



ESTRUTURA DO LÉXICO: MODELO LINGÜÍSTICO-COMPUTACIONAL DE
REPRESENTAÇÃO DAS RELAÇÕES SEMÂNTICAS
(THE LEXICON STRUCTURE: A LINGUISTIC-COMPUTACIONAL MODEL OF
REPRESENTATION OF THE SEMANTIC RELATIONS)

Bento Carlos DIAS-DA-SILVA (Universidade Estadual Paulista)
Mirna Fernanda de OLIVEIRA (Universidade Estadual Paulista - PG)

ABSTRACT: In this paper an attempt is made to investigate the linguistic and computational representation issues of the main sense relations responsible for the lexicon semantic structure: the conceptual and lexicosemantic relations. After a brief review of the relevant theories, a machine-tractable model of such relations is presented.

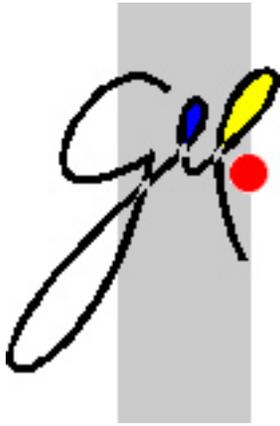
KEYWORDS: lexical semantics; sense relations; lexicon structure; natural language processing.

0. Introdução

Seja pela mente humana, seja pela máquina, o processamento das informações lingüisticamente codificadas por uma língua natural é uma atividade surpreendentemente complexa. Vários níveis e subprocessos são simultaneamente ativados durante os processos de compreensão e produção de enunciados. Desde processos elementares - análise da onda sonora, identificação de formas simbólicas ou estímulo de certos músculos, até processos mais complexos - análise e síntese de propriedades lingüísticas e conceituais. Do ponto de vista da teoria lingüística, seu estudo busca regras e princípios subjacentes ao processamento lingüístico; do ponto de vista da psicolingüística, o centro das investigações são os processos mentais envolvidos nesse processamento; do ponto de vista do processamento automático das línguas naturais (PLN), coloca-se a questão de simular esses processos, investigando se há fundamentação psicológica para as hipóteses levantadas, ou as possibilidades de se enriquecer os programas com “conhecimento lingüístico” e executá-los para realizar tarefas como tradução automática, análise gramatical, sumarização automática de textos, entre outras (Allen, 1987).

Nesse contexto, o léxico desempenha um papel essencial. Nele estão armazenados todos os lexemas que entram no “jogo da sintaxe”, associados a um complexo de informações de naturezas diversas: fonológicas, morfológicas, sintáticas, semânticas e, até mesmo, pragmáticas (Andrews, 1989). Em particular, os programas de PLN necessariamente pressupõem uma “base de dados” lexicais, cuja organização, mesmo motivada por princípios lingüísticos e psicolingüísticos, precisa conformar-se a uma série de restrições: a rígida estrutura dos programas que a codifica, o modo mecânico de gerenciamento de bases de dados e a busca de eficiência de processamento, que se traduz como a rapidez com que a máquina é capaz de executar tarefas. Visto dessa forma, o léxico computacional ideal é concebido como uma base de dados gigantesca, contendo milhares de entradas, cujo modo de acesso é imediato e rápido, e incorporando princípios de armazenamento eficientes e técnicas de manipulação de dados altamente flexíveis.

Para atingir esse ideal, são inúmeras as questões que se colocam: descoberta dos princípios de organização do léxico, especificação do conteúdo e da estrutura das entradas e das



informações a elas associadas e a proposição de formalismos suficientemente expressivos para representar essas especificações.

Este trabalho, que dá continuidade à tarefa que vimos desenvolvendo nos últimos anos (Dias-da-Silva, 1998a; 1998b), tecer pontes de contato entre lingüistas e pesquisadores do PLN, investiga uma pequena parte dessas questões, relacionada à organização do léxico: a representação lingüístico-computacional das principais relações semânticas responsáveis pela estruturação semântica do léxico. Dentre elas, destacam-se as relações de natureza conceptual – hiponímia, meronímia, holonímia, e as relações de natureza léxico-semântica – sinonímia e antonímia (Cruse, 1986). Após revisar as principais propostas teórico-metodológicas apontadas na literatura – análise componencial, campos semânticos, redes semânticas e redes relacionais – apresenta-se um modelo lingüístico-computacional que possibilita a montagem de uma base relacional de dados lexicais para o português brasileiro, cuja configuração espelha a organização do léxico em função dessas relações.

Em (1), apresenta-se uma breve revisão dos principais modelos de representação do significado lexical: *teoria de traços*, *redes semânticas*, *teoria de protótipos*, *teoria de frames* e *teoria da dependência conceptual* e, em (2), delinea-se o modelo lingüístico-computacional de representação, baseado em redes relacionais. Conclui-se, em (3), com o esboço de um protótipo de base relacional de dados lexicais.

1. Modelos de representação do significado lexical

Dentre as teorias que focalizam a descrição semântica do léxico, Handke (1995: 91) destaca as que mais têm contribuído para o estudo do PLN: *teoria de traços*, *redes semânticas*, *teoria de protótipos*, *teoria de frames* e *teoria da dependência conceptual*.

O ponto de partida para a descrição do significado lexical é a teoria de traços, proposta por Katz e Fodor (1963) e Katz e Postal (1964). Essa teoria propõe técnicas que consistem na “quebra” dos conceitos expressos pelas unidades lexicais em partes menores. Várias denominações são aplicadas a essas subpartes do significado: traços lexicais ou semânticos, átomos de significado, primitivos semânticos ou componentes semânticos. Esses componentes não são lexemas, mas elementos abstratos de uma meta-linguagem de descrição do significado lexical: [mosca: (ANIMAL), (INSETO), (ASAS), (PERNAS), (FERRÃO)]. Para fins de representação, pode-se interpretar esses traços como atributos, e para cada atributo associar valores específicos, que podem ser binários, numéricos ou descritivos: [mosca: (ANIMAL +), (INSETO +), (ASAS:2), (PERNAS:6), (FERRÃO -)]; [menino: ... (NOME: JOÃO) ...]. Essa teoria possibilita definir formalmente as relações de sentido. Por exemplo, um lexema A é um hipônimo de B se todos os traços de B estão contidos na especificação de traços de A. A antonímia, considerada uma relação entre conjuntos de lexemas, define-se simultaneamente pelo compartilhamento e contraste de traços.

Um outro modelo de representação, que tem por base a psicologia e a ciência da computação, trazido por Quilian (1968), é conhecido como *redes semânticas*. Numa rede semântica simples, os conceitos (objetos e eventos) são representados por nós e as inter-relações entre conceitos, por arcos. As redes semânticas são populares pela forma elegante de representação do significado e pela facilidade de se inferir deduções. Por exemplo, para deduzir que uma mosca possui cabeça, é só traçar a hierarquia É-UM (*IS-A*), assumindo que os atributos associados aos nós mais altos também são válidos para nós situados mais abaixo. É o que se chama de propriedade de herança, que permite inferir que o conceito MOSCA herda os atributos de INSETO, e que ambos INSETO e MOSCA herdam os atributos de ANIMAL. Por serem enriquecidas com o



procedimento de herança, oferecem uma técnica precisa de implementação de hierarquias de conceitos e dos diferentes tipos de relações que se estabelecem entre eles, evitando redundâncias, posto que os traços semânticos são especificados uma única vez para os conceitos hierarquicamente superiores, e permitindo o estabelecimento de uma estrutura de conceitos para ancorar o significado lexical.

Na *teoria de protótipos*, a relação de pertença (*membership*) está centrada na representação do membro prototípico da classe que um certo item denota (Rosch e Lloyd, 1978). Experiências baseadas nessa abordagem trouxeram dois tipos de resultados: prototípicos e básicos. Os resultados prototípicos demonstram que PARDAL, por exemplo, é mais representativo de PÁSSAROS do que AVESTRUZ ou PINGÜIM. Assim, como PARDAL é o membro mais representativo da categoria AVE, ele é o protótipo dessa categoria. Observe-se que o protótipo não é o significado de um lexema. Ele está localizado num espaço multidimensional, com dimensões correspondentes aos atributos, que podem ser diferentes através de uma categoria. Os resultados básicos demonstram que certas categorias são mais básicas que outras: porque são reconhecidas mais rapidamente, aprendidas mais cedo, usadas mais frequentemente ou processadas mais facilmente que outras categorias. Por exemplo, se MOBÍLIA é o conceito superordenado, CADEIRA é o nível básico e CADEIRA DE BALANÇO, uma categoria subordinada. A idéia por trás da postulação de categorias básicas é que elas: servem de base para a interação entre membros de uma categoria, servem como âncoras perceptuais e proporcionam imagens mentais que refletem a categoria toda.

Originalmente proposta por Minsky (1975), como uma base para a percepção visual, a *teoria de frames* visa ao equacionamento da compreensão de diálogos das línguas naturais e outras atividades cognitivas igualmente complexas. A idéia principal de *frame* é a integração de conceitos novos a conceitos adquiridos por experiência prévia. Em termos formais, um *frame* é uma estrutura de atributo-valor. O *frame* sempre contém um atributo GRUPO, que especifica a qual classe geral um conceito pertence. O atributo “Especialização de” estabelece o conceito superordenado do qual o *frame* é hipônimo. Ambos os atributos são essenciais para estabelecer propriedades de herança entre *frames*. Considerados uma variação notacional das redes semânticas, os *frames* são também representações suficientemente ricas para especificar relações de sentido como a hiponímia, meronímia e sinonímia.

Por fim, a combinação de redes semânticas e *frames*, deu origem à teoria de dependência conceitual, que visa representar conceitos predicativos e seus argumentos (Schank e Riesbeck, 1981). Essa teoria pressupõe que todo evento comporta os seguintes elementos: um ATOR, uma AÇÃO desencadeada pelo ATOR, um OBJETO como objetivo da AÇÃO, uma DIREÇÃO como orientação básica para a AÇÃO. Como consequência, estabelece-se um inventário de ações primitivas como, por exemplo, PTRANS (transferência física de objetos, por exemplo, *ir*), ou MOVE (o ato de mover uma parte [do corpo]). No caso da frase “John foi para Londres”, o primitivo da ação é o PTRANS, ou conceito de transferência física, e teríamos a representação (PTRANS (ATOR X) (OBJETO X) (DIREÇÃO (DE Y) (PARA Z))), que codifica o seguinte: um ATOR X transferiu fisicamente um OBJETO X do local Y para o local Z. Note que ATOR e OBJETO (ambos X, ligado a John) são idênticos, indicando que o ATOR se transfere fisicamente. A variável Z liga-se a Londres e Y é uma variável livre, porque não se sabe a origem da ação. Observe-se que, com esse mecanismo, é possível traduzir um conceito primitivo em uma estrutura de dependências conceituais e armazená-la no léxico. Permite ainda que se estabeleçam relações como hiponímia e antonímia.

Todos esses modelos, entretanto, apresentam problemas resistentes a soluções transparentes: a delimitação e explicitação de valores semânticos, sua seleção e quantificação. Note-se que, apesar das denominações distintas - “traços”, “atributos”, “valores”, “protótipo” ou



“elementos primitivos”, todos os modelos apresentados pressupõem a existência de “átomos conceituais” para proceder à decomposição do significado lexical.

No sentido de evitar o recurso a primitivos semânticos, e com isso evitar o enfrentamento dos problemas apontados, Miller e Fellbaum (1991) propõem um modelo alternativo de representação de conceitos lexicalizados em uma língua e das relações semânticas que se estabelecem entre as lexicalizações (lexemas) desses conceitos. Suas características principais são: (i) a adoção do “método diferencial”, que pressupõe o princípio de ativação de conceitos por meio de um conjunto de formas lexicais relacionadas pela relação de sinonímia, eliminando a necessidade de se especificar o valor semântico associado às entradas do léxico; (ii) a noção constitutiva básica de *synset*, isto é, um conjunto de sinônimos; (iii) a noção de “matriz lexical”, que especifica uma correspondência biunívoca entre sentido e *synset*. A matriz lexical pode ser visualizada por meio de um plano cartesiano: no eixo das abscissas representam-se as formas lexicais (F1,...,Fn) e no eixo das ordenadas, os significados (S1,...,Sm). Nesse plano, um ponto de coordenadas (F1,S1) representa o fato de que a forma lexical F1 é usada para expressar o significado S1.

2. Modelo de representação computacional: a rede *Wordnet*

Em termos formais, o sistema *Wordnet 1.6* pode ser entendido como uma base relacional de lexemas.¹ O construto básico dessa base, responsável pela organização do léxico, é o *synset*, um conjunto de sinônimos ou quase sinônimos. Por exemplo, para um falante do inglês que saiba que *board* pode ser *tábua velha* e *grupo de pessoas reunidas com algum propósito* (entre outros), o conjunto de sinônimos {*board, plank*} e {*board, committee*} servem como designadores não ambíguos dos dois significados de *board* (Miller e Fellbaum, op. cit.:201). Embora esses conjuntos não “expliquem” o conceito, servem de índice para sinalizar a existência de um conceito lexicalizado. As exigências de uma teoria diferencial são bastante modestas, porém suficientes para a construção de uma base relacional de dados lexicais, que refletem as intuições do falante de uma língua no que diz respeito às relações de sentido.

Do ponto de vista lógico, os *synsets* são conjuntos munidos de dois tipos de ponteiros que representam dois tipos de relação (R) entre os conjuntos: ponteiros que especificam relações léxico-semânticas e ponteiros que especificam relações conceptuais. Os ponteiros do primeiro tipo especificam relações entre as formas (vocábulos e expressões) e os do segundo tipo entre os conceitos atualizados pelas formas. As relações de sinonímia e antonímia enquadram-se no primeiro tipo. Entretanto, como a relação de sinonímia é a relação constitutiva básica dos *synsets*, ela não é especificada por meio de ponteiros, mas pela relação de pertença. As demais relações incluem-se no segundo tipo e especificam relações entre *synsets*. Dentre elas, estão, por exemplo, as relações de hiperonímia, hiponímia, holonímia, meronímia, implicação.

3. Protótipo de uma base relacional de dados lexicais

Do ponto de vista da implementação, o sistema *Wordnet* poderia ser composto de arquivos preparados por lexicógrafos (ALs), um programa que converte esses arquivos em uma base de dados (DB), rotinas de busca e interfaces para a apresentação da informação a partir da base de dados. Nos ALs substantivos, verbos, adjetivos e

¹ Disponível *on line* no endereço <http://www.cogsci.princeton.edu/~wn/>.



advérbios seriam organizados em conjuntos de sinônimos; as demais relações acima nomeadas seriam também especificadas nesses arquivos. Prevê-se também um programa que converte os ALs na DB, responsável pela codificação dessas relações. As diferentes interfaces de acesso à DB utilizariam uma biblioteca comum de rotinas, criadas para exibir os diversos tipos de relação.

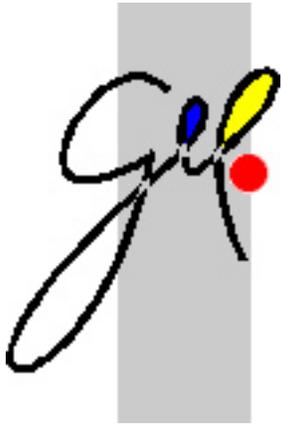
Os *synsets* formados por substantivos e verbos seriam organizados em hierarquias baseadas nas relações de hiperonímia/hiponímia. Já os *synsets* formados por adjetivos seriam estruturados de modo diferente. Conforme previsto na rede *Wordnet*, seriam construídas configurações em forma de constelações de *synsets* nucleares e *synsets* satélites. Cada constelação teria como âncora pares de antônimos, especificados nos *synsets* nucleares. Em torno dos *synsets* nucleares teríamos um ou mais *synsets* satélites, cada um dos quais representando um conceito semelhante ao conceito representado pelo *synset* nuclear. A relação de antonímia seria metaforicamente comparada a um eixo que liga duas rodas raiadas, em cujos centros estariam localizados os *synsets* nucleares (definidos como contendo pelo menos um lexema associado a seu antônimo), e na extremidade de cada raio, os *synsets* satélites. Nos adjetivos deverbais, formados a partir do particípio passado de verbos, seriam incluídos ponteiros, apontando para os verbos dos quais são derivados. Por fim, como os advérbios são em sua grande maioria derivados de adjetivos, os *synsets* formados por advérbios também conteriam ponteiros, apontando para os adjetivos dos quais foram derivados, herdando destes a relação de antonímia correspondente.

RESUMO: *Este artigo examina a questão da representação lingüístico-computacional das principais relações responsáveis pela estruturação semântica do léxico: relações de natureza conceptual e léxico-semântica. Após uma breve revisão das principais propostas teórico-metodológicas apontadas na literatura, apresenta-se um modelo de representação computacionalmente tratável que espelha essas relações.*

PALAVRAS-CHAVE: *semântica lexical; relações de sentido; estrutura do léxico; processamento automático das línguas naturais.*

REFERÊNCIAS BIBLIOGRÁFICAS

- ALLEN, J.F. *Natural language understanding*. Menlo Park: Benjamin Cummings, 1987.
- ANDREWS, A.D. Lexical Structure. In: F. Newmeyer (ed.) *Linguistics: the Cambridge survey I*. Cambridge: Cambridge University Press, p. 60-88, 1989.
- CRUSE, D. A. *Lexical semantics*. Cambridge: Cambridge University Press, 1986.
- HANDKE, J. *The structure of the lexicon: human versus machine*. Berlin-New York: Mouton de Gruyter, 1995.
- MILLER, G.A. & FELLBAUM, C. Semantic networks of English. In: *Cognition*, Amsterdam, v. 41, n. 1 e 2, p. 197-229, 1991.
- DIAS-DA-SILVA, B. C. Bridging the gap between linguistic theory and natural language processing. In: Bernard CARON (éd.) *Actes du 16^e Congrès International des Linguistes*. Oxford: Elsevier Sciences, Paper 0425, 1998a.
- _____. Os domínios lingüístico e tecnológico do estudo do processamento automático das línguas naturais. *Estudos Lingüísticos*, Jaú, v. 26, p. 612-617, 1998b.
- KATZ, G.G. e FODOR, J.A. The structure of a semantic theory. In: *Language*, Pittsburgh, v. 39, p. 170-210, 1963.



- KATZ, G.G. e POSTAL, P.M. *An integrated theory of linguistic descriptions*. Cambridge: MIT Press, 1964.
- MINSKY, M. A framework for representing knowledge. In: J. Haugeland (ed.) *Mind design*. Cambridge: MIT Press, p. 95-128, 1975.
- QUILLIAN, M.R. Semantic memory. In: M. Minsky. *Semantic information processing*. Cambridge: MIT Press, p. 227-70, 1968.
- ROSCH, E. e LLOYD, B.B. (eds.) *Categorization and cognition*. Hillsdale: Lawrence Erlbaum Ass., 1978.
- SCHANK, R.C. e C.K. RIESBECK (eds.). *Inside computer understanding*. Hillsdale: Lawrence Erlbaum, 1981.