

Corpus psicolinguístico Léxico do Português Brasileiro

Gustavo Lopez Estivalet¹ Fanny Meunier²

Resumo: O Léxico do Português Brasileiro foi desenvolvido com o objetivo de oferecer um *corpus* baseado em palavras para a pesquisa em psicolinguística no português brasileiro. Ele foi criado a partir de um corpus com mais de 32 milhões de palavras. Assim, o Léxico do Português Brasileiro contém mais de 215 mil entradas lexicais e apresenta 21 colunas com informações metalinguísticas e psicolinguísticas relevantes, como categoria gramatical, frequência ortográfica, número de letras, vizinhos ortográficos, entre outras. Ele é um corpus aberto e de livre acesso na internet, possuindo uma plataforma amigável e dinâmica para pesquisas simples e complexas. O Léxico do Português Brasileiro ainda disponibiliza uma série de dados já computados, oferece um motor de geração de pseudopalavras do português brasileiro e um conjunto de ferramentas de linguística e estatística. Sendo assim, o presente artigo tem como objetivo introduzir e apresentar o Léxico do Português Brasileiro, e servir como seu manual de utilização. Ainda, é realizada uma descrição do desenvolvimento e criação do *corpus*. Enfim, o Léxico do Português Brasileiro preenche uma enorme lacuna na pesquisa em psicolinguística e linguística computacional, oferecendo um *corpus* baseado em palavras com valiosas informações metalinguísticas e psicolinguísticas do português brasileiro.

Palavras-chave: Psicolinguística. Linguística computacional. *Corpus*. Lexicografia. Linguística. Português brasileiro.

Introdução

O principal objetivo do Léxico do Português Brasileiro (LexPorBR)³ é oferecer um *corpus* baseado em palavras do português brasileiro (PB) que disponibilize o máximo de informações metalinguísticas e psicolinguísticas sobre as palavras do PB. O Léxico do Português Brasileiro é um *corpus* livre e aberto, consultado em uma plataforma simples e dinâmica através da internet. A partir de uma pesquisa, os resultados são apresentados de forma organizada e hierárquica, contendo dados metalinguísticos e psicolinguísticos das palavras ou grupos de palavras pesquisados.

Corpora psicolinguísticos são utilizados: (1) no controle, seleção e manipulação de palavras e critérios específicos para a criação de experiências psicolinguísticas; e (2) em

¹ Doutor em Neurociências e Ciências Cognitivas na École Doctorale Neurosciences et Cognition (NSCo) da Université Claude Bernard Lyon 1 (UCBL) no Laboratoire sur le Langage, le Cerveau et la Cognition (L2C2), Lyon, França, com bolsa do Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq). E-mail: gustavoestivalet@hotmail.com.

Doutorado em Psicologia Cognitiva pela Universidade René Descartes (1997, Paris, França). Diretora de pesquisa do Centre National de la Recherche Scientifique (CNRS). E-mail: fanny.meunier@isc.cnrs.fr.

³ http://www.lexicodoportugues.com/.



análises em linguística computacional da distribuição e do comportamento lexical (BAAYEN, 2001). Alguns exemplos de *corpora* psicolinguísticos baseados em palavra são: francês - *Lexique*⁴ (NEW et al., 2001, 2004), espanhol ó *Busca Palabras* (DAVIS; PEREA, 2005), inglês ó MRC ⁵ (COLTHEART, 1981), alemão, espanhol, francês, holandês e inglês - ClearPON⁶ (MARIAN et al., 2012), alemão, cirílico, holandês e inglês - CELEX (BAAYEN; PIEPENBROCK; VAN RIJN, 1995).

Esses *corpora* foram utilizados, por exemplo, em megaestudos que investigam o comportamento psicolinguístico no processamento de palavras e pseudopalavras: *English Lexicon Project* (BALOTA et al., 2007), *French Lexicon Project* (FERRAND et al., 2010), *Dutch Lexicon Project* (KEULEERS; DIEPENDAELE; BRYSBAERT, 2010), e *British Lexicon Project* (KEULEERS et al., 2012). Ainda, eles são utilizados na seleção, controle e manipulação de palavras para criação de experiências psicolinguísticas em inúmeros estudos e pesquisas específicas(GIMENES; NEW, 2015), assim como no desenvolvimento e simulação de modelizações linguísticas (SCHREUDER; BAAYEN, 1995).

NILC/São Carlos e Linguateca

O Léxico do Português Brasileiro foi desenvolvido a partir do *corpus* do Núcleo Interinstitucional de Linguística Computacional de São Carlos (NILC) ⁸ (PINHEIRO; ALUÍSIO, 2003) sediado no Instituto de Ciências Matemáticas e de Computação de São Carlos (ICMC/São Carlos) ⁹, da Universidade de São Paulo em São Carlos (USP/São Carlos) ¹⁰. As listas de formas e lemas divididas em categorias gramaticais foram baixadas do site do Linguateca ¹¹(SANTOS; BICK, 2000), onde se encontram informações do NILC, como

⁴http://www.lexique.org/.

⁵http://www.psych.rl.ac.uk/.

⁶http://clearpond.northwestern.edu/.

⁷http://celex.mpi.nl/.

⁸http://www.nilc.icmc.usp.br/nilc/index.php.

⁹http://www.icmc.usp.br/Portal/.

¹⁰http://www.saocarlos.usp.br/.

¹¹http://www.linguateca.pt/.



dados quantitativos e estatísticos¹², descendência do corpus¹³ e os arquivos de formas¹⁴ e lemas¹⁵ no formato .txt, separados por categorias gramaticais.

Lexique

A criação e o desenvolvimento do Léxico do Português Brasileiro foram inspirados no corpus psicolinguístico do francês Lexique (NEW et al., 2001, 2004). O Lexique tem oferecido dados sobre as palavras do francês a uma série de estudos e pesquisas, sendo um ótimo exemplo de corpus psicolinguístico simples e eficaz. Ele exemplifica as funcionalidades e utilidades que um corpus psicolinguístico deve e pode oferecer como recursos para a pesquisa em psicolinguística e linguística computacional. Oferece, ainda, uma série de informações indispensáveis para criação das experiências e análise dos resultados (categoria gramatical, frequência, número de letras, vizinhos ortográficos, entre outras), motores para criação de pseudopalavras, links e referências relacionadas ao corpus, assim como listas com dados já computados (FERRAND et al., 2010). Uma descrição detalhada desse corpus é encontrada no manual do Lexique¹⁶.

Programa e pacotes R

O Léxico do Português Brasileiro foi desenvolvido com o programa R¹⁷, com os dados originais importados a partir de arquivos .txt e cada coluna sendo criada e computada através de determinadas funções e algoritmos. O número de vizinhos ortográficos (Coltheartøs N) (COLTHEART et al., 1977) e a distância de Levenshtein ortográfica das 20 palavras mais próximas (OLD20) (YARKONI; BALOTA; YAP, 2008) foram calculados a partir das funções õcoltheart.Nö e õold20ö disponibilizadas no pacote õvwrö¹⁸ (KEULEERS, 2013).

¹²http://www.linguateca.pt/acesso/desc corpus.php?corpus=SAOCARLOS.

¹³http://www.linguateca.pt/acesso/NILCsaocarlos.html.

¹⁴http://www.linguateca.pt/acesso/contabilizacao.php#listaPosSAOCARLOS.

¹⁵http://www.linguateca.pt/acesso/contabilizacao.php#listaLemasSAOCARLOS.

¹⁶http://www.lexique.org/docLexique.php.

¹⁷http://www.r-project.org/.

¹⁸http://cran.r-project.org/web/packages/vwr/index.html.



Uma série de funções do pacote õlanguageRö¹⁹(BAAYEN, 2013) também foram utilizadas no desenvolvimento do Léxico do Português Brasileiro.

Léxico do Português Brasileiro

O projeto de criar o Léxico do Português Brasileiro nasceu de uma necessidade, em 2013, quando começamos a investigar a representação e o processamento morfológico flexional verbal no PB, no francês e em bilíngues com PB como língua materna e francês como língua estrangeira. Para as experiências em francês, os estímulos foram selecionados a partir do *corpus Lexique* (NEW et al., 2004), quando começamos a preparar as experiências em PB, deparamo-nos com a completa falta de um *corpus* psicolinguístico do PB. Procurando suprir nossas necessidades, tivemos acesso ao site do Linguateca (SANTOS; BICK, 2000) que reúne vários *corpora* do português europeu e brasileiro. Entretanto, não encontramos nenhum corpus do PB com dados metalinguísticos e psicolinguísticos apropriados para a criação rigorosa de experiências psicolinguísticas em PB. Foi nesse momento que decidimos fazer o Léxico do Português Brasileiro, que apresenta a página principal conforme a **Figura 1**.



Figura 1. Página principal do Léxico do Português Brasileiro.

¹⁹https://cran.r-project.org/web/packages/languageR/index.html.



Categorias de informações

No início de 2014, o Léxico do Português Brasileiro começou a ser desenvolvido em quatro etapas: (1) construção do *corpus* com palavras e informações metalinguísticas e psicolinguísticas, (2) construção das páginas na internet em HTML, (3) importação do *corpus* para um banco de dados MySQL na internet e (4) programação em PHP do funcionamento do corpus. Além disso, foram criadas as demais páginas do site: atualizações, downloads, ferramentas, créditos. Em seguida, foi desenvolvido o motor de geração de pseudopalavras do PB e as ferramentas de linguística estatística.

Para tanto, foi desenvolvida uma série de conhecimentos de programação computacional em R, HTML²⁰, MySQL²¹, PHP²², Java²³ e CSS²⁴. Selecionou-se no Linguateca o corpus do Núcleo Interdisciplinar de Linguística Computacional de São Carlos (NILC)²⁵como o mais pertinente para a criação do Léxico do Português Brasileiro. Essa seleção foi baseada nos seguintes critérios: (1) número total de palavras (32 milhões) condizente com outros corpora psicolinguísticos (Lexique, CELEX. ClearPOND)(BRYSBAERT; NEW, 2009), (2) quantidade e tamanho dos arquivos (13 arquivos, tamanho total 49 MB), (3) organização do corpus em arquivos .txt separados por categorias gramaticais, (4) organização dos arquivos em duas colunas (ortografia e frequência) separadas por tabulação e (5) recursos e publicações já desenvolvidos pelo NILC (PINHEIRO; ALUÍSIO, 2003).

Para primeira etapa foi realizado o download dos 13 arquivos em formato .txt do corpus do NILC no site do Linguateca²⁶ separados por categorias gramaticais (6 arquivos de formas: adjetivos, advérbios, gramaticais, nomes, numerais e verbos; 7 arquivos de lemas: adjetivos, advérbios, gramaticais, nomes, nomes próprios, numerais e verbos). Em seguida, utilizou-se o programa R para a criação das categorias de informações de todas as palavras com menos de 30 letras. Criaram-se diferentes colunas com: (1) ortografia (orto) e

²⁰http://pt.wikipedia.org/wiki/HTML.

²¹http://www.mysql.com/.

²²http://www.php.net/.

²³http://www.java.com/pt_BR/.

²⁴http://pt.wikipedia.org/wiki/Cascading_Style_Sheets.

²⁵http://www.linguateca.pt/acesso/corpus.php?corpus=SAOCARLOS.

²⁶http://www.linguateca.pt/acesso/contabilizacao.php.



(2) categoria gramatical (cat_gram), além de uma coluna com o tipo de palavra (forma ou lema). Todas as palavras foram padronizadas em letras minúsculas e as formas repetidas foram somadas. Criou-se uma coluna com um (3) número de identificação (id) da palavra de acordo com a organização do corpus por frequência em ordem decrescente e ordem alfabéticaa-z.

Em seguida, as seguintes colunas com informações sobre as palavras foram contabilizadas: (4) frequência ortográfica (freq_orto), (5) frequência ortográfica por milhão de palavras (freq_orto/M, [1000000*freq_orto/freq_total]), (6) logaritmo natural da freq_orto/M (log10_freq_orto), (7) número de letras (nb_letras) (BRYSBAERT; NEW, 2009). Logo após, foram criadas colunas com: (8) número de formas homógrafas (nb_homogr) e (9) categorias gramaticais das formas homógrafas (homografas). Ainda, foram criadas colunas com: (10) informações gramaticais (inf_gram), (11) forma ortográfica invertida (inv_orto), (12) estrutura CVCV(CVCV_orto), (13) estrutura CVCV invertida (inv_CVCV_orto), (14) bigramas (bigramas), (15) bigramas invertidos (inv_bigra), (16) trigramas (trigramas), (17) trigramas invertidos (inv_trigra) e (18) número aleatório entre 0 e 1 com oito dígitos de precisão (aleatorio). Enfim, foi calculado: (19) ponto de unicidade ortográfico (pu_orto), (20) número de vizinhos ortográficos (viz_orto) (COLTHEART et al., 1977) e (21) distância de Levenshtein ortográfica (old20) (YARKONI; BALOTA; YAP, 2008) com a utilização do pacote õvwrö (KEULEERS, 2013) e õlanguageRö (BAAYEN, 2013) para o programa R.

Sendo assim, o Léxico do Português Brasileiro versão Alfa conta com 21 colunas de informações metalinguísticas e psicolinguísticas conforme o **Quadro 1**. Ele possui 215.175 linhas com diferentes palavras do PB. Portanto, cada linha do Léxico do Português Brasileiro contém uma palavra e cada coluna uma determinada informação sobre esta palavra. O corpus completo é disponibilizado em uma tabela em formato .csv com codificação UTF-8 com um tamanho de 45 MB.

1	Ortografia (orto): forma ortográfica da palavra em letras minúsculas (com exceção dos nomes próprios),
	respeitando os acentos específicos de cada palavra ²⁷ .
2	Categoria gramatical (cat_gram): categorial gramatical da palavra (adj, adv, gram, nom, num, prop, ver).
3	Informação gramatical (inf_gram): informações gramaticais sobre a palavra (e.g. singular/plural,
	masculino/feminino, passado/presente/futuro, 1/2/3 pessoas, etc.).
4	Frequência ortográfica (freq_orto): número de vezes que a palavra aparece no NILC.
5	Frequência ortográfica por milhão (freq orto/M): número de vezes que a palavra aparece entre 1 milhão

²⁷ O corpus do NILC foi realizado em 1999, antes da reforma ortográfica do português.

_



	de palavras. Valor padrão para frequência de palavras(BRYSBAERT; NEW, 2009).
6	Logaritmo natural da frequência ortográfica (log10_freq_orto): logarítmico natural da frequência
	ortográfica. Utilizado para linearizar-se o comportamento da frequência das palavras no corpus(BAAYEN,
	2001).
7	Número de letras (nb_letras): número de letras da palavra.
8	Número de homógrafas (nb_homogr): número de palavras homógrafas. Palavras que possuem a mesma
	ortografia ou diferenças de acentos, mas pertencem a categorias gramaticais diferentes.
9	Homógrafas (homografas): categorias gramaticais das palavras homógrafas.
10	Ponto de unicidade ortográfico (pu_orto): letra a partir da qual a palavra se dissocia das outras, ou seja,
	letra a partir da qual a palavra é única. Sentido da esquerda para direita.
11	Vizinhos ortográficos (viz_orto): número de vizinhos ortográficos a partir do N de Coltheart, ou seja,
	alterando-se apenas uma letra por vez (COLTHEART et al., 1977).
12	Distância de Leveinshtein ortográfica (old20): distância ortográfica de Leveinshtein das 20 palavras mais
	próximas calculadas a partir de regressões lineares (YARKONI; BALOTA; YAP, 2008).
13	Estrutura CVCV (CVCV_orto): estrutura CVCV da palavra, onde C para consoantes C e V para vogais.
	Ainda, A para acentos, P para pontuação, N para números e S para símbolos.
14	Bigramas (bigramas): bigramas que constituem a palavra separados por õ_ö e limitados por õ#ö. O
	número de bigramas é igual ao número de letras da palavra mais 1.
15	Trigramas (trigramas): trigramas que constituem a palavra separados por õ_ö e limitados por õ#ö. O
	número de trigramas é igual ao número de letras da palavra.
16	Ortografia invertida (inv_orto): forma invertida da ortografia (orto).
17	Estrutura CVCV invertida (inv_CVCV_orto): estrutura CVCV da palavra invertida a partir de
10	(CVCV_orto).
18	Bigramas invertidos (inv_bigra): bigramas que constituem a palavra separados por õ_ö e limitados por
10	õ#ö invertidos a partir de (bigramas).
19	Trigramas invertidos (inv_trigra): trigramas que constituem a palavra separados por õ_ö e limitados por
20	õ#ö invertidos a partir de (trigramas).
20	Número aleatório entre 0 e 1 (aleatório): número aleatório entre 0 e 1 com oito algarismos de precisão.
21	Número de identificação (id): número de identificação da palavra designado a partir da organização do
	corpus por frequência decrescente e ordem alfabética a-z. O número de identificação é a posição da palavra
\Box	no corpus.

Quadro 1. Nome e descrição das 21 columas de informações metalinguísticas e psicolinguísticas do Léxico do Português Brasileiro versão Alfa.

Páginas e funcionamento

Para segunda etapa, utilizou-se o programa Notepad++²⁸para o desenvolvimento de toda programação visual em HTML e CSS, e programação lógica em PHP e MySQL do site do Léxico do Português Brasileiro em um servidor local com o programa XAMPP²⁹, contendo os módulos Apache, MySQL, PHP e Perls pré-instalados.

A página principal do site do Léxico do Português Brasileiro possui dois motores de pesquisa: (1) pesquisa simples e (2) pesquisa complexa. A pesquisa simples contém uma área de texto onde se podem inserir uma palavra específica ou múltiplas palavras em forma de lista. Pode-se copiar e colar uma lista de palavras de uma planilha ou editor de texto. A

_

²⁸http://notepad-plus-plus.org/.

²⁹http://www.apachefriends.org/pt_br/index.html.



pesquisa complexa contém quatro campos de inserção de critérios específicos das palavras a serem pesquisadas. No primeiro campo, o usuário deve escolher a coluna de informação pela qual deseja realizar a pesquisa. No segundo campo, deve escolher se deseja considerar õsimö ou desconsiderar õnãoö o critério. E no terceiro campo, o usuário deve inserir os critérios específicos de sua pesquisa.

Os símbolos coringas õ_ö para uma letra e õ%ö para uma cadeia de letras são aceitos pelo MySQL e podem ser utilizados em ambos os motores de pesquisa. Ainda, os símbolos maior que õ>ö e menor que õ<ö podem ser utilizados para pesquisas numéricas de grupos de palavras na pesquisa complexa. Inicialmente, a pesquisa complexa apresenta quatro campos de critérios para pesquisa, clicando-se no botão õ+Critériosö, o usuário é enviado a uma página que apresenta oito campos de critérios para a pesquisa. Cada motor de pesquisa possui um botão õProcurarö para iniciar a pesquisa e apresentar os resultados e um botão õLimparö para apagar os dados presentes nos campos. O usuário pode escolher a categoria utilizada para a organização e apresentação das palavras e o sentido de organização crescente ou decrescente, conforme a Figura 2:



Figura 2. Motores de pesquisa simples e pesquisa complexa do Léxico do Português Brasileiro.

Além da página principal do Léxico do Português Brasileiro: Léxico³⁰, as seguintes páginas ainda foram criadas para complementar o site: Pseudopalavras³¹, Downloads³², Ferramentas³³, Atualizações³⁴, Créditos³⁵e Linguística Estatística³⁶. õPseudopalavrasö acessa o motor de geração de pseudopalavras do PB, conforme descrito abaixo. õDownloadsö

³⁰http://www.lexicodoportugues.com/index.php.

³¹http://www.lexicodoportugues.com/pseudowords.php.

³²http://www.lexicodoportugues.com/downloads.php.

³³http://www.lexicodoportugues.com/tools.php.

³⁴http://www.lexicodoportugues.com/updates.php.

³⁵http://www.lexicodoportugues.com/credits.php.

³⁶http://www.lexicodoportugues.com/stat_ling.php.



disponibiliza uma série de arquivos pertinentes do Léxico do Português Brasileiro para downloads (corpus.txt, manuais, listas, convenções, bigramas, trigramas, scripts em R, entre outros). õFerramentasö disponibiliza uma série links de corpora, programas e literatura em psicolinguística e linguística computacional. õAtualizaçõesø descreve o desenvolvimento do Léxico do Português Brasileiro e as atualizações realizadas. õCréditosö apresenta o objetivo, a origem do Léxico do Português Brasileiro, assim como descreve as referências e pertinência do corpus do NILC, do Linguateca, do Lexique, do programa e dos pacotes R e da licença Creative Commons, finalizando com os agradecimentos. Enfim, Linguística Estatística é uma página que disponibiliza diversos recursos e ferramentas de livre acesso, conforme descrito abaixo.

Finalmente, todas as páginas e informações do Léxico do Português Brasileiro foram traduzidas para o inglês (*Brazilian Portuguese Lexicon*)³⁷. Implementou-se também o Google Tradutor em todas as páginas do Léxico do Português Brasileiro para a tradução do site para as línguas disponibilizadas nesse mecanismo. Sugere-se que o Google Tradutor seja utilizado a partir da versão inglês do *Brazilian Portuguese Lexicon*, pois assim não traduzirá os resultados das pesquisas, que por sua vez são sempre apresentados em PB.

Para terceira etapa, foi realizado o registro do domínio próprio do Léxico do Português Brasileiro (www.lexicodoportugues.com) junto ao HostGator ³⁸ e redirecionamento deste domínio para o servidor onde o corpus foi hospedado (http://www.biz.nf/). Esse servidor foi escolhido a partir dos critérios: (1) espaço de 250 MB, (2) banco de dados MySQL 5, (3) suporte à PHP 4/5, (4) 5000 MB de transferência, (5) hospedagem gratuita, (6) domínio gratuito do tipo http://portugueselexicon.co.nf, (7) webmail POP3/SMTP e (8) controle de arquivos por FTP. Assim, os dados do Léxico do Português Brasileiro foram importados em formato .csv e configurados com a utilização do phpMyAdmin ³⁹ para um banco de dados MySQL no servidor acima.

Na quarta etapa, foi realizada a programação de algoritmos em Java e PHP para: (1) manter os dados preenchidos nos campos da página HTML após pesquisa, (2) inserção de dois campos para organização dos resultados, um para seleção do critério de organização e outro para ordem crescente ou decrescente, (3) inserção do botão õ+Critériosö na pesquisa

³⁷http://www.lexicodoportugues.com/index_en.php.

³⁸http://hostgator.com.br/.

³⁹http://www.phpmyadmin.net/home_page/index.php.



complexa para disponibilização de oito campos de pesquisa, (4) reconhecimento dos símbolos maior que õ>ö e menor que õ<ö para as pesquisas numéricas, (5) desenvolvimento de um módulo de limitação e navegação dos resultados apresentados com o número de palavras a serem apresentadas (50, 100, 200 ou 500) e dois botões (õAnteriorö e õPosteriorö) para navegar entre as páginas de resultados, (6) apresentação de quatro informações gerais da pesquisa: i. total de palavras encontradas, ii. total de páginas de resultados, iii. intervalo das palavras apresentadas e iv. página apresentada, e (7) desenvolvimento do botão õExportar .csvö para exportar o resultado da pesquisa realizada em um arquivo .csv disponibilizado para download do usuário.

Convenções

Para a utilização do Léxico do Português Brasileiro, algumas convenções foram determinadas para realização das pesquisas e apresentação dos resultados.

- Categorias gramaticais: adj adjetivo, adv advérbio, gram gramatical, nom substantivo, num numeral, prop nome próprio, ver verbo.
- Estruturas CVCV das palavras possuem: V vogais, C consoantes, P pontuação, N
 números, A acentos, S símbolos.
- Símbolos coringas utilizados: õ<ö menor que, õ>ö maior que, õ_ö substitui uma letra,
 õ%ö substitui uma cadeia de letras.
- Ordem de apresentação dos resultados: crescente apresenta os resultados na ordem crescente, decrescente - apresenta os resultados na ordem decrescente.
- Botões: Procurar realiza a pesquisa e apresenta os resultados, Limpar -limpa os dados dos campo do formulário, +Critérios - direciona o usuário para uma página com mais critérios para a pesquisa complexa.
- Escolha sim/não: sim considera o critério, não desconsidera o critério.

Versão Alfa

Tendo em vista a enorme quantidade de informações metalinguísticas e psicolinguísticas que podem e serão computados, implementados e disponibilizados no



Léxico do Português Brasileiro, seu desenvolvimento foi dividido em três versões: 1) Alfa (2014), 2) Beta (2017) e 3) Delta (2019). O Léxico do Português Brasileiro versão Alfa foi inaugurada em 25 de março de 2014 e o surgimento do primeiro corpus psicolinguístico baseado em palavra do PB. A principal característica do Léxico do Português Brasileiro versão Alfa é que ele disponibiliza um *corpus* ortográfico em que as informações foram computadas a partir de dados ortográficos das palavras do PB do NILC.

A versão Beta contará com as informações: (1) fonológicas, (2) silábicas e (3) dos lemas associados às formas. A versão Delta contará com uma série de: (1) informações morfológicas, (2) informações sintáticas e (3) medidas de tempo de reação do reconhecimento de um grande número de palavras e pseudopalavras do PB, conforme os *Lexicon Projects* (BALOTA et al., 2007; FERRAND et al., 2010; KEULEERS et al., 2012; KEULEERS; DIEPENDAELE; BRYSBAERT, 2010).

Resultados

Na seção de resultados (Figuras 3 e 4), o usuário encontra os resultados da pesquisa organizada em diferentes linhas e as informações metalinguísticas e psicolinguísticas nas diferentes colunas. Encontram-se ainda uma série de informações pertinentes à pesquisa, conforme o lado esquerdo da Figura 4: (1) número total de palavras encontradas na pesquisa, (2) intervalo de palavras apresentados, (3) número total de páginas da pesquisa e (4) número da página apresentada. Pode-se escolher no campo superior à esquerda o número de palavras apresentadas em cada página e o usuário pode navegar entre os resultados e as páginas da pesquisa através dos botões õAnteriorö e õPróximoö. Um exemplo de pesquisa que apresenta palavras que possuem a categoria gramatical definida como õverboö pode ser visualizada a partir da pesquisa complexa com o critério **cat_gram - sim - ver**, conforme a Figura 3.



Vere 2 2 2 2 2 2 2 2 2	0,0057 0,0286 1,0441 0,0956 0,0956 0,0956 0,0019 0,0119 0,0019 0,	1,4150 1,4150 1,4150 1,4150 1,4150 1,5164 1,4150 1,5164 1,4150 1,5164 1			41141	11 21 21 11 11		1,00 SCVCCV #_bba_al_id_da_s# 1.20 SCVCCV #_bba_an_nj_a_s# 4	#b ba bal aid ida da# adl #b ba ban anj nja ja# ajn	adlab' VCCVCS ajnab' VCCVCS	#a_ad_dl_la_ab_b'#	#ed_adl_dla_lab_ab' b#	0,35822656 128902 0,08662146 34980 0,88272030 21495
Vee		1,17833 6 1,17833 6 1,17833 6 1,17833 6 1,17833 6 1,17833 6 1,17833 6 1,17833 6 1,17834 6 1,			41			SCVCCV # 'b ba an nj ja a#				with the steer and real title	
Vee		1,7833 6 6,2010 6 6 6,2010 6 6,2010 6 6 6,2010 6 6 6,2010 6 6 6,2010 6 6 6 6 6 6 6 6 6				[A] = [A				ř		-o de ne ne ne o	
Vee		0,3010 0,4771 0,1		1 2 1 1				1,00 SCVCCV # bbaar mraaf	#b ba bar arr ma ran arrab'		#a_ar_nra_ab_b' #	#ar_anr_nra_rab_ab'_b#	
Vee		0,4771 61 1,4472 51 1,4472 51 1,4472 51 1,0000 1,00000 1,4314 61 60 60 60 60 60 60 60		2 1 1			11 1,	1,20 SCVCCV # ba_ar_n_re_e#	#b ba bar an ne re# enab'	ab VCCVCS	#e_er_n_ra_ab_b'#	#er_err_rra_rab_ab'_b#	0,88779777 129002
Vec 13 Vec 15 Vec 15		1,4472 5			nom,	7	22 1	SCVCCV # 'b ba ar nr ro o#	≖ъ ъа ъаг ап по го≐	orrab' VCCVCS	#o or m ra ab b' '#	#ог оп па гар аb' b'#	0,05634216 102039
Vee 13 14 15 15 15 15 15 15 15		1,0000 1,5441 1,5				2	25 1,	1,00 SCVCV # bboocca at	#b_bococa_ca#	ob' VCVCS	#a_ac_co_ob_b' '#	#,q_,qo_qo5_o5=_5s#	0,46764344 33552
Vee		1,544			1	60	-1	.55 SCVCVC # bbooc ca am m#	#1b_bo_boç_oça_çam_am# ma	maçob' CVCVCS	#m ma aç ço ob b' "	# #ma_mac_aco_cob_ob' b'#	0,61274448 57308
Vere 1 1 1 1 1 1 1 1 1		0,0000 0,3010 1,4314 0,0000 0,0000 0,5342 1,5325 0,0000 0,0000 0,8431 0,6990		-1	-1	60	+	SCVCVC # b bo og ga ar r#	#1b bo boç oça çar ar# raç	raçob' CVCVCS	#t ra as so ob b' "#	#ra rac ace cob ob b#	0,68417173 29533
Vee Vee		0,3010 1,4314 0,0000 1,8325 0,9542 0,0000 0,0000 0,8431 0,6990				1	7	1.10 SCVCV #_bbo_cc_ce_e#	#b_bo_boc_oce_ce#	ob' VCVCS	#= ec_co_ob_b' #	#ec_eco_cob_ob" b#	0,59514092 200969
Vee Vee		1,4314 0,0000 1,8325 0,9542 0,0000 0,0000 0,8431 0,6990		1	1	3	1	1,80 SCVCVV # 'b bo oc ce ei #	#b bo boc oce cei ei# iec	iecob' VVCVCS	#i ie ec co op p, #	#ie iec eco cob ob' b'#	0,38859450 129371
Ver		0,0000 1,8325 0,9542 2,5703 0,0000 0,8451 0,6990	2 2 9 9 2 7	1		Per	44	1,60 SCVCVV # b bo og go ou u#	on #no nos oso soq oq q.#	noçob' VVCVCS	#, ,q qo o5 50 on n#	#4 'do do coo cou eu#	0,52948808 34251
Vee		1,8325 0,9542 0,9542 0,0000 0,0000 0,8451	5 9 9 5 7	1	4	en.	I.	1.80 SCVCVV = b_bo_or_ro_oe_e#		eorob' VVCVCS	#e_eo_cr_ro_ob_b' '#		0,16018517 201255
Vee 9 172 172 173 174 175		0,9542 2,5705 0,0000 0,8451 0,6990	2 9 9 7	2 10	nem,	95	59 1.	1,00 SCVCV # 'c ca al la a#	#'c 'ca cal ala la# alac'	c VCVCs	#a al la ac c' '#	#al ala lac ac' c'#	0,46227256 20129
Vee		0,0000 0,0000 0,8451 0,6990	,	1		2	1,	相	#'c_'ca_cal_ala_lam_am# ma	malac' CVCVCS	#m_ma_al_la_ac_c'_#	#ma_mal_ala_lac_ac_c'#	0,56367418 60356
Vee		0,0000	2 2	2 8	adj, 5	2	26 1,	1,00 SCVCVC # 'c_ca_al_la_ar_r#	#'c_'ca_cal_ala_lar_ar# ralac	ac CVCVCS	#r_ra_al_la_ac_c'_#	#ra_ral_ala_lac_ac'_c'#	0,03561097 6324
Vee		0,8451	2	1	-1	9		1,55 SCVCCV # 'c ca al 1d de ett	#'c 'ca cal aid ide de# ediac'	lac' VCCVCS	#e ed di la ac c' #	#ed edl dla lac ac' c'#	0,16646319 202411
Vee		0669'0	, v	1	4	10	39 1.	1,00 SCVCV #_c_ca_al_le_e#	#c_ca_cal_ale_le# elac	c VCVCS	#e_el_la_ac_c' '#	#el_ela_lac_ac'_c'#	0,47199025 68022
Vee 1 106				1		89	1	,45 SCVCVV # 'c_ca_al_le_ei_i#	#c ca cal ale lei ei# ielac	rc, VVCVCS	#i ie el la ac c' #	#ie iel ela lac ac' c#	0,11981471 79837
Veer 106 1 106 106 106 106 106 106 106 106 106 106 107 1	0,0319	0,000,0	2	1		-1	12 1,	1,15 SCVCVC # 'c_ca_al_le_em_m#	#'c 'ca cal ale lem em# me	melac' CVCVCS	#m me el la ac c' #	#me mel ela lac ac' c'#	0,40868094 202431
Veer 106 Veer 1 1 1 1 1 1 1 1 1	0,0637	0,3010	2	1	4	4	48	1,00 SCVCV #_c_ca_al_lo_o#	#c_ca_cal_alo_lo# olac	ie VCVCS	#0_ol_la_ac_c'_#	#ol_ola_lac_ac'_c'#	0,81120979 129890
Veer 1 1 1 1 1 1 1 1 1	3,3782	2,0253	9	1		9-4	15 1,	,00 SCVCVV = c ca al lo ou u#	#'c 'ca cal ale lou ou# uo	uotac' VVCVCS	#u uo ol la ac c' #	#uo uol ola lac ac' c'#	0,39501775 15185
Ver	0,0319	0,000,0	9	1	4	1	17 1,	1,00 SCVCCV #_'c_ca_al_lp_pa_a#	#c_ca_cal_alp_lpa_pa# aplac	lac' VCCVCS	#a_ap_pl_la_ac_c' "#	#ap_apl_pla_lac_ac'_c'#	0,10601483 202510
	0,0319	0,000	2	2 12	nom,	व	43 1,	1,00 SCVCV #_'c_ca_am_ma_a#	#c_ca_cam_ama_ma# amac	ac' VCVCS	#a_am_ma_ac_c'#	#am_ama_mac_ac_c'#	0,73092336 202537
Veer 15.2	6150'0	0,000	9	1	- 1		1	,25 SCVCCV # 'c ca an nd de e#	#'c 'ca can and nde de# edi	ednac' VCCVCS	#e ed dn na ac c' #	#ed edn dna nac ac' c'#	0,80485197 202811
Veer	4,8442	2,1818	\$	1	4	S.	50 1.	1,00 SCVCV # 'c_ca_ap_pa_a#	#c_cap apa pa# apac	ac VCVCS	#a_ap_pa_ac_c_#	#ap_apa_pac_ac_c'#	0,13694648 11994
Vee	2,3902	1,8751	9	1		5 21		1,00 SCVCVC # 'c_ca_ap_pa_am_m#	#"c_'ca_cap_apa_pam_am# ma	mapac' CVCVCS	#m_ma_ap_pa_ac_c' #	#ma_map_apa_pac_ac_c'#	0,52825827 18958
Veer	16,4448	2,7126	9	1		++	19 1,	,00 SCVCVC = 'c_ca_ap_pa_ar_r#		rapac' CVCVCS	#r ra ap pa ac c' #	#ra rap apa pac ac' c'#	0,54493440 4915
Ver	0,2868	0,9542	2	2 12	nom,	2	29 1.	1,00 SCVCV #_'c_ca_ap_pe_e#	#c_ca_cap_ape_pe# epac	ac VCVCS	#e_ep_pa_ac_c' #	#ep_epa_pac_ac_c'#	0,51854128 60381
Ver 9 193 Ver 5 Ver 5 Ver 14 Ver	0,1593	06690	9	1	1		10 1,	1,25 SCVCVV # 'c_ca_ap_pe_ei_#	#c ca cap ape per er iepac	ac VVCVCS	#i ie ep pa ac c' #	#ie iep epa pac ac' c'#	0,00668694 79890
Ver 193 Ver 5 Ver 4 4	0,2868	0,9542	9	1	1	1	10 1,	35 SCVCVC # c.ca.ap pe_em_m#	#'c 'ca cap ape pem em# me	mepac' CVCVCS	#m me ep pa ac c' '#	#me mep epa pac ac' c'#	0,60624337 60382
Ver Ver	6,1509	2,2856	9	1			11 [1,	1,35 SCVCVV # 'c_ca_ap_po_ou_u#	#c_ca_cap_apo_pou_ou#	uopac VVCVCS	#u_uo_op_pa_ac_c'_"#	#no_nop_opa_pac_ac_c'#	0,51369837 10274
ver ver		0669'0	2	-1	1	MF.	47 1,	,00 SCVCV = 'c_ca_av_va_a#	#'c_'ca_cav_ava_va# avac	ac VCVCS	#a av va ac c' #	#av ava vac ac' c'#	0,89154389 79988
1		0,6021	9	1	-1	ped .	18 1.	1,00 SCVCVC # 'c_ca_av_va_am_m#	#'c 'ca cav ava vam am# ma	mayac' CVCVCS	#m_ma_av_va_ac_c'_#	#ma mav ava vac ac c'#	0,49069928 89155
cavar ver	0,4462	1,1461	9	1		100	19 11	1,00 SCVCVC # 'c ca av va ar r#	#'c_'ca_cav_ava_var_ar# rav	ravac' CVCVCS	#r ra_av_va_ac_c' #	#ra_rav_ava_vac_ac_c'#	0,35418832 48578
cave ver 1	0,0319	0,000,0	2	1		(3)	32 1,	,00 SCVCV # 'c ca av ve em	#c_ca_cav_ave_ve# evac	ac VCVCS	#e_ev_va_ac_c' #	#ev_eva_vac_ac'_c'#	0,15725172 204013
clar ver 1 0	0,0319	0	\$	1	1	9		SCCVC #_'c_cl_la_ar_r#	#c_'cl_ cla_lar_ar# ralc'	c CVCCS	#r_ra_al_lc_c' #	#m_ral_alc_lc'_c'#	0,0194298 205621
coa ver 6 0	0,1912	0,7782	1	1	. 4	C1	27 1,	1,00 SCVV # 'c co oa a#	#c co cos co#	s Avcs	#a_ao_oc_c' '#	#20 20c 0c c'#	0,46102706 73412
coada ver 3	9560'0	0,4771	9	1	-1	and .	12 1	SCVVCV # 'c co oa ad da a#	#c_co_coa_oad_ada# ad	adaoc' VCVVCS	#a ad da ao oc c' #	#2d ada dao aoc oc c'#	0,80872251 103022
coam ver 3 0	0,1593	0,6990	2	1		1	13 1,	1,00 SCVVC #_'c_co_oa_am_m#	#c_co_coa_oam_am# maoc		#m ma ao oc c' #	#ma_mao_soc_oc" c'#	0,28047141 80122
coar ver 25 0	1961,0	1,3979	2	1			1	# 'c co oa ar r#	#c 'co coa oar ar# raoc'		#r ra ao oc c' '#	≓ra rao aoc oc' c'≅	0,98648495 35839
coars ver 1	0,0319	0	9	1		10	-	1,6 SCVVCV # 'c_co_os_ar_rs_s#	#'c_'co_coa_oar_arâ_râ# âra	áraoc' VCVVCS	#3_34_12_30_00_C_#	#sr_sta_rao_aoc_oc_c'#	0,09491235 205910
coava ver 1 0	0,0319	0	9	1		and .	1	SCVVCV # 'c_co_oa_av_va_a#	#c co cos oav ava va# ava	avaec' VCVVCS	#a av va ao oc c' #	#av ava vao aoc oc' c'#	0,2472736 205924

Figura 3. Exemplo de resultado de pesquisa do Léxico do Português Brasileiro.



Na parte superior a direita dos resultados, conforme a **Figura 4**, apresenta-se uma série de dados estatísticos estabelecidos e calculados a partir da pesquisa realizada (DAVIS, 2005; DAVIS; PEREA, 2005): 1) média, 2) valor máximo e 3) valor mínimo, das seguintes categorias: 1) freq_orto, 2) log10_freq_orto, 3) nb_letras, d) viz_orto e 4) old20. Futuramente, mais dados estatísticos serão inseridos neste módulo. Por fim, o botão õExportar .csvö exporta todos os dados da pesquisa para um arquivo .csv disponibilizado para *download* do usuário.

Resultados	Estatísticas						
50 • Anterior Posterior Exportar csv	categoria	freq_orto	log10_freq_orto	nb_letras	viz_orto	old20	
The Control of the Co	Média	48.6683	0.307745434371568	9.3912	1.9714	1.8379357585227	
Página 1 de 1767	Minimo	1	0	1	0	1	
0 - 50 palavras de um total de 88323 palavras encontradas		239218	5,3788	24	167	9,95	

Figura 4. Informações dos resultados e estatísticas básicas.

Pseudopalavras do PB

O motor gerador de pseudopalavras do PB foi desenvolvido para a criação de pseudopalavras baseadas na estrutura e frequência das palavras do PB. Diferentemente de outros motores de geração de pseudopalavras que se baseiam na estrutura silábica das palavras existentes da língua (KEULEERS; BRYSBAERT, 2010; MOTA; RESENDE, 2013), o motor de geração de pseudopalavras do PB do Léxico do Português Brasileiro utiliza os bigramas e trigramas(NEW et al., 2001). Todos os bigramas e trigramas foram contabilizados a partir de todas as palavras do Léxico do Português Brasileiro. As pseudopalavras são geradas a partir da frequência e combinação dos bigramas ou trigramas. Contabilizaram-se a (1) frequência geral dos bigramas e trigramas, (2) frequência dos bigramas e trigramas de acordo a posição na palavra e (3) frequência dos bigramas e trigramas por categoria gramatical.

No motor de geração de pseudopalavras do PB, o usuário deve inserir quatro campos: (1) número de palavras a serem geradas, (2) número de letras das palavras a serem geradas, (3) categoria gramatical que estas palavras devem pertencer (todas, adj, adv, gram, nom, num, ver) e (4) tipo de critério para a construção das palavras (bigramas ou trigramas). O motor de geração de pseudopalavras do PB constrói as palavras simultaneamente nos dois sentidos, da



esquerda para a direita e da direita para a esquerda, começando com um bigrama ou trigrama do tipo õ#xxö ou õxx#ö. De acordo com o número de letras, o motor vai concatenando novos bigramas ou trigramas que dividam o máximo de informação ortográfica com bigrama ou trigrama anterior (1 letra para os bigramas e 2 letras para os trigramas). O motor apresenta dois botões: õEnviarö para gerar e apresentar os resultados das pseudopalavras e õLimparö para limpar os dados dos campos, conforme a Figura 5.

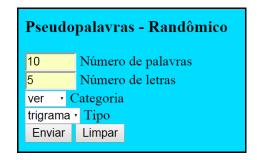


Figura 5. Motor de geração de pseudopalavras do PB.

Na tabela de resultados da geração de pseudopalavras do PB (10 pseudopalavras baseadas em verbos de 5 letras a partir de trigramas), conforme a Figura 6, quatro colunas com dados sobre as pseudopalavras são apresentadas: (1) categoria gramatical definida pelo usuário, (2) frequência da pseudopalavras calculada a partir da soma das frequências dos bigramas ou trigramas que compõem a pseudopalavra, (3) log10 da frequência calculada da pseudopalavra e (4) número de letras da pseudopalavra. Nos resultados, ainda é disponibilizado o botão õExportar .csvö para exportar os resultados da geração de pseudopalavras do PB para um arquivo .csv disponibilizado para download do usuário.

Resultados Exportar .csv									
pseudo esq-dir	categoria esq-dir	freq esq-dir	log_freq esq-dir	nb_letras esq-dir	pseudo dir-esq	categoria dir-esq	freq dir-esq	log_freq dir-esq	nb_letras dir-esq
cosas	ver	10522	9.2612	5	cados	ver	11736	9.3704	5
desta	ver	7338	8.9008	5	atica	ver	8879	9.0914	5
presa	ver	4743	8.4644	5	atico	ver	8662	9.0667	5
resta	ver	6972	8.8497	5	inais	ver	5958	8.6925	5
antes	ver	8512	9.0492	5	apres	ver	4182	8.3385	5
cassa	ver	4555	8.424	5	prado	ver	7894	8.9739	5
supes	ver	4215	8.3464	5	dista	ver	7424	8.9125	5
manta	ver	7666	8.9446	5	prada	ver	6615	8.7971	5
extra	ver	2909	7.9756	5	dinal	ver	4247	8.354	5
estas	ver	10813	9.2885	5	cante	ver	8074	8.9964	5

Figura 6. Exemplo dos resultados da geração de pseudopalavras do PB.



Linguística Estatística

A página Linguística Estatística do Léxico do Português Brasileiro disponibiliza livremente e abertamente recursos e ferramentas psicolinguísticas e de estatística linguística que podem ser consultadas diretamente na página através da internet. Esses recursos e ferramentas foram desenvolvidos em HTML/PHP, sendo eles: (1) Føe minFø- MS, (2) minFø - F1.F2, (3) teste de Hartley, (4) normalização entre 0 e 1, (5) inversor de palavras, (6) distância de Hamming, (7) distância de Levenshtein, (8) vizinhos ortográficos (Coltheart's N), (9) média das distâncias de Levenshtein, (10) entropia relativa, (11) frequência de palavras e (12) distribuição de Zipf, entre outros (BAAYEN, 2001, 2013; BRYSBAERT; NEW, 2009; DAVIS, 2005; DAVIS; PEREA, 2005; KEULEERS, 2013).

Licença Creative Commons

O Léxico do Português Brasileiro está licenciado com uma Licença Creative Commons - Atribuição-NãoComercial-CompartilhaIgual 4.0 Internacional ⁴⁰. Baseado no trabalho disponível em http://www.linguateca.pt/acesso/contabilizacao.php. Podem estar disponíveis autorizações adicionais às concedidas no âmbito desta licença em http://www.lexicodoportugues.com/creditos.php.

Agradecimentos

Agradecemos ao Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) pela bolsa de Doutorado Pleno no Exterior (GDE) e ao *Centre National de La Recherche Scientifique* (CNRS) pela estrutura para o desenvolvimento do Léxico do Português Brasileiro. Agradecemos às pesquisadoras do NILC Profa. Dra. Sandra M. Aluísio e Profa. Dra. Maria das Graças Volpe Nunes pelos valiosos materiais, informações e auxílio sobre o NILC, assim como o apoio na realização deste trabalho. Agradecemos aos pesquisadores Dr. Léo Varnet e Dr. Emmanuel Trouche pelas discussões sobre os scripts e

⁴⁰http://creativecommons.org/licenses/by-nc-sa/4.0/.



algoritmos para o desenvolvimento do Léxico do Português Brasileiro. Agradeço aos usuários dos fóruns de discussão e tutoriais da internet sobre o desenvolvimento de páginas e bancos de dados.

Referências

BAAYEN, R. H. Word Frequency Distributions. Dodrecht; Boston; London: Kluwer Academic Publishers, 2001.

BAAYEN, R. H. languageR: Data sets and functions with õAnalyzing Linguistic Data: A practical introduction to statisticsö. *R Package*, p. 133, 2013.

BAAYEN, R. H.; PIEPENBROCK, R; VAN RIJN, H. *The CELEX lexical database. Release 2 [CD-ROM]*. Philadelphia, PA: Linguistic Data Consortium, University of Pennsylvania, 1995.

BALOTA, D. A. et al. The English Lexicon Project. *Behavior Research Methods*, v. 39, n. 3, p. 4456459, 2007.

BRYSBAERT, M.; NEW, B. Moving beyond Ku era and Francis: A critical evaluation of current word frequency norms and the introduction of a new and improved word frequency measure for American English. *Behavior Research Methods*, v. 41, n. 4, p. 9776990, 2009.

COLTHEART, M. et al. Access to the internal lexicon. In: DORNIC, S. (Ed.). *Attention and Performance VI*. Hillsdale, NJ: Lawrence Erlbaum Associates, p. 5356555, 1977.

COLTHEART, M. The MRC psycholinguistic database. *The Quarterly Journal of Experimental Psychology Section A*, v. 33, n. 4, p. 4976505, 1981.

DAVIS, C. J. N-Watch: A program for deriving neighborhood size and other psycholinguistic statistics. *Behavior Research Methods*, v. 37, n. 1, p. 65670, 2005.

DAVIS, C. J; PEREA, M. BuscaPalabras: A program for deriving orthographic and phonological neighborhood statistics and other psycholinguistic indices in Spanish. *Behavior Research Methods*, v. 37, n. 4, p. 6656671, 2005.

FERRAND, L. et al. The French Lexicon Project: Lexical decision data for 38,840 French words and 38,840 pseudowords. *Behavior Research Methods*, v. 42, n. 2, p. 4886496, 2010.

GIMENES, M; NEW, B. Worldlex: Twitter and blog word frequencies for 66 languages. *Behavior Research Methods*, v. 48, n. 3, p. 963-972, 2015.

KEULEERS, E. et al. The British Lexicon Project: Lexical decision data for 28,730 monosyllabic and disyllabic English words. *Behavior Research Methods*, v. 44, n. 1, p. 2876 304, 2012.



KEULEERS, E. vwr: Useful functions for visual word recognition reserach. *R Package*, p. 19, 2013.

KEULEERS, E; BRYSBAERT, M. Wuggy: A multilingual pseudoword generator. *Behavior Research Methods*, v. 42, n. 3, p. 6276633, 2010.

KEULEERS, E; DIEPENDAELE, K.; BRYSBAERT, M. Practice Effects in Large-Scale Visual Word Recognition Studies: A Lexical Decision Study on 14,000 Dutch Mono- and Disyllabic Words and Nonwords. *Frontiers in Psychology*, v. 1, 2010.

MARIAN, V. et al. CLEARPOND: Cross-Linguistic Easy-Access Resource for Phonological and Orthographic Neighborhood Densities. *PLoS ONE*, v. 7, n. 8, p. e43230, 2012.

MOTA, M. B.; RESENDE, N. Metodologia da pesquisa em psicolinguística: desenvolvimento de uma ferramenta para a geração automática de pseudoverbos. *Letras de Hoje*, v. 48, n. 1, p. 100ó107, 2013.

NEW, B. et al. Une base de données lexicales du français contemporain sur internet : LEXIQUETM//A lexical database for contemporary french : LEXIQUETM. *Løannée psychologique*, v. 101, n. 3, p. 4476462, 2001.

NEW, B. et al. Lexique 2: A new French lexical database. *Behavior Research Methods, Instruments*, & *Computers*, v. 36, n. 3, p. 516ó524, 2004.

PINHEIRO, G. M.; ALUÍSIO, S. M. Corpus NILC: descrição e análise crítica com vistas ao projeto Lacio - WebSérie de Relatórios do Núcleo Interinstitucional de Lingüística Computacional NILC - ICMC - USP. São Carlos, SP: Universidade Federal de São Carlos - UFSCar, 2003.

SANTOS, D.; BICK, E. *Providing internet access to Portuguese corpora: the AC/DC project.* (M. Gavrilidou et al., Eds.)Proceedings of the Second International Conference on Language Resources and Evaluation (LREC2000). *Anais.*..Athens, Greece: 2000.

SCHREUDER, R; BAAYEN, R. H. Modeling Morphological Processing. In: FELDMAN, L. B. (Ed.). *Morphological Aspects of Language Processing*. Hillsdale, New Jersey: Lawrence Erlbaum Associates, Inc., Publishers, p. 1316154, 1995.

YARKONI, T; BALOTA, D.; YAP, M. Moving beyond Coltheart N: A new measure of orthographic similarity. *Psychonomic Bulletin & Review*, v. 15, n. 5, p. 9716979, 2008.



The Brazilian Portuguese Lexicon Psycholinguistic Corpus

Abstract: The Brazilian Portuguese Lexicon was developed to offer a word-based corpus for psycholinguistic research in Brazilian Portuguese language. It was created from a corpus with more than 32 million words. Thus, the Brazilian Portuguese Lexicon contains more than 215 thousand lexical entries and presents 21 columns with relevant metalinguistic and psycholinguistic information, such as grammatical category, orthographic frequency, word length, and orthographic neighbors, among others. It is an open free-access corpus consulted on the Internet; it has a friendly and dynamic interface for simple and complex searches. The Brazilian Portuguese Lexicon still offers a series of data computed, a Brazilian Portuguese pseudo word generation engine, and a collection of statistical linguistic tools. Therefore, the present article aims to introduce and present the Brazilian Portuguese Lexicon, as well as to serve as a manual for its use. Further, a description of its development and creation is given. Finally, the Brazilian Portuguese Lexicon fills a huge gap in the research in psycholinguistics and computational linguistics, offering a word-based corpus with valuable metalinguistic and psycholinguistic information from Brazilian Portuguese language.

Keywords: Psycholinguistics. Computational Linguistics. Corpus. Lexicography. Linguistics. Brazilian Portuguese.

Recebido em: 29 de maio de 2017.

Aprovado em: 07 de junho de 2017.