

Comparador: forma de auscultar corpos no AC/DC

Alberto Simões

19 de Julho de 2014

Resumo

Com a quantidade de diferentes corpos e a grande diversidade entre o que está anotado e marcado em cada um, é necessário desenvolver ferramentas que possam permitir ao linguista comparar fenómenos de várias perspetivas diferentes.

O Comparador, já mencionado em vários textos sobre o AC/DC ou a Linguateca em geral [2, 3], mas nunca ainda documentado independentemente, é uma destas ferramentas, que agora com a Gramateca urge tornar disponível para uma comunidade maior de interessados.

Os requisitos do Comparador foram discutidos e acordados entre Diana Santos e Alberto Simões, e a documentação foi feita em conjunto, mas o desenvolvimento é apenas da responsabilidade de Alberto Simões.

1 Objetivo

Embora os dados e os resultados existam acessíveis para todos os que usam o AC/DC na sua pesquisa, a quantidade de informação associada aos corpos e a diversidade do material podem tornar as pesquisas complicadas para um linguista desprevenido.

Poder comparar diferentes pesquisas com uma única interface, e obter os resultados diretamente num único lugar, torna muito mais fácil o reaproveitamento do material e a própria interação com o AC/DC, da mesma forma que o Ensinador¹ [4] tornou possível usar o mesmo material num contexto didático de sala de aula.

Outras ferramentas existem ou foram desenvolvidas para reuso do material do AC/DC para outros fins ou por outros tipos de utilizadores, como é o caso do Ordenador, do Distribuidor, do Castor, do Ensinador e do VARRA.

A inspiração original do Comparador, para que conste, foi um programa que comparava a frequência na rede (Internet) no final dos anos 90, usado para ensinar inglês para estrangeiros. (NOTA: estamos à procura da referência exata!)

¹<http://www.linguateca.pt/Ensinador/>

2 Como utilizar

O Comparador, como todas as ferramentas associadas ao AC/DC, resume-se, do ponto de vista do utilizador, a uma interface na Rede em que se escolhe um conjunto de parâmetros e se obtém o resultado em HTML.

Do ponto de vista de interface, as duas pesquisas a comparar são especificadas em paralelo, e o resultado aparece em paralelo no écran, ou numa tabela única.

Seguem-se alguns exemplos com a sua motivação linguística. Para exemplos práticos, ver a página de Ajuda do próprio Comparador².

2.1 A mesma pesquisa em dois corpos diferentes

Uma das questões mais básicas da linguística com corpos é a da importância do corpo para uma determinada conclusão: será que ter feito uma pesquisa num dado corpo é suficiente, ou se tivesse usado outro não poderia ter chegado à mesma conclusão?

Sendo o AC/DC um serviço que disponibiliza o acesso a todos os corpos cujos donos nos dêem autorização (sem impedir ou prejudicar que os mesmos corpos sejam interrogados e aprimorado noutros locais), o AC/DC é o local privilegiado para a comparação de corpos (ou melhor de pesquisas em corpos).

2.2 A mesma pesquisa em duas variantes diferentes

Outra das áreas que tem importância fulcral para a política linguística da língua portuguesa, é conhecer a diversidade da nossa língua, e em particular saber se as diferenças entre variantes são absolutas (categóricas) ou uma questão de grau e de variabilidade intrínseca a todas as línguas.

Para isso nada como poder comparar todas as análises em variantes diferentes.

2.3 A mesma pesquisa em dois géneros diferentes

Outra questão essencial na linguística com corpos é a de saber quais as generalizações a fazer para a língua toda, a língua geral, ou apenas para géneros específicos. Para isso é preciso que se possa recuperar o género dos dados a que recorreremos, e que se possa pedir distribuição por género e tipo de texto, pese embora a grande dificuldade de concordância na delimitação dos mesmos.

2.4 A mesma pesquisa com dois analisadores diferentes

É preciso também reconhecer que a análise subjacente à maior parte do material presente no AC/DC tem dois vieses: a teoria linguística implícita no PALAVRAS[1], e o desempenho (prático) do mesmo analisador.

²<http://www.linguateca.pt/Comparador/>, escolha Ajuda

Embora estejamos gratos pela simples existência do PALAVRAS e pela possibilidade de desenvolvimento conjunto que muitos projetos com a Linguatca proporcionaram, nada obsta a que a totalidade do mesmo material seja também analisado por outros sistemas, e isso já acontece nos seguintes corpos: o FrasesPP e o AmostRA, e estamos de momento a trabalhar na mesma vertente com o ReLI e o COLONIA.

2.5 Comparando duas cores no mesmo corpo

Obviamente que a comparação não se esgota com a variação entre parâmetros externos, mas que até é provavelmente mais interessante do ponto de vista linguístico. Embora partes dessa comparação já se possam fazer através da interface padrão do AC/DC, é certamente interessante poder comparar em paralelo, por exemplo, diferentes distribuições.

Quais as palavras mais correntes associadas com a cor verde, e com a cor azul? Uma forma muito simples (que tem de ser naturalmente refinada) é a descrita na figura seguinte:

[Ajuda](#)

Comparador

Procurar: [pos="N.*"] [lema="verde"]

Corpo: VERCIAL

Distribuir por: lema

Procurar: [pos="N.*"] [lema="azul"]

Corpo: VERCIAL

Distribuir por: lema

olho	66	olho	73
vinho	34	céu	59
veludo	24	túnica	19
seda	21	seda	19
ramo	21	tinta	19
erva	15	olhar	18
pano	15	óculo	16
repe	14	vestido	16
folha	14	pano	15
palma	13	cetim	14
campo	13	veludo	14
mar	12	água	14
cana	11	repe	13
baeta	10	sangue	12
portinha	9	mar	12
vara	9	casaca	11
monte	9	espaço	11
persiana	8	manto	10
caldo	8	sobrecasaca	8
pó	8	luneta	8
água	8	poeira	8
ano	8	papel	8
batente	7	fita	8
cetim	7	firmamento	8
óculo	7	cor	8
cor	7	algodão	7
tabuinha	6	campo	7

2.6 Comparando duas emoções no mesmo género

Da mesma forma, pode ser interessante comparar quais os fenómenos que provocam diferentes emoções.

Na próxima figura temos uma simples comparação entre aquilo que inspira amor ou ódio em texto jornalístico.

[Ajuda](#)

Comparador

Procurar: [sema="emo:amor" & pos="V.*"] [pos!="[VN].*"]*

Corpo: CHAVE

Distribuir por: word

Procurar: [sema="emo:odio" & pos="V.*"] [pos!="[VN].*"]*

Corpo: CHAVE

Distribuir por: word

música	62	mulheres	14
trabalho	45	peçoas	11
futebol	35	idéia	10
país	34	vida	10
vida	32	publicidade	9
caso	32	ideia	8
mulher	31	música	7
coisas	27	homens	7
cinema	23	tipo	7
mulheres	23	política	7
projecto	20	piscinas	6
obras	20	filme	6
idéia	19	filmes	6
paisagem	18	palavra	6
crianças	18	mundo	6
proposta	18	reprise	5
situação	18	imagem	5
homem	17	judeus	5
liberdade	16	futebol	5
filmes	16	sistema	5
propostas	16	gente	5

Agradecimentos A Linguateca agradece à Universidade de Oslo apoio financeiro para o desenvolvimento desta e outras ferramentas associadas à Gramateca.

Referências

- [1] Eckhard Bick. *The Parsing System "Palavras": Automatic Grammatical Analysis of Portuguese in a Constraint Grammar Framework*. PhD thesis, Aarhus University, Aarhus, Denmark, November 2000.
- [2] Diana Santos. Linguateca's infrastructure for Portuguese and how it allows the detailed study of language varieties. In J.B. Johannessen, editor, *Language Variation Infrastructure: Papers on selected projects*, volume 3, pages 113–128, 2011.
- [3] Diana Santos. Corpora at Linguateca: Vision and Roads Taken. In Tony Berber Sardinha and Telma São Bento Ferreira, editors, *Working with Portuguese corpora*, pages 219–236. Bloomsbury, 2014.
- [4] Alberto Simões and Diana Santos. Ensinador: corpus-based Portuguese grammar exercises. *Procesamiento del Lenguaje Natural*, 47:301–309, Setembro 2011.