

Resolución de la anáfora: estructura del diálogo y conocimiento lingüístico*

Patricio Martínez-Barco y Manuel Palomar

Departamento de Lenguajes y Sistemas Informáticos

Universidad de Alicante

Carretera de San Vicente del Raspeig - Alicante - España

Tel. 965903653 Fax. 965909326

patricio@dlsi.ua.es, mpalomar@dlsi.ua.es

Resumen En este artículo presentamos un esquema de anotación sobre la estructura de los diálogos que puede ser aplicado a un algoritmo de resolución de anáfora basado en restricciones y preferencias. Este algoritmo será capaz de combinar la información lingüística tradicionalmente usada en algoritmos de resolución de anáfora junto con la información que proporciona la estructura del diálogo. Esta información ha sido aplicada al módulo de resolución de la anáfora de un sistema de Procesamiento de Lenguaje Natural al que se le ha introducido un corpus de diálogos en castellano obteniéndose un porcentaje de acierto del 73.8% para el tratamiento de la anáfora pronominal y un 78.9% para el tratamiento de la anáfora adjetiva. Para demostrar la influencia de esta información en la resolución de la anáfora se ha comparado nuestra propuesta con otras propuestas de sistemas en los que se usaban distintas fuentes de conocimiento que alcanzaban siempre resultados inferiores.

1 Introducción

La resolución de la anáfora es una de las tareas básicas a realizar en la mayoría de los sistemas de Procesamiento del Lenguaje Natural. De acuerdo a la definición de Hirst [Hir81], la anáfora en el discurso es un mecanismo para formular una referencia abreviada (esto quiere decir que contiene sólo alguna información de desambiguación más que sea léxica o fonéticamente corta) sobre alguna entidad esperando que el receptor sea capaz de desabreviar esa referencia y determinar, por tanto, la identidad de esa entidad. Así, la

referencia a una entidad es lo que se conoce como anáfora y la entidad se conoce como referente o antecedente¹.

De esta forma, los participantes en un diálogo usan la anáfora como una herramienta para indicar al oyente que aquello de lo que está hablando hace referencia a algo previamente nombrado. Y por otra parte, el oyente que detecta una anáfora debe desambiguar entre todas las entidades nombradas para obtener aquella a la que se ha hecho referencia. Así podrá añadir en su esquema mental la nueva información que ha sido vertida sobre esta entidad antigua. La anáfora por tanto, se convierte en un mecanismo que se usa en el diálogo para que tanto hablante como oyente centren su atención en las mismas entidades.

Los mecanismos de resolución de anáfora en el Procesamiento del Lenguaje Natural intentan emular el comportamiento de los humanos en conversación poniendo en marcha los mismos conocimientos que usan éstos para poder desambiguar el antecedente correcto. En este sentido se puede hacer una distinción entre aquellos sistemas de resolución de anáfora que se basan en el empleo de conocimiento lingüístico (como información léxica, morfológica, sintáctica o semántica) y aquellos que intentan modelar la estructura del discurso para extraer la información de desambiguación. Entre las primeras aproximaciones podemos destacar los algoritmos planteados por Hobbs [Hob86], Lappin y Leass [LL94], Mitkov [Mit98], y Baldwin [Bal97], quienes resuelven la anáfora discursiva con buenos resultados. Entre las segundas apro-

¹Destacaremos en este punto, la existencia de un tipo de anáfora conocida como anáfora abstracta en la que el referente no es una entidad como tal sino un concepto abstracto como puede ser una acción, una oración completa, etc. Sin embargo, en este artículo trabajaremos únicamente con anáforas individuales, es decir, aquellas cuyo referente es siempre una entidad

* Este artículo ha sido subvencionado por la CICYT (Comisión Interministerial de Ciencia y Tecnología) bajo el proyecto TIC97-0671-C02-01/02.

ximaciones destacamos el trabajo de Grosz *et al* [GJW95] quienes presentan el *centering* como un modelo para explicar la coherencia en el discurso local². Posteriormente, Strube y Hahn [SH99] realizan una adaptación del *centering* adaptándola a lenguajes de orden libre de palabras como el alemán o el castellano. Así mientras las aproximaciones lingüísticas centran su atención en establecer restricciones que les permitan definir claramente el antecedente, las aproximaciones que modelan la estructura del discurso se centran en la definición de un espacio de accesibilidad anafórico claro donde se pueden encontrar los antecedentes.

De esta forma, nuestra propuesta pretende aunar ambas aproximaciones de tal forma que el algoritmo de resolución de la anáfora sea capaz de manejar un espacio de accesibilidad anafórico basado en la estructura del diálogo, y además use la información lingüística para restringir el antecedente. Por ello, presentamos un algoritmo de resolución de la anáfora basado en restricciones (definidas mediante conocimiento lingüístico) y preferencias (basadas en la estructura del discurso junto con un conjunto de reglas heurísticas).

Puesto que otros trabajos ya han tratado abiertamente el uso de la información lingüística en la resolución de la anáfora para el castellano [FPM98] [FPM99], en este trabajo nos centraremos básicamente en la definición de un esquema de anotación apropiado para aportar la información sobre la estructura del diálogo que ayuda a la resolución de la anáfora de la anáfora pronominal (introducida a través de pronombre) y adjetiva³ (introducida mediante la omisión del núcleo en un sintagma nominal con modificadores, como en *la blanca*, *el primer*, etc.) . Destacaremos que el estudio de la anáfora en diálogos escritos en castellano constituye una línea de investigación novedosa actualmente.

²Según Beaugrande y Dressler [BD72], todo texto se caracteriza por tener dos propiedades: la coherencia y la cohesión. La coherencia mide la ligadura entre constituyentes a nivel semántico. Por otra parte la cohesión contribuye a la coherencia en niveles inferiores, como el sintáctico, léxico o morfológico. Así la anáfora es un mecanismo de cohesión que contribuye a la coherencia del discurso

³Nos centramos en estos dos tipos de anáfora puesto que son los más frecuentes en los diálogos analizados. Otros tipos de anáfora como los que genera la omisión del sujeto pronominal se abordarán en el futuro

Para ello, el artículo ha sido organizado de la siguiente forma: en la sección 2 se presentará el esquema de anotación propuesto, en la sección 3 expondremos brevemente el sistema de restricciones y preferencias usado y finalmente se demostrará la importancia de la fuente de conocimiento presentando una evaluación del sistema completo usando para ello el corpus de diálogos proporcionado por el proyecto *Basurde*⁴. Este corpus contiene una serie de diálogos telefónicos grabados entre el operador de una compañía de ferrocarriles y algunos usuarios de la misma que llamaban para solicitar información sobre los servicios de la compañía.

2 Esquema de anotación para la estructura del diálogo

De acuerdo a la exposición anterior, consideramos necesario el uso de una anotación de la estructura del diálogo para la resolución de la anáfora. Desde este punto de vista, proponemos un esquema de anotación para diálogos basado en el trabajo llevado a cabo por Gallardo [Gal96], quien aplica al castellano las teorías pragmáticas de Sacks *et al* [SSJ74] sobre los sistemas conversacionales. De acuerdo a estas teorías, la unidad básica del diálogo es el *turno*.

Puesto que nuestro trabajo parte de diálogos hablados que han sido transcritos mediante un transcriptor automático, el turno ya aparece anotado en los textos que vamos a procesar junto con el hablante que formula cada uno de los turnos. Además, uno de los principales problemas que tiene el tratamiento de diálogos que es la falta de signos de puntuación, ya ha sido resuelto por el transcriptor puesto que cada turno contiene oraciones más o menos completas pero al menos puntuadas. De esta forma nuestra anotación manual se centra en la clasificación de los turnos, y cómo éstos se agrupan para formar pares adjacientes, cuya etiqueta podrá ser procesada por nuestro sistema.

Como conclusión proponemos el siguiente esquema de anotación para la estructura del diálogo:

Turno (T) identificado por el cambio de hablante. De acuerdo a la clasificación propuesta por Gallardo, hay dos clases de

⁴BASURDE: Sistema de diálogo de habla espontánea en dominios limitados. CICYT (TIC98-423-C06).

turnos diferentes: Se considera **Turno de Intervención (IT)** aquel que añade información al flujo del diálogo. Estos turnos constituyen el *sistema primario de conversación*. Los hablantes hacen uso de sus intervenciones para proporcionar información que facilita el progreso de la conversación. Además las Intervenciones pueden ser clasificadas como **iniciaciones (IT_I)** cuando formulan invitaciones, requerimientos, ofertas, etc., o **reacciones (IT_R)** cuando responden o evalúan una iniciativa anterior. Por otra parte, un **Turno de Continuación (CT)** representa un turno vacío, cuyo uso por parte del hablante se limita al refuerzo y ratificación formal de los roles conversacionales. Básicamente son pequeñas interrupciones sin información que serán desechadas por el oyente.

Par Adyacente o Intercambio (AP) es un grupo de turnos T encabezados por un turno iniciación (IT_I) y finalizado por un turno reacción (IT_R). Una forma de anáfora que suele ser muy común en diálogos es la referencia cuyo espacio de accesibilidad corresponde a un par adyacente [Fox87]. El par adyacente es usado por los hablantes para definir un tema de conversación local, por lo tanto, muchas anáforas generadas hacen referencia a este tema local.

Segmento de tema (TOPIC) es un segmento del diálogo formado por un conjunto de pares adyacentes en los que los hablantes tratan un tema común. Definen, por tanto, el tema global. Otro tipo de anáforas muy corriente hará referencia siempre a este tema global.

Según el esquema anterior, conociendo el ámbito de los pares adyacentes así como el del tema global, podremos definir el espacio de accesibilidad anafórico donde los oyentes suelen encontrar los antecedentes de las anáforas, tanto a nivel local como global. Para ello proponemos etiquetar el siguiente conjunto de etiquetas: IT_I, IT_R, CT and AP.

En la figura 1 presentamos un ejemplo de texto anotado según estos criterios. La entrada para el esquema de anotación ya contiene las marcas de los turnos, donde la etiqueta (OP) indica el turno del operador de la compañía de ferrocarriles y (US) indica el

INPUT
<US> ¿me puedes decir algún tren que salga mañana para Monzón?
<OP> ¿por la mañana o por la tarde?
<US> por la tarde
<OP> si, hay un Talgo a las tres y media
<US> si
<OP> y un Intercity a las cinco y media

OUTPUT
<TOPIC> algún tren que salga mañana para Monzón
<AP1><IT _I ><US> ¿me puedes decir algún tren que salga mañana para Monzón?
<AP2><IT _R ><OP> ¿por la mañana o por la tarde?
<IT _R ><OP> por la tarde
<\AP2>
<IT _R ><OP> si, hay un Talgo a las tres y media
<CT><US> si
<IT _R ><OP> y un Intercity a las cinco y media
<\AP1>
<\TOPIC>

Figura 1: *Ejemplo de esquema de anotación*

turno del usuario, y como salida se proporciona una clasificación de turnos junto con la correspondiente generación de pares adyacentes, así como la definición del segmento del tema global y el propio tema global. Destacaremos en este punto que si bien nuestra labor de anotación de la estructura del diálogo ha sido realizada manualmente, hay actualmente anotadores automáticos de pares adyacentes como los que se están desarrollando en el proyecto BASURDE [BAS01] así como segmentadores automáticos de temas como los desarrollados en Reynar [Rey94] y [Rey99], o la propuesta de segmentación automática aplicada a la resolución de la anáfora desarrollada por Martínez-Barco y Palomar [MB99].

Para garantizar la fiabilidad de los resultados obtenidos con este corpus, el proceso de anotación manual debe ser realizado por al menos dos anotadores y posteriormente realizar una prueba que confirme la fiabilidad de la anotación. En nuestro caso confirmaremos la anotación usando el test de fiabilidad propuesto por Carletta *et al* [CII⁺97].

3 *Algoritmo de resolución de la anáfora basado en restricciones y preferencias*

El algoritmo de resolución de la anáfora que proponemos está basado en el uso de restricciones y preferencias. De acuerdo a este planteamiento, el algoritmo tiene tres fases básicas:

3.1 Identificación del espacio de accesibilidad anafórico

En la primera fase el algoritmo debe identificar el espacio donde se encontrarán los antecedentes. De acuerdo con las ideas anteriormente expuestas, utilizaremos el par adyacente como unidad básica para la definición del tema de conversación local. Así consideramos que el antecedente de una anáfora se encontrará fundamentalmente en el par adyacente donde aparece la anáfora, o bien en el par adyacente anterior (esto es debido a que las anáforas que aparecen al principio del par adyacente suelen encontrar el antecedente en el anterior). También, en aquellos pares adyacentes que contengan internamente el par adyacente de la anáfora (pares adyacentes anidados) pueden contener el antecedente. Por último, para aquellos casos en los que la anáfora haga referencia al tema del segmento tomaremos la entidad del tema también como parte del espacio de accesibilidad. De esta forma el espacio de accesibilidad define una lista con todas las entidades aparecidas en estos lugares hasta la aparición de la anáfora.

3.2 Aplicación de restricciones

Tras obtener la lista de posibles antecedentes, se debe filtrar la existencia de incompatibilidades entre la anáfora y el antecedente. Para ello se propone usar un conjunto de restricciones como la concordancia morfológica, el cumplimiento de las condiciones sintácticas tal y como proponen Lappin y Leass [LL94] en el caso de la anáfora pronominal, y la condición de no ser un nombre propio para la anáfora adjetiva. De esta forma, aquellos candidatos que no cumplan alguna de las restricciones serían automáticamente desechados.

3.3 Aplicación de preferencias

Finalmente, para garantizar una única solución al final del proceso, se aplicará un conjunto de preferencias que sea capaz de obtener el candidato más probable. De nuevo, la estructura del diálogo junto con un pequeño conjunto de reglas heurísticas nos permitirá obtener el resultado de la anáfora. Así definimos una lista de reglas ordenada de mayor a menor grado de preferencia:

1. Preferencia por candidatos que se encuentran en el mismo par adyacente que la anáfora.

2. Preferencia por candidatos que se encuentran en el par adyacente anterior a la anáfora.
3. Preferencia por candidatos que se encuentran en algún par que contenga al de la anáfora.
4. Preferencia por candidatos que se encuentran el tema global.
5. Preferencia por candidatos que se encuentran en la misma posición relativa respecto al verbo que la anáfora, es decir, antes o después (para la pronominal) o preferencia por candidatos que comparten el mismo tipo de modificador (para la adjetiva).
6. Preferencia por candidatos que comparten el mismo número de constituyente dentro de la oración (para la pronominal) o por los que comparten el mismo modificador (para la adjetiva).
7. Preferencia por el candidato más cercano.

De esta forma, nuestro algoritmo para la resolución de la anáfora usa y combina dos clases de conocimientos: a) el conocimiento lingüístico para definir el conjunto de restricciones y algunas reglas heurísticas del conjunto de preferencias; y b) el conocimiento de la estructura del diálogo que se usa para definir el espacio de accesibilidad anafórico y las preferencias derivadas del uso de este espacio.

4 Evaluación de la propuesta

Para llevar a cabo la evaluación del algoritmo se seleccionaron 40 de los diálogos proporcionados por el proyecto *Basurde* (5 de los cuales sirvieron para el propio entrenamiento del proceso), que fueron previamente transcritos, y posteriormente anotados usando un etiquetador automático de categorías gramaticales junto con información morfológica. Después se realizó la anotación manual de la estructura del diálogo asegurando su fiabilidad, y finalmente se les aplicó un analizador sintáctico parcial para extraer la información sintáctica necesaria para el algoritmo de resolución de la anáfora.

Una descripción completa del preproceso de nuestra propuesta junto con el módulo de resolución de la anáfora se presenta en la figura 2.

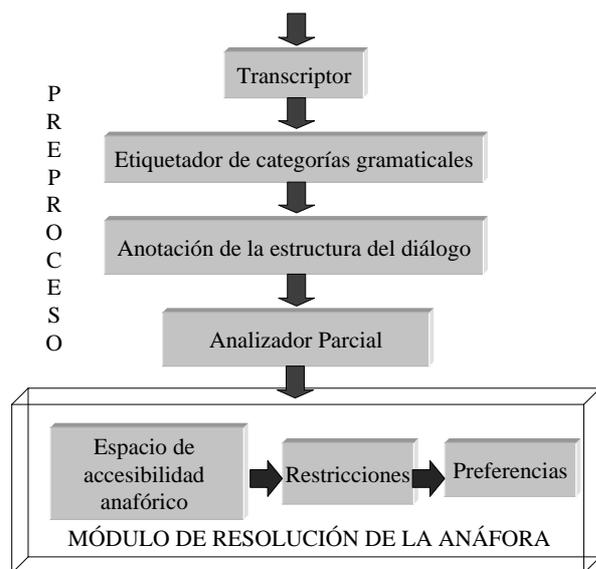


Figura 2: Descripción del proceso y preproceso del diálogo

Para demostrar la importancia de la inclusión de la información sobre la estructura del diálogo en el algoritmo de resolución de anáfora, se realizaron varios experimentos estudiando diferentes combinaciones de conocimientos. Estos experimentos pueden quedar resumidos en:

Experimento 1. Basado en conocimiento lingüístico. Para este experimento se tomó como base el algoritmo de restricciones y preferencias definido en Ferrández *et al* [FPM99] quien reporta resultados del 82.3% de precisión en el caso de la anáfora pronominal en textos no dialogados aplicando únicamente conocimiento léxico, morfológico y sintáctico (sin incluir información sobre la estructura del diálogo). Sin embargo, cuando el sistema se aplica a los textos de diálogos seleccionados, la precisión del sistema desciende hasta el 59% de precisión en la anáfora pronominal y el 23.7% de precisión para la anáfora adjetiva.

Experimento 2. Basado en conocimiento lingüístico y conocimiento sobre pares adyacentes. Partiendo del mismo conjunto de restricciones definido en el experimento 1 y sustituyendo el espacio de accesibilidad y el conjunto de preferencias por los definidos en la sección anterior (exceptuando la información sobre el tema global que no se incluye en este experimento) se obtiene una mejora en las prestaciones del algoritmo que alcanza un 70.2% de precisión en la anáfora

pronominal y un 68.4% para la adjetiva.

Experimento 3. Basado en conocimiento lingüístico y conocimiento sobre pares adyacentes y el tema global. Finalmente, la inclusión de la información sobre el tema global mejora el resultado del sistema hasta obtener 73.8% de precisión para la anáfora pronominal y 78.9% para la adjetiva.

Discusión. Un estudio de los casos de en los que nuestro sistema del experimento 3 falló, nos llevó a estimar las siguientes causas: un 16.4% del total de anáforas pronominales resueltas fueron erróneas debido a la falta de información semántica, un 6.6% debido a fallos en el sistema de restricciones, 1.6% debido a fallos en el sistema de preferencias y 1.6% a errores no clasificados. En la anáfora adjetiva, el 13% fue erróneo por la falta de información semántica, el 2.7% por el sistema de preferencias y el 5.4% por errores no clasificados.

5 Conclusiones

En este artículo hemos demostrado la necesidad de combinar distintas clases de conocimientos para resolver la anáfora pronominal y adjetiva en diálogos, usando por una parte la información de la estructura del diálogo a través de pares adyacentes y el tema global, y por otra parte el conocimiento lingüístico que proporciona la información léxica, morfológica y sintáctica. Además, tal y como demuestra la estimación realizada, la incorporación de información semántica al algoritmo podría llevar a alcanzar un 90.2% de precisión en el caso de la anáfora pronominal y un 91.9% en la adjetiva.

6 Agradecimientos

Agradecemos a Ferrán Plá (Universitat Politècnica de València) su ayuda prestada en el etiquetado de los diálogos así como a Encarna Segarra (Universitat Politècnica de València) por proporcionarnos el corpus de diálogos de proyecto Basurde, sin cuya colaboración no hubiera sido posible la evaluación de nuestro sistema.

Referencias

- [Bal97] B. Baldwin. CogNIAC: high precision coreference with limited knowledge and linguistic resources. In *Proceedings of ACL/EACL workshop on Operational factors in prac-*

- tical, robust anaphora resolution*, Madrid (Spain), July 1997.
- [BAS01] Proyecto BASURDE. *Spontaneous-Speech Dialogue System in Limited Domains*. CICYT (TIC98-423-C06), 1998-2001.
- [BD72] R. Beaugrande and W.U. Dressler. *Einführung in die Textlinguistik*. Max Niemeyer Verlag, Tübingen (Germany), 1972.
- [CII⁺97] J. Carletta, A. Isard, S. Isard, J.C. Kowtko, G. Doherty-Sneddon, and A.H. Anderson. The Reliability of a Dialogue Structure Coding Scheme. *Computational Linguistics*, 23(1):13–32, 1997.
- [col98] *36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics (COLING-ACL'98)*, Montreal (Canada), August 1998.
- [Fox87] B. Fox. *Discourse Structure and Anaphora*. Written and conversational English. Cambridge Studies in Linguistics. Cambridge University Press, Cambridge, 1987.
- [FPM98] A. Ferrández, M. Palomar, and L. Moreno. Anaphora resolution in unrestricted texts with partial parsing. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics (COLING-ACL'98)* [col98], pages 385–391.
- [FPM99] A. Ferrández, M. Palomar, and L. Moreno. An empirical approach to Spanish anaphora resolution. *Machine Translation*, (14), 1999.
- [Gal96] B. Gallardo. *Análisis conversacional y pragmática del receptor*. Colección Sinapsis. Ediciones Episteme, S.L., Valencia, 1996.
- [GJW95] B. Grosz, A. Joshi, and S. Weinstein. Centering: a framework for modeling the local coherence of discourse. *Computational Linguistics*, 21(2):203–225, 1995.
- [Hir81] G. Hirst. *Anaphora in Natural Language Understanding*. Springer-Verlag, Berlin, 1981.
- [Hob86] J. Hobbs. Resolving pronoun references. In B. Grosz B.L. Webber and K. Jones, editors, *Readings in Natural Language Processing*. Morgan Kaufmann, Palo Alto, CA, 1986.
- [LL94] S. Lappin and H.J. Leass. An algorithm for pronominal anaphora resolution. *Computational Linguistics*, 20(4):535–561, 1994.
- [MB99] P. Martínez-Barco. Algoritmo de resolución de la anáfora pronominal en diálogos. *Procesamiento del Lenguaje Natural*, (24):75–82, 1999.
- [Mit98] R. Mitkov. Robust pronoun resolution with limited knowledge. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics (COLING-ACL'98)* [col98].
- [Rey94] Jeffrey C. Reynar. An automatic method of finding topic boundaries. In *Proceedings of 32nd Annual Meeting of the Association for Computational Linguistics*, pages 331–333, Las Cruces, New Mexico, 1994.
- [Rey99] Jeffrey C. Reynar. Statistical Models for Topic Segmentation. In *Proceedings of 37th Annual Meeting of the Association for Computational Linguistics*, pages 357–364, Maryland, USA, June 1999.
- [SH99] M. Strube and U. Hahn. Functional Centering - Grounding Referential Coherence in Information Structure. *Computational Linguistics*, 25(5):309–344, 1999.
- [SSJ74] H. Sacks, E. Schegloff, and G. Jefferson. A simplest systematics for the organization of turn taking for conversation. *Language*, 50(4):696–735, 1974.