n, Argentina

RESUMEN

A partir de la expansi n y consolidaci n de Internet, como medio principal de comunicaci n electr nica de datos, se ha puesto a disposici n de casi toda la humanidad una importante cantidad de informaci n de todo tipo. A los efectos de aprovechar todo este potencial de informaci n, es necesario poseer accesos que permitan que la tarea de recuperaci n sea efectiva y eficiente en t rminos de recursos invertidos por los usuarios. Este art culo plantea cu l es el objeto de estudio del  rea denominada "recuperaci n de informaci n", en que estado se encuentra y cuales son sus principales l neas de trabajo.

Palabras clave: Recuperaci n de informaci n, web, motores de b squeda

ABSTRACT

As a result of the expansion and consolidation of the Internet as the main medium for the transmission of electronic data, a huge amount of information of all kinds has become readily available to humanity. For the purpose of exploiting this information potential it is necessary to have ways of access that would make the information retrieval task both effective and efficient in terms of user resources. This article describes the object of study of "Information Retrieval", its state of the art and main lines of research.

Keywords: Information retrieval, web, search engines

EL ENTORNO DE LA RECUPERACI N DE INFORMACI N

Hist ricamente, el hombre ha necesitado de medios sobre los cuales representar todo acerca del mundo que lo rodea y de reflejar – de alguna



manera – su evoluci n. La escritura ha sido el mecanismo “tradicional” y fundamental que soporta su conocimiento en el tiempo.

Esta misma evoluci n ha facilitado la existencia de diferentes medios de representaci n de la escritura, llegando hasta nuestros d as donde la informaci n se representa digitalmente y es posible su almacenamiento y distribuci n masiva en forma simple y r pida, a trav s de redes de computadoras. La digitalizaci n abri  nuevos horizontes en las formas que el hombre puede tratar con la informaci n que produce.

De igual manera, el volumen de informaci n existente crece permanentemente y adquiere diferentes formas de representaci n, desde simples archivos de texto en una computadora personal o un peri dico electr nico hasta librer as digitales y espacios mucho m s grandes y complejos como la web. Algunos investigadores han planteado que – desde hace varios a os – existe un fen meno denominado “sobrecarga de informaci n” [MAES] debido a que el volumen y la disponibilidad hacen que los usuarios no cuenten con suficiente tiempo f sico para “procesar” todo el c mulo de medios a su alcance [CARLSON].

Entonces, resulta importante tratar con toda esa informaci n disponible electr nicamente para que pueda servir a diferentes personas (usuarios) en diferentes situaciones. Esto plantea un desaf o interesante: hay importantes vol menes de informaci n y hay usuarios que se pueden beneficiar de alguna manera con la posibilidad de acceder a  sta, por lo tanto, c mo poder unir preguntas con respuestas, necesidades de informaci n con documentos, consultas con resultados. Bien, en las ciencias de la computaci n existe un  rea, la Recuperaci n de Informaci n (*Information Retrieval*), que estudia y propone soluciones al escenario presentado, planteando modelos, algoritmos y heur sticas.

La Recuperaci n de Informaci n (RI) no es un  rea nueva, sino que se viene desarrollando desde finales de la d cada de 1950. Sin embargo, en la actualidad adquiere un rol m s importante debido al valor que tiene la misma. Se puede plantear que disponer o no de la informaci n justa en tiempo y forma puede resultar en el  xito o fracaso de una operaci n. De aqu , la importancia de los Sistemas de Recuperaci n de Informaci n (SRI) que pueden manejar – con ciertas limitaciones – estas situaciones de manera eficaz y eficiente.

Pero,  Qu  se entiende concretamente por “Recuperaci n de Informaci n”? Para Ricardo Baeza-Yates y otros [BAEZA-YATES] “la



Recuperaci n de Informaci n trata con la representaci n, el almacenamiento, la organizaci n y el acceso a  tems de informaci n”.

A os antes, Salton [SALTON_a] propuso una definici n amplia que plantea que el  rea de RI “*es un campo relacionado con la estructura, an lisis, organizaci n, almacenamiento, b squeda y recuperaci n de informaci n*”.

Cabe aclarar que en las definiciones anteriores los elementos de informaci n son no estructurados, tales como documentos de texto libre (por ejemplo, un archivo de texto que contenga La Biblia)   semi-estructurados, como lo son las p ginas web.

Croft [CROFT] estima que la recuperaci n de informaci n es “el conjunto de tareas mediante las cuales el usuario localiza y accede a los recursos de informaci n que son pertinentes para la resoluci n del problema planteado. En estas tareas desempe an un papel fundamental los lenguajes documentales, las t cnicas de resumen, la descripci n del objeto documental, entre otras”. Por otro lado, Korfhage [KORFHAGE] defini  la RI como “la localizaci n y presentaci n a un usuario de informaci n relevante a una necesidad de informaci n expresada como una pregunta”

Ciertamente, es un  rea amplia, donde se abarcan diferentes t picos, algunos computacionales como el almacenamiento y la organizaci n; y otros relacionados con el lenguaje y los usuarios como la representaci n y la recuperaci n propiamente dicha.

N tese que Croft y Korfhage plantean expl citamente el rol del usuario como fuente de consultas y destinatario de las respuestas. Por lo tanto, de manera m s gen rica, podemos plantear que la recuperaci n de informaci n intenta resolver el problema de “**encontrar y rankear documentos relevantes que satisfagan la necesidad de informaci n de un usuario, expresada en un determinado lenguaje de consulta**”. Sin embargo, existe un problema que dificulta sobremanera esta tarea y consiste en poder “compatibilizar” y comparar el lenguaje en que est  expresada tal necesidad de informaci n y el lenguaje de los documentos.

LA PROBLEM TICA DE LA RI

De forma general – seg n Baeza-Yates [BAEZA-YATES] – el problema de la RI puede ser estudiado desde dos puntos de vista: el computacional y el humano. El primer caso tiene que ver con la construcci n de estructuras

de datos y algoritmos eficientes que mejoren la calidad de las respuestas. El segundo caso corresponde al estudio del comportamiento y de las necesidades de los usuarios.

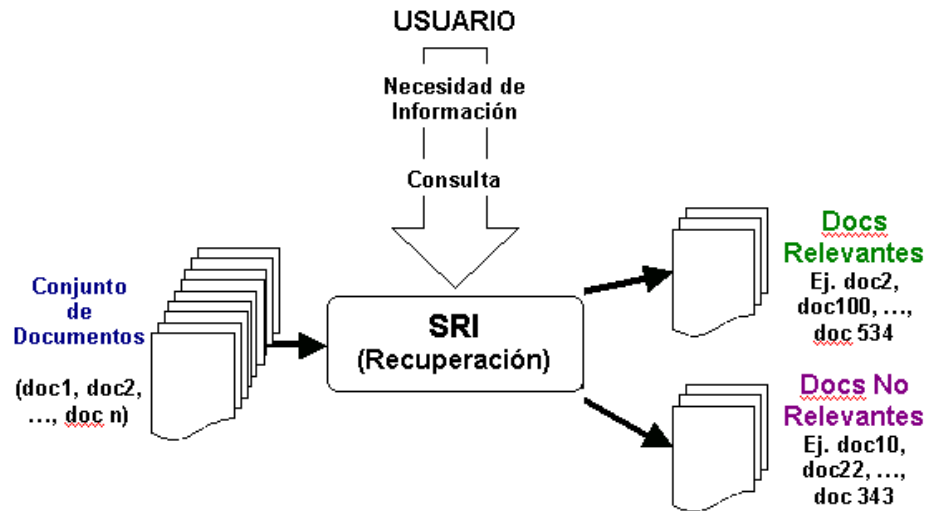


Figura 1 – La problemática de la RI

Si se analiza la problemática de la RI desde un alto nivel de abstracción (Figura 1) podemos establecer que:

- Existe una colección de documentos que contienen información de interés (sobre uno o varios temas)
- Existen usuarios con necesidades de información, quienes las plantean al SRI en forma de una consulta (en inglés, *query*. En adelante, ambas palabras se utilizarán indistintamente)
- Como respuesta, el sistema retorna – de forma ideal – referencias a documentos “relevantes”, es decir aquellos que satisfacen la necesidad expresada, generalmente en forma de una lista rankeada.

Planteamos que la respuesta “ideal” de un SRI está formada solamente por documentos relevantes a la consulta, pero – en la práctica – esta no es aún alcanzable. Esto se debe a que – entre otros motivos – existe el problema de compatibilizar la expresión de la necesidad de información y el lenguaje y de los documentos. Además, hay una carga de subjetividad subyacente y depende de los usuarios. Entonces, el SRI recupera la mayor cantidad posible de documentos relevantes, minimizando la cantidad de documentos no relevantes (ruido) en la respuesta. En términos de eficiencia, se plantea la idea de **precisión** de la respuesta, es decir, cuando más

documentos relevantes contengan el conjunto soluci n (para una consulta dada), m s preciso ser .

Para cumplir con sus objetivos, un SRI debe realizar algunas tareas b sicas, las cuales se encuentran – fundamentalmente – planteadas en cuestiones computacionales, a saber:

- Representaci n l gica de los documentos y – opcionalmente – almacenamiento del original. Algunos sistemas solo almacenan porciones de los documentos y otros lo hacen de manera completa.
- Representaci n de la necesidad de informaci n del usuario en forma de consulta.
- Evaluaci n de los documentos respecto de una consulta para establecer la relevancia de cada uno.
- Ranqueo de los documentos considerados relevantes para formar el “conjunto soluci n” o respuesta.
- Presentaci n de la respuesta al usuario.
- Retroalimentaci n o refinamiento de las consultas (para aumentar la calidad de la respuesta)

En la figura 2 se puede apreciar con mayor detalle la arquitectura b sica de un SRI, el tratamiento de los documentos y la interacci n con el usuario. Aqu  se ven algunos componentes que no se hab an mencionado hasta el momento.

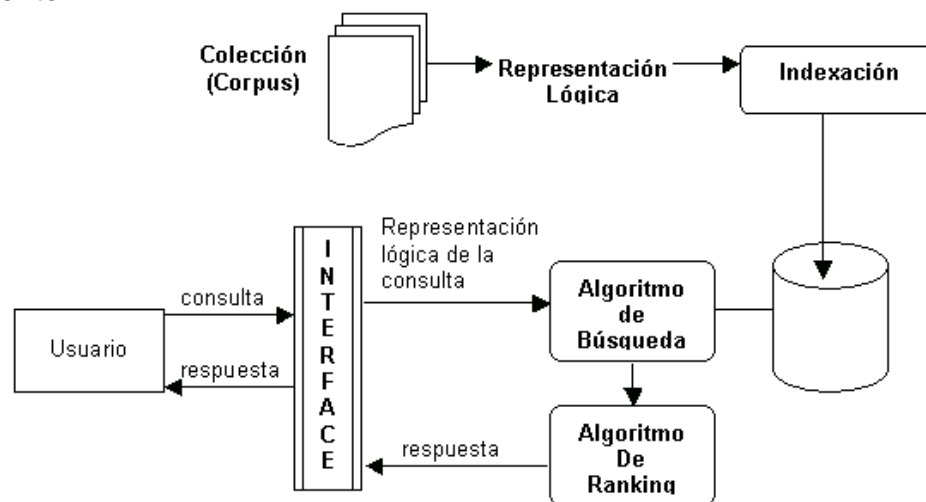


Figura 2 – Arquitectura b sica de un SRI



Como podemos observar, se inicia desde un conjunto de documentos de texto, los cuales est n compuestos por sucesiones de palabras que forman estructuras gramaticales (por ejemplo, oraciones y p rrafos). Tales documentos est n escritos en lenguaje natural y expresan ideas de su autor sobre un determinado tema. El conjunto de todos los documentos con los que se trata y sobre los que se deben realizar operaciones de RI se denomina **corpus**, **colecci n** o **base de datos textual o documental**. Para poder realizar operaciones sobre un corpus, es necesario obtener primero una **representaci n l gica** de todos sus documentos, la cual puede consistir en un conjunto de t rminos, frases u otras unidades (sint cticas o sem nticas) que permitan – de alguna manera – caracterizarlos. Por ejemplo, la representaci n de los documentos mediante un conjunto de sus t rminos se la conoce como “bolsa de palabras” (*bag of words*).

A partir de la representaci n l gica existe un proceso (**indexaci n**) que llevar  a cabo la construcci n de estructuras de datos (normalmente denominadas ** ndices**) que la almacene y soporte b squedas eficientes. Es importante destacar que una vez construidos los  ndices, los documentos del corpus pueden ser eliminados del sistema ya que  ste retornar  las referencias a los mismos porque cuenta con la informaci n necesaria para hacerlo. En tal caso, el usuario ser  el encargado de localizar el documento para consultarlo. A los sistemas que funcionan bajo este modelo se los denomina “sistemas referenciales”, en contraste con los que s  almacenan y mantienen los documentos denominados “sistemas documentales” [PE A]. Un ejemplo de sistemas referenciales son algunos de los motores de b squeda web, que retornan una lista de *urls* a los documentos, como – por ejemplo – Altavista¹. Un caso particular es el motor de b squeda Google² el cual – en algunos casos – almacena en memoria cach  el documento completo, el cual puede ser consultado durante cierto tiempo, incluso si ha desaparecido del sitio original.

El **algoritmo de b squeda** acepta como entrada una expresi n de consulta o *query* de un usuario y verificar  en el  ndice cu les documentos pueden satisfacerlo. Luego, un algoritmo de ranking determinar  la relevancia de cada documento y retornar  una lista con la respuesta. Se establece que el primer  tem de dicha lista corresponde al documento m s relevante respecto a la consulta y as  sucesivamente en orden decreciente.

¹ <http://www.altavista.com/>

² <http://www.google.com/>



La **interface** de usuario permite que  ste especifique la consulta mediante una expresi n escrita en un lenguaje preestablecido y – adem s – sirve para mostrar las respuestas retornadas por el sistema.

Si bien hasta aqu  se plante  la tarea b sica de la RI y la arquitectura general de un SRI, el  rea es muy amplia y abarca diferentes t picos. En general, un SRI no entrega una respuesta directa a una consulta, sino que permite localizar referencias a documentos que pueden contener informaci n  til. Pero  ste es s lo uno de los aspectos del  rea de RI en la actualidad, ya que se ha atacado el problema con una perspectiva m s amplia, proponiendo y desarrollando estrategias y modelos para mejorar y aumentar la funcionalidad de los SRI. Entre otras, la RI abarca t picos como:

- Modelos de Recuperaci n: La tarea de la recuperaci n puede ser modelada desde distintos enfoques, por ejemplo la estad stica, el  lgebra de Boole, el  lgebra de vectores, la l gica difusa, el procesamiento del lenguaje natural y dem s.
- Filtrado y Ruteo: Es un  rea que permite definir perfiles de necesidades informativas por parte de usuarios y ante el ingreso de nuevos documentos al SRI, se los analiza y reenv a a quienes se estime a que van a ser relevantes. [RIGOUTSOS]
- Clasificaci n: Aqu  se realiza la rotulaci n autom tica de documentos de un corpus en base a clases previamente definidas. [BEKKERMAN]
- Agrupamiento (*Clustering*): Es una tarea similar a la clasificaci n pero no existen clases predefinidas. El proceso autom ticamente determinar  cu les son las particiones.
- Sumarizaci n:  rea que entiende sobre t cnicas de extracci n de aquellas partes (palabras, frases, oraciones, p rrafos) que contienen la sem ntica que determina la esencia de un documento. [SMEDT]
- Detecci n de novedades (*Novelty Detection*): Se basa en la determinaci n de la introducci n de nuevos t picos o temas a un SRI. [LI]
- Respuestas a Preguntas (*Question Answering*): Consiste en hallar aquellas porciones de texto de un documento que satisfacen expresamente a una consulta, es decir, la respuesta concreta a una pregunta dada. [CORRADA] [KISUH]



- Extracci n de Informaci n: Extraer aquellas porciones de texto con una alta carga sem ntica y establecer relaciones entre los t rminos o pasajes extra dos. [MCCALLUM] [WADE]
- Recuperaci n cross-language: Hallar documentos escritos en cualquier lenguaje que son relevantes a una consulta expresada en otro lenguaje (b squeda multilingual). [CLOUGH]
- B squedas Web: Se refiere a los SRI que operan sobre un corpus web privado (intranet) o p blico (Internet). La web ha planteado nuevos desaf os al  rea de RI, debido a sus caracter sticas particulares como – por ejemplo – dinamismo y tama o. [CASTILLO]
- Recuperaci n de Informaci n Distribuida: A diferencia de los SRI cl sicos donde el corpus y las estructuras de datos que auxilian a la b squeda est n centralizadas, aqu  se plantea la tarea sobre los mismos elementos pero distribuidos sobre una red de computadoras.
- Modelado de Usuarios: Esta  rea – a partir de la interacci n de los usuarios con un SRI – estudia c mo se generan de forma autom tica perfiles que definan las necesidades de informaci n de  stos. [LIU] [JOACHIMS]
- Recuperaci n de Informaci n Multimedia: M s all  que los SRI tradicionales operan sobre corpus de documentos textuales, la recuperaci n de informaci n tiene que tratar con otras formas alternativas de representaci n como im genes, registro de conversaciones y video. [CLOUGH]
- Desarrollo de Conjuntos (data-sets) de Prueba: A los efectos de evaluar SRI completos o nuevos m todos y t cnicas es necesario disponer de juegos de prueba normalizados (corpus con preguntas y respuestas predefinidas, corpus clasificados, etc.). Esta  rea tiene que ver con la producci n tales conjuntos, a partir de diferentes estrategias que permitan reducir la complejidad de la tarea, manejando la dificultad inherente a la carga de subjetividad existente. [GUSTMAN] [SANDERSON]

 RECUPERACI N DE INFORMACI N O RECUPERACI N DE DATOS?

Muchos usuarios se encuentran familiarizados con el concepto de recuperaci n de datos (RD), especialmente aquellos que – a menudo – interact an con sistemas de consulta en bases de datos relacionales   en registros de alguna naturaleza, como por ejemplo, un registro de los



empleados de una organizaci n. Sin embargo, hay diferencias significativas en los conceptos que definen que el tratamiento de las unidades (datos o informaci n) en cada caso sean completamente diferentes.

B sicamente, existen diferencias sustanciales en cuanto a los objetos con que se trata y su representaci n, la especificaci n de las consultas y los resultados.

En el  rea de RD los objetos que se tratan son estructuras de datos conocidas. Su representaci n se basa en un formato previo definido y con un significado impl cito (hay una sintaxis y sem ntica no ambigua) para cada elemento. Por ejemplo, una tabla en una base de datos que almacena instancias de clientes de una organizaci n posee un conjunto de columnas que definen los atributos de todos los clientes y cada fila corresponde a uno en particular. N tese que cada elemento (atributo) tiene un dominio conocido y su sem ntica est  claramente establecida. Por otro lado, en el  rea de RI la unidad u objeto de tratamiento es b sicamente un documento de texto – en general – sin estructura.

En cuanto a la especificaci n de las consultas, en el  rea de RD se cuenta con una estructura bien definida dada por un lenguaje de consulta que permite su especificaci n de manera exacta. Las consultas no son ambiguas y consisten en un conjunto de condiciones que deben cumplir los  tems a evaluar para que la misma se satisfaga. Por ejemplo, en el modelo de bases de datos, las consultas especifican – entre otros – utilizando el lenguaje SQL (Structured Query Language) cuya sem ntica es precisa:

SQL

```
SELECT *  
FROM Clientes  
WHERE Localidad = "Chivilcoy"  
AND Saldo_Cuenta > 10000
```

En lenguaje natural

Seleccionar todos los clientes de Chivilcoy que deban m s de 10000 pesos (se sabe, por definici n, que lo que deben es su saldo de cuenta)

En este ejemplo, se puede ver la clara sem ntica de la consulta en SQL a partir de conocer que existe un atributo Localidad y otro Saldo_Cuenta y lo que cada uno representa. Sin embargo, esto no es tan directo ni tan simple cuando se trata de recuperar documentos en el contexto de la RI. En primer lugar, debido a que la necesidad de un usuario puede ser dif cil de expresar. Por ejemplo, sup ngase que se desea encontrar:



“Documentos que contengan información biográfica de los entrenadores de los equipos de fútbol de Argentina que ganaron más torneos en los últimos 10 años”

La primera dificultad consiste en construir una expresión de consulta que refleje exactamente esta necesidad de información del usuario. Especialmente, si se tiene en cuenta que para resolverla completamente quizá primero se requiera de conocer información parcial, por ejemplo, “ganaron más torneos en los últimos 10 años”. ¿Qué significa “ganaron más torneos”? Esta es una situación subjetiva y – en muchos casos – el sistema debe manejar estas cuestiones, junto con ambigüedades (por ejemplo, palabras cuyo significado está determinado por el contexto) e incompletitud de la mejor manera posible. De hecho, los documentos y las expresiones de consulta se interpretan de forma que el proceso de recuperación determine un grado de similitud entre éstos.

En un sistema de RD los resultados consisten en el conjunto completo de elementos que satisfacen todas las condiciones del *query*. Como la consulta no admite errores, el resultado es exacto, ni uno más, ni uno menos. Y el orden de aparición es simplemente casual (a menos que específicamente se desee ordenar por alguna columna), pero en todos los casos este orden es irrelevante respecto de la consulta y no significa nada, es decir, no se puede implementar sistema de ranqueo alguno. En el área de RI, aparece el concepto de relevancia y la salida (respuesta) se encuentra confeccionada de acuerdo a algún criterio que evalúa la “similitud” que existe entre la consulta y cada documento. Por lo tanto, el resultado es un ranking (que no es sinónimo de “orden”, tal como se lo entiende habitualmente en RD), donde la primera posición corresponde al documento más relevante a la consulta y así decrece sucesivamente. El proceso de recuperación de información puede retornar documentos que no sean relevantes para el usuario, es decir, que el conjunto de respuesta no es exacto.

Otros autores también establecieron las diferencias entre ambos conceptos: Grossman [GROSSMAN] claramente muestra la diferencia cuando enuncia que *“la recuperación de información es encontrar documentos relevantes, no encontrar simples correspondencias a unos patrones de bits”*. Nótese la diferencia sustancial que existe en tratar de encontrar documentos “relevantes” a una consulta o – simplemente – encontrar aquellos donde “coinciden” patrones de términos o se cumplen ciertas condiciones. En el caso de la RD, la tarea es relativamente sencilla, mientras que en área de RI es extremadamente compleja y no existe aún una solución definitiva al problema.



LA INTERACCI N DEL USUARIO CON EL SRI

La tarea de recuperar informaci n puede ser planteada de diversas formas, de acuerdo a c mo el usuario interact a con el sistema o bien qu  facilidades  ste le brinda. B sicamente, la tarea se la puede dividir en:

a. Recuperaci n inmediata: El usuario plantea su necesidad de informaci n y – a continuaci n – obtiene referencias a los documentos que el sistema eval a como relevantes. Existen dos modalidades:

1. B squeda (propriadamente dicha) o recuperaci n “ad-hoc”, donde el usuario formula una consulta en un lenguaje y el sistema la eval a y responde. En este caso, el usuario tiene suficiente comprensi n de su necesidad y sabe c mo expresar una consulta al sistema. Un ejemplo cl sico son los buscadores de Internet como Google³, Altavista⁴ o AllTheWeb⁵.

2. Navegaci n o *browsing*: En este caso, el usuario utiliza un enfoque diferente al anterior. El sistema ofrece una interface con temas donde el usuario “navega” por dicha estructura y obtiene referencias a documentos a relacionados. Esto facilita la b squeda a usuarios que no pueden definir claramente c mo comenzar con su consulta e – inclusive – van definiendo su necesidad a medida que observan diferentes documentos. Es este enfoque no se formula consulta expl cita. Un ejemplo t pico es el proyecto Open Directory⁶.

En ambos casos, la colecci n es relativamente est tica, es decir, se parte de un conjunto de documentos y la aparici n de nuevos no es muy significativa. Por otro lado, las consultas son las que se van modificando ya que este proceso es proactivo por parte del usuario.

b. Recuperaci n diferida: El usuario especifica sus necesidades y el sistema entregar  de forma continua los nuevos documentos que le lleguen y concuerden con  sta. Esta modalidad recibe el nombre de **filtrado y ruteo** y la necesidad del usuario – generalmente – define un “perfil” (*profile*) de los documentos buscados. N tese que un “perfil” es – de alguna forma – un *query* y puede ser tratado como tal. Cada vez que un nuevo documento arriba al sistema se compara con el perfil y – si es relevante – se env a al

³ <http://www.google.com/>

⁴ <http://www.altavista.com/>

⁵ <http://www.alltheweb.com/>

⁶ <http://www.dmoz.org/>



usuario. Un ejemplo, es el servicio provisto por la empresa Indigo Stream Technologies denominado GoogleAlert⁷.

En esta modalidad la consulta es relativamente estática (corresponde al *profile*) y el usuario tiene un rol pasivo. El dinamismo está dado por la aparición de nuevos documentos y es lo que determina mas resultados para el usuario.

En algunos casos, se plantea que documentos y consultas son objetos de la misma clase por lo que estos enfoques son – de alguna manera – visiones diferentes de una misma problemática. Bajo este punto de vista, documentos y consultas se pueden intercambiar. Sin embargo, esto no es siempre posible debido al tratamiento que se aplica a cada uno en diferentes sistemas. Algunos sistemas representan *queries* y documentos de diferente manera. Es más, existe una diferencia obvia en cuanto a la longitud de uno y otro que se tiene en cuenta bajo ciertos modelos de recuperación. Finalmente, en un sistema de búsqueda, existe el concepto de ranking de los documentos respuesta, mientras que en un modelo de filtrado, cada nuevo objeto es relevante a un perfil o no.

EL CONCEPTO DE RELEVANCIA

Como mencionamos, la recuperación de información intenta resolver el problema de encontrar documentos relevantes que satisfagan la necesidad de información del usuario. Sin embargo, se ha planteado la dificultad para llevar a cabo esta tarea debido a la imposibilidad de expresar exactamente tal necesidad. Además, la noción de relevancia es un juicio subjetivo [RIJSBERGEN] y depende de diferentes factores relacionados más cercanamente con el usuario. La relevancia de un documento respecto a un *query* se refiere a cuánto el primero responde al segundo. De igual manera, luego el usuario evalúa qué tanto, es decir, en qué medida, se satisface su necesidad de información [KORFHAGE].

Es por ello, que se plantea la relevancia como similitud, para poder comparar documentos con consultas y – bajo ciertos criterios – definir una medida de distancia entre ambos. Por lo tanto, se puede plantear la idea que “un documento es relevante a una consulta si son similares”, donde la medida de similitud puede estar basada en diferentes criterios (coincidencias de términos, significado de éstos, frecuencia de aparición de términos y distribución del vocabulario, entre otros).

⁷ <http://www.googlealert.com/>

Martínez Méndez y otros [MARTINEZ] resaltan la dificultad para determinar la relevancia o no de un documento respecto de una consulta. Plantean – por ejemplo – que dos personas pueden juzgar un mismo documento de diferente manera y que es difícil establecer los criterios para la evaluación de la relevancia. Finalmente, mencionan la idea de relevancia parcial, es decir, cuando solo una parte del documento se considera relevante.

Por otro lado, como el *query* no describe exactamente la necesidad de información del usuario, algunos autores [KORFHAGE] definen el concepto de “pertinencia”, donde se incluyen las restricciones impuestas por el SRI. Este concepto está relacionado con la utilidad del documento para el usuario [MARTINEZ], de acuerdo a la necesidad de información original que guió su búsqueda, independientemente si es en parte o todo el documento.

Sin embargo – y a pesar de las dificultades para determinarla – el concepto genérico de relevancia es aceptado ampliamente por la comunidad de RI para evaluar la respuesta de un SRI respecto de una consulta de un usuario, la cual – como ya mencionamos – surge a partir de una necesidad de información.

MODELOS DE RI

Los SRI toman un conjunto de documentos (colección) para procesar y luego poder responder consultas. De forma básica, podemos clasificar los documentos en estructurados y no estructurados. Los primeros son aquellos en los que se pueden reconocer elementos estructurales con una semántica bien definida, mientras que los segundos corresponden a texto libre, sin formato. La diferencia fundamental de un SRI que procese documentos estructurados se encuentra en que puede extraer información adicional al contenido textual, la cual utiliza en la etapa de recuperación para facilitar la tarea y aumentar las prestaciones.

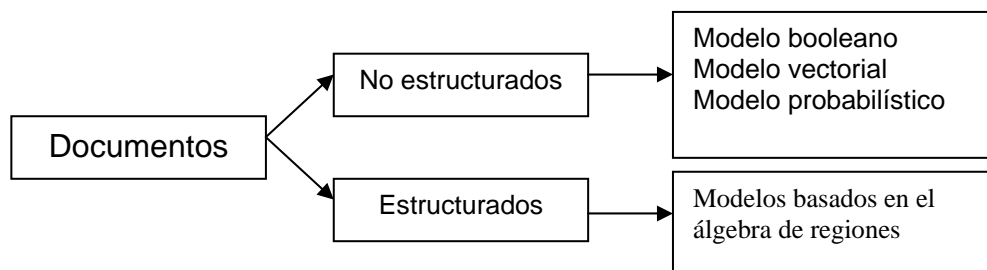


Figura 3 – Clasificación de modelos de RI



A partir de lo expresado anteriormente en la figura 3 se presenta una posible clasificación de modelos de RI – la cual no es exhaustiva – de acuerdo a características estructurales de los documentos. A continuación se describen – de forma somera – los modelos clásicos y el álgebra de regiones.

a) MODELO BOOLEANO

En el modelo booleano el representar la colección de documentos se realiza sobre una matriz binaria documento–término, donde los términos han sido extraídos manualmente o automáticamente de los documentos y representan el contenido de los mismos.

Las consultas se arman con términos vinculados por operadores lógicos (AND, OR, NOT) y los resultados son referencias a documentos donde cuya representación satisface las restricciones lógicas de la expresión de búsqueda. En el modelo original no hay ranking de relevancia sobre el conjunto de respuestas a una consulta, todos los documentos poseen la misma relevancia.

Si bien es el primer modelo desarrollado y aún se lo utiliza, no es el preferido por los ingenieros de software para sus desarrollos. Existen diversos puntos en contra que hacen que cada día se lo utilice menos y – además – se han desarrollado algunas extensiones, bajo el nombre modelo booleano extendido [WALLER] [SALTON_b], que tratan de mejorar algunos puntos débiles.

b) MODELO VECTORIAL

Este modelo fue planteado y desarrollado por Gerard Salton [SALTON_c] y – originalmente – se implementó en un SRI llamado SMART. Aunque el modelo posee más de treinta años, actualmente se sigue utilizando debido a su buena performance en la recuperación de documentos.

Conceptualmente, este modelo utiliza una matriz documento–término que contiene el vocabulario de la colección de referencia y los documentos existentes. En la intersección de un término t y un documento d se almacena un valor numérico de importancia del término t en el documento d ; tal valor representa su *poder de discriminación*. Así, cada documento puede ser visto como un vector que pertenece a un espacio n -dimensional, donde n es la cantidad de términos que componen el vocabulario de la colección. En teoría, los documentos que contengan términos similares estarán a muy poca

distancia entre sí sobre tal espacio. De igual forma se trata a la consulta, es un documento más y se la mapea sobre el espacio de documentos. Luego, a partir de una consulta dada es posible devolver una lista de documentos ordenados por distancia (los más relevantes primero). Para calcular la semejanza entre el vector consulta y los vectores que representan los documentos se utilizan diferentes fórmulas de distancia, siendo la más común la del coseno.

Obsérvese el siguiente ejemplo donde se representa a un documento d y a una consulta c :

Documento: “*La República Argentina ha sido nominada para la realización del X Congreso Americano de Epidemiología en Zonas de Desastre. El encuentro se realizará...*”

Consulta: “*argentina congreso epidemiología*”

	argentina	...	congreso	epidemiología	...
d_1	0.5		0.3	0.2	
...					
d_n					
Consulta	0.4		0.3	0.3	

Matriz término-documento con pesos normalizados entre 0 y 1

c) MODELO PROBABILÍSTICO

Fue propuesto por Robertson y Spark-Jones [ROBERTSON] con el objetivo de representar el proceso de recuperación de información desde el punto de vista de las probabilidades. A partir de una expresión de consulta se puede dividir una colección de N documentos (figura 4) en cuatro subconjuntos distintos: REL conjunto de documentos relevantes, REC conjunto de documentos recuperados, RR conjunto de documentos relevantes recuperados y NN el conjunto de documentos no relevantes no recuperados.

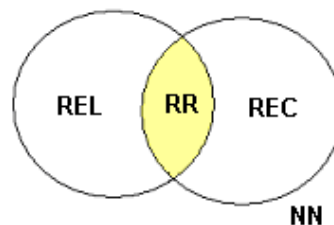


Figura 4. División de la colección



El resultado ideal de una consulta se da cuando el conjunto REL es igual REC. Como resulta dif cil lograrlo en primera intenci n, el usuario genera una descripci n probabil stica del conjunto REL y a trav s de sucesivas interacciones con el SRI se trata de mejorar la performance de recuperaci n. Dado que una recuperaci n no es inmediata dado que involucra varias interacciones con el usuario y que estudios han demostrado que su performance es inferior al modelo vectorial, su uso es bastante limitado.

d) MODELOS PARA DOCUMENTOS ESTRUCTURADOS

Los modelos cl sicos responden a consultas, buscando sobre una estructura de datos que representa el contenido de los documentos de una colecci n,  nicamente como listas de t rminos significativos. Un modelo de recuperaci n de documentos estructurados utiliza la organizaci n de los mismos a los efectos de mejorar la performance y brindar servicios alternativos al usuario (por ejemplo, uso de memoria visual, recuperaci n de elementos multimedia, mayor precisi n sobre el  mbito de la consulta y dem s).

La estructura de los documentos a indexar est  dada por marcas o etiquetas, siendo los est ndares m s utilizados el SGML (*Standard General Markup Language*), el HTML (*HyperText Markup Language*), el XML (*eXtensible Markup Language*) y LATEX.

Al poseer la descripci n de parte de la estructura de un documento es posible generar un grafo sobre el que se navegue y se respondan consultas de distinto tipo, por ejemplo:

- Por estructura: *  Cu les son las secciones del segundo cap tulo?*
- Por metadatos o campos: *Documentos de "Editorial UNLu" editados en 1998*
- Por contenido: *T rmino "agua" en t tulos de secciones*
- Por elementos multimedia: *Im genes cercanas a p rrafos que contengan "Bosch"*

Para Baeza-Yates existen dos modelos en esta categor a "nodos proximales" [NAVARRO] y "listas no superpuestas" [BURKOWSKI]. Ambos modelos se basan en almacenar las ocurrencias de los t rminos a indexar en estructuras de datos diferentes, seg n aparezcan en alg n elemento de estructura (regi n) o en otro como cap tulos, secciones, subsecciones y dem s. En general, las regiones de una misma estructura de datos no



poseen superposici n, pero regiones en diferentes estructuras s  se pueden superponer. Los tipos de consultas soportados son simples:

- Seleccione una regi n que contenga una palabra dada
- Seleccione una regi n X que no contenga una regi n Y
- Seleccione una regi n contenida en otra regi n

Sobre una estructura tipo libro un ejemplo de consulta ser a:

[subsecci n[+] CONTIENE "tambo"]

Como respuesta el SRI buscar a subsecciones y sub-subsecciones que contengan el t rmino "tambo".

Cabe mencionar que algunos motores de b squeda de Internet ya utilizan ciertos elementos de la estructura de un documento – por ejemplo, los t tulos – a los efectos de realizar tareas de ranqueo, resumen autom tico, clasificaci n y otras.

La expansi n de estos lenguajes de demarcaci n, especialmente en servicios sobre Internet, hace que se generen y publiquen cada vez m s documentos semiestructurados. Es necesario – entonces – desarrollar t cnicas que aprovechen el valor agregado de los nuevos documentos. Si bien – en la actualidad –  stas no se encuentran tan desarrolladas como los modelos tradicionales, consideramos su evoluci n como una cuesti n importante en el  rea de RI, especialmente a partir de investigaciones con enfoques diferentes que abordan la problem tica [EGNOR] [OGILVIE] [RAGHAVAN].

LA RI EN LA ERA DE LA WEB

Con la aparici n de la web surgieron nuevos desaf os para resolver en el  rea de recuperaci n de informaci n como consecuencia – principalmente – a sus caracter sticas y tama o. La web puede ser vista como un gran repositorio de informaci n, completamente distribuido sobre Internet y accesible por gran cantidad de usuarios. Por sus  rdenes como un espacio p blico existen millones de organizaciones y usuarios particulares que incorporan, quitan   modifican contenido continuamente, por lo que su estructura no es est tica.

Su contenido no respeta est ndares de calidad, ni estilos ni organizaci n. Como medio de publicaci n de informaci n de naturaleza diversa se ha



convertido en un servicio de permanente crecimiento. Una de las caracter  sticas de la informaci  n publicada en la web es su dinamismo, dado que pueden variar en el tiempo tanto los contenidos como su ubicaci  n [BREWINGTON] [LAWRENCE].

El tama  o de la web es imposible de medir exactamente y muy dif  cil de estimar. Sin embargo, se calcula que son decenas de terabytes de informaci  n, y crece permanentemente. Est   formada por documentos de diferente naturaleza y formato, desde p  ginas HTML hasta archivos de im  genes pasando por gran cantidad de formatos est  ndar y propietarios, no solamente con contenido textual, sino tambi  n con contenido multimedial.

La b  squeda de informaci  n en la web es una pr  ctica com  n para los usuarios de Internet y los sistemas de recuperaci  n de informaci  n web (conocidos como motores de b  squeda) se han convertido en herramientas indispensables para los usuarios. Su arquitectura y modo de operaci  n se basan en poder recolectar mediante un mecanismo adecuado los documentos existentes en los sitios web. Una vez obtenidos, se llevan a cabo tareas de procesamiento que permiten extraer t  rminos significativos contenidos dentro de los mismos, junto con otra informaci  n, a los efectos de construir estructuras de datos (  ndices) que permitan realizar b  squedas de manera eficiente.

Luego, a partir de una consulta realizada por un usuario, un motor de b  squeda extraer   de los   ndices las referencias que satisfagan la consulta y se retornar   una respuesta, acomodada en el ranking por diversos criterios al usuario. El modo de funcionamiento de los diferentes motores de b  squeda puede diferir en diversas implementaciones de los mecanismos de recolecci  n de datos, los m  todos de indexaci  n y los algoritmos de b  squeda y ranqueo.

Sin embargo, esta tarea no es sencilla y se ha convertido en un desaf  o para los SRI debido las caracter  sticas propias de la web. Baeza-Yates [BAEZA-YATES] plantea que hay desaf  os de dos tipos:

a) Respecto de los datos

- Distribuidos: La web es un sistema distribuido, donde cada proveedor publica su informaci  n en computadoras pertenecientes a redes conectadas a Internet, sin una estructura    topolog  a predefinida.



- Volátiles: El dinamismo del sistema hace que exista información nueva a cada momento ó bien que cambie su contenido ó inclusive desaparezca otra que se encontraba disponible.
- No estructurados y redundantes: Básicamente, la web está formada de páginas HTML, las cuales no cuentan con una estructura única ni fija. Además, mucho del contenido se encuentra duplicado (por ejemplo, espejado).
- Calidad: En general, la calidad de la información publicada en la web es altamente variable, tanto en escritura como en actualización (existe información que puede considerarse obsoleta), e inclusive existe información con errores sintácticos, ortográficos y demás.
- Heterogeneidad: La información se puede encontrar publicada en diferentes tipos de medios (texto, audio, gráficos) con diferentes formatos para cada uno de éstos. Además, hay que contemplar los diferentes idiomas y diferentes alfabetos (por ejemplo, árabe ó chino).

b) Respecto de los usuarios

- Especificación de la consulta: Los usuarios encuentran dificultades para precisar – en el lenguaje de consulta – su necesidad de información.
- Manejo de las respuestas: Cuando un usuario realiza una consulta se ve sobrecargado de respuestas, siendo una parte irrelevante.

Estas características – sumadas al tamaño de la web – imponen restricciones a las herramientas de búsqueda en cuanto a la cobertura y acceso a los documentos, exigiendo cada vez mayores recursos computacionales (espacio de almacenamiento, ancho de banda de las redes, ciclos de CPU) y diferentes estrategias para mejorar la calidad de las respuestas.

RECOMENDACIONES

Se plantean dos recomendaciones relativas al tema recuperación de información, su enseñanza e investigación por parte de instituciones hispanoamericanas.

- Se sugiere incorporar contenidos mínimos relativos a recuperación de información en carreras de grado relacionadas con la informática. Deberían



implementarse como cursos regulares, que se impartan luego que el alumno haya aprendido conceptos de bases de datos. Tal capacitación le permitiría tener una visión de la realidad existente en los sectores de las organizaciones donde la materia prima sean los documentos.

- Se recomienda crear y apoyar grupos de investigación interdisciplinarios que puedan estudiar y proponer nuevas técnicas de RI adaptadas a nuestra lengua. Uno de los principales problemas para la región es que hay pocos investigadores dedicados a este tema y, en general, cuando hacen sus investigaciones en países no hispanos el lenguaje principal foco de sus estudios es el inglés.

REFERENCIAS BIBLIOGRÁFICAS

[BAEZA-YATES] Baeza-Yates, R. y Ribeiro-Neto, B. (1999) **Modern Information Retrieval**. Ed. ACM Press. Addison Wesley.

[BEKKERMAN] Bekkerman, R. y Allan, J., (2005) **Using Bigrams in Text Categorization**, CIIR Technical Report.

[BREWINGTON] Brewington, B. E. y Cybenko Thayer, G. (2000) **How Dynamic is the Web?** Proceedings of the Ninth International World Wide Web Conference. 2000.

[BURKOWSKI] Burkowski, F. (1992) **Retrieval activities in a database consisting of heterogeneous collections of structured texts**. Belkin, N., Ingwersen, P., Pejtersen, A. M., and Fox, E., editors, Proceedings of the 15th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval. New York. ACM Press. 112–125.

[CARLSON] Carlson, C. (2003) **Information overload, retrieval strategies and Internet user empowerment**. Haddon, Leslie, Eds. Proceedings The Good, the Bad and the Irrelevant (COST 269), Helsinki (Finland). 1(1): 169-173.

[CASTILLO] Castillo, C y Baeza-Yates. R. (2005) **WIRE: an Open Source Web Information Retrieval Environment**. Workshop on Open Source Web Information Retrieval (OSWIR).

[CLOUGH] Clough, P, Sanderson, M. y Muller, H. (2004) **The CLEF Cross Language Image Retrieval Track (ImageCLEF)**. CIVR 2004: 243-251



- [CORRADA] Corrada E. A. y Croft, W.B., (2004) **Answer Models for Question Answering Passage Retrieval**, Poster en Proceedings of the 27th Annual International ACM SIGIR Conference, pp. 516-517.
- [CROFT] Croft, W. B. (1987) **Approaches to intelligent information retrieval**. Information Processing & Management, 23(4): 249-254
- [EGNOR] Egnor, D. y Lord, R. (2000) **Structured information retrieval using XML**. Proceedings of the ACM SIGIR 2000 Workshop on XML and Information Retrieval.
- [GROSSMAN] Grossman, D. y Frieder, O. (1998) **Information Retrieval. Algorithms and Heuristics**. Ed. Kluwer Academic Publishers.
- [GUSTMAN] Gustman, James Mayfield, Liliya Kharevych y Stephanie Strassel, (2004) **Building an information retrieval test collection for spontaneous conversational speech**. Proceedings of ACM SIGIR 2004, pp. 41-48.
- [JOACHIMS] Joachims, T. Granka, L. Pang, B. Hembrooke, H. y Gay, G. (2005) **Accurately Interpreting Clickthrough Data as Implicit Feedback**, Proceedings of the ACM Conference on Research and Development on Information Retrieval (SIGIR)
- [KISUH] Kisuh Ahn, Johan Bos, James R. Curran, Dave Kor, Malvina Nissim, Bonnie Webber (2005): **Question Answering with QED at TREC-2005. In Voorhees and Buckland** (eds.): The Fourteenth Text REtrieval Conference, TREC 2005.
- [KORFHAGE] Korfhage, R. R. (1997) **Information Storage and Retrieval**. New York. Ed. Wiley Computer Publishing.
- [LAWRENCE] Lawrence, S. y Giles, L. (1999) **Accessibility and Distribution of Information on the Web**. Nature, 400(6740): 107-109. 1999.
- [LI] Li, X. and Croft, W.B., (2005) **Novelty Detection Based on Sentence Level Patterns**, En Proceedings of ACM Fourteenth Conference on Information and Knowledge Management (CIKM), Bremen, Germany, October, pp. 744-751



- [LIU] Liu, X., Croft, W.B., Oh, P. and Hart, D., (2004) **Automatic Recognition of Reading Levels from User Queries**, En Proceedings of SIGIR '04, pp. 548-549.
- [MAES] Maes, P. (1994). **Agents that Reduce Work and Information Overload**. Communications of the ACM, 37(7): 30-40.
- [MARTINEZ] Mart nez M endez, F. J. y Rodr guez Mu oz, J. V. (2004) **Reflexiones sobre la Evaluaci n de los Sistemas de Recuperaci n de Informaci n: Necesidad, Utilidad y Viabilidad**". Anales de Documentaci n, 7:153-170.
- [MCCALLUM] McCallum, A., (2005) **"Information Extraction: Distilling Structured Data from Unstructured Text," in ACM Queue**. Vol 3, No 9, November 2005, pp. 48-57.
- [NAVARRO] Navarro, G. y Baeza-Yates, R. (1995) **A language for queries on structure and contents of textual databases**. Proceedings of the 18th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval. New York. ACM Press. 93-101.
- [OGILVIE] Ogilvie, P. y Callan, J. (2003) **Language Models and Structured Document Retrieval**. Proceedings of the First INEX Workshop.
- [PE A] Pe a, R., Baeza-Yates, R., Rodr guez, J. V. (2003) **Gesti n Digital de la Informaci n**. Ed. Alfaomega Grupo Editor.
- [RAGHAVAN] Raghavan, S. y Garcia-Molina, H. (2001) **Integrating diverse information management systems: A brief survey**. IEEE Data Engineering Bulletin, 24(4):44-52.
- [RIGOUTSOS] Rigoutsos, I. y Chung-kwei, T. (2004) **A pattern-discovery-based system for the automatic identification of unsolicited e-mail messages (spam)**. En Proceedings of the First Conference on Email and Anti-Spam (CEAS'04).
- [RIJSBERGEN] van Rijsbergen, C.J. (1979) **Information Retrieval. Department of Computing Science**. Ed. University of Glasgow.
- [ROBERTSON] Robertson, S.E y Spark-Jones, K. (1976) **Relevance Weighting of Search terms**. Journal of Documentation. 33:126-148.



- [SANDERSON] Sanderson, M. y Hideo Joho (2004) **Forming test collections with no system pooling**. SIGIR 2004: 33-40
- [SALTON_a] Salton, G. Y Mc Gill, M.J. (1983) **Introduction to Modern Information Retrieval**. New York. Ed. Mc Graw-Hill Computer Series.
- [SALTON_b] Salton, G.; Fox, E.A. y Wu, H. (1983) **Extended Boolean information retrieval**. Communications of the ACM, 26(11):1022-1036. Noviembre,
- [SALTON] Salton, G. (1971) (editor). **The SMART Retrieval System – Experiments in Automatic Document Processing**. Ed. Prentice Hall Inc. Englewood Cliffs, NJ.
- [SMEDT] Smedt, K. Liseth, A. Hassel, M. y Dalianis, H. (2005). **How short is good? An evaluation of automatic summarization**. En Holmboe, H. (ed.) **Nordisk Sprogteknologi 2004**. Årbog for Nordisk Språkteknologisk Forskningsprogram 2000-2004, pp 267-287.
- [WADE] Wade, C. y Allan, J., (2005) **Passage Retrieval and Evaluation**, CIIR Technical Report.
- [WALLER] Waller, W. G. y Kraft, D. H. (1979) **A mathematical model for a weighted Boolean retrieval system**". Information Processing and Management, 15(5):235-245. 1979.