



Leitura Distante: o que podemos com corpora anotados?

Cláudia Freitas

Linguateca



Preâmbulo

Leitura Distante (Moretti, 2003)

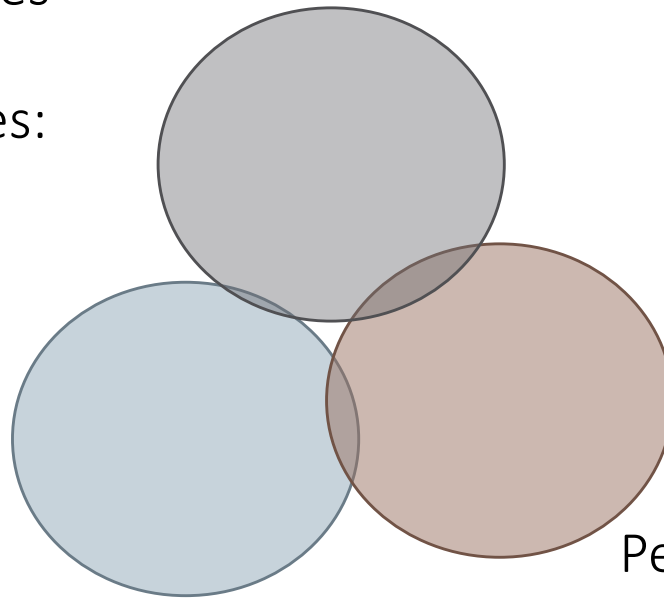
A distância como aliada → padrões

Dados que dão forma aos padrões:
metadados (em geral)
Estão *fora* do texto

Leitura não-linear (Freitas 2017)

Padrões, mas também
singularidades

Rearranjos tornados possíveis
apenas com o texto
descorporificado (Paixão de
Souza, 2013)



Corpus, anotação e ferramentas (Santos, 2009)

Um corpus é um trio: o texto, a
anotação, a ferramenta de busca

Permite alternar lentes entre a leitura
distante e a leitura aproximada

Estratégias e abordagens



Analisar dados

Atribuir sentido aos
dados
(para analisá-los)

Um exemplo - gênero de personagens como operador analítico

- Como são caracterizadas personagens masculinas e femininas em obras literárias brasileiras dos séculos XIX e XX ?

Corpus OBras

248 obras da
literatura brasileira –
v 5.3

32 autores séc
XIX-XX

3 obras de autoria
feminina

Obras escritas por quem e
para quem?
Censo de 1872: 16% de
população brasileira era
alfabetizada....

Álvares de Azevedo
Manuel Antonio de Almeida
Tomás Antônio Gonzaga
Visconde de Taunay Joaquim Nabuco
Cláudio Manoel da Costa Gregório de Matos
Maria Firmina dos Reis Inglês de Sousa Mário de Andrade
Franklin Távora Lima Barreto
Júlia Lopes de Almeida Castro Alves Aluísio Azevedo
Paulo Setúbal Alvarenga Peixoto João do Rio
Machado de Assis
Júlio Ribeiro Manuel de Oliveira Paiva
Raul Pompéia Adolfo Caminha
Bernardo Guimarães Olavo Bilac Euclides da Cunha
Joaquim Manuel de Macedo
Humberto de Campos
Basílio da Gama José de Alencar
Coelho Neto
Artur Azevedo

Um corpus é um trio....

Projeto **AC/DC**: **corpo OBras**

[AC/DC : Linguateca](#)

O corpo **OBras** (Obras Brasileiras) é um corpo de textos brasileiros que já alcançaram o domínio público, criado numa colaboração entre a Linguateca, a Universidade de Oslo, a PUC-Rio, a Universidade Estadual do Maranhão (UEMA) e Anya Campos. Para mais informações veja-se a [página do projeto](#).

Procurar:

Resultado:

- Concordância
- Distribuição das formas (*word*)
- Distribuição dos lemas ([lema](#))
- Distribuição da categoria gramatical (PoS) ([pos](#))
- Distribuição do tempo verbal e/ou do caso pronominal ([temcagr](#))
- Distribuição de pessoa e/ou número ([pessnum](#))
- Distribuição do género morfológico ([gen](#))
- Distribuição da função sintáctica ([func](#))
- Distribuição pelas obras (*obra*)
- Distribuição por autores (*autor*)
- Distribuição por género de texto (*classe*)
- Distribuição pela corrente literária (*escola*)
- Distribuição pelo sexo do entrevistado, do biografado ou do autor (*sexo*)
- Distribuição por original ou traduzido (*trad*)
- Distribuição por campo semântico (*sema*)
- Distribuição por grupo (de cor, roupa, etc.) (*grupo*)
- Distribuição das dependências (*dependencias*)

Tipo	Literário
Variante(s)	BR
Tamanho (unidades)	8.9 milhões
Tamanho (palavras)	6.3 milhões

Carateres úteis: | { } []

[Página principal](#)

Procure noutros corpos:

[AmostRA-NILC](#) [ANCIB](#) [Avante!](#)
[Corpus Brasileiro](#) [CD HAREM](#)
[CETEMPúblico](#) [CHAVE](#) [Ciência Viva](#)
[Colonia CONDIVport](#) [CONDIVport2](#)
[CoNE](#) [C-Oral-Brasil](#) [DHBB](#) [DiaCLAV](#)
[Diáspora TL-PT](#) [ECI-EBR](#) [ECI-EE](#)
[ENPCPUB](#) ([parte em português](#))
[Floresta](#) [FrasesPB](#) [FrasesPP](#) [Mariano](#)
[Gago](#) [Literateca](#) [Marielle, presente!](#)
[Moçambula](#) [Museu da Pessoa](#)
[Natura/Minho](#) [NOBRE](#) [OBras](#) [P'lo](#)

Metodologia

1. Mineração de textos: Busca por estruturas predicadoras

NOMES COMUNS

NOMES PRÓPRIOS

PRONOMES PESSOAIS

+

função

APOSTO

ADJUNTO ADNOMINAL

PREDICATIVO DO SUJEITO

Concretamente....

id="Os_irmãos_Leme Prosa:romance PS 1933 histórico masc ":

Rodrigo César, iradíssimo, deliberou invadir as atribuições do Ouvidor.

id="O_Coruja Prosa:romance AA 1889 realismo masc ":

é feio, mas enfim, sempre há **homens sérios**, cujo nome o público não ignora;

id="O_Cabeleira Prosa:romance FT 1876 realismo_regionalismo_romantismo masc ":

Mas **José César era ativo**, enérgico, esforçado e de grandes espíritos.

id="Uma_lágrima_de_mulher prosa:romance AA 1880":

E **Rosalina, meiga**, encarava com chorosa ternura o olhar sombrio de Miguel.

id="A_serpente_de_bronze prosa:conto HC 1921":

Criemos as meninas com decoro, vestindo-as com discrição, e teremos **moças discretas, pudicas, decorosas, ciosas** do seu corpo e dos seus encantos.

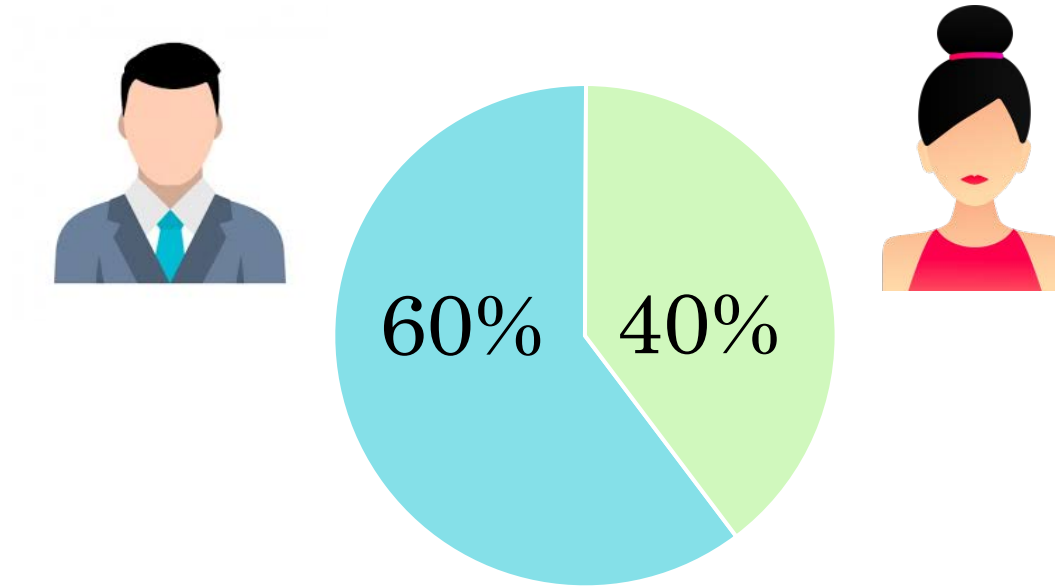
id="O_seminarista prosa:romance BG 1872":

A **mãe** de Eugênio era **fanática** e **supersticiosa**.

Mais concretamente...

```
[pos="PROP.*" & func="SUBJ>"] [lema="ser|estar"] [pos="ADV.*"]* @[temcagr!=".*PASS.*" & pos="ADJ|N|V" & gen="M" & func="<SC"]
[pos="PROP.*"] "era" [pos="ADV.*"]* [pos="ADJ.*" & gen="M" & func!=">N"] [word="e"] @[pos="ADJ.*" & gen="M" & func!=">N"]
[lema="ele" & func="SUBJ>"] [lema="ser|estar"] [pos="ADV.*"]* @[temcagr!=".*PASS.*" & pos="ADJ|N|V" & gen="M" & func="<SC"]
[pos="PROP.*" & func!="P<"] ", " [pos="ADV.*"]* @[func="N<PRED|.APP.*" & gen="M" & pos="ADJ"]
[pos="PROP.*" & func!="P<"] ", " [pos="ADV.*"]* [func="N<PRED|.APP.*" & gen="M" & pos="ADJ"] "e" @[gen="M" & pos="ADJ"]
[lema="homem|moço|rapaz|marido"] @[pos="N|ADJ|V" & func="<PRED|<OC|N<"]
[lema="homem|moço|rapaz|marido"] [pos="N|ADJ|V" & func="<PRED|<OC|N<"] ", " @[pos="ADJ" & gen="M"]
[lema="homem|moço|rapaz|marido"] [pos="N|ADJ|V" & func="<PRED|<OC|N<"] "e" @[pos="ADJ" & gen="M"]
```

Primeiros resultados: o que vemos?



3.862 Predicações



O que vemos?



sentada perdida
doente linda alta feia
solteira casada boa **bela**
pobre idosa velha livre moça
capaz rica
bonita
filha amada loura pálida feliz
gorda alegre viúva honesta
formosa
encantadora

<500 predicacões diferentes



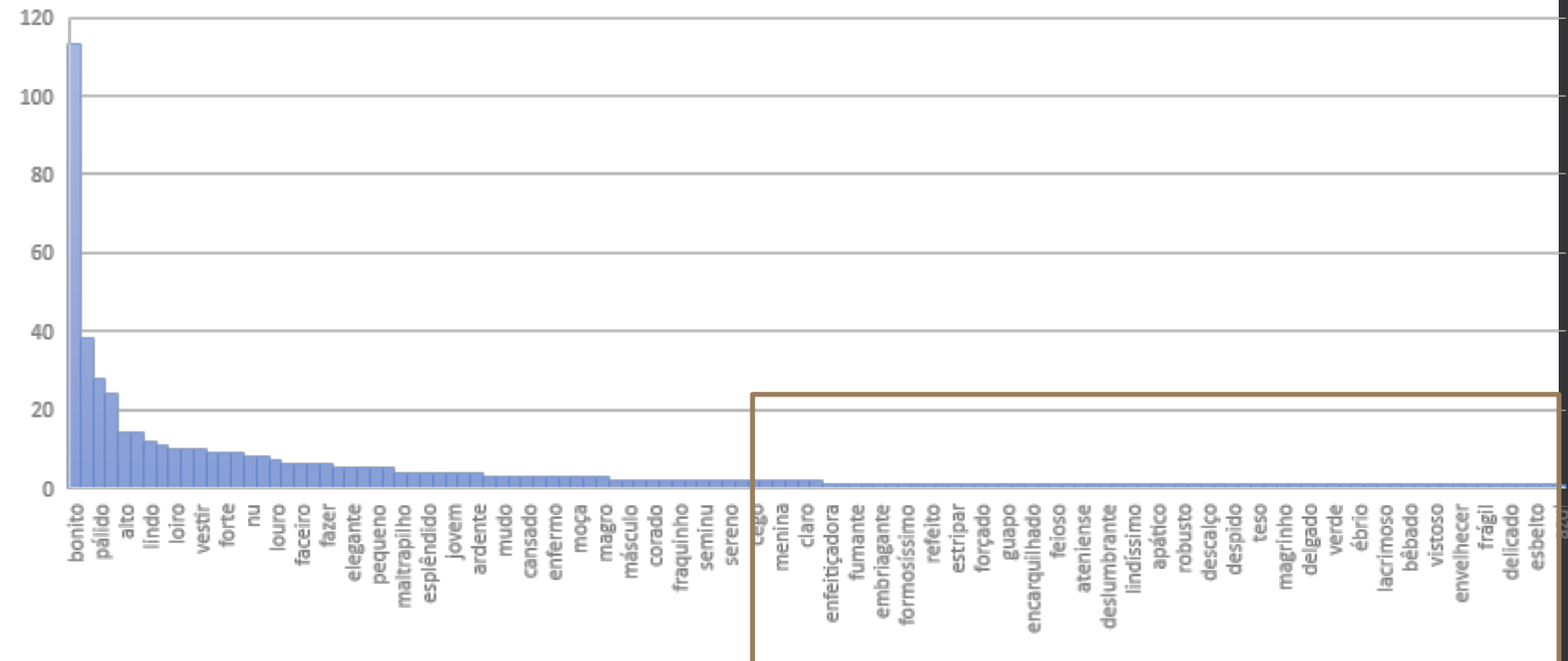
alto inteligente
feliz gordo
pobre solteiro capaz
moço bom
cândido mau
doente digno
forte sério honrado
robusto doido feito magro
pálido público
honesto alegre grave rico
velho baixo
bonito branco

< 800 predicacões diferentes



Primeiros resultados

Cerca de 60% dos predicadores masculinos e femininos aparece uma única vez



Como generalizar?

Metodologia

2. Categorização: Criação de eixos para *dar sentido* à cauda longa.

Análise manual

APARÊNCIA /
CORPO

CARÁTER

EMOÇÃO / AFETO

LUGAR SOCIAL

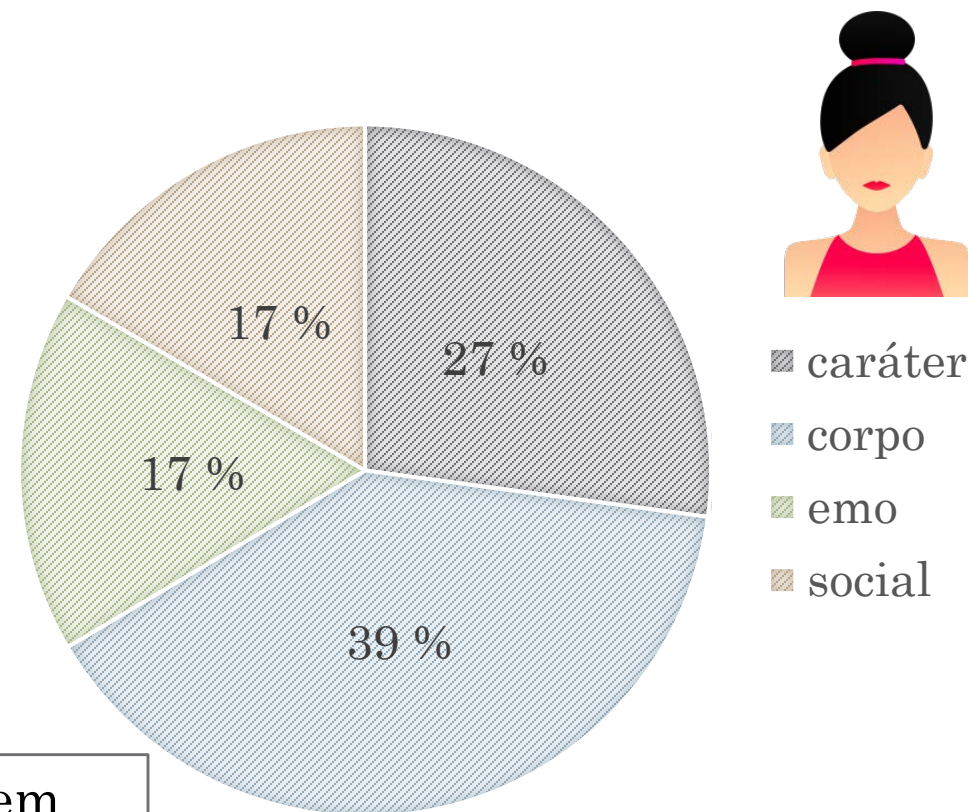
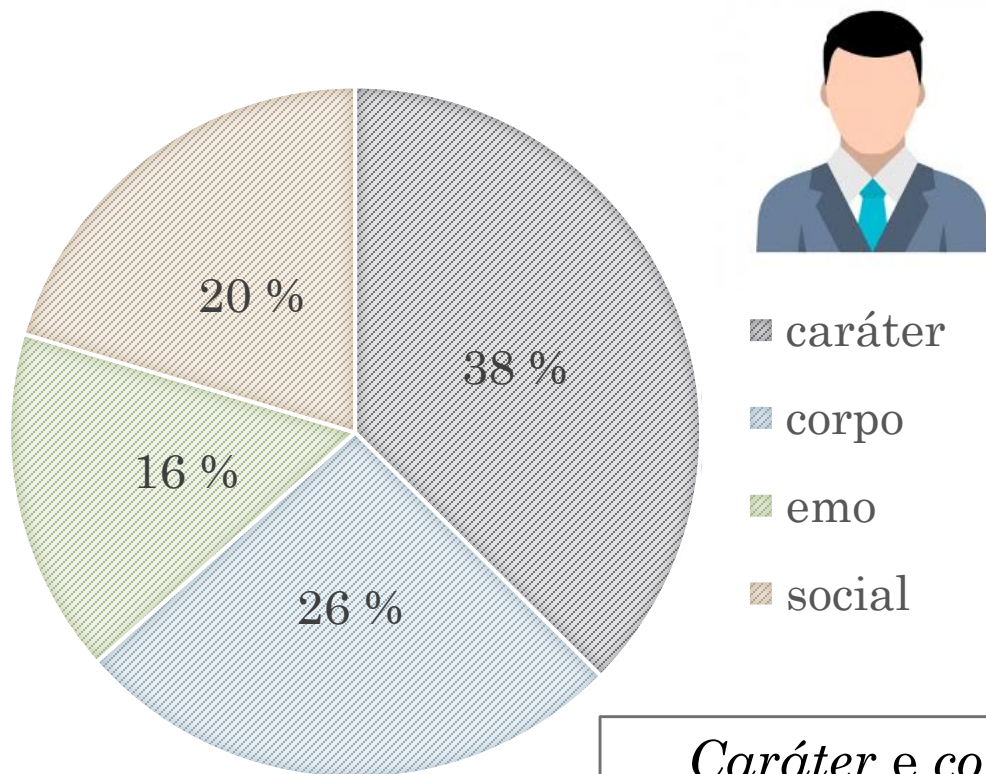
Classificação múltipla

- “Nem os moços **valentes**_[emo & caráter], nem os senhores respeitáveis, nem os jornalistas vão sequer à delegacia.” (*A Alma Encantadora das Ruas*)
- “Eugênio, **trêmulo**_[físico & emo], confuso e de olhos no chão deixou cair sobre sua cabeça toda esta tremenda trovoadas.” (*O Seminarista*)

Metodologia

predicador	qtd	eixo
bonito	146	corpo
belo	66	corpo
amado	43	emo
casado	43	social
formoso	38	corpo
pálido	37	corpo
honesto	36	caráter
velho	32	corpo
alto	30	corpo
doente	25	corpo
rico	24	social
bom	24	caráter
idoso	22	corpo
capaz	22	caráter
solteiro	21	social
moreno	21	corpo
moço	21	corpo
pobre	21	social
sentado	20	erro
feliz	19	emo
lindo	18	corpo
encantador	17	corpo

O que vemos?



Caráter e corpo tem distribuição invertida

Emo e social têm distribuição parecida





O que vemos?



	Predicações masculinas	Predicações femininas
Emoção/afeto	Feliz, amigo, contente, alegre, valente*, atônito, amado, triste, espantado, pensativo, comovido, apaixonado, respeitado*, radiante, ansioso, trêmulo, entusiasmado*, aborrecido, indignado, orgulhoso, furioso*	Amada, feliz, alegre, amante*, triste, faceira*, trêmula, impaciente, amiga, assustada, agitada*, chorosa*, radiante, espantada, tranquila, nervosa, admirada, pensativa, risonha, coitada, ardente*
Caráter	Bom, sério, grande, honesto, digno, capaz, inteligente, mau, honrado*, doido, grave*, valente*, generoso, frio, vulgar, direito, calmo, ilustre*, livre, simples, hábil	Honesta, boa, perdida*, capaz, livre, vulgar, de verdade*, digna, impaciente*, fácil, indiferente, discreta, fresca, infame, inteligente, ingênu*, inocente, fria, pura, santa*, namoradeira*
Aparência/corpo	Alto, magro, bonito, pálido, gordo, forte, velho, (homem)feito, moço, branco, doente, baixo, elegante, robusto, vestido, morto, fraco, novo, maduro*, moreno	Bonita, bela, pálida, formosa*, alta, feia, linda*, velha, loira, morena, vestida, encantadora*, forte, moça, nua, doente, loura, gorda, faceira*, trêmula, (mulher)feita
Lugar social	Rico, pobre, solteiro, público*, amigo, casado, filho, político*, capitão*, feiticeiro, notável*, português*, sertanejo, primitivo*, importante*, nobre, ativo, desconhecido*, sozinho, abastado*, civilizado	Rica, casada, pobre, perdida, filha, viúva, solteira, mãe*, sozinha, amiga, distinta, empregada, guerreira*, feiticeira, virgem, suprema*, íntima*, irmã, cozinheira*, brasileira

Quadro 3 – Predicações masculinas e femininas, listadas por frequência. As palavras com * referem-se a predicações usadas



O que vemos?



	Predicações exclusivamente masculinas	Predicações exclusivamente femininas
emoção/afeto	Valente, respeitável, entusiasmado, furioso, terrível, rude, indiferente, deslumbrado, feroz, cruel, apressado, hediondo, corajoso, destemido, triunfante, abatido, maravilhado, ruim, desconfiado, pacato	Amante; faceira, agitada, chorosa, ardente, ciumenta, adorável, zangadíssima, corada, sôfrega, espavorida, afortunada, desgraçada, adorada, soberbíssima, <u>brabinha</u> , tiriricas, <u>desafortunada</u> , agastada
caráter/personalidade	Mau, honrado, grave, valente, ilustre, extraordinário, ébrio, bobo, rude, sábio, perdido, seguro, franco, religioso, singular, polido, <u>circunspeto</u> , teimoso, leal, desgraçado	Perdida, de verdade, impaciente, ingênua, santa, namoradeira, obediente, ímpia, travessa, má, carinhosa, bondosa, insensível, angélica, morfética, frívola, trabalhadeira, rebelde, curiosa
aparência/corpo	Maduro, baixinho, barbado, sadio, limpo, galhardo, calvo, miúdo, rijo, quadragenário, encolhido, vesgo, míope, idoso, suado, tismado, espadaúdo, bonitinho, macilento, preto	Formosa, linda, encantadora, faceira, maltrapilha, tísica, esplêndida, ardente, muda, grávida, cansada, moça, coxa, máscula, cardíaca, corada, fraquinha, sã, seminua, deliciosa
lugar social	Público, político, capitão, notável, português, primitivo, importante, desconhecido, abastado, célebre, ilustrado, útil, formado, oficial, negociante, poderoso, pa, engenheiro	Mãe, livre, viúva, guerreira, suprema, íntima, cozinheira, noiva, esposa, donzela, núbil, troiana

Lugares públicos

Parecidos em quantidade; diferentes em qualidade

Lugares domésticos



O que ainda poderíamos ver?



CRUZAMENTO COM OUTRAS VARIÁVEIS



- Diferenças conforme período/tempo
- Diferenças conforme a escola literária
- Diferenças conforme os autores
- Diferenças conforme o gênero (M/F) dos autores
- Diferenças conforme o país dos autores

- Diferenças conforme *quem predica*
 - Personagens...
 - Narradores...



Outro exemplo – incompleto...

Discurso relatado (*reported speech*) como operador analítico

- Quem fala (mais)?
- O quanto fala?
- Fala sobre o quê?

```
([sema="dizer_relato.*"] [func="<SUBJ" & gen="M" & pos="PROP.* | PERS.*"] | [func="SUBJ.*" & gen="M" & pos="PROP.* | PERS.*"] [pos!="V"] {0,6} @[sema="dizer_relato.*"]) within s
```



O AC/DC já tem anotação semântica dos verbos de elocução :)

Procura: ([sema="dizer_relato.*"] [func="<SUBJ" & gen="F" & pos="PROP.*|PERS.*"]|[func="SUBJ.*" & gen="F" & pos="PROP.*|PERS.*"] [pos!="V"] {0,6} @[sema="dizer_relato.*"]) within s
Pedido de uma concordância em contexto
Corpo: OBras v. 7.12

814 ocorrências.

Número de ocorrências excessivo! Tente restringir a sua procura a menos de 500 casos.

Concordância

Procura: ([sema="dizer_relato.*"] [func="<SUBJ" & gen="F" & pos="PROP.*|PERS.*"]|[func="SUBJ.*" & gen="F" & pos="PROP.*|PERS.*"] [pos!="V"] {0,6} @[sema="dizer_relato.*"]) within s.

Apresenta-se uma amostra aleatória de 500 das **814** ocorrências encontradas.

id="Histórias_sem_Data Prosa:conto Mda 1884 masc ": -- Nunca fui janeleira, **dizia ela**, e acho muito feio que uma moça viva com o sentido na rua .

id="Dona_Guidinha_do_Poço Prosa:romance MdOP 1891 naturalismo_regionalismo masc ": **Guida** **noticiara-lhe** que logo que o Silveira estivesse de volta mandava-o tomar conta das bestas e outros animais .

id="Casa,_não_casa Prosa:conto Mda 1975 masc ": -- Bem, **disse ela**, já sei que me despreza .

id="Iaiá_Garcia Prosa:romance Mda 1878 romantismo masc ": **Estela supunha** que o amor de Jorge, ao fim de tão longo período, estaria acabado de todo, como produto da primeira estação .

id="Memorial_de_Aires Prosa:romance Mda 1908 realismo masc ": como se falasse da morte do Barão de Santa-Pia e da situação da filha, **D. Cesária perguntou** se ela realmente não casava .

id="O_Sacrifício Prosa:romance FT 1879 romantismo_realismo masc ": **perguntou Virgínia**, tanto que por entre árvores e sombras reconheceu Maurícia .

id="Clara_dos_anjos Prosa:romance LB 1948 naturalismo_realismo masc ": -- Penso que sim -- **disse Clara**, e acrescentou: -- olhe, papai, não pude passar a limpo a música .

id="A_Escrava_Isaura Prosa:romance BG 1875 romantismo masc ": Leôncio! ... onde vais! -- **exclamou Malvina** precipitando-se para ele; mal, porém, havia ela chegado à porta, ouviu-se a explosão atoadora de um tiro .

QUEM fala? *para além do gênero...*

- Caracterização dos personagens de Dickens por meio dos verbos de elocução

-

Ruano San Segundo, Pablo. 2016. “**A Corpus-Stylistic Approach to Dickens’ Use of Speech Verbs: Beyond Mere Reporting.**” *Language and Literature* 25 (2): 113-129.

Utilização de corpus para recuperar verbos de elocução e revelar como Dickens usa certos verbos para reportar a fala de certas personagens, contribuindo para a sua caracterização. A análise de 14 obras completas mostra que a prática não é um fenômeno isolado, mas um recurso estilístico importante em suas obras.

Cuidados - porque nossa língua é o Português...

Tabela 1. Distribuição de sujeitos ocultos por corpus

Corpus	Frases com sujeito oculto
Bosque-UD (v.2.4)	15.8%
DHBB	39.5%
Machado de Assis	28.42%

Desafios (para o português também):
Automatic Attribution of Quoted Speech in Literary Narrative
(Elson & McKeown, 2010)

Columbia Quoted
Speech Attribution
Corpus (2017)

AC/DC:
QUEMDISSE
(Freitas et al., 2016)

Quantos desses casos
estão envolvidos em
predicações?

Quantos desses casos
referem-se a verbos de
elocução?

Considerações finais (?)

MUITOS caminhos para explorar!

- Uma gramática das predicções humanas?
- Mas quem são os humanos (para além dos PROP)?
 - Humanos são personagens...
- Que categorias de análise são relevantes para enriquecer a leitura distante?
- Faz sentido pensar em categorias novas?
- O que já temos (para o português)?
- O que o AC/DC nos dá hoje?
- O que podemos fazer para deixar ainda melhor?
- O quão boa precisa ser a análise automática (morfológica, sintática, semântica) para nos oferecer resultados confiáveis?

Desafios do pré-processamento (ainda):
Do pdf ao txt (bem tratado), o caminho é longo...
: (

Referências

- ELSON, D. & MCKEOWN K. 2010. Automatic attribution of quoted speech in literary narrative. Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence (AAAI'10) , p. 1013-1019.
- FREITAS, C. Estudos linguísticos e Humanidades Digitais: corpus e descorporificação. *Gragoatá*, vol. 22, n. 44, p.1207-1227.
- FREITAS, B., FREITAS, C. & SANTOS, D. "QUEMDISSE?: Reported speech in Portuguese". In: Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)
- FREITAS, C., DE SOUZA, E. & ROCHA, L. *Quantificando e qualificando o sujeito oculto em Português*. Jornada de descrição do Português, Salvador, BA, 16-17 Outubro 2019.
- SABURI COSTA, B. & FREITAS, C. "Verbos de Elocução em Português: um estudo descritivo com base em grandes corpora e motivado pela Linguística Computacional". Revista Fórum Linguístico, vol. 14, n. 3, p. 2266-2285, 2017.
- PAIXÃO de SOUZA, M. C. P. A Filologia Digital em Língua Portuguesa: alguns caminhos. In: GONÇALVES e BANZA, Ana Paula Banza (Eds.). Património Textual e Humanidades Digitais: da antiga à nova Filologia. Évora: CIDEHUS, 2013. Disponível em <<http://dspace.uevora.pt/rdpc/bitstream/10174/10468/1/e-book.pdf>>. Acessado em: 30 out. 2017.
- RUANO SAN SEGUNDO, Pablo. 2016. "A Corpus-Stylistic Approach to Dickens' Use of Speech Verbs: Beyond Mere Reporting." *Language and Literature* 25 (2): 113-129.
- SANTOS, D., FREITAS, C. & BICK, E. OBRas: a fully annotated and partially human-revised corpus of Brazilian literary works in the public domain, OpenCor, Canela, RGS, Brasil, 24 de setembro de 2018.
- SANTOS, D.; BICK, E. Providing Internet access to Portuguese corpora: the AC/DC project. In Maria Gavrilidou, George Carayannis, Stella Markantonatou, Stelios Piperidis & Gregory Stainhauer (eds.), Proceedings of the Second International Conference on Language Resources and Evaluation (LREC 2000) (Atenas, Grécia, 31 de Maio a 2 de Junho de 2000), pp. 205-210.
- SANTOS, D. "Pesquisando corpos: Foi você que pediu um Corpo Ferreira?". Escola de Verão Belinda Maia (Edv 2009; Edv 2009) (FLUP, Porto, Portugal, 29 de Junho - 3 de Julho 2009)

Obrigada!

claudiafreitas@puc-rio.br