



Universidade do Minho



UNIVERSIDADE
DE VIGO

*lingua*MÁTICA

Volume 13, Número 2 (2021)

ISSN: 1647-0818

lingua

Volume 13, Número 2 – 2021

LinguaMÁTICA

ISSN: 1647-0818

Editores

Alberto Simões

José João Almeida

Xavier Gómez Guinovart

Conteúdo

Artigos de Investigação

Uso de tecnologias linguísticas para estudar a evolução dos sufixos -ÇOM e -VEL no galego-português medieval a partir de *corpora* históricos

Pablo Gamallo, José Ramon Pichel, José Martinho Montero Santalha & Marco Neves 3

AIA-BDE: um Corpo de Perguntas, Variações e outras Anotações

Hugo Gonçalo Oliveira & Ana Alves 19

Comissão Científica

Alberto Álvarez Lugrís,
Universidade de Vigo

Alberto Simões,
Instituto Politécnico do Cávado e Ave

Aline Villavicencio,
Universidade Federal do
Rio Grande do Sul

Álvaro Iriarte Sanroman,
Universidade do Minho

Anselmo Peñas,
Universidad Nacional de
Educación a Distancia

Antoni Oliver González,
Universitat Oberta de Catalunya

Antonio Moreno Sandoval,
Universidad Autónoma de Madrid

António Teixeira,
Universidade de Aveiro

Arkaitz Zubiaga,
Dublin Institute of Technology

Bruno Martins,
Instituto Superior Técnico

Carmen García Mateo,
Universidade de Vigo

Diana Santos,
Linguatca/Universidade de Oslo

Fernando Batista,
Instituto Universitário de Lisboa

Ferran Pla,
Universitat Politècnica de València

Gael Harry Dias,
Université de Caen Basse-Normandie

Gerardo Sierra,
Universidad Nacional
Autónoma de México

German Rigau,
Euskal Herriko Unibertsitatea

Helena de Medeiros Caseli,
Universidade Federal de São Carlos

Horacio Saggion,
University of Sheffield

Hugo Gonçalo Oliveira,
Universidade de Coimbra

Irene Castellón Masalles,
Universitat de Barcelona

Iria da Cunha,
Universidad Nacional de
Educación a Distancia

Itziar Gonzalez-Dios,
Euskal Herriko Unibertsitatea

Joaquim Llisterri,
Universitat Autònoma de Barcelona

José João Almeida,
Universidade do Minho

José Paulo Leal,
Universidade do Porto

Juan-Manuel Torres-Moreno,
Université d'Avignon et
des Pays du Vaucluse

Kepa Sarasola,
Euskal Herriko Unibertsitatea

Laura Plaza,
Complutense University of Madrid

Liliana Ferreira,
Fraunhofer Portugal AICOS & FEUP

Lluís Padró,
Universitat Politècnica de Catalunya

Luis Morgado da Costa,
Nanyang Technological University

Manex Agirrezabal,
University of Copenhagen

Marcos Garcia,
Universidade da Corunha

María Inés Torres,
Euskal Herriko Unibertsitatea

Mário Rodrigues,
Universidade de Aveiro

Mercè Lorente Casafont,
Universitat Pompeu Fabra

Miguel Solla Portela,
Universidade de Vigo

Mikel Forcada,
Universitat d'Alacant

Pablo Gamallo Otero,
Universidade de Santiago de
Compostela

Patrícia Cunha França,
Universidade do Minho

Patricia Martin Rodilla
Universidade de Santiago de
Compostela

Ricardo Rodrigues
Instituto Politécnico de Coimbra

Rogelio Nazar
Pontificia Universidad Católica de Val-
paraíso

Rui Pedro Marques,
Universidade de Lisboa

Sebastião Pais,
Universidade da Beira Interior

Susana Afonso Cavadas,
University of Exeter

Xavier Gómez Guinovart,
Universidade de Vigo

Artigos de Investigação



Uso de tecnologias linguísticas para estudar a evolução dos sufixos -ÇOM e -VEL no galego-português medieval a partir de *corpora* históricos

Use of Linguistic Technologies for Analysing the Evolution of Suffixes -ÇOM and -VEL in Medieval Galician-Portuguese from Historical Corpora

Pablo Gamallo  
CiTIUS

Univ. de Santiago de Compostela

José Martinho Montero Santalha 
Universidade de Vigo

José Ramom Pichel  
CiTIUS, USC / imaxin software
Santiago de Compostela

Marco Neves  
Universidade Nova de Lisboa

Resumo

O trabalho apresentado neste artigo tem dois objectivos. Por um lado, descreve a adaptação de duas ferramentas de processamento da língua natural ao galego-português medieval, nomeadamente um analisador morfosintático e um reconhededor de variedades medievais, e por outro, visa testar hipóteses linguísticas sobre a evolução de sufixos medievais mediante o uso dessas ferramentas em *corpora* históricos. Apesar de o desempenho das ferramentas ser inferior do que quando utilizadas para variedades modernas mais estandardizadas e com menos variabilidade formal, mostramos que é possível usá-las com grande fiabilidade para estudos quantitativos baseados em *corpus*. O estudo linguístico baseado em *corpus* permite-nos conferir que, pela sua distribuição de frequências, a presença dos sufixos -CION e -BLE nos textos medievais da Galiza foi provavelmente influenciada pelo castelhano baixo medieval.

Keywords

etiquetagem morfosintática, reconhededor de línguas, linguística histórica, humanidades digitais

Abstract

The work presented in this paper has two objectives. On the one hand, it describes how to adapt two natural language processing tools to medieval Galician-Portuguese, namely a morphosyntactic analyzer and a medieval language recognizer, and on the other hand, it verifies linguistic hypotheses about the evolution of medieval suffixes by using these tools by using historical corpora. Although the performance of the tools is inferior to those used for more standardized modern varieties with less formal variability, we show that it is possible to use them with great reliability for quantitative corpus-based studies.

The corpus-based linguistic study allows us to verify that, on the basis of their frequency distribution, the presence of the suffixes -CION and -BLE in medieval Galician texts is probably influenced by medieval Castilian

Keywords

part-of-speech tagging, language recognizer, historical linguistics, digital humanities

1. Introdução

A principal dificuldade para analisar automaticamente os textos medievais que se conservam é a sua grande variabilidade e falta de estandardização, o que faz com que ainda não se disponha de mecanismos para lematizar ou normalizar os textos mediante formas canónicas. Um exemplo claro desta variabilidade aparece refletido na Tabela 1, onde são registadas 15 variantes gráficas diferentes do termo *exceção*, encontradas em textos medievais da Galiza entre os séculos XIII e XV.

Esta extrema variabilidade formal está por trás da falta de recursos lexicais ou dicionários relativos ao período medieval, exceto algumas tentativas apresentadas em glossários ou dicionários etimológicos incompletos (Fillo & Lopes, 2013), ou dicionários parciais em fase de construção, como o Dicionário de Verbos do Português Medieval — DVPM (Xavier, 2005). Também há escassez de *corpora* anotados para serem usados como treino de sistemas de etiquetagem morfosintática e lematização. O Corpus Informatizado de Textos Portugueses Medievais (CIPM) (Xavier, 2016) é um interessante recurso textual com

exçeçon	2
exceiçon	2
exçeçon	7
excepcion	1
exçepcion	1
exçeçom	3
exçeçon	132
exceçon	2
exçeçon	1
exepçon	19
exeçon	4

Tabela 1: Variantes de *exceção* nos textos medievais da Galiza.

etiquetagem morfossintática revisada, mas não fornece nenhum tipo de lematização dos tokens etiquetados.¹

Em datas recentes, foi desenvolvido um novo módulo da ferramenta *LinguaKit* para a etiquetagem morfossintática de textos do galego-português medieval (Canosa et al., 2019). Para garantir uma maior abrangência do sistema e minimizar as limitações do módulo estatístico e do modelo de língua, o sistema foi enriquecido com regras e heurísticas que resultaram num sistema híbrido simbólico-estatístico. Como se mostrará na experimentação, a prestação deste sistema híbrido resulta em valores de exatidão modestos, bastante inferiores aos valores obtidos por sistemas treinados para a análise sincrónica de línguas modernas, já normalizadas e estandardizadas.

O presente artigo tem dois objetivos perfeitamente entrelaçados, um focado na área do processamento da linguagem natural e um outro filológico, orientado para realizar um estudo diacrónico baseado na análise de corpus. Desde o ponto de vista do processamento da linguagem, o artigo tem como objetivo descrever a adequação de duas ferramentas linguísticas ao galego-português medieval, nomeadamente a adaptação do módulo de análise morfossintática da *suite* linguística *LinguaKit*, e o treino e adição de novos modelos medievais ao reconhecedor de línguas *QueLingua*. Tanto a *LinguaKit* como o *QueLingua* são ferramentas que desenhamos e implementámos há vários anos e nas quais continuamos a trabalhar, para adaptá-las a novos domínios. Desde a perspectiva da análise de corpus, tencionamos estudar a evolução da distribuição da frequência dos nomes e adjectivos terminados em *-ÇOM* ou *-VEL*, em conjunto com as variantes terminadas em *-CION* e *-BLE*, tomando como fonte de estudo um corpus representativo do galego medieval da Galiza.

A partir das experiências realizadas, chegámos a dois tipos de conclusões, umas relacionadas com as ferramentas linguísticas e outras com o estudo filológico.

No tocante às ferramentas, concluímos que o módulo de etiquetagem morfossintática do galego-português medieval e o reconhecedor de variedades medievais são úteis e fiáveis para serem usados em tarefas de análise de corpus, mesmo com limitações importantes na prestação devido à grande variabilidade da língua. Chegamos a esta conclusão mediante uma avaliação extrínseca, que é a principal contribuição do presente trabalho. Além das avaliações intrínsecas com valores de exatidão obtidos a partir dum corpus de teste, neste artigo também realizámos uma avaliação extrínseca das ferramentas, ou seja, medimos o seu desempenho a partir de outra tarefa externa, que é a análise quantitativa da distribuições dos sufixos medievais. Mais concretamente, usámos primeiro um método semi-manual sem ferramentas de PLN para calcularmos a distribuição dos sufixos, muito custoso em tempo e trabalho. E a seguir realizámos a mesma tarefa com o uso exclusivo das ferramentas PLN adaptadas à língua medieval e sem revisão manual. Os resultados desta comparação mostram-nos que as duas ferramentas avaliadas são válidas para automatizar a compilação e análise de dados quantitativos fiáveis extraídos de textos medievais.

No que diz respeito ao estudo baseado em corpus, mostramos que a análise automática destes *corpora* medievais permite mesmo validar ou refutar hipóteses linguísticas sobre a *prolificidade* (Viaro, 2012) dos sufixos objeto de estudo.² O presente trabalho insere-se, desta forma, no âmbito das Humanidades Digitais, sendo a primeira vez que se pretende verificar mediante técnicas de linguística de corpus e processamento da língua as hipóteses filológicas formuladas sobre a evolução diacrónica dos sufixos objeto de estudo.

O que resta do artigo está organizado do seguinte jeito. As hipóteses linguísticas sobre a evolução dos sufixos são introduzidas na Secção 2. A seguir, na Secção 3, descrevem-se as duas ferramentas usadas: identificador de línguas medievais e analizador morfossintático do galego-português medieval. Estas ferramentas são utilizadas para a análise quantitativa dos dois pares de sufixos, *-ÇOM/-CION* e *-VEL/-BLE*, que é descrita na Secção 4. A avaliação das ferramen-

¹<https://cipm.fcsh.unl.pt/>

²Segundo Viaro (2012), a prolificidade refere-se à produtividade do ponto de vista diacrónico, e aponta, portanto, para o passado.

tas e da análise quantitativa é levada a cabo na Secção 5, para finalizarmos na Secção 6 com a enumeração das principais conclusões tiradas do estudo e de algumas novas ideias sobre trabalho futuro a desenvolver.

2. Hipóteses linguísticas sobre a evolução dos sufixos

2.1. Os nomes derivados de -TION

Segundo investigadores como [Mariño \(1998\)](#) e anteriormente [Lorenzo \(1985\)](#)³, no tocante às vozes derivadas dos nomes latinos terminados em -TION, no galego sempre houve duas tendências, uma mais popular ou patrimonial, com a remoção do iode (-ÇOM), e uma outra mais culta, com a conservação ou reposição deste: -CION. No galego moderno, à diferença do português, com a penetração dos cultimos, triunfou sempre a forma culta. Trata-se portanto, segundo estes autores, duma evolução interna da língua da Galiza que se produziu independentemente de interferências externas.

Existem, porém, outros pesquisadores que consideram que o triunfo do sufixo culto -CION no galego moderno é devido às interferências do castelhano sobre a língua da Galiza ([Freixeiro Mato, 1997](#); [Ferreiro, 1997](#)). No entanto, este processo não aconteceu na língua usada em Portugal, onde triunfaram as formas patrimoniais em -ÇOM.⁴

No presente estudo, iremos à procura de evidências no corpus diacrónico do galego-português medieval que nos permitam fortalecer ou rechaçar uma das duas hipóteses formuladas pela filologia: evolução interna ou interferência externa.

2.2. Os adjetivos derivados de -BĪLIS

Os sufixos descendentes de -BĪLIS sofrem dois tipos de evolução: a que mantém a vogal pós-tónica dando lugar a -vel(e) -uel(e), -uil(e), e a que se forma por síncope da vogal pós-tónica não final, dando lugar a -ble, -ule ou -bre. Segundo [Mariño Paz \(2005\)](#), as variantes sem síncope (variações de -VEL) são maioritárias até 1450 e so-

frem uma rápida descida no seu uso a partir dessa data, a favor das variantes de -BLE. O próprio autor conclui que este fenómeno pode explicar-se por influência do castelhano, tal e como afirma em [Mariño Paz \(2005, p.111\)](#):

[...] en todos os xéneros se percibe o predominio da opción -uel / -ueles (ou variantes) ata o ecuador do século XV. Paréceme fundada, por tanto, a sospeita de que a fulgurante expansión de -ble(s) na prosa notarial posterior a 1451 estivo en relación directa co aumento da familiaridade co castelán que se daría na actividade profesional de notarios e escribáns a partir do ecuador do século XV.

Como no caso anterior, procuraremos evidências baseadas em corpus que permitam conferir ou não esta hipótese filológica.

No entanto, devemos ter em conta uma limitação derivada das características das fontes analisadas. Como bem explica [Santalha & José-Martinho \(2005\)](#), é preciso distinguir entre “edição paleográfica” e “edição filológica.” A primeira é substancialmente fiel aos manuscritos, enquanto que a segunda procura reproduzir, de maneira regular, a língua que o escriba queria representar. No presente estudo, limitamos a analisar os textos digitalizados das edições filológicas pois há muito poucas edições paleográficas digitalizadas.

No que resta do artigo, usamos as maiúsculas para representar a forma gráfica normalizada dos sufixos, enquanto escrevemos em minúsculas as variantes de cada um deles. Por exemplo -çon e -zom são variantes de -ÇOM, e -uel e -vil são variantes de -VEL.

3. Ferramentas PLN para a língua medieval

Para levar a cabo a análise linguística alvejada no presente trabalho, foram adaptadas e treinadas duas ferramentas: um etiquetador morfosintático e um identificador de línguas.

3.1. Etiquetador morfosintático para o galego-português medieval

Em [Canosa et al. \(2019\)](#), foi descrito o módulo de classificação e reconhecimento de entidades mencionadas para textos do galego-português medieval. Este módulo para textos históricos contém também um tokenizador, um lematizador e um

³Na sua edição da *Crónica Troiana*, Ramón Lorenzo escreve: “No galego, [...] sempre apareceron en xogo dúas tendencias: unha máis popular, coa supresión do iode, e outra máis culta, coa conservación ou reposición. Na lingua moderna, coa penetración dos cultimos, triunfou case sempre a forma culta” ([Lorenzo, 1985, p.96](#)).

⁴A nível lexical, pelo contrário, houve influência castelhana no Português ([Messner, 2007](#); [Silvestre, 2008](#); [Venâncio, 2019](#)).

etiquetador morfossintático, todos eles adaptados para o galego-português medieval, e integrados na suite de ferramentas linguísticas, *LinguaKit* (Gamallo & Garcia, 2017; Gamallo et al., 2018), com licença livre GPLv3.⁵ *Linguakit* está também disponível como serviço web.⁶

O etiquetador e lematizador do galego-português medieval, ainda em fase de protótipo, contém os seguintes três elementos, dos quais o terceiro foi parcialmente adaptado para o presente trabalho:

Léxico medieval: O etiquetador morfossintático de *LinguaKit* contém um léxico computacional de formas, onde cada forma é associada a um ou vários lemas e às correspondentes etiquetas morfossintáticas. O *tagset* empregado tem 255 etiquetas diferentes e baseia-se nas recomendações do Grupo EAGLES (Leach & Wilson, 1996). Este *tagset* é comum a outros sistemas de análise morfossintática, nomeadamente o *FreeLing* (Padró, 2012) e o *TreeTagger* para português.⁷

Em Canosa et al. (2019), descreve-se como foi construído um léxico específico medieval a partir dos termos mais frequentes dum corpus medieval de desenvolvimento. Este léxico foi inserido noutra maior, constituído pela reunião dos léxicos pertencentes aos módulos de galego e português contemporâneo de *LinguaKit*.

Modelo de língua: O classificador estatístico de *LinguaKit* está baseado num simples algoritmo bayesiano que trabalha com um modelo de língua constituído por bigramas de pares $\langle \text{forma}, \text{etiqueta} \rangle$ e as suas probabilidades. O modelo usado para o processamento dos textos medievais foi desenvolvido a partir de duas fontes de dados: utilizou-se, por um lado, um modelo de galego moderno previamente treinado para o módulo correspondente do etiquetador (Garcia & Gamallo, 2015) e, por outro, treinou-se um novo modelo em base a um pequeno conjunto de textos medievais anotados automaticamente e revisados manualmente. Os dois modelos foram agrupados de tal maneira que se adicionaram ao modelo de galego moderno os novos pares do modelo medieval.

Regras linguísticas: O algoritmo de etiquetagem é um sistema híbrido, linguístico-estadístico, que contém, além dum modelo

de probabilidades e um desambiguador bayesiano, um conjunto de regras que podem alterar a escolha do classificador. Dois tipos de regras são consideradas pelo sistema:

- Regras usadas para corrigir erros do classificador devido ao excesso de ambiguidade. Por exemplo, uma regra específica impede etiquetar as formas *as* ou *os* como pronome pessoal seguido de nome comum. Este tipo de regras atuam diretamente sobre o classificador bayesiano.
- Regras morfológicas pensadas para garantir uma maior abrangência do sistema, que alargam a cobertura do léxico. Estas regras, que se aplicam sobre o léxico, permitem associar etiquetas morfossintáticas a tokens desconhecidos quando estes contêm determinados afixos. Por exemplo, se um token desconhecido ou OOV (*out of vocabulary*) finaliza com a sequência *-ados*, é etiquetado como verbo em modo participio, com os traços flexivos masculino e plural. As regras morfológicas são especialmente relevantes para melhorar a prestação de sistemas aplicados a línguas e variedades, como as dos textos medievais, de grande variabilidade formal e nenhuma estandardização. Para o presente trabalho, foram desenvolvidas novas regras morfológicas com diferentes variantes de afixos do galego-português medieval. Nomeadamente, as 4 regras adicionadas ao módulo para levar a cabo as experiências do presente trabalho estão descritas na Tabela 2. Repare-se que as regras tomam em conta, não só informação morfológica, mas também informação sobre a presença ou não do token a etiquetar no dicionário de formas (OOV) e o facto de ser identificado ou não como nome próprio pelo módulo NER (Named Entity Recognition - Reconhecimento de Entidades Mencionadas).

3.2. Identificador de línguas medievais

Nos textos medievais galegos, nomeadamente os de tipo notarial, encontramos um grande número de parágrafos escritos em castelhano. Isto tem uma influência dupla: por um lado restringe a efetividade das ferramentas de processamento da língua e, por outro, distorce as conclusões que se consigam tirar dos dados quantitativos extraídos. É portanto necessário efetuar um processo automático de identificação da língua parágrafo a parágrafo.

⁵<https://github.com/citiususc/Linguakit>

⁶<https://www.linguakit.com>

⁷(<https://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/data/Portuguese-Tagset.html>)

Regra 1	Se T é um token OOV não identificado como nome próprio e cuja forma termina numa variante do sufixo -ÇOM ou -CION, então T é um nome comum masculino singular: NCMS000
Regra 2	Se T é um token OOV não identificado como nome próprio e cuja forma termina numa variante do sufixo -ÇONS ou -CIONS, então T é um nome comum masculino plural: NCMP000
Regra 3	Se T é um token OOV não identificado como nome próprio e cuja forma termina numa variante do sufixo -VEL ou -BLE, então T é um adjetivo singular com género neutro: AQ0CS0
Regra 4	Se T é um token OOV não identificado como nome próprio e cuja forma termina numa variante do sufixo -VEIS ou -BLES, então T é um adjetivo plural com género neutro: AQ0CP0

Tabela 2: Regras morfológicas adicionadas ao módulo de etiquetagem morfossintática medieval. A notação morfossintática é a requerida pelo tagset da *LinguaKit*.

Com esse objetivo, elaborámos o primeiro identificador de idioma entre galego-português medieval e castelhano medieval. Treinamos a ferramenta *QueLingua* (Gamallo et al., 2014) com o corpus diacrónico Carvalho (Pichel et al., 2020),⁸ escolhendo para treino as variantes do castelhano entre os séculos XIII e XV, e textos de galego-português de Portugal dos mesmo séculos. O corpus de treino foi também enriquecido com textos do TMILG - Tesouro Medieval Informatizado da Língua Galega (Varela Barreiro et al., 2016). Em total, escolhemos 1 milhão e meios de tokens para o castelhano medieval e um número semelhante para o galego-português medieval. O identificador está baseado num método simples que toma como fonte de informação um dicionário ordenado com as formas mais frequentes de cada idioma a identificar. Dado que a distribuição das frequências segue a lei de Zipf, um pequeno conjunto de palavras dum língua representa uma ampla proporção do total de ocorrências de tokens em qualquer corpus textual dessa língua. Na configuração do nosso treino, usámos as 200 formas mais frequentes por cada idioma, que tendem a ser palavras gramaticais. A maioria são palavras muito curtas (entre 1 e 3 caracteres), pois são principalmente palavras gramaticais: artigos, preposições, pronomes e conjunções. Foi escolhido este número como resultado das experiências realizadas no trabalho descrito em Gamallo et al. (2014). Nesse trabalho, foram feitas avaliações com diferentes partições dum dicionário com as 5000 formas mais frequentes e observou-se que o F1-Score deixava de crescer significativamente a partir dum tamanho pequeno: 200. Este reduzido número tem a vantagem de

permitir uma revisão manual rápida e eficiente para limpar malformações frequentes derivadas dos textos de treino.⁹

4. Corpus histórico e estudo quantitativo

Nesta secção, o objetivo é levarmos a cabo um estudo quantitativo da distribuição dos sufixos -ÇOM/-CION e -VEL/-BLE mediante o uso das ferramentas descritas na secção precedente (3) aplicadas a um corpus histórico. Consideramos essencial que, para levar a cabo estudos filológicos sobre estes ou outros fenómenos linguísticos, é importante termos a nosso dispor dados quantitativos que espelhem a distribuição e evolução temporal das formas alvejadas. No presente estudo, compararemos a distribuição no corpus das formas terminadas em -ÇOM e -VEL com a distribuição das variantes em -CION e -BLE. A análise quantitativa ajudará a mostrar qual das hipóteses formuladas na Secção 2 condiz melhor com os dados quantitativos extraídos. Começamos a secção descrevendo a coleção de textos históricos.

4.1. O corpus histórico

O corpus medieval da Galiza utilizado para este estudo forma parte do TMILG (Varela Barreiro et al., 2016) e consta de 24 documentos datados entre os séculos XIII e XV com 1,5 milhões de tokens. Foi inicialmente construído para estudos sobre distâncias entre variedades diacrónicas (Pichel et al., 2020). Os documentos abrangem diferentes géneros textuais: líricos, ensaístas e no-

⁸O corpus Carvalho está disponível em fegalaz.usc.es/~gamallo/resources/Carvalho.tgz

⁹O reconhecedor de línguas, *QueLingua*, e os novos dicionários medievais construídos estão disponíveis em <https://github.com/gamallo/QueLingua>.

tariais. Entre as obras compiladas, incluem-se as Cantigas de Santa Maria, Crónica Geral de Castela, Cancioneiro de Ajuda, cantigas de Airas Nunes e numerosas atas notariais. A lista completa dos textos incluídos no corpus está no Anexo A. Para dar a conhecer a língua empregada nos textos deste corpus, deixamos uma breve amostra na Tabela 3.

“et este nombre munene quiere dezir en arauigo tanto como enel nuestro language de castiella. lo que desseamos. et cuenta aquel sabio que esta duenna era de buen seso et de grand conseio.
(General Estoria. Alfonso X)”

Tabela 3: Amostra de texto medieval em galego-português na parte galega.

O corpus medieval de Portugal inclui textos dos séculos XIII ao XV e forma parte do corpus diacrónico e multilingue Carvalho com licença livre (Pichel et al., 2019, 2020). Pode ver-se uma breve amostra dum texto deste corpus na Tabela 4. O sub-corpus português de Carvalho consta de 1,7 milhões de tokens (ligeiramente maior que o corpus medieval da Galiza), e inclui também textos de diferentes géneros, como Chronica de Dom João I, Cantigas de Dom Dinis e documentos notariais.

“A quantos esta carta uiren faço saber que Domingos perez filho de Maria. martjz dicta Daynha mj mostrou hũa mha carta que de mjn ten pola qual eu enpraizei a el. e aa primeira molher con que fosse casado dous casaaes e hũu Moynho que eu ei na quintaa de Maceeira.”
(Chancelaria de Dom Afonso. Volume I)

Tabela 4: Amostra de texto medieval em galego-português na parte portuguesa.

O identificador de idioma foi aplicado aos dois *corpora* para separar os parágrafos em galego-português, que são a imensa maioria, dos parágrafos em castelhano medieval, presentes sobretudo nos textos notariais do corpus da Galiza. Uma vez realizada a selecção dos parágrafos, o texto em galego-português medieval foi processado com o etiquetador morfossintático medieval de *LinguaKit*.

No resto da secção, analisamos primeiro a distribuição dos sufixos -ÇOM/-CION, e a seguir apresentamos a mesma análise com o par -VEL/-BLE.

4.2. Distribuição dos sufixos -ÇOM/-CION

Para levar a cabo o estudo sobre este par de sufixos, derivados do latino -TION, que constrói nomes de ação a partir de verbos, foram listadas primeiro todas as variantes gráficas identificadas nos textos e mostradas na Tabela 5. Foi tomada em conta a listagem publicada em Diéguez (2018). A variação é devida principalmente a três fenómenos de natureza morfofonológica: alomorfia da vogal temática, haplogogia e fusão de vogais semelhantes (Cristine Prado & Massini-Cagliari, 2014).

-ÇOM:	-çom, -som, -zom, -çon, -son, -zon, -çãõ, -sãõ, -sao, -çao, -çõ, -çón, -són, -zón, -zóm
-CION:	-cion, -sion, -siom, -çiom, -çion, -ción, -syon, -sióm, sión, -çión

Tabela 5: Variantes de -ÇOM e -CION em textos medievais.

Devido ao grande número de variantes para um mesmo sufixo, e para simplificar a tarefa, só as formas no singular foram consideradas.

4.2.1. Distribuição dos sufixos no corpus da Galiza

Uma vez definidas as variantes de cada sufixo, iniciamos o processo de extração a partir dos textos da Galiza etiquetados morfossintaticamente e identificados como sendo escritos em galego-português pelo reconhecedor de idioma. Deste jeito, foram extraídas todas as ocorrências das palavras etiquetadas mediante o tag NC (nome comum) e contendo alguma das variantes listadas na Tabela 5. Não foram considerados os nomes com iode desaparecido em estágios iniciais, nomeadamente *coraçom*, *razom* e *sazom*, cuja evolução é coincidente com a do castelhano.

Uma vez extraídas as palavras com os sufixos alvejados, contamos o número de ocorrências totais e número de palavras diferentes por variante e agregamos os resultados para obtermos o valor total de cada sufixo. As tabelas 6 e 7 apresentam a distribuição dos dados quantitativos de -ÇOM e -CION, respetivamente, em relação às suas variantes no corpus de textos medievais da Galiza.

Na última coluna das tabelas 6 e 7 mostra-se a probabilidade de cada variante do sufixo. Por exemplo, a probabilidade P de -çon, como variante de -ÇOM, é calculada na equação 1, onde f é a função que devolve a frequência de cada variante, v , pertencente ao conjunto {-ÇOM}.

-ÇOM	freq	formas	P
-sao	11	4	0,003
-som	80	18	0,022
-son	769	74	0,220
-são	1	1	0,0002
-són	24	8	0,006
-zom	19	4	0,005
-zon	61	21	0,017
-zóm	7	1	0,002
-zón	22	5	0,006
-çao	2	2	0,0005
-çom	67	23	0,019
-çon	2119	226	0,606
-ção	14	3	0,004
-çón	61	29	0,017
-çõ	236	66	0,067
Total	3493	485	1

Tabela 6: Ocorrências (freq) das variantes de -ÇOM no corpus medieval da Galiza, junto com o número de formas de palavras diferentes por variante (formas) e probabilidade de cada variante (P).

-CION	freq	formas	P
-cion	15	11	0,087
-siom	2	2	0,011
-sion	33	16	0,191
-sióm	4	2	0,023
-sión	16	8	0,093
-çiom	1	1	0,005
çion	78	40	0,453
-çión	23	19	0,133
Total	172	99	1

Tabela 7: Ocorrências (freq) das variantes de -CION no corpus medieval da Galiza, junto com o número de formas de palavras diferentes por variante (formas) e a probabilidade da variante (P).

$$P(-çon) = \frac{f(-çon)}{\sum_{v \in \{-ÇOM\}} f(v)} = \frac{2119}{3493} = 0.606 \quad (1)$$

A variante mais frequente do sufixo -ÇOM é -çon, com uma probabilidade de ocorrência de 0.606. É de salientar que a probabilidade de ocorrência desta variante é muito superior ao resto. No tocante, ao sufixo -CION, a variante mais frequente é çion, com uma probabilidade de 0.453. Segundo Santalha & José-Martinho (2005), nas edições paleográficas com os textos originais, o uso de 'õ' final era maioritário, mas

este foi transcrito como 'n' final nas edições filológicas dos textos da Galiza, provavelmente por influência do castelhano.

-TION	freq	P	R ₁	R ₂
-ÇOM	3493	0,953	20,308	7,202
-CION	172	0,047	0,049	1,737

Tabela 8: Comparativa dos sufixos derivados de -TION (-ÇOM e -CION) no corpus medieval da Galiza, no tocante às ocorrências totais (freq), a probabilidade de ocorrência (P) e às duas ratios, R_1 e R_2 .

A probabilidade dum sufixo dadas as duas alternativas derivadas do sufixo latino -TION é calculada dividindo a frequência agregada de todas as variantes do sufixo pelo total de ocorrências das duas alternativas. A equação 2 mostra o cálculo da probabilidade do sufixo -ÇOM dadas as alternativas (patrimoniais e cultas) derivadas do conjunto {-TION}, que reúne todas as variantes de -ÇOM e -CION.

$$P(-ÇOM) = \frac{\sum_{v \in \{-ÇOM\}} f(v)}{\sum_{z \in \{-TION\}} f(z)} = \frac{3493}{3665} = 0.953 \quad (2)$$

A ratio R_1 devolve a razão entre os dois sufixos. As duas equações em 3 mostram a razão de -ÇOM para -CION e vice-versa.

$$\begin{aligned} R_{1(-ÇOM, -CION)} &= \frac{\sum_{v \in \{-ÇOM\}} f(v)}{\sum_{z \in \{-CION\}} f(z)} \quad (3) \\ &= \frac{3493}{172} = 20,308 \end{aligned}$$

$$\begin{aligned} R_{1(-CION, -ÇOM)} &= \frac{\sum_{v \in \{-CION\}} f(v)}{\sum_{z \in \{-ÇOM\}} f(z)} \\ &= \frac{172}{3493} = 0,049 \end{aligned}$$

A ratio R_2 devolve a razão entre a frequência total dum sufixo e o número de formas diferentes. É a ratio inversa à fórmula da *token/type ratio*, conhecida como TTR (Kettunen, 2014). Nas duas equações em 4, calcula-se a R_2 dos dois sufixos, onde T é a função que devolve o número de formas diferentes (ou *types*) associadas a todas as variantes dum sufixo.

$$R_{2(-ÇOM)} = \frac{\sum_{v \in \{-ÇOM\}} f(v)}{T(-ÇOM)} = \frac{3493}{485} = 7,202 \quad (4)$$

$$R_{2(-CION)} = \frac{\sum_{v \in \{-CION\}} f(v)}{T(-CION)} = \frac{172}{99} = 1,737$$

Em resumo: o sufixo -ÇOM tem uma frequência total de 3493 ocorrências em 485 palavras ou formas diferentes, enquanto -CION ocorre só 172 vezes em 99 formas diferentes. A probabilidade de aparição do sufixo patrimonial é muito mais alta (0,953) que a do sufixo culto (0,047). No tocante à ratio de uso entre eles (R_1), o primeiro é 20 vezes mais usado do que o segundo. No que diz respeito à ratio R_2 , cada forma diferente em -CION ocorre de média menos de 2 vezes no corpus, enquanto que as formas com variantes de -ÇOM tendem a ter uma frequência média muito superior (7,2) e, portanto, têm mais produtividade, é dizer *prolificidade* (vd. nota supra 2), e uso na época medieval. Esta informação comparativa pode consultar-se na Tabela 8, que mostra os resultados agregados de frequência e as medidas relativas: P , R_1 e R_2 . Pode-se concluir, a partir destes cálculos, que o uso do sufixo patrimonial -ÇOM é quase hegemónico nos textos medievais da Galiza, sendo o uso do sufixo culto muito minoritário.

4.2.2. Distribuição dos sufixos no corpus de Portugal

Realizamos a mesma análise sobre os textos de Portugal etiquetados morfossintaticamente. Mostramos os resultados nas tabelas 9 e 10

-ÇOM	freq	formas	P
-sao	28	1	0.007
-som	470	77	0.124
-son	38	13	0.010
-são	38	1	0.018
-són	1	1	0.0002
-zom	275	13	0,072
-zon	37	14	0.009
-çao	14	13	0.003
-çom	2239	386	0.591
-çon	155	51	0.040
-çãõ	282	102	0.074
-çón	1	1	0.0002
-çõ	175	90	0.046
Total	3784	770	1

Tabela 9: Ocorrências (freq) das variantes de -ÇOM no corpus medieval de Portugal, formas diferentes e probabilidade de cada variante (P).

A variante mais frequente do sufixo -ÇOM é -çom (terminada em ‘m’), com uma probabilidade de ocorrência de 0,59 (ver Tabela 9). Repare-se que nos textos da Galiza a variante mais frequente termina em ‘n’, nomeadamente -çon.

A Tabela 11 mostra os resultados agregados

-CION	freq	formas	P
cion	8	5	0,307
siom	9	3	0,346
sion	6	3	0,230
çiom	2	2	0,076
çion	1	1	0,038
Total	26	14	1

Tabela 10: Ocorrências (freq) das variantes de -CION no corpus medieval de Portugal, formas diferentes e a probabilidade da variante (P).

de frequência e as medidas relativas: P , R_1 e R_2 . Observamos que não há grandes diferenças entre os dois *corpora* no tocante aos resultados agregados, além do uso oposto do ‘m/n’ final e um uso de -CION ainda mais residual no corpus de Portugal. De resto, existe uma grande simetria no uso das variantes de -ÇOM nos textos da Galiza e de Portugal: encontramos quase as mesmas variantes gráficas e um número muito próximo de ocorrências de palavras com este sufixo. Além de mais, a probabilidade das duas variantes mais frequentes nos textos da Galiza, -çon e -son, é de 0,606 e 0,220, respetivamente, enquanto nos textos de Portugal as duas variantes mais frequentes são -çom e -som, com uma probabilidade de 0,591 e 0,124. Uma vez identificadas as diferenças no ‘m/n’ final, as distribuições no uso das variantes são muito semelhantes.

Apesar destes fortes paralelos na língua medieval, o galego moderno padronizou a forma -CIÓN (como em espanhol) enquanto o português moderno reuniu em -ÇÃO as formas derivadas de -TION junto com as doutros sufixos latinos: -ANT, -UNT, -AN, -ON.

-TION	freq	P	R_1	R_2
-ÇOM	3784	0.993	145,53	4,914
-CION	26	0.007	0,006	1,857

Tabela 11: Comparativa dos sufixos derivados de -TION (-ÇOM e -CION) no corpus medieval de Portugal, no tocante às ocorrências totais (freq), à probabilidade de ocorrência (P) e a duas ratios, R_1 e R_2 .

4.2.3. Índícios de castelhanização nos textos da Galiza

Dois novos testes foram conduzidos com o intuito de conferir se o uso do sufixo -CION está relacionado com a crescente castelhanização da língua medieval na Galiza.

Em primeiro lugar, foi realizada uma análise diacrónica, século a século, da distribuição dos dois sufixos na Galiza. Para este propósito, o corpus da Galiza foi dividido em três partições: textos do século XIII, do XIV e do XV. A Figura 1 mostra a evolução da ratio R_1 de -ÇOM para -CION ao longo dos três séculos. A figura revela que a proporção de palavras com o sufixo -ÇOM ao respeito de -CION é claramente menor no século XV que em séculos anteriores. Enquanto nos séculos XIII e XIV a ratio de palavras com sufixo em -ÇOM é por volta de 27 vezes superior, este valor baixa a 13 no século XV. Isto parece ser um indício de que a presença de -CION está relacionada com uma maior castelhanização, pois é no século XV onde a influência do castelhano sobre o galego é maior (Ferreiro, 1997).

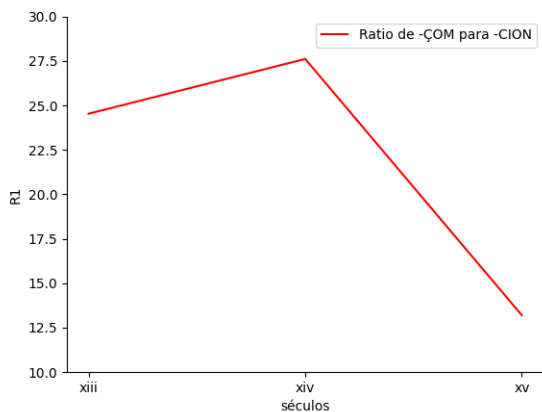


Figura 1: Evolução diacrónica ao longo de três séculos da ratio R_1 de -ÇOM para -CION, no corpus medieval da Galiza.

Em segundo lugar, foi levada a cabo a análise de distribuição dos dois sufixos utilizando só os excertos dos textos da Galiza identificados automaticamente como sendo escritos em castelhano pelo reconhecedor de língua (ver Tabela 12). Mesmo se a quantidade de texto é muito menor, as medidas relativas P , R_1 e R_2 mostram uma tendência contrária à encontrada nos textos em galego-português: maior presença de -CION que de -ÇOM e com maior ratio de uso.

-TION	freq	P	R_1	R_2
-ÇOM	62	0.319	0,469	1,589
-CION	132	0.681	2,129	2,425

Tabela 12: Comparativa dos sufixos derivados de -TION (-ÇOM e -CION) no corpus medieval da Galiza em excertos em castelhano, no tocante às ocorrências totais (freq), à probabilidade de ocorrência (P) e às duas ratios, R_1 e R_2 .

4.3. Distribuição dos sufixos -VEL e -BLE

O mesmo tipo de experiência foi levada a cabo para analisar a distribuição do par -VEL/-BLE, derivados do sufixo latino -BĪLIS, que constrói adjetivos a partir de verbos. Baseamo-nos no trabalho de Mariño Paz (2005) para listar as variantes possíveis dos dois sufixos (ver Tabela 13). Foram extraídas todas as ocorrências das palavras etiquetadas mediante o tag AQ (adjetivo comum) e contendo alguma das variantes listadas na Tabela 13. Como no caso descrito na secção anterior, só as formas no singular foram consideradas.

-VEL:	-vel, -bel, -uel, -vil, -uil, -uel, -uele, -vele, -velle, -uelle
-BLE:	-ble, -vle, -ule

Tabela 13: Variantes de -VEL e -BLE em textos medievais.

4.3.1. Distribuição dos sufixos no corpus da Galiza

As tabelas 14, 15 e 16 mostram os resultados obtidos a partir dos textos da Galiza. Embora haja muitos mais casos de -VEL que de -BLE, as ratios são menores do que no par -ÇON/-CION.

Podemos observar também que as palavras com sufixos derivados de -BĪLIS são muito menos numerosas que as derivadas de -TION. Trata-se portanto dum sufixo pouco produtivo na Idade Média, aplicável a muito poucas formas lexicais.

-VEL	freq	formas	P
-bel	1	1	0,006
-uel	49	17	0,304
-uele	5	3	0,031
-uelle	6	4	0,037
-uil	31	14	0,192
-vel	52	8	0,322
-vele	3	1	0,018
-vil	14	9	0,086
Total	161	57	1

Tabela 14: Ocorrências (freq) das variantes de -VEL no corpus medieval da Galiza, formas diferentes por variante e probabilidade de cada variante (P).

-CION	freq	formas	P
ble	16	7	0,80
ule	2	1	0,10
vle	2	11	0,10
Total	20	9	1

Tabela 15: Ocorrências (freq) das variantes de -BLE no corpus medieval da Galiza, formas diferentes por variante e a probabilidade de cada variante (P).

-BĪLIS	freq	P	R₁	R₂
-VEL	161	0,889	8,050	2,824
-BLE	20	0,111	0,124	2,222

Tabela 16: Comparativa das ocorrências totais (freq), da probabilidade de ocorrência (P) e de duas ratios (R_1 e R_2) dos sufixos derivados de **-BĪLIS** (-VEL e -BLE) no corpus medieval da Galiza.

4.3.2. Distribuição dos sufixos no corpus de Portugal

As tabelas 17, 18 e 19 mostram os resultados obtidos a partir dos textos de Portugal. Como no caso de -CION, as variantes de -BLE são marginais ao respeito das de -VEL. Observamos que o número de variantes diferentes dos sufixos -VEL e -BLE são menores que nos textos da Galiza. Parece, portanto, que há um maior grau de estandarização nos textos de Portugal em relação a este tipo de adjetivos.

-VEL	freq	formas	P
-bel	3	3	0,008
-uel	120	44	0,335
-uil	94	16	0,262
-vel	118	54	0,329
-vil	23	12	0,064
Total	358	129	1

Tabela 17: Ocorrências (freq) das variantes de -VEL no corpus medieval de Portugal, formas diferentes e probabilidade de cada variante (P).

-BLE	freq	formas	P
ule	2	2	1
Total	2	2	1

Tabela 18: Ocorrências (freq) das variantes de -BLE no corpus medieval de Portugal, junto com o número de formas de palavras diferentes por variante (formas) e a probabilidade da variante (P).

-BĪLIS	freq	P	R₁	R₂
-VEL	358	0,994	179.0	2,0
-BLE	2	0,006	0,005	1.0

Tabela 19: Comparativa das ocorrências totais (freq), da probabilidade de ocorrência (P) e de duas ratios (R_1 e R_2) dos sufixos derivados de **-BĪLIS** (-VEL e -BLE) no corpus medieval de Portugal.

4.3.3. Indícios de castelhanização nos textos da Galiza

Como no caso do par anterior, foram realizadas experiências suplementares em busca de indícios que demostrem ou não que o uso do sufixo -BLE está relacionado com a castelhanização da língua medieval na Galiza.

A Figura 2 mostra uma maior proporção (ratio R_1) de variantes de -VEL sobre variantes de -BLE nos séculos XIII e XIV, frente ao século XV. Esta mesma tendência também se pode observar no trabalho de Mariño Paz (2005), que defende, para este par de sufixos, a hipótese da influência castelhana.

Por outro lado, não há maior número de ocorrências de -VEL frente a -BLE nos parágrafos identificados automaticamente como sendo escritos em castelhano. Como no caso do par anterior, encontramos a tendência inversa (embora com muito poucos casos e valores mais iguais), tal e como se mostra na Tabela 20.

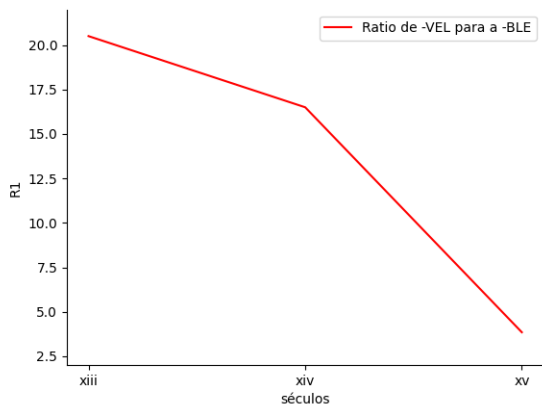


Figura 2: Evolução diacrónica ao longo de três séculos da ratio R_1 de -VEL para -BLE, no corpus medieval da Galiza.

-BĪLIS	freq	P	R_1	R_2
-VEL	3	0,429	0,750	1,00
-BLE	4	0,571	1,333	1,333

Tabela 20: Comparativa das ocorrências totais (freq), da probabilidade de ocorrência (P) e de duas ratios (R_1 e R_2) dos sufixos derivados de -BĪLIS (-VEL e -BLE) no corpus medieval da Galiza em excertos em castelhano.

4.4. Discussão

As experiências descritas parecem demonstrar que as vozes cultas das palavras com sufixos derivados de -TION e -BĪLIS, é dizer, -CION e -BLE, respetivamente, são muito minoritárias nos textos da Galiza e residuais nos de Portugal. Os testes suplementares mostram que a maior presença na Galiza destas formas cultas correlaciona com indicadores de castelhanização, o que contradiz a hipótese de Lorenzo (1985) e Mariño (1998) sobre o par -ÇON/-CION, que sustentam que se trata duma evolução interna da língua, alheia a influências externas. Aliás, os indícios de castelhanização são muito mais conclusivos no caso de -CION que no de -BLE, mesmo se no segundo caso é admitida a influência castelhana (Mariño Paz, 2005).

Por último, se houve desde a Idade Media um registo culto interno e próprio da língua medieval da Galiza que favoreceu as formas em -CION, então é difícil explicar como é possível que a forma “popular,” é dizer, patrimonial, apareça em muitas vozes de registo culto, tais como: *absoluçon, alegaçõn, anexaçõn, comutaçõn, disençõn, diçõn, encarnaçõn, incorporaçõn, interposiçõn, precisson, procuraçõn, restituyçõn, va-*

lidaçõn, etc. Os dados achados no corpus seem chamados a contradizer a principal hipótese filológica na que se baseia a norma moderna do galego para defender o uso de -CION frente a -ÇOM.

5. Avaliações das ferramentas

Nesta secção, o objetivo é avaliarmos as ferramentas PLN de duas maneiras diferentes:

Avaliação intrínseca: mediante o uso de corpus de teste anotados.

Avaliação extrínseca: mediante a comparação dos resultados obtidos automaticamente com as ferramentas e os obtidos mediante uma análise tradicional de corpus sem ferramentas PLN, só com expressões regulares e revisão manual. Para levar a cabo esta avaliação, os resultados da análise tradicional deram lugar à criação dum *gold standard* composto pelos dados de referência revisados manualmente.

5.1. Avaliação intrínseca

A Tabela 21 mostra os valores de exatidão (*accuracy*) obtidos na avaliação do reconhecedor de língua e do etiquetador de galego-português medieval.

O reconhecedor de língua foi avaliado a partir dum conjunto de 300 parágrafos extraídos de textos dos três séculos: 100 parágrafos do século XIII, 100 do XIV e 100 do XV. Foram extraídos aleatoriamente e alguns deles representam só uma pequena frase de poucas palavras, o que dificulta a identificação da língua. O sistema atingiu um 94,3% de exatidão (283 parágrafos identificados corretamente de 300). Os resultados foram avaliados manualmente. Devemos ter em conta que, apesar de não se poder fazer uma comparação direta, uma vez que os dados são diferentes, as experiências reportadas por Zampieri et al. obtiveram menos de 95% de exatidão para variantes e línguas similares.

	exatidão
Reconhecedor de língua	94,3%
Etiquetador morfossintático	87,4%

Tabela 21: Valores de exatidão do reconhecedor de línguas medievais e do etiquetador morfossintático do galego-português medieval.

Para avaliarmos o etiquetador, utilizamos como corpus de teste o documento *Pedro Rodriguez notario publico del Rey en Trasanços* datado

de 1257, que consta de 564 tokens. O valor de exatidão do sistema é de 87,4%, considerando só a etiqueta principal e sem tomar em conta o lema nem a informação morfológica. Portanto, a avaliação foi levada a cabo tomando em conta a primeira letra de cada etiqueta, pois é a que codifica a classe de palavra (nome, adjetivo, verbo, etc). O número total de classes de palavra do *tagset* é 11.

O valor de exatidão obtido é baixo, mesmo se comparado com outros etiquetadores adaptados às variedades históricas. No trabalho descrito em Sánchez-Marco et al. (2011) para o espanhol antigo, a melhor configuração do etiquetador atinge o 94,5% de exatidão. Em Rögnvaldsson & Helgadóttir (2008), para textos medievais do antigo nórdico (Old Norse), língua da que derivam o norueguês, islandês, danês e sueco, o sistema atinge um 92% de exatidão.

5.2. Avaliação extrínseca

Mesmo se os resultados do etiquetador não são muito encorajadores, na seguinte avaliação medimos a qualidade dos resultados obtidos automaticamente com o etiquetador, contrastando os erros e acertos com um conjunto de dados de referência revisados manualmente e criados a partir do corpus de textos da Galiza.

O conjunto de dados de referência é a lista correta de formas com os sufixos estudados (-ÇOM/-CION e -VEL/-BLE). Esta lista foi construída selecionando todos os tokens possíveis cuja terminação coincide com as variantes dos sufixos, gerando assim uma lista ampla de candidatos. Para gerar esta lista de candidatos, utilizamos um simples tokenizador e expressões regulares. Posteriormente, a lista de candidatos foi revisada manualmente e só os casos corretos foram selecionados. A partir desta revisão manual foram criados dois léxicos de frequências, um com todas as variantes terminadas em -COM/-CION e outro com as terminadas em -VEL/-BLE, que estão disponíveis para descarga livre.¹⁰

A Tabela 22 mostra os valores de precisão, abrangência (*recall*) e F1 obtidos comparando a lista de referência ou *gold standard* com as listas extraídas mediante o uso do etiquetador morfosintático. Estamos, de facto, a avaliar o desempenho da estratégia automática (sem revisão manual) de extração das formas com os sufixos objeto de estudo e descrita na Secção 4.

Nesta tarefa, o método automático baseado na etiquetagem atinge um F1 muito aceitável

(95,54%) e muito superior ao desempenho do etiquetador (87,4%). Isto é devido a que o etiquetador tende a ser mais preciso na etiquetagem de nomes e adjetivos, as etiquetas pertinentes na nossa tarefa filológica, que na etiquetagem doutro tipo de formas com significado gramatical. Estas formas tendem a ser mais ambíguas e frequentes, como é o caso dalguns determinantes, pronomes e conjunções, que as categorias lexicais.

Na última coluna da Tabela 22 mostramos o valor da correlação de *Pearson*, calculado tomando em conta a relação de valores da coluna de frequências da lista extraída automaticamente e a da lista de referência. A correlação média é muito alta: 0.9838. Repare-se que o valor para -VEL/-BLE é menor (mesmo tendo maior F1) porque a frequência de casos é muito mais baixa e este é um parâmetro determinante no cálculo da correlação.

5.3. Análise de erros

A seguir, após analisar os falsos positivos da última avaliação, apresentamos o tipo de erro mais comum que comete este método de extração automática:

– Erros de -ÇOM:

- Nomes próprios terminados em variantes de -ÇOM, que podem ir em maiúscula ou não. Por exemplo: *Jaason, iaason, ia-som*.
- Variantes do verbo ser, por exemplo: *sao, ssom, sson*.
- Variantes não esperadas dos termos *razom, coraçom* e *sazom*, com iode que desaparece também em castelhano, e que não foram considerados na contagem. Por exemplo: *rrason, rrazon, rrazóm*.

– Erros de -VEL:

- Nomes próprios terminados em variantes de -VEL: *Çentule*.
- Verbos em forma terminada em 'u', seguido pelo pronome *le* enclítico, por exemplo: *doule, deule, outorgoule*.

No tocante aos falsos negativos, o mais habitual é encontrarmos substantivos etiquetados como adjetivos, e viceversa. Como estes casos são mais frequentes que os falsos positivos, a abrangência do sistema é menor que a sua precisão.

¹⁰http://fegalaz.usc.es/~gamallo/resources/sufixos_medievais.zip

	prec.	abrang.	F1	ρ
-ÇOM/-CION	94,37%	92,54%	93,45%	0,9972
-VEL/-BLE	99,39	95,93	97,63	0,9705
Média	96,85	94,23	95,54	0,9838

Tabela 22: Precisão, abrangência e F1 dos resultados obtidos pelo sistema (etiquetador), junto com a correlação *Pearson* entre os resultados do sistema e o *gold standard*.

6. Conclusões

No presente artigo, descrevemos a adaptação, uso e avaliação de ferramentas de etiquetagem e de identificação de variedades medievais, cujo desempenho é inferior às ferramentas utilizadas para variedades modernas mais estandardizadas. As dificuldades para etiquetar textos medievais derivam principalmente da grande variabilidade de formas não estandardizadas.

Usamos estas ferramentas para um estudo filológico concreto cujo objetivo é verificar hipóteses linguísticas sobre a evolução de sufixos medievais. Conferimos que, pela sua distribuição de frequências, o uso dos sufixos -CION e -BLE nos textos da Galiza já está influenciado pelo castelhano baixo medieval. Os indícios de castelhanização são claros: a ratio de -ÇOM para -CION e de -VEL para -BLE desce consideravelmente nos textos do século XV, quando a influência do castelhano é mais evidente, e esta ratio inverte-se nos excertos escritos em castelhano ou mais castelhanizados. É pertinente considerar que o triunfo das soluções castelhanas dos dois sufixos, -CION e -BLE, no galego moderno viu-se favorecido pelo facto de se aplicarem, em geral, a palavras cultas, muitas delas herdadas diretamente do castelhano a partir do XV e séculos posteriores com a necessidade de incorporarmos na língua termos técnicos e cultismos.

No estudo quantitativo, observamos que as variantes dos sufixos patrimoniais, -ÇOM e -VEL, são claramente maioritárias tanto nos textos da Galiza como de Portugal. É também importante sublinhar que a produtividade (ou *prolificidade*) de -ÇOM é muito maior que a de -VEL, tanto nos textos medievais de Galiza como de Portugal. Também mostramos como a proporção do uso de -VEL na Galiza frente ao castelhanismo -BLE (7 vezes maior) é claramente menor que a de -ÇOM frente a -CION (20 vezes maior). Isto põe em questão a polémica decisão de recuperarmos para o galego moderno o uso de formas em -VEL mas não em -ÇOM.

Os resultados obtidos automaticamente (etiquetagem e identificação de língua) foram avaliados mediante a construção dum *gold standard*


com revisão manual. Verificamos que os resultados da análise automática estão muito próximos do *gold standard*, mostrando que a abordagem automática é fiável e pode ser utilizada para este ou outros estudos filológicos que precissem de informação quantitativa extraída de corpus.

Como trabalho futuro, desenharemos um método automático de normalização de variantes para a língua medieval galego-portuguesa tomando em conta regras fonológicas, como já se fez para uma parte do léxico do Inglês Antigo (Sáenz, 2015). O normalizador/lematizador assim construído permitirá alargar o modelo de língua do nosso etiquetador, com base no corpus anotado CIPM previamente enriquecido com os lemas automaticamente gerados. Finalmente, continuaremos desenvolvendo o modelo híbrido do etiquetador com a definição de mais regras morfológicas e de correção, que enriquecerão o classificador estatístico.

Agradecimentos

Este trabalho foi financiado pelo projeto NÓS, da Xunta de Galicia, pelos projetos DOMINO (PGC2018-102041-B-I00, MCIU/AEI/FEDER, UE) e eRisk (RTI2018-093336-B-C21), e também pela Consellería de Cultura, Educación e Ordenación Universitaria (acreditação 2016-2019, ED431G/08 e Programa de Formación Posdoctoral da Xunta de Galicia 2016) e European Regional Development Fund (ERDF).

Referências

- Canosa, Xavier, Pablo Gamallo, Xavier Canosa, José Ángel Taboada, Paulo Martínez Lema & Marcos Garcia. 2019. Uma utilidade para o reconhecimento de topónimos em documentos medievais. *Linguamática* 2(11). 3–15.  [10.21814/lm.11.1.291](https://doi.org/10.21814/lm.11.1.291).
- Cristine Prado, Natália & Gladis Massini-Cagliari. 2014. Formação de nomes deverbais nas cantigas de Santa Maria: Um estudo morfológico. *Revista Do GEL* 11(2). 71–96.

- Diéguez, Ignacio Vázquez. 2018. Sobre algúns sufixos galegos medievais. *Estudios de Lingüística del Español* 39. 241–277.
- Ferreiro, Manuel. 1997. *Gramática histórica da lingua galega. ii. lexicología*. Santiago de Compostela: Lailovento.
- Fillo, Machado & Américo Venâncio Lopes. 2013. *Dicionário etimológico do português arcaico: Projeto DEPARC*. Salvador: Edufba.
- Freixeiro Mato, Xosé Ramón. 1997. *Lingua galega: normalidade e conflito*. Santiago de Compostela: Lailovento.
- Gamallo, Pablo & Marcos Garcia. 2017. LinguaKit: uma ferramenta multilingue para a análise linguística e a extração de informação. *Linguamática* 9(1). 19–28. doi 10.21814/lm.9.1.243.
- Gamallo, Pablo, Marcos Garcia, Cesar Pineiro, Rodrigo Martínez-Castano & Juan C. Pichel. 2018. LinguaKit: a big data-based multilingual tool for linguistic analysis and information extraction. Em *5th International Conference on Social Networks Analysis, Management and Security (SNAMS)*, 239–244. doi 10.1109/SNAMS.2018.8554689.
- Gamallo, Pablo, Susana Sotelo & José Ramom Pichel. 2014. Comparing ranking-based and naive bayes approaches to language detection on tweets. Em *Workshop TweetLID: Twitter Language Identification Workshop at SEPLN 2014*, n/p.
- Garcia, Marcos & Pablo Gamallo. 2015. Yet another suite of multilingual NLP tools. Em *Languages, Applications and Technologies*, 65–75. doi 10.1007/978-3-319-27653-3_7.
- Kettunen, Kimmo. 2014. Can type-token ratio be used to show morphological complexity of languages? *Journal of Quantitative Linguistics* 21. 223–245. doi 10.1080/09296174.2014.911506.
- Leach, Geoffrey & Andrew Wilson. 1996. Recommendations for the morphosyntactic annotation of corpora. Em *Technical Report, Expert Advisory Group on Language Engineering Standard (EAGLES)*.
- Lorenzo, Ramón. 1985. *Crónica troiana. introducción e texto*. A Coruña: Fundación Pedro Barrié de la Maza, Conde de Fenosa.
- Mariño, Ramón. 1998. Notas sobre a historia das terminacións -ión / -ón en galego. Em D. Kremer (ed.), *Homenaxe a Ramón Lorenzo*, 735–760. Vigo, Galaxia, vol. 2.
- Mariño Paz, Ramón. 2005. Forma e función do sufixo -uel no galego medieval. *Cadernos de Lingua* 27. 155–193.
- Messner, Dieter. 2007. Os dicionários portugueses, devedores da lexicografia espanhola. *Península, Revista de Estudos Ibéricos* 4. 141–151.
- Padró, Lluís. 2012. Analizadores multilingües en FreeLing. *Linguamática* 3(2). 13–20.
- Pichel, José Ramom, Pablo Gamallo, Iñaki Alegria & Marco Neves. 2020. A methodology to measure the diachronic language distance between three languages based on perplexity. *Journal of Quantitative Linguistics* 28(4). 306–336. doi 10.1080/09296174.2020.1732177.
- Pichel, José Ramom, Pablo Gamallo & Inaki Alegria. 2019. Measuring diachronic language distance using perplexity: Application to english, portuguese, and spanish. *Natural Language Engineering* 26(4). 433–454. doi 10.1017/S1351324919000378.
- Rögnvaldsson, Eiríkur & Sigrún Helgadóttir. 2008. Morphological tagging of old norse texts and its use in studying syntactic variation and change. Em *2nd Workshop on Language Technology for Cultural Heritage Data*, 40–46.
- Sánchez-Marco, Cristina, Gemma Boleda & Lluís Padró. 2011. Extending the tool, or how to annotate historical language varieties. Em *5th ACL-HLT Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities*, 1–9.
- Santalha, Montero & José-Martinho. 2005. Documentos medievais galegos (3). *Agália* 81–82. 255–264.
- Silvestre, João Paulo. 2008. *Bluteau e a origens da lexicografia moderna*. Lisbon: Imprensa Nacional – Casa da Moeda: Coleção filologia portuguesa.
- Sáenz, Marta. 2015. The lemmatization of Old English verbs from the second weak class on a lexical database. *Journal of English Studies* 13. 135. doi 10.18172/jes.2861.
- Varela Barreiro, Xavier, Maria Francisca Xavier & Charlotte Galves. 2016. Corpus informatizado Galego-Português antigo. Instituto da Lingua Galega / Centro de Lingüística da Universidade Nova de Lisboa / Universidade de Campinas. <http://ilg.usc.gal/tmilg>.
- Venâncio, Fernando. 2019. *Assim nasceu uma língua. sobre as origens do português*. Lisbon: Guerra e Paz.

- Viaro, Mário Eduardo. 2012. A produtividade dos sufixos do ponto de vista diacrônico. Em T. Lobo, Z. Carneiro, J. Soleidade, A. Almeida & S. Ribeiro (eds.), *Rosae: linguística histórica, história das línguas e outras histórias*, 275–292. SciELO Books.
- Xavier, Maria Francisca. 2005. A caminho de um dicionário do português medieval. Em *Des(a)fiando discursos: Homenagem a Maria Emília Ricardo Marques*, 667–686. Lisboa: Universidade Aberta, Língua, Literatura e Cultura Portuguesas.
- Xavier, Maria Francisca. 2016. O CIPM — corpus informatizado do português medieval, fonte de um dicionário exaustivo. Em Carlota de Benito Moreno Johannes Kabatek (ed.), *Lingüística de corpus y lingüística histórica iberorrománica*, 137–156. De Gruyter.
- Zampieri, Marcos, Shervin Malmasi, Nikola Ljubešić, Preslav Nakov, Ahmed Ali, Jörg Tiedemann, Yves Scherrer & Noëmi Aepli. ????. Findings of the VarDial evaluation campaign 2017. Em *4th Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial)*, 1–15.

Appendices



A. Lista de obras do corpus galego

Nome	Tamanho
Cantigas de Afonso Eanes do Coton	7,1K
Cantiga de Afonso Gomes de Sarria	0,706K
Cantigas de Airas Nunes	5,9K
Arte de Trovar	11K
BMSEH (Tui século XV)	334K
Cancioneiro de Ajuda	77K
Cantigas de Amigo	20K
CDMO (Santa María de Oseira)	1,3M
CI (Crónica de Santa María de Iria - Rui Vasques)	79K
CSMp (Cantigas de Santa María: Pauta - Afonso X)	89K
CSMr (Cantigas de Santa María: Rúbricas - Afonso X)	45K
CT (Crónica troiana - Benoît de Saint-Maure)	1,2M
DAG (Documentos antigos de Galicia)	179K
FDUSC (Fontes documentais da Universidade de Santiago de Compostela)	1,2M
FR (fragmento do Foro Real - Afonso X)	4,7K
HGPg (Historia do galego-português - Galicia)	134K
HT (Historia troiana - Î de Saint-Maure)	528K
LCS (Libro do Concello de Santiago)	680K
LNAP (Libro de notas de Álvaro Pérez)	237K
MSCDR (San Clodio do Ribeiro)	1,2M
MSPT (San Salvador de Pedroso)	260K
OMOM (San Martiño de Vilourente)	1,1M

AIA-BDE: um Corpo de Perguntas, Variações e outras Anotações

AIA-BDE: a corpus of Portuguese Questions, Variations and other Annotations

Hugo Gonçalo Oliveira  
CISUC, DEI, Universidade de Coimbra

Ana Alves  
CISUC, Universidade de Coimbra
ISEC, Instituto Politécnico de Coimbra

Resumo

Apresentamos neste artigo o corpo AIA-BDE, que tem como principal objetivo a avaliação de sistemas que procuram associar necessidades de informação expressas em linguagem natural a perguntas com resposta conhecida (i.e., FAQ). Este corpo inclui várias perguntas no domínio da Administração Pública em Portugal e respetivas respostas. A 855 dessas perguntas foram adicionadas, manual e automaticamente, formas alternativas de as fazer, a que chamamos variações, e que podem ser utilizadas para simular interações de humanos. Essas perguntas encontram-se classificadas de acordo com a sua origem, com quatro valores possíveis, e têm ainda associado um tipo, atribuído com base na opinião de cinco anotadores. Para além de apresentar o AIA-BDE, ilustramos como pode ser utilizado através de três experiências, com resultados que podem ser vistos como base para melhorias futuras: associação de variações às respetivas perguntas; identificação automática da origem das variações; e classificação automática das perguntas quanto ao seu tipo.

Palavras chave

corpora, FAQs, resposta a perguntas, paráfrases, similaridade semântica, classificação de texto

Abstract

We present the AIA-BDE corpus, which has as main goal the evaluation of computational systems that attempt at assigning questions with known answers (i.e., FAQs) to information needs, expressed in natural language. This corpus includes several questions in the domain of the Portuguese Public Administration and their answers. To 855 of those questions, alternative ways of making them were manually and automatically added. We call them variations and they can be used in the simulation of human user interactions. Such questions are classified according to their source, with four possible values, and have also a question type, based on the opinion of five human annotators. Besides presenting AIA-

BDE, we illustrate how it can be used through three experiments, with results that might be seen as the baselines for future improvements, namely: variation assignment to the original questions; automatic automatic identification of the questions according to their source; and automatic classification of the questions according to their type.

Keywords

corpora, FAQs, question answering, paraphrases, semantic similarity, text classification

1. Introdução

O projeto *AIA: Apoio Inteligente a Empreendedores* decorreu entre os anos de 2018 e 2020 numa colaboração entre o Centro de Informática e Sistemas da Universidade de Coimbra (CISUC), o INESC-ID, e a Agência para a Modernização Administrativa (AMA), e teve como objetivo inicial o estudo de mecanismos de interação em linguagem natural para assistir automaticamente empreendedores em Portugal. Nesse contexto, foram exploradas diferentes técnicas para agentes artificiais com a capacidade de resposta automática a perguntas colocadas em português (Gonçalo Oliveira et al., 2019; Santos et al., 2020b,a). A maioria destes agentes usa como base listas de Perguntas Já Respondidas (em inglês, *Frequently Asked Questions*) e procura associar as perguntas do utilizador a perguntas conhecidas.

Para avaliar progressos, ao longo do projeto foi feita uma recolha de dados, focada em FAQ no domínio do problema, inicialmente a partir de conteúdos do Balcão do Empreendedor (BDE), desde 2019 integrado no portal e-Portugal¹. O resultado desta recolha de perguntas, da sua expansão e da sua organização, é o corpo AIA-BDE, cuja versão 2.1 apresentamos aqui.

Apesar desta versão incluir um conjunto adicional de perguntas e respostas, a principal con-

¹<https://eportugal.gov.pt>



tribuição do trabalho está num conjunto de 855 perguntas e respetivas respostas, obtidas a partir de quatro fontes distintas, e num total de 5.298 variações, algumas produzidas automaticamente, outras manualmente, por diferentes grupos de pessoas. As variações são outras formas de fazer as mesmas perguntas, utilizando outras palavras ou construções sintáticas alternativas, na maior parte dos casos paráfrases das perguntas originais. Considerando que um agente que dá respostas com base em perguntas conhecidas (i.e., FAQ) deve ter a capacidade de lidar com perguntas feitas de diferentes formas, a principal utilidade destas variações está na simulação de interações humanas e, conseqüentemente, na avaliação de agentes deste tipo.

Para além das variações, numa fase final do projeto, foi atribuída uma (ou mais) de nove classes a cada uma das 855 perguntas, de acordo com o seu tipo. Ainda que, até à data, esta anotação tenha sido pouco explorada, a identificação automática do tipo de pergunta pode ser vista como um problema de classificação, com utilidade para sistemas de resposta automática a perguntas, dado que esse tipo vai condicionar a resposta dada.

Este artigo descreve o corpo AIA-BDE, que acreditamos tratar-se de um contributo importante para o desenvolvimento e avaliação de sistemas de resposta automática a perguntas, deteção de paráfrases ou similaridade semântica textual, entre outros, na língua portuguesa. É também por isso que disponibilizamos o corpo a todos os potenciais interessados.

Antes de descrever o corpo, apresentamos algum trabalho relacionado, mais propriamente outros corpos com características semelhantes, para português e para outras línguas, bem como a sua utilização. Depois de descrever o conteúdo do AIA-BDE, relatamos três experiências realizadas com o corpo, nomeadamente na associação automática de variações às perguntas originais; classificação automática de perguntas de acordo com a sua origem; e classificação automática do tipo de pergunta. Os resultados obtidos podem ser vistos como base para experiências futuras com o corpo, onde outras técnicas podem vir a ser exploradas.

2. Trabalho Relacionado

Resposta Automática a Perguntas (RAP) é uma tarefa que tem como objetivo obter, de forma automática, respostas para perguntas colocadas em linguagem natural. Ainda que esse não seja um requisito, sistemas de RAP são normalmente

focados em perguntas factuais (Voorhees, 2008) e baseiam-se em técnicas de Recuperação de Informação (IR, do inglês *Information Retrieval*) (Kolomiyets & Moens, 2011) para encontrar a resposta em coleções de documentos.

O desenvolvimento de sistemas de RAP em português foi impulsionado pelo fórum de avaliação CLEF, onde esta língua esteve presente de 2004 (Magnini et al., 2004) a 2008 (Forner et al., 2008). Neste âmbito, destaca-se a criação da coleção CHAVE (Santos & Rocha, 2004), que inclui textos noticiosos publicados entre 1994 e 1995 e que foi usada para avaliações de IR e RAP (QA@CLEF).

Para avaliação de IR, foram definidos tópicos apropriados, considerando que podiam ser usados em pesquisas noutras línguas, e documentos foram marcados como relevantes ou não para cada tópico. A avaliação de RAP é semelhante, mas com perguntas factuais em vez de tópicos, incluindo a indicação do tipo de resposta esperada, e com respostas em vez de documentos. Ao longo das várias edições onde o português esteve incluído, foi criado um total de 4.380 perguntas em português e respetivas respostas na coleção CHAVE, quando essa resposta existe.

Outra avaliação relacionada foi o Págico (Mota et al., 2012), para a qual foi criada uma lista de 150 tópicos acerca da cultura lusófona e, para cada um, indicadas as páginas da Wikipédia que lhes respondiam, sendo que os tópicos correspondiam a perguntas cuja resposta, regra geral, não se encontrava em apenas uma página da Wikipédia.

A resposta automática a Perguntas Já Respondidas (doravante, FAQ) é uma tarefa específica no âmbito de IR e RAP, em que para cada pergunta há uma resposta bem definida. Devido à sua natureza e estrutura, listas de FAQ têm sido exploradas no desenvolvimento de sistemas de esclarecimento de dúvidas sobre um determinado domínio, em alguns casos culminando com a disponibilização das coleções usadas.

Um dos primeiros sistemas deste tipo, para a língua inglesa, terá sido o FAQFinder (Burke et al., 1997), que para além de uma abordagem baseada em IR, tira partido de conhecimento semântico na base de conhecimento lexical WordNet (Fellbaum, 1998). Para a língua croata, um sistema do mesmo tipo foi treinado e testado para, com base em técnicas de similaridade semântica textual (Agirre et al., 2012), fornecer respostas com base em 1.222 perguntas recolhidas a partir do sítio de uma operadora móvel na Web (Karan et al., 2013). As perguntas e respostas foram ainda enriquecidas por

10 voluntários que, sem conhecer à partida as perguntas recolhidas, criaram questões-exemplo que poderiam ser colocadas por utilizadores reais, assim como as suas paráfrases. O último passo consistiu na aplicação de modelos tradicionais de IR (e.g., pesquisa com base em palavras-chave, TF-IDF, modelação de linguagem) para, por cada questão colocada pelos voluntários, recuperar FAQ cuja relevância binária foi finalmente atribuída de forma manual.

A coleção FAQIR (Karan & Šnajder, 2016) tem 4.313 FAQ acerca de “manutenção e reparações”, obtidas a partir do portal *Yahoo! Answers*, e 1.233 perguntas de utilizadores acerca desse domínio, obtidas através do parafraseamento de 50 necessidades de informação base. A cada uma das perguntas anteriores foi posteriormente associado um de quatro valores de relevância para cada FAQ, com base em métodos de IR não supervisionados: Relevante (associação perfeita), Útil (tópico relacionado e útil), Inútil (tópico relacionado, mas sem utilidade), e Irrelevante (tópico não relacionado). Contudo, apesar do elevado número de FAQ, apenas uma pequena parte ($\approx 22\%$, 779) tem pelo menos uma pergunta considerada relevante ou útil.

Os mesmos autores criaram uma nova coleção de FAQ, agora no domínio das aplicações web (Karan & Šnajder, 2018), obtida a partir do portal StackExchange². A recolha focou-se nas perguntas mais populares neste domínio, correspondentes a 125 necessidades de informação. Ao associar cada uma dessas perguntas a diferentes respostas dadas por utilizadores e consideradas boas, de acordo com os votos dados, foi possível obter 719 FAQ. Ainda que as perguntas fossem constituídas por um título e uma descrição mais detalhada, foram usados apenas os títulos, após uma revisão manual que procurou garantir que continham toda a informação relevante. Depois desta recolha, dois anotadores produziram, cada um, cinco formas diferentes de exprimir a necessidade de informação associada a cada FAQ. Este processo resultou em 1.250 variações e à associação da relevância binária de cada forma para cada FAQ.

Ainda a este respeito, a QA4FAQ (Caputo et al., 2016) foi um tarefa de avaliação conjunta para resposta automática a perguntas baseada em FAQ escritas em italiano, organizada no contexto da avaliação EVALITA 2016. O objetivo consistia em recuperar FAQ relevantes para perguntas feitas por utilizadores. Neste contexto, foram disponibilizadas 406 FAQ (perguntas, respostas, etiquetas), 1.232 perguntas de utilizado-

res recolhidas a partir de registos de um sistema IR, e um conjunto de mapeamentos entre perguntas e FAQ.

Uma tarefa relacionada com a resposta a FAQ é a resposta a perguntas da comunidade (em inglês, *Community Question Answering*) (Nakov et al., 2015, 2016, 2017), onde o objetivo é ordenar pares de perguntas-perguntas e perguntas-comentários em fóruns de discussão na Web, de acordo com a sua similaridade. O fórum do portal *Qatar Living*³ tem sido utilizado como fonte de dados para esta tarefa. A partir de uma lista de perguntas originais, outras perguntas relacionadas são obtidas para cada, assim como os primeiros comentários nos seus tópicos. A relevância de cada pergunta relacionada para a pergunta original foi atribuída, e o mesmo foi feito para a relevância dos comentários em relação às perguntas originais e à pergunta original, embora diferentes anotações tenham sido utilizadas em diferentes sub-tarefas. As perguntas foram finalmente geradas a partir do assunto de cada pergunta original e o motor de busca Google foi usado para recolher 200 comentários de perguntas no sítio do fórum. Os resultados com dez ou mais comentários e perguntas com menos de 2.000 caracteres foram consideradas como perguntas relacionadas válidas.

No que concerne especificamente ao português, e tanto quanto conhecemos, o mais próximo com uma coleção de FAQ e respetivas respostas será a coleção MilkQA (Criscuolo et al., 2017), que inclui perguntas colocadas de forma mais densa, no domínio dos laticínios, seguidas pelas suas respostas. A disponibilização do corpo AIA-BDE é mais uma contribuição neste sentido, já que inclui FAQ em português, no domínio da administração pública, e está preparada para avaliar um conjunto de tarefas relevantes para sistemas de IR, RAP, e mesmo diálogo, com foco na associação de interações em linguagem natural com perguntas conhecidas.

Por se abordar normalmente como uma tarefa de IR, o desempenho de sistema de RAP é avaliado com recurso a medidas como a precisão, a abrangência e a medida-F. Já na recuperação de FAQ é normal assumir que existe uma e uma só resposta para cada pergunta (Burke et al., 1997), e, nesse caso, o desempenho é dado pela proporção de perguntas de uma lista para as quais foi dada a resposta correta (Burke et al., 1997; Caputo et al., 2016). Como a mesma pergunta pode ser feita de diferentes formas, o principal desafio é associar qualquer pergunta do utilizador a uma das perguntas conhecidas. Aqui é co-

²<http://www.stackexchange.com/>

³<https://www.qatarliving.com/>

num recorrer a técnicas de IR, podendo passar-se pelo reconhecimento automático de paráfrases ou cálculo de similaridade semântica textual. Para as últimas tarefas, foram também organizadas avaliações conjuntas em português (Fonseca et al., 2016; Real et al., 2020), que resultaram na disponibilização de coleções que incluem pares de frases e respetivos valores de similaridade e classificação de paráfrases.

Um outro tipo de RAP em que o interesse tem aumentado recentemente é a resposta extractiva a perguntas (em inglês, *Extractive Question Answering*). Este interesse tem sido impulsionado pela sua inclusão na avaliação de modelos de compreensão de linguagem natural (em inglês, *Machine Reading Comprehension*) (Devlin et al., 2019), recorrendo a coleções como a SQuAD (Rajpurkar et al., 2016), que inclui mais de 100 mil perguntas colocadas acerca de artigos da Wikipédia, e respetivas respostas. Apesar de originalmente em inglês, foram recentemente disponibilizadas duas traduções do SQuAD para português (Carvalho et al., 2021; Guillou, 2021).

Tal como acontece nas coleções para RAP e de FAQ, no SQuAD, a ordem das perguntas para o mesmo documento não é importante. Por outro lado, há corpos com uma estrutura semelhante ao SQuAD que procuram simular conversas, e por isso mais adequados, por exemplo, ao desenvolvimento e avaliação de sistemas de diálogo. Aqui destacamos o QuAC (Choi et al., 2018) e o CoQA (Reddy et al., 2019), ambos criados com base numa tarefa em que uma pessoa realiza perguntas acerca de um dado assunto e outra responde, tão naturalmente quanto possível, tendo por base um texto acerca do mesmo assunto. Já no âmbito dos sistemas de diálogo orientados à resolução de tarefas, há corpos como o MultiWOZ (Budzianowski et al., 2018) que podem ser utilizados para treino e avaliação. Este corpo tem vários diálogos com mais do que um turno, estruturados em atos de diálogo (em inglês, *Dialog Acts*) atribuídos manualmente, que incluem uma intenção (em inglês *intent*) e um conjunto de pares atributo-valor de acordo com a tarefa a realizar. Um exemplo será o ato de *informar*(domínio=hotel, preço=moderado), que representa a intenção de obter nomes de hotéis com um preço moderado.

Ainda no âmbito do diálogo, vários corpos têm sido utilizados no desenvolvimento de sistemas que seguem uma abordagem orientada a dados (em inglês, *data-driven*), mais propriamente, para aprenderem a traduzir interações do utilizador em respostas do sistema. Aqui destacam-se diálogos obtidos a partir de plataformas de con-

versação ou redes sociais, respetivamente compilados no Ubuntu Dialogue Corpus (Lowe et al., 2015) ou no Twitter Conversation Corpus (Ritter et al., 2011), ou diálogos obtidos a partir de legendas de filmes, compilados, por exemplo, nos corpos Open Subtitles (Lison & Tiedemann, 2016).

Ao contrário dos sistemas de recuperação de FAQ, a avaliação de sistemas de diálogo aberto (e.g., Vinyals & Le (2015)) é um desafio, e passa muitas vezes por comparar as respostas de um dado sistema com as respostas dadas por humanos, no mesmo contexto (Gunasekara et al., 2020). Na medida em que não existem objetivos claramente definidos para esses sistemas, o problema mantém-se mesmo quando são usadas métricas recentes, como a BertScore (Zhang et al., 2020), que considera representações semânticas e não se limita a comparar a sobreposição de sequências.

Uma alternativa mais demorada e dispendiosa passa por colocar utilizadores humanos a interagir com um dado sistema e avaliarem o quão natural e fluído decorreu o diálogo (Gunasekara et al., 2020). Quando estão a ser avaliados sistemas orientados à realização de tarefas, o sucesso na concretização da tarefa poderia ser acrescentado a esta avaliação (Wen et al., 2017). O recrutamento de utilizadores, em qualquer um dos casos, pode ser feito via *crowdsourcing*. Porém, para medir o progresso do sistema, a avaliação deveria ser realizada para cada nova atualização.

Especificamente para o português, os corpos de legendas de filmes têm sido usados por agentes conversacionais, não para responder a perguntas ou resolver tarefas em concreto, mas para melhor lidar com interações fora do domínio (Magarreiro et al., 2014). Considerando que, para decidir que tipo de resposta dar e como a obter é importante identificar o tipo de interação, relacionado com o ato de diálogo, para o português há também a coleção UC-PT (Fernandes et al., 2019), que inclui interações em diálogos com um conjunto de anotações: pergunta ou não; pergunta com resposta sim/não; pergunta de cariz pessoal.

3. Conteúdo do Corpo

Esta secção descreve o formato e o conteúdo do corpo AIA-BDE. Começa por explicar a estrutura dos ficheiros em que é distribuído, e depois apresenta e quantifica a origem das perguntas, as variações e respetiva produção e, finalmente, a atribuição do tipo de pergunta. Nas secções seguintes, ilustram-se algumas das possibilidades que o AIA-BDE oferece através de um conjunto de experiências realizadas.

S:Espaço Empresa
 ...
 P:Como pedir o Cartão Provisório de Identificação de Pessoa Coletiva?
 VG1:Como solicitar o cartão de identidade provisório?
 VG2:Como solicitar o cartão de identidade provisório?
 VUC:Como posso obter o cartão provisório de identificação de pessoa coletiva?
 VUC:Onde posso pedir o cartão provisório de pessoa coletiva?
 VIN:Como posso pedir o Cartão provisório de identificação de pessoa coletiva?
 VIN:Qual o procedimento para obter o Cartão Provisório de Identificação de Pessoa Coletiva?
 R:O Cartão Provisório de Identificação de Pessoa Coletiva deixou de ser emitido, (...) Atualmente, existe apenas o Cartão da Empresa e o Cartão de Pessoa Coletiva, que são emitidos para entidades definitivamente registadas ou inscritas.

Figura 1: Sequência de linhas do AIA-BDE 2.1

3.1. Estrutura e Organização

O corpo AIA-BDE é distribuído a partir do repositório <https://github.com/NLP-CISUC/AIA-BDE>, num ficheiro principal, `AIA-BDE.v2.1.txt`, e de outros ficheiros complementares. Todos os ficheiros são constituídos por perguntas e respetivas respostas. Mais propriamente, cada linha começa com um marcador que pode indicar se se trata de uma pergunta (P:), da resposta à pergunta anterior (R:), ou ainda a origem de todas as perguntas que se seguem (S:).

Especificamente no ficheiro principal, entre as 855 perguntas e as respetivas respostas, encontram-se variações da pergunta, cujos marcadores começam por V. A título de exemplo, a Figura 1 mostra uma sequência de linhas deste ficheiro. Na versão 2.1, este ficheiro inclui perguntas obtidas a partir de quatro fontes, mais propriamente:

- 625 perguntas do *Espaço Empresa* (EE), que cobrem informações relacionadas com o exercício de uma atividade económica e com o ciclo de vida de uma empresa;
- 118 do *Guia de Aplicação do Regime Jurídico de Acesso e Exercício de Atividades de Comércio, Serviços e Restauração* (RJACSR);
- 56 acerca da legislação de *Alojamento Local* (AL);
- 56 obtidas a partir de um conjunto de guias relacionados com *Apoios sociais* (AS), as últimas acrescentadas.

Ao nível das origens, o corpo tem uma distribuição altamente desequilibrada, com mais de três quartos das perguntas originárias do EE. De forma a não perder a informação relativa à origem mais específica destes documentos, foram introduzidos dois marcadores com a indicação da

zona (SS) e do ficheiro (SSS) onde estas perguntas se encontravam dentro do EE.

O ficheiro `AIA-BDE.tipo_pergunta.txt` tem exatamente as mesmas perguntas e respostas que o anterior, mas: (i) não inclui variações; e, (ii) depois de cada resposta, tem uma linha iniciada por um marcador F: seguida do nome do tipo ou tipos da pergunta anterior, separados por vírgulas, com a respetiva confiança associada a cada, separada por um # (e.g., `Condição#0.6,Procedimento#0.4`).

Para além dos dois ficheiros anteriores, há um conjunto de ficheiros com mais perguntas, respostas e, em alguns casos, variações, todas do mesmo tipo e produzidas automaticamente a partir de documentos estruturados. Estas perguntas foram geradas no âmbito do projeto AIA e, mesmo nos casos em que existem variações, elas seguem todas uma estrutura padrão. A sua geração teve como principal objetivo uma rápida ampliação do número de perguntas a que os agentes desenvolvidos poderiam responder. No entanto, devido à sua simplicidade e falta de revisão manual, a sua utilidade para fins de avaliação é limitada, e não foram por isso utilizados em nenhuma das experiências descritas neste artigo.

Entre os ficheiros anteriores destacamos o `AIA_atividades.txt` e o `AIA_licencas.txt`. O primeiro inclui 844 perguntas que usam uma de três formas para perguntar o que é determinada atividade económica, as outras duas formas como variação (marcada por VAU), e como resposta uma definição da atividade. O segundo ficheiro inclui 1.281 perguntas de um de dois tipos: (i) o que permite determinada licença; ou (ii) licença necessária para fazer qualquer coisa. Cada uma também pode ser feita de uma de duas formas, sendo a outra usada como variação. A Figura 2 ilustra estes dois ficheiros com uma seleção de linhas de cada um.

P:O que é um Café?
 VAU:O que faz um Café?
 VAU:Para que serve um Café?
 R:Estabelecimentos de bebidas que servem, através de pagamento, bebidas e cafetaria ...

P:De que preciso para circular, parar e estacionar veículos de tração animal?
 VAU:O que permite circular, parar e estacionar veículos de tração animal?
 R:Veículo de tração animal - licença de circulação
 P:O que permite a licença Grua - licença de ocupação do espaço público ?
 VAU:O que posso fazer com uma Grua - licença de ocupação do espaço público ?
 R:Permite a instalação de uma grua (aparelho para levantar e deslocar corpos pesados), ...

Figura 2: Linhas dos ficheiros AIA_actividades.txt (cima) e AIA_licencas.txt (baixo).

3.2. Perguntas e Variações

Para cada uma das 855 perguntas no ficheiro principal do AIA-BDE foram produzidas variações, isto é, reformulações da pergunta original utilizando outras palavras ou construções, mas mantendo o significado original ou um suficientemente próximo, de tal forma que a resposta original continue a ser válida.

A sua produção teve em conta que, na maioria dos casos, os utilizadores não escrevem uma pergunta exatamente da mesma forma que ela se encontra numa lista de perguntas já respondidas (FAQ). Isto implica que, para ter sucesso, um sistema computacional que procure encontrar respostas com base numa lista de perguntas, terá de conseguir associar perguntas com base na sua proximidade semântica, ainda que feitas por outras palavras. Ou seja, terá de lidar com fenómenos como a sinonímia e o relacionamento semântico ou, ao nível da frase, similaridade semântica textual (Agirre et al., 2012), parafraseamento e inferência (Bowman et al., 2015). Assim, as variações têm como objetivo principal permitir a avaliação de sistemas de resposta automática a perguntas com resposta conhecida, ou simplesmente de sistemas focados nas tarefas anteriores, mas em contexto interrogativo.

Ainda que produzidas de diferentes formas, para cada pergunta do AIA-BDE há pelo menos cinco variações. Por não haver uma forma ideal de produzir variações, e porque a sua criação manual é um processo moroso, as variações foram sendo produzidas ao longo do tempo, por diferentes pessoas, e seguindo abordagens diferentes. Assim, optamos por marcá-las consoante o processo de criação, em alguns casos automático e noutros manual. Para as primeiras, utilizamos a API do Google Translate API⁴ como uma abordagem rápida e de baixo custo para gerar paráfrases da pergunta original. Mais propriamente, cada pergunta tem duas variações

deste tipo: tradução do texto em português para inglês e novamente para português (VG1), e tradução do resultado anterior novamente para inglês e para português (VG2). Dada a simplicidade da abordagem, algumas das variações acabam por ser muito próximas, ou até iguais, à pergunta original. Mais propriamente, há 51 variações VG1 e 41 VG2 iguais à pergunta original. Numa minoria de casos, a semântica acaba mesmo por sofrer alterações, devido a problemas na tradução e conseqüente introdução de termos incorretos.

Devido às limitações da abordagem anterior, foram também produzidas variações de forma manual. Neste caso, por terem sido criadas por diferentes grupos de pessoas, optamos por separá-las em três tipos:

- Variações produzidas pela equipa de investigadores do projeto AIA na Universidade de Coimbra (VUC);
- Variações produzidas por alunos da unidade curricular de Língua Natural, leccionada em mestrados do Instituto Superior Técnico (VIN);
- Variações produzidas com recurso à plataforma de *crowdsourcing* Amazon Mechanical Turk⁵ (VMT).

Para qualquer um dos tipos, foi pedido aos voluntários que lessem tanto a pergunta original como a sua resposta, e para reescrever a pergunta usando outras palavras, ainda que mantendo o significado original ou, pelo menos, um significado próximo ou implicado pelo original.

A Tabela 1 contabiliza, para cada origem, o número de perguntas e variações de cada tipo disponível. Apresenta ainda o número médio de átomos nas perguntas de cada origem (**Comprimento**), que permite verificar que as perguntas de AL são, regra geral, mais longas,

⁴<https://cloud.google.com/translate/docs/>

⁵<https://www.mturk.com/>

Origem	Perguntas	Comprimento	Variações					Total
			VG1	VG2	VUC	VIN	VMT	
EE	625	11,6±5,8	625	625	430	2.279	0	4.584
RJACSR	118	14,6±9,1	118	118	380	0	0	734
AL	56	20,1±12,5	56	56	122	0	0	290
AS	56	12,5±5,0	56	56	0	0	168	336
Total	855	12,6±7,3	855	855	932	2.279	168	5.944

Tabela 1: Distribuição de perguntas e variações por origem.

e as do EE mais curtas. Para cada pergunta original há uma variação do tipo VG1 e VG2, no entanto, a existência dos restantes tipos é variável, perfazendo pouco mais de 5.000 variações. O tipo mais predominante são as variações VIN, no entanto, estas foram realizadas apenas para perguntas do EE. Por outro lado, as variações VMT foram produzidas apenas para as perguntas AS, mais propriamente, três para cada pergunta, e não há mais nenhum tipo de variação manual para estas perguntas.

A Tabela 2 ilustra o conteúdo do corpo AIA-BDE com uma pergunta para cada origem, seguida de um conjunto de variações e, finalmente, da resposta. No primeiro exemplo, do EE, as variações VIN aparentam ser mais “conservadoras” do que as VUC, o que se pode ver neste e nos restantes exemplos. Isto não acontece apenas aqui e é uma das razões para termos decidido marcar as variações de acordo com a forma de criação. Podemos considerar que na criação das VUC terá havido mais criatividade, com maior utilização de sinónimos e variações a omitir partes da pergunta original. Finalmente, e apesar das guias serem as mesmas, nas VMT houve um menor controlo no processo de criação, o que torna as suas diferenças para a pergunta original mais variáveis. Veja-se o último exemplo da tabela.

3.3. Tipo de pergunta

A última anotação adicionada às perguntas do corpo AIA-BDE foi o tipo de pergunta. Apesar de não ter sido explorada no contexto do projeto AIA, a sua identificação automática pode ser útil para um sistema computacional saber como responder ou onde procurar pela resposta. Num sistema de diálogo, o tipo de pergunta estará relacionado com os atos do diálogo e, por essa razão, ser útil para o seu reconhecimento.

Depois de olhar para as várias perguntas do AIA-BDE, foram definidas nove tipos em que as perguntas poderiam ser classificadas, todas elas com exemplos identificados. Esses tipos foram:

Binário, Condição, Custo, Definição, Local, Pré-requisito, Procedimento, Tempo, Vantagem.

De forma a agilizar o processo de anotação dos tipos e a considerar mais do que uma opinião nesta escolha, optámos por recorrer à plataforma de *crowdsourcing* Amazon Mechanical Turk. Mais propriamente, o tipo de cada pergunta foi atribuído de forma independente por cinco trabalhadores.

Antes de realizar a tarefa, foram apresentadas diretivas onde se descrevia cada um dos tipos e se incluía um exemplo para cada um. A combinação das cinco anotações permitiu calcular a confiança relativamente à adequação de cada tipo a cada pergunta. Mais propriamente, para uma pergunta, a confiança num tipo resulta da divisão do número de vezes que esse tipo foi atribuído pelo total de anotações obtidas (5).

A Tabela 3 apresenta a distribuição de perguntas por tipo e confiança. No ficheiro disponibilizado, optou-se por omitir os tipos com menor confiança nos casos em que o tipo com maior confiança tinha sido atribuído por mais de dois trabalhadores. Ou seja, quando há um tipo atribuído por três ou quatro trabalhadores e os outros por apenas um, apenas se apresenta o primeiro. Verifica-se que há tipos frequentemente atribuídos, nomeadamente a Definição e o Procedimento, enquanto outros são mais raros. Por exemplo, não há qualquer pergunta dos tipos Binário e Local com confiança 100%.

Na Tabela 4 apresentamos exemplos de perguntas do EE com os seus tipos e confianças calculadas. Tal como em alguns exemplos da tabela, há várias perguntas que acabam por ter mais do que um tipo e, na maior parte das vezes, estão ambos corretos e nem fará muito sentido escolher apenas um. Quando definimos os tipos preocupamo-nos mais em abranger todas as perguntas do que propriamente garantir que os tipos eram mutuamente exclusivos. A nossa opção por recorrer a cinco trabalhadores por pergunta e de incluir a confiança também está relacionada com esta situação. E assim, potenciais interessados poderão utilizar essa informação da forma que

Origem	Var	Texto
EE	P	<i>É necessário submeter documentos para efectivar o Registo por Depósito?</i>
	VG1	<i>É necessário enviar documentos para efetuar o registo por depósito?</i>
	VG2	<i>Preciso enviar documentos para registrar por depósito?</i>
	VIN	<i>Para efectivar o Registo por Depósito é necessário submeter documentos?</i>
	VIN	<i>Tenho que submeter documentos para efectivar o Registo por Depósito?</i>
	VUC	<i>De que documentos preciso para realizar um registo por depósito?</i>
	VUC	<i>Que documentos é necessário enviar para fazer um registo de depósito?</i>
	R	<i>Sim, deverá submeter os documentos que titulem o acto requerido.</i>
RJACSR	P	<i>Qual a coima aplicável às contraordenações graves?</i>
	VG1	<i>Qual é a multa aplicável à falta grave?</i>
	VG2	<i>Qual é a multa aplicável à falta grave?</i>
	VUC	<i>coima para contraordenação grave</i>
	VUC	<i>Qual o valor da multa para contraordenações graves?</i>
	R	<i>As contraordenações graves são sancionáveis com coima: ...</i>
AL	P	<i>No alojamento local é obrigatória a certificação energética? Em que termos deve ser efetuada?</i>
	VG1	<i>No alojamento local é obrigatório a certificação energética? Em que condições deveria ser feito?</i>
	VG2	<i>A certificação energética é necessária em alojamento local? Em que condições deve ser feito?</i>
	VUC	<i>Como deve ser feita certificação energética do meu alojamento local?</i>
	VUC	<i>Qual o procedimento para certificar energeticamente o meu alojamento local?</i>
	R	<i>De acordo com esclarecimento da DGEG (Direção-Geral de Energia e Geologia) ...</i>
AS	P	<i>Quando é que me dão uma resposta sobre o apoio social a crianças e jovens?</i>
	VG1	<i>Quando recebo uma resposta sobre apoio social para crianças e jovens?</i>
	VG2	<i>Quando recebo uma resposta sobre apoio social a crianças e jovens?</i>
	VMT	<i>Quando receberei a resposta do apoio social a crianças e jovens?</i>
	VMT	<i>Tenho que esperar muito para ter uma resposta sobre o apoio social a crianças e jovens?</i>
	R	<i>Depois de fazer a sua inscrição na instituição que lhe interessa, pode acontecer ter de ficar em lista de espera...</i>

Tabela 2: Exemplos de perguntas, variações e respostas no AIA-BDE.

Tipo	Confiança			
	100%	80%	60%	40%
Binário	0	8	35	66
Condição	1	12	45	124
Custo	15	11	5	8
Definição	115	64	38	58
Local	0	4	3	11
Pré-requisito	2	10	24	58
Procedimento	40	48	59	96
Tempo	23	14	9	15
Vantagem	2	9	8	5

Tabela 3: Distribuição de perguntas por tipo.

4. Associação de variações a perguntas

Esta e as próximas secções ilustram algumas das experiências que o AIA-BDE permite realizar. A primeira foca-se nas variações e, em trabalhos anteriores (Burke et al., 1997; Karan et al., 2013; Karan & Šnajder, 2018), foi apelidada de Recuperação de FAQ (*FAQ retrieval*). O objetivo passa por associar perguntas feitas de diferentes formas a perguntas conhecidas e já com uma resposta conhecida. Isto permite simular a capacidade de um sistema identificar paráfrases ou calcular a similaridade semântica entre frases interrogativas, um cenário com aplicabilidade, por exemplo, nos sistemas de diálogo ou de resposta a perguntas. Se a base de conhecimento de um sistema deste tipo for uma lista de FAQ, ao associar a pergunta recebida a uma conhecida, ele pode imediatamente retornar a sua resposta.

lhes for mais conveniente.

Pergunta	Tipo
<i>O que é o Cartão da Empresa e o Cartão de Pessoa Coletiva?</i>	Definição#0.8
<i>Para a constituição de uma empresa através do serviço de criação de Empresa Online é necessária a presença de todos os sócios?</i>	Procedimento#0.4, Pré-requisito#0.4
<i>Como pedir o Cartão Provisório de Identificação de Pessoa Coletiva?</i>	Procedimento#1
<i>O Cartão de Identificação de Pessoa Coletiva ou entidade equiparada, emitido pelo RNPC e de que sou titular, continua a ser válido?</i>	Procedimento#0.4, Binário#0.4
<i>Quando é possível o levantamento do capital social da Empresa Online?</i>	Procedimento#0.4, Tempo#0.4
<i>Onde posso adquirir um certificado digital qualificado?</i>	Procedimento#0.4, Local#0.6
<i>O que acontece se o trabalhador adoecer durante as férias?</i>	Condição#0.6
<i>Quais as vantagens de aderir a um centro de arbitragem?</i>	Vantagem#1
<i>Onde posso pedir uma certidão permanente?</i>	Local#0.8
<i>Qual o custo do Cartão da Empresa e do Cartão de Pessoa Coletiva?</i>	Custo#0.8
<i>O Cartão da Empresa e o Cartão de Pessoa Coletiva podem ser cancelados?</i>	Binário#0.8

Tabela 4: Seleção de perguntas do Espaço Empresa, respetivo tipo e confiança.

O que acabamos por avaliar com esta primeira experiência é um conjunto de métodos para o cálculo da similaridade semântica. Mais propriamente, para cada variação, utilizamos cada um dos métodos para calcular a similaridade com cada uma das 855 perguntas, e analisamos quantas vezes ele atribuiu a maior similaridade à pergunta original. Como esta análise pode ser limitada, olhamos ainda para o número de vezes em que a pergunta original está nas três (top-3) e nas cinco (top-5) mais similares. Isto porque, principalmente em casos de dúvida, pode ser mais útil fornecer três, ou até cinco respostas, em que uma delas é a pretendida, do que não fornecer nenhuma ou fornecer uma que não a desejada.

Foram testados diferentes métodos, todos eles não-supervisionados. Esta opção prende-se, por um lado, com a escassez de dados para treino, e, por outro, tem em vista a flexibilidade dos métodos e sua aplicabilidade a diferentes quantidades de dados, com diferentes origens (e.g., FAQ noutros domínios). A principal diferença entre os vários métodos testados é a forma adoptada para representar os textos. Aqui incluímos métodos mais simples, baseados em técnicas tradicionais de IR, e outros baseados na representação de texto com recurso a modelos distribucionais, nomeadamente *word embeddings* e *sentence embeddings*. Mais precisamente, testamos:

- Um método baseado numa biblioteca de IR para Python, Whoosh⁶, com a configuração base (Whoosh base), ou com a análise de

stems (StemmingAnalyzer) e tratamento de acentos (Charset Filter) ativados para português (Whoosh+). Para ambas as configurações, cada pergunta do AIA-BDE foi representada por um documento com dois campos —pergunta e resposta, com a pesquisa feita apenas no primeiro. O parâmetro *group* usou o valor *orGroup*, para evitar que todos os termos da pergunta fossem obrigatórios, e o método de ordenamento utilizado foi o BM25F.

- Métodos baseados na representação das perguntas através de *word embeddings* estáticos, nomeadamente word2vec CBOW (CBOW), GloVe e FastText. Mais propriamente, considerou-se que cada pergunta seria representada pelo vetor médio dos vetores dos seus átomos no modelo de *embeddings*. Para cada modelo, considerou-se também uma variação com uma média pesada, usando como peso o valor do TF-IDF⁷ de cada átomo na pergunta em relação ao corpo constituído por todas as perguntas originais do AIA-BDE. Foram usados modelos com vetores de 300 dimensões, pré-treinados para Português, obtidos a partir do repositório do NILC (Hartmann et al., 2017) (CBOW e GloVe) e do fastText⁸ (FastText).
- Métodos baseados em *sentence embeddings* obtidos a partir de modelos de linguagem neurais BERT (Devlin et al., 2019). Mais pro-

⁷Term Frequency - Inverse Document Frequency

⁸<https://fasttext.cc/>

⁶<https://whoosh.readthedocs.io/>

priamente, cada pergunta foi codificada com recurso a um modelo BERT, carregado e disponibilizado localmente a partir da plataforma `bert-as-a-service`⁹, com todos os parâmetros por defeito, à excepção do parâmetro *maximum length of sequences*, usado com o valor NONE para que o tamanho máximo das sequências fosse igual ao tamanho da pergunta mais longa no AIA-BDE. Foram testados dois modelos BERT pré-treinados, nomeadamente: BERT-Base, Multilingual Cased (Multi-BERT), treinado e disponibilizado pelos criadores do BERT¹⁰ para 104 línguas, que permite representar o texto em vetores de dimensão 768; BERTimbau-large-portuguese-cased (Souza et al., 2020), especificamente para português, que codifica o texto em vetores de 1.024 dimensões.

Para cada um dos métodos anteriores, e para cada tipo de variação, a Tabela 5 mostra a proporção de variações corretamente associadas às perguntas originais. Confirma-se uma dificuldade variável, dependendo do tipo de variação, o que suporta a nossa opção por identificar esse tipo. Como seria de esperar, é mais fácil associar as variações geradas automaticamente (VG1 e VG2) à pergunta original, o que fica claro com a observação de que todos os métodos têm um melhor desempenho nessas duas variações. Relembramos que estas variações são geradas com recurso ao Google Translate e é normal apresentarem poucas alterações relativamente à forma das perguntas originais. Assim, não é de admirar que vários métodos apresentem uma taxa de acerto superior a 85% para a primeira resposta, e 90% considerando a presença no top-5. Ainda que inferior, o melhor desempenho nas variações VIN é claramente superior ao desempenho nas VMT e, principalmente, nas VUC. Isto sugere que as variações VIN são mais fáceis de associar automaticamente à pergunta original. Por outro lado, as variações VUC parecem ser as mais desafiantes, o que estará relacionado com o seu maior nível de criatividade, já referido na Secção 3.2.

Também não é fácil de identificar o melhor método para esta tarefa, já que este varia para diferentes tipos de variação. Por exemplo, nas variações geradas automaticamente, o BERT multilingue tem um dos melhores desempenhos, apenas batido pelo word2vec-CBOW no top-3 e top-5 das variações VG1. No entanto, nas restantes variações há vários métodos com um desempenho superior. Nas variações VIN, o BERTimbau tem o melhor desempenho (83.5%

para a primeira), mas não muito superior ao word2vec-CBOW (82.2%) ou a um método mais simples, como o Whoosh+ (82.3%). Nas VUC, apesar do desempenho médio ser inferior, os três melhores métodos são os mesmos que para as VIN, ainda que numa ordem diferente, nomeadamente Whoosh+ (62.5%), BERTimbau (60.4%), word2vec-CBOW (59.1%). Finalmente, para as VMT, o cenário é um pouco diferente, com o word2vec-CBOW a ser claramente o melhor (70.8%), o que mostra que este talvez seja o método mais equilibrado, já que o desempenho do BERTimbau nestas variações foi bastante inferior (47.6%).

Por se tratar de um modelo recente, responsável por avanços significativos em várias tarefas do Processamento de Linguagem Natural, seria expectável que o melhor desempenho fosse alcançado pelos modelos BERT, o que não acontece. Contudo, chamamos a atenção de que estes modelos são mais complexos que os restantes e podem fornecer diferentes representações para o texto, considerando, por exemplo, a representação em diferentes camadas ou a sua combinação. Para além de não termos explorado todas essas possibilidades, utilizamos as versões pré-treinadas destes modelos, que não afinamos (em inglês, *fine-tuned*) para a tarefa em questão. Apesar de não termos uma quantidade suficiente de dados do domínio alvo para este fim, uma possibilidade seria afinar os modelos para calcular a similaridade semântica textual em português, tal como alguns investigadores já fizeram (Rodrigues et al., 2020b) (Rodrigues et al., 2020a).

5. Classificação da Origem das Variações

Como referido na Secção 3.1, as perguntas do AIA-BDE foram obtidas a partir de quatro fontes principais: Espaço Empresa (EE), Regime Jurídico de Acesso e Exercício de Atividades de Comércio, Serviços e Restauração (RJACSR), Alojamento Local (AL) e Apoios Sociais (AS). A origem das perguntas pode ser vista como o seu assunto de alto nível, isto é, todas as perguntas com a mesma origem estarão relacionadas e serão acerca do mesmo tipo de serviços. Só por si, esta informação pode ser importante porque, em alguns casos, um agente de pesquisa ou resposta a perguntas pode começar por identificar o assunto ou domínio da pergunta do utilizador, e assim diminuir o conjunto onde procurar a resposta. No limite, tal agente poderá nem encontrar uma resposta adequada, mas pelo menos

⁹<https://github.com/hanxiao/bert-as-service>

¹⁰<https://github.com/google-research/bert>

Método	Variações														
	VG1 (855)			VG2 (855)			VIN (2,279)			VUC (816)			VMT (168)		
	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5
Whoosh base	83.2	90.9	93.1	80.2	88.3	90.5	73.8	85.7	88.2	50.9	65.6	69.5	59.5	73.8	79.2
Whoosh+	88.1	95.6	96.6	85.4	93.7	95.6	82.3	91.4	94.0	62.5	78.9	83.0	54.8	70.2	81.5
CBOw	88.4	95.8	97.1	86.5	94.7	96.4	82.2	91.0	93.3	59.1	75.5	79.8	70.8	86.3	89.9
FastText	57.2	68.5	72.6	51.5	63.4	68.0	40.0	51.0	55.2	21.0	28.8	34.2	39.9	57.7	64.9
GloVe	85.1	92.0	93.9	82.5	90.3	92.2	70.0	79.2	81.7	46.3	58.0	64.6	63.7	79.8	83.3
Multi-BERT	90.6	95.3	96.6	90.6	96.0	97.5	73.7	84.0	87.1	46.4	59.6	65.8	39.9	53.0	56.5
BERTimbau	86.1	94.9	96.1	83.6	93.1	94.5	83.5	92.4	94.1	60.4	75.2	79.4	47.6	57.1	62.5

Tabela 5: Proporção de variações de diferentes tipos corretamente mapeadas com as perguntas originais (Top1), nas Top-3 e nas Top-5 mais similares, utilizando diferentes representações vetoriais.

conseguir indicar ao utilizador uma página web ou uma lista de perguntas sobre o assunto identificado.

De forma a simular este processo, outra experiência realizada com o AIA-BDE procurou, de forma automática, identificar a fonte das variações. Para tal, adoptamos uma abordagem de Aprendizagem Automática Supervisionada (em inglês, *Supervised Machine Learning*), em que classificadores foram treinados com as perguntas originais, e testados na identificação da origem das variações. Entre os classificadores testados, destacou-se um classificador baseado numa *Support Vector Machine* (SVM) linear, que utilizamos com duas representações diferentes do texto:

- Vetores TF-IDF com dimensão máxima 750 (i.e., comparável com a dos vetores do BERT-base), utilizando como *features* os átomos que ocorriam em pelo menos duas perguntas e, de forma a não considerar palavras demasiado frequentes, no máximo, em 50% de todas as perguntas originais.
- Vetores de dimensão 768, resultantes da codificação pelo modelo pré-treinado BERTimbau-base (Souza et al., 2020)¹¹, desta vez recorrendo à biblioteca `transformers`, da HuggingFace¹², e à pipeline *feature-extraction*.

A experiência foi realizada com recurso à biblioteca Python `scikit-learn` (Pedregosa et al., 2011) e respetivas implementações dos classificadores (i.e., LinearSVC com os parâmetros por omissão, para o classificador baseado em SVM), vetorização TF-IDF (`TfidfVectorizer`), e cálculo de métricas. Ao optar pelo BERT, os vetores deste último eram utilizados em alternativa aos vetores TF-IDF.

As Tabelas 6 e 7 apresentam o desempenho nesta experiência usando os vetores TF-IDF e BERT. Em cada uma, incluímos a precisão (P),

abrangência (A) e medida-F1 (F1) do modelo treinado com todas as perguntas e testado com as variações de cada tipo. Estas métricas são apresentadas para cada classe (i.e., origem) e também para o total, através de uma macro-média (Macro- μ), onde a proporção de instâncias de cada classe não é considerada, e de uma média pesada (μ -pesada). Para as variações que não cobrem alguma classe, as classes em falta não foram consideradas no cálculo das médias.

Uma vez mais, verifica-se que, dependendo da variação, pode ser mais ou menos difícil identificar a origem. Com a representação mais tradicional, TF-IDF, as variações onde o classificador teve mais dificuldades foram as VUC, o que volta a sugerir que são aquelas que mais se desviam das perguntas originais. Foi também para estas variações que o desempenho foi claramente superior com as representações baseadas no BERT, o que mostra a capacidade deste modelo lidar com diferenças lexicais em textos com o mesmo significado. Ainda que o desempenho para as variações VMT seja superior ao desempenho para as VUC, o melhor desempenho para ambas é obtido com os vetores baseados no BERT. Como referido anteriormente, e passando a redundância, estas são nada mais nada menos do que as variações onde há maior variação lexical relativamente à pergunta original. Por outro lado, para os outros três tipos de variação, a utilização do BERT parece não trazer grandes benefícios, apresentando mesmo um desempenho ligeiramente inferior.

Importa ainda destacar o desempenho superior na classificação de variações VIN, que, utilizando os vetores TF-IDF, atinge uma medida-F1 próxima do 100%, quando para as VG1 e VG2 este valor se situa em torno dos 90%. Acreditamos que isto também seja uma consequência do método adoptado para a geração automática de variações. Mais propriamente, a uma pequena proporção de resultados que, devido a problemas de tradução, nomeadamente do nome de serviços, introduzem termos fora do esperado e, conseqüentemente, fora do domínio de cada classe.

¹¹Neste caso optámos pela versão do modelo *base*, mais simples, porque, com a versão *large*, a SVM tinha dificuldade em convergir, e por isso a melhorar os resultados.

¹²<https://huggingface.co/transformers/>

Origem	VG1			VG2			VUC			VIN			VMT		
	P	A	F1	P	A	F1	P	A	F1	P	A	F1	P	A	F1
EE	91%	99%	95%	91%	99%	95%	65%	99%	78%	100%	98%	99%	N/A	N/A	N/A
RJACSR	92%	67%	77%	90%	69%	78%	96%	49%	65%	N/A	N/A	N/A	N/A	N/A	N/A
AL	93%	75%	83%	98%	71%	82%	99%	69%	81%	N/A	N/A	N/A	N/A	N/A	N/A
AS	100%	86%	92%	100%	84%	91%	N/A	N/A	N/A	N/A	N/A	N/A	100%	70%	83%
Macro- μ	94%	82%	87%	95%	81%	87%	87%	72%	75%	100%	98%	99%	100%	70%	83%
μ -pesada	92%	92%	92%	92%	92%	92%	82%	74%	73%	100%	98%	99%	100%	70%	83%

Tabela 6: Classificação da origem das variações, com SVM e representação TF-IDF.

Origem	VG1			VG2			VUC			VIN			VMT		
	P	A	F1	P	A	F1	P	A	F1	P	A	F1	P	A	F1
EE	93%	92%	93%	92%	95%	94%	87%	93%	89%	100%	91%	95%	N/A	N/A	N/A
RJACSR	61%	68%	64%	70%	63%	66%	87%	85%	86%	N/A	N/A	N/A	N/A	N/A	N/A
AL	91%	95%	93%	91%	89%	90%	86%	66%	74%	N/A	N/A	N/A	N/A	N/A	N/A
AS	96%	89%	93%	89%	89%	89%	N/A	N/A	N/A	N/A	N/A	N/A	100%	81%	89%
Macro- μ	85%	86%	86%	87%	81%	83%	90%	86%	87%	100%	91%	95%	100%	81%	89%
μ -pesada	89%	88%	89%	89%	89%	89%	87%	86%	86%	100%	91%	95%	100%	81%	89%

Tabela 7: Classificação da origem das variações, com SVM e representação BERT.

Vejam-se alguns exemplos onde isto acontece:

P: *Há sociedades que não podem ser constituídas nos balcões “Empresa na Hora”?*

VG1: *Existem empresas que não podem ser configuradas nos contadores “Empresa no Horário”?*

VG2: *Existem empresas que não podem ser configuradas nos contadores “Business on Time”?*

P: *A que balcão de atendimento “Empresa na Hora” me devo dirigir?*

VG1: *Qual service desk devo entrar em contato?*

Apesar do AIA-BDE sofrer de desequilíbrio ao nível da origem, o desempenho individual em cada origem não parece ser muito afetado pelo número de perguntas originais com essa origem. É de notar, aliás, que o desempenho inferior é claramente para as variações RJACSR, quando o número de perguntas com esta origem (118) é mais do dobro das perguntas de AL e AS (56). Uma possível explicação seria o facto de estas perguntas serem mais longas que as demais, mas como a Tabela 1 mostra, as perguntas de AL são as mais longas, com uma diferença considerável. Uma última possibilidade será um maior vocabulário utilizado por estas perguntas. No entanto, identificar claramente a razão para este desempenho inferior implicaria uma análise mais profunda das perguntas.

Após esta experiência, não fica claro se seria benéfica a utilização de um classificador inicial que fizesse a triagem das perguntas de acordo com a sua origem. Com a exceção das variações VIN, haveria ainda uma proporção considerável

de perguntas mal classificadas, o que impossibilitaria à partida a identificação da sua resposta correta. No entanto, tal como acontece para a experiência anterior, estes resultados devem ser vistos como ilustrativos daquilo que pode ser feito com o AIA-BDE, e apenas uma base (*baseline*) com grande margem de melhoria. Por exemplo, à semelhança da experiência anterior, poderiam ter sido obtidas representações alternativas a partir do modelo BERT, ou utilizada uma versão deste modelo afinada (*fine-tuned*) para a classificação automática, sem recurso ao classificador SVM.

6. Classificação do Tipo de Pergunta

Considerando que o tipo de uma pergunta pode condicionar a resposta a dar e o processo de a obter, numa última experiência procuramos identificar automaticamente esse tipo. Contudo, aqui deparámo-nos com dois problemas, também referidos na Secção 3.3: (i) as anotações do tipo de pergunta foram muito variáveis, com apenas cerca de um quarto das perguntas em que os cinco anotadores concordaram; (ii) considerando estas anotações, o corpo AIA-BDE é altamente desequilibrado relativamente ao tipo de pergunta com maior confiança.

Por se tratar apenas de uma experiência ilustrativa do que é possível fazer com o AIA-BDE, optamos por realizar algumas simplificações. Mais propriamente, consideramos que cada pergunta só podia ter um tipo, e que esse seria o tipo em que tínhamos maior confiança. Isto eliminou automaticamente da nossa experiência as per-

guntas em que havia tipos empatados, com 40% ou mesmo 20% de confiança. Para lidar com o segundo problema, decidimos focar-nos apenas nos tipos para os quais, após a aplicação da condição anterior, restava um número aceitável de exemplos.

A Tabela 8 apresenta os quatro tipos com mais de 30 perguntas (Definição, Procedimento, Tempo, Condição), e a respetiva quantidade de perguntas desse tipo, após a simplificação anterior. Ao verificarmos que para dois dos quatro tipos anteriores, Tempo e Condição, há pouco mais de 40 perguntas, decidimos realizar a experiência de duas formas: uma considerando os quatro tipos, outra considerando apenas os dois majoritários, Definição e Procedimento.

Tipo	Quantidade
Definição	204
Procedimento	132
Tempo	44
Condição	43

Tabela 8: Distribuição após remoção de instâncias com menor confiança.

A experiência inicial consistiu em treinar um classificador para prever o tipo da pergunta, dada o texto da pergunta. No entanto, acabamos também por experimentar até que ponto seria possível fazer o mesmo, mas considerando apenas o texto da resposta.

Voltamos a experimentar um conjunto de classificadores, incluídos na biblioteca scikit-learn, e baseados nas mesmas representações vetoriais da experiência anterior (Secção 5), TF-IDF e BERT. Uma vez mais, o classificador baseado numa SVM linear voltou a destacar-se. Foi também utilizado com os parâmetros por omissão, desta vez com a exceção do número máximo de iterações. Ao verificarmos, através de uma mensagem de *warning*, que nem sempre havia convergência, decidimos aumentar o número máximo de iterações de 1.000, o valor por omissão, para 3.000, minimizando desta forma o problema.

Outra diferença desta experiência foi a existência de menos dados e inexistência de uma separação clara em dados a usar para treino e para teste. Assim, nas Tabelas 9 e 10, optamos por apresentar os resultados de uma validação cruzada em 10 subconjuntos (em inglês, *10-fold cross validation*), respetivamente para as experiências com quatro e dois tipos. Os resultados podem ser analisados com base na exatidão (em inglês, *accuracy*), i.e., a proporção de tipos corretamente identificados, mas também

através das macro médias da precisão (Macro-P), abrangência (Macro-A) e medida-F1 (Macro-F1). Enquanto que na *accuracy* os tipos mais frequentes terão um maior peso para a média, nas macro médias o desempenho em cada tipo tem o mesmo peso.

Os desempenhos reportados mostram que, mesmo com as simplificações realizadas, a identificação do tipo de pergunta é desafiante. Esta situação é mais evidente quando se consideram quatro tipos, e ainda mais quando a identificação se baseia na resposta e não na pergunta. O impacto do desequilíbrio observa-se nos valores da medida-F1, sempre mais baixos do que a *accuracy*, principalmente quando são considerados quatro tipos. Ou seja, como seria de esperar, o desempenho será melhor para os tipos mais representados.

De notar ainda os desvios padrão elevados, que mostram que o desempenho depende muito da escolha do conjunto de treino, mesmo quando esse conjunto inclui 90% dos dados (i.e., validação cruzada em 10 subconjuntos). Estes desvios dificultam a análise da melhor forma de representação e não é possível tirar grandes conclusões. Mesmo que, por exemplo, no cenário com quatro tipos, as médias sugiram que a representação BERT funcione melhor quando se usa a pergunta, e que a representação TF-IDF seja preferível quando se usa a resposta, os desvios padrão, respetivamente 10% e 8%, mostram que também pode acontecer o contrário.

7. Conclusões

Apresentamos neste artigo o corpo AIA-BDE, focado em FAQ que cobrem um pequeno conjunto de domínios da Administração Pública de Portugal. Apesar de outros ficheiros complementares, focamo-nos em 855 perguntas, para as quais variações foram produzidas manual e automaticamente, permitindo assim a avaliação de sistemas de recuperação de FAQ, que podem estar integrados em sistemas de resposta automática a perguntas, ou até de diálogo. Para além das variações, a cada uma das 855 perguntas anteriores está associada uma lista de tipos de pergunta e respetiva confiança.

Utilizamos ainda o AIA-BDE em três experiências, com vista à avaliação de: (i) utilização das variações como perguntas de um utilizador, e respetiva associação a perguntas conhecidas; (ii) classificação de variações de acordo com a origem da respetiva pergunta original; (iii) identificação do tipo de pergunta. As experiências realizadas ajudaram-nos a compreender melhor

Entrada	Modelo	Accuracy	Macro-P	Macro-A	Macro-F1
Pergunta	SVM + TFIDF	86±5%	85±10%	78±7%	79±8%
	SVM + BERT	87±8%	84±10%	83±11%	82±10%
Resposta	SVM + TFIDF	73±6%	66±12%	60±9%	61±10%
	SVM + BERT	69±9%	61±14%	58±11%	56±11%

Tabela 9: Validaao cruzada para a identificaao do tipo de pergunta para as quatro classes com mais de 40 instâncias: Definiao, Procedimento, Condiao, Tempo

Entrada	Modelo	Accuracy	Macro-P	Macro-A	Macro-F1
Pergunta	SVM + TFIDF	93±2%	92±2%	92±3%	92±3%
	SVM + BERT	94±4%	94±4%	93±5%	93±4%
Resposta	SVM + TFIDF	82±7%	82±8%	80±7%	80±7%
	SVM + BERT	82±7%	84±7%	81±7%	81±7%

Tabela 10: Validaao cruzada da identificaao do tipo de pergunta para as duas classes com mais de 100 instâncias: Definiao e Procedimento

o seu cont eudo, e demonstraram a sua utilidade, ao mesmo tempo que estabeleceram resultados base (*baselines*), com margem para melhoria no futuro. Verificamos que diferentes tipos de variaao trazem desafios diferentes e que, regra geral, n o parece haver um m todo que se adapte bem a todos os tipos. O mesmo acontece com as perguntas de diferentes origens. Por exemplo, verificamos que, apesar das representaoes obtidas a partir de modelos de linguagem baseados em *transformers* (BERT) levarem a resultados interessantes, isso n o se verifica para todo o tipo de variaao, e uma abordagem baseada em IR tradicional, bem mais simples, pode ser bastante competitiva. Ainda assim, acreditamos que as abordagens baseadas em *transformers* tenham maior margem de progresso, por exemplo, se forem afinadas para o c culo da similaridade sem ntica; ou se forem pr -treinados em dados espec fico do dom nio. O mesmo para a abordagem baseada em *word embeddings*, onde poderia adicionalmente ser ben fico considerar expressões ou entidades multipalavra (e.g., *Cart o Provis rio de Identificaao de Pessoa Colectiva* ou *Empresa na Hora*).

Depois da terceira experi ncia, verificamos tamb m que seria importante rever os tipos de pergunta atualmente considerados, tentando diminuir ou uniformizar as sobreposioes poss veis. Paralelamente, poderia ser interessante abordar a classificaao autom tica do tipo como um problema de multi-classificaao.

O AIA-BDE   disponibilizado   comunidade atrav s de <https://github.com/NLP-CISUC/AIA-BDE>, para que possa ser utilizado em experi ncias como aquelas aqui apresentadas, com

vista   melhoria dos resultados base, ou em experi ncias alternativas. Apesar de focado num dom nio, acreditamos que as variaoes podem ser usadas como dados de treino para a identificaao mais generalizada de par frases em contexto interrogativo, ou como base para a criaao de uma coleao para a avaliaao da similaridade sem ntica textual no mesmo contexto. O AIA-BDE pode ainda servir de base a um agente que responde a perguntas acerca dos dom nios cobertos. Aqui, as variaoes podem servir apenas para avaliaao do sistema, como feito recentemente (Santos et al., 2020a), mas tamb m usadas como variaao das intenoes no processo de compreens o de linguagem natural (em ingl s, *Natural Language Understanding*).

Ainda que n o seja uma prioridade, no futuro poder o ser inclu das mais perguntas e mais variaoes ao corpo, aumentando assim a sua dimens o e tornando-o mais apto para o treino de modelos mais poderosos. Para al m do AIA-BDE, mostramos que um sistema de recuperaao de FAQ pode ser avaliado com base num conjunto de variaoes para cada pergunta, que simulem as necessidades de informaao dos utilizadores. Estas necessidades podem, em alguns casos, ser expressas de uma forma semelhante  quela em que as perguntas est o armazenadas mas, devido   variabilidade lingu stica, podem tamb m ser colocadas de formas radicalmente diferentes, ao n vel lexical ou sint tico. Apesar do trabalho manual envolvido, acreditamos que   uma vantagem ter um recurso deste tipo, que permita ir avaliando progressos desta forma. O processo de criaao poder  ser replicado para outros dom nios e, ainda que possa

ser benéfico recorrer a especialistas do domínio para o fazer, o mais importante é cobrir diferentes necessidades de informação, incluindo aquelas de utilizadores menos experientes, obtidas, por exemplo, com recurso a *crowdsourcing*.

Agradecimentos

Parte deste trabalho foi realizado no âmbito do projeto demonstrador AIA, “Apoio Inteligente a empreendedores (chatbots)”, financiado pela FCT, através da iniciativa INCoDe 2030.

Gostaríamos também de agradecer: ao João Ferreira, pelo seu envolvimento na integração das variações mais recentes no corpo e na definição dos tipos de pergunta; à Luísa Coheur e aos seus alunos, pela criação das variações VIN; à AMA, em especial ao Jorge Cabrita de Sousa, por nos ceder uma grande parte dos materiais de onde foram extraídas as perguntas originais.

Referências

- Agirre, Eneko, Mona Diab, Daniel Cer & Aitor Gonzalez-Agirre. 2012. SemEval-2012 task 6: A pilot on semantic textual similarity. Em *1st Joint Conference on Lexical and Computational Semantics: 6th International Workshop on Semantic Evaluation*, 385–393.
- Bowman, Samuel R., Gabor Angeli, Christopher Potts & Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. Em *Conference on Empirical Methods in Natural Language Processing*, 632–642. doi 10.18653/v1/D15-1075.
- Budzianowski, Paweł, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan & Milica Gašić. 2018. MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling. Em *Conference on Empirical Methods in Natural Language Processing*, 5016–5026. doi 10.18653/v1/D18-1547.
- Burke, Robin D, Kristian J Hammond, Vladimir Kulyukin, Steven L Lytinen, Noriko Tomuro & Scott Schoenberg. 1997. Question answering from frequently asked question files: Experiences with the FAQ finder system. *AI magazine* 18(2). 57–57.
- Caputo, Annalina, Marco Degemmis, Pasquale Lops, Francesco Lovecchio & Vito Manzari. 2016. Overview of the EVALITA 2016 question answering for frequently asked questions (QA4FAQ) task. Em *3rd Italian Conference on Computational Linguistics (CLiC-it): 5th Evaluation Campaign of Natural Language Processing and Speech Tools for Italian (EVALITA)*, CEUR-WS.
- Carvalho, Nuno Ramos, Alberto Simões & José João Almeida. 2021. Bootstrapping a data-set and model for question-answering in portuguese (short paper). Em *10th Symposium on Languages, Applications and Technologies (SLATE)*, 18:1–18:5. doi 10.4230/OASICS.SLATE.2021.18.
- Choi, Eunsol, He He, Mohit Iyyer, Mark Yatskar, Wen-tau Yih, Yejin Choi, Percy Liang & Luke Zettlemoyer. 2018. QuAC: Question answering in context. Em *Conference on Empirical Methods in Natural Language Processing*, 2174–2184. doi 10.18653/v1/D18-1241.
- Criscuolo, Marcelo, Erick Rocha Fonseca, Sandra Maria Aluísio & Ana Carolina Sperança-Criscuolo. 2017. MilkQA: a dataset of consumer questions for the task of answer selection. Em *6th Brazilian Conference on Intelligent Systems (BRACIS)*, 354–359. doi 10.1109/BRACIS.2017.12.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee & Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. Em *Conference of the North American Chapter of the Association for Computational Linguistics*, 4171–4186. doi 10.18653/v1/N19-1423.
- Fellbaum, Christiane (ed.). 1998. *WordNet: An Electronic Lexical Database (language, speech, and communication)*. The MIT Press.
- Fernandes, Mariana Gaspar, Cátia Dias & Luísa Coheur. 2019. Distinguishing different classes of utterances - the UC-PT corpus. Em *8th Symposium on Languages, Applications and Technologies (SLATE)*, 14:1–14:8. doi 10.4230/OASICS.SLATE.2019.14.
- Fonseca, Erick, Leandro Santos, Marcelo Criscuolo & Sandra Aluísio. 2016. Visão geral da avaliação de similaridade semântica e inferência textual. *Linguamática* 8(2). 3–13.
- Forner, Pamela, Anselmo Peñas, Eneko Agirre, Iñaki Alegria, Corina Forăscu, Nicolas Moreau, Petya Osenova, Prokopis Prokopidis, Paulo Rocha, Bogdan Sacaleanu et al. 2008. Overview of the CLEF 2008 multilingual Question Answering track. Em *Workshop of the Cross-Language Evaluation Forum for European Languages*, 262–295.

- Gonalo Oliveira, Hugo, Ricardo Filipe, Ricardo Rodrigues & Ana Alves. 2019. Using Lucene for developing question-answering agent in Portuguese. Em *8th Symposium on Languages, Applications and Technologies (SLATE)*, 2:1–2:14. doi 10.4230/OASICS.SLATE.2019.2.
- Guillou, Pierre. 2021. Portuguese BERT base cased QA (Question Answering), finetuned on SQUAD v1.1. <https://huggingface.co/pierreguillou/bert-base-cased-squad-v1.1-portuguese>.
- Gunasekara, Chulaka, Seokhwan Kim, Luis Fernando D’Haro, Abhinav Rastogi, Yun-Nung Chen, Mihail Eric, Behnam Hedayatnia, Karthik Gopalakrishnan, Yang Liu, Chao-Wei Huang, Dilek Hakkani-Tür, Jinchao Li, Qi Zhu, Lingxiao Luo, Lars Liden, Kaili Huang, Shahin Shayandeh, Runze Liang, Baolin Peng, Zheng Zhang, Swadheen Shukla, Minlie Huang, Jianfeng Gao, Shikib Mehri, Yulan Feng, Carla Gordon, Seyed Hossein Alavi, David Traum, Maxine Eskenazi, Ahmad Beirami, Eunjoon, Cho, Paul A. Crook, Ankita De, Alborz Geramifard, Satwik Kottur, Seungwhan Moon, Shivani Poddar & Rajen Subba. 2020. Overview of the 9th dialog system technology challenge: DSTC9. ArXiv:2011.06486 [cs.CL].
- Hartmann, Nathan S., Erick R. Fonseca, Christopher D. Shulby, Marcos V. Treviso, Jéssica S. Rodrigues & Sandra M. Aluísio. 2017. Portuguese word embeddings: Evaluating on word analogies and natural language tasks. Em *11th Brazilian Symposium in Information and Human Language Technology (STIL)*, 122–131.
- Karan, Mladen & Jan Šnajder. 2016. FAQIR—a frequently asked questions retrieval test collection. Em *International Conference on Text, Speech, and Dialogue*, 74–81.
- Karan, Mladen & Jan Šnajder. 2018. Paraphrase-focused learning to rank for domain-specific frequently asked questions retrieval. *Expert Systems with Applications* 91. 418–433. doi 10.1016/j.eswa.2017.09.031.
- Karan, Mladen, Lovro Źmak & Jan Šnajder. 2013. Frequently asked questions retrieval for Croatian based on semantic textual similarity. Em *4th Biennial International Workshop on Balto-Slavic Natural Language Processing*, 24–33.
- Kolomiyets, Oleksandr & Marie-Francine Moens. 2011. A survey on question answering technology from an information retrieval perspective. *Information Sciences* 181(24). 5412–5434. doi 10.1016/j.ins.2011.07.047.
- Lison, Pierre & Jörg Tiedemann. 2016. Open-Subtitles2016: Extracting large parallel corpora from movie and TV subtitles. Em *10th International Conference on Language Resources and Evaluation (LREC)*, 923–929.
- Lowe, Ryan, Nissan Pow, Iulian Serban & Jolene Pineau. 2015. The Ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems. Em *16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 285–294. doi 10.18653/v1/W15-4640.
- Magarreiro, Daniel, Luísa Coheur & Francisco S. Melour. 2014. Using subtitles to deal with out-of-domain interactions. Em *18th Workshop on the Semantics and Pragmatics of Dialogue (SemDial)*, 98–106.
- Magnini, Bernardo, Alessandro Vallin, Christelle Ayache, Gregor Erbach, Anselmo Peñas, Marten de Rijke, Paulo Rocha, Kiril Ivanov Simov & Richard F. E. Sutcliffe. 2004. Overview of the CLEF 2004 Multilingual Question Answering track. Em *Multilingual Information Access for Text, Speech and Images, 5th Workshop of the Cross-Language Evaluation Forum (CLEF), Revised Selected Papers*, 371–391.
- Mota, Cristina, Alberto Simões, Cláudia Freitas, Luís Costa & Diana Santos. 2012. Páxico: Evaluating Wikipedia-based information retrieval in Portuguese. Em *8th International Conference on Language Resources and Evaluation (LREC)*, 2015–2022.
- Nakov, Preslav, Doris Hoogeveen, Lluís Màrquez, Alessandro Moschitti, Hamdy Mubarak, Timothy Baldwin & Karin Verspoor. 2017. SemEval-2017 task 3: Community question answering. Em *11th International Workshop on Semantic Evaluation (SemEval)*, 27–48. doi 10.18653/v1/S17-2003.
- Nakov, Preslav, Lluís Màrquez, Walid Magdy, Alessandro Moschitti, Jim Glass & Bilal Randeree. 2015. SemEval-2015 task 3: Answer selection in community question answering. Em *9th International Workshop on Semantic Evaluation (SemEval)*, 269–281. doi 10.18653/v1/S15-2047.
- Nakov, Preslav, Lluís Màrquez, Alessandro Moschitti, Walid Magdy, Hamdy Mubarak, Abed Alhakim Freihat, Jim Glass & Bilal Randeree. 2016. SemEval-2016 task 3: Community question answering. Em *10th International Workshop on Semantic Evaluation*, 525–545. doi 10.18653/v1/S16-1083.

- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot & E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12. 2825–2830.
- Rajpurkar, Pranav, Jian Zhang, Konstantin Lopyrev & Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. Em *Conference on Empirical Methods in Natural Language Processing*, 2383–2392.
- Real, Livy, Erick Fonseca & Hugo Gonalo Oliveira. 2020. The ASSIN 2 shared task: a quick overview. Em *13th International Conference on Computational Processing of the Portuguese Language (PROPOR)*, 406–412. doi 10.1007/978-3-030-41505-1_39.
- Reddy, Siva, Danqi Chen & Christopher D Manning. 2019. CoQA: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics* 7. 249–266.
- Ritter, Alan, Colin Cherry & William B Dolan. 2011. Data-driven response generation in social media. Em *Conference on empirical methods in natural language processing*, 583–593.
- Rodrigues, Ruan Chaves, Jessica Rodrigues da Silva, Pedro Vitor Quinta de Castro, Nadia Felix Felipe da Silva & Anderson da Silva Soares. 2020a. Multilingual transformer ensembles for Portuguese natural language tasks. Em *ASSIN 2 Shared Task: Evaluating Semantic Textual Similarity and Textual Entailment in Portuguese*, CEUR–WS.
- Rodrigues, Rui, Paula Couto & Irene Rodrigues. 2020b. IPR: The semantic textual similarity and recognizing textual entailment systems. Em *ASSIN 2 Shared Task: Evaluating Semantic Textual Similarity and Textual Entailment in Portuguese*, CEUR–WS.
- Santos, Diana & Paulo Rocha. 2004. The key to the first CLEF with Portuguese: topics, questions and answers in CHAVE. Em *Workshop of the Cross-Language Evaluation Forum for European Languages*, 821–832.
- Santos, Jose, Luıs Duarte, Joao Ferreira, Ana Alves & Hugo Gonalo Oliveira. 2020a. Developing Amaia: A conversational agent for helping Portuguese entrepreneurs – An extensive exploration of question-matching approaches for Portuguese. *Information* 11(9). doi 10.3390/info11090428.
- Santos, Jose, Ana Alves & Hugo Gonalo Oliveira. 2020b. Leveraging on Semantic Textual Similarity for developing a Portuguese dialogue system. Em *13th International Conference on Computational Processing of the Portuguese Language (PROPOR)*, 131–142. doi 10.1007/978-3-030-41505-1_13.
- Souza, Fabio, Rodrigo Nogueira & Roberto Lotufo. 2020. BERTimbau: Pretrained BERT models for Brazilian Portuguese. Em *Brazilian Conference on Intelligent Systems (BRACIS)*, 403–417. doi 10.1007/978-3-030-61377-8_28.
- Vinyals, Oriol & Quoc V. Le. 2015. A neural conversational model. Em *International Conference on Machine Learning, Deep Learning Workshop*, arXiv:1506.05869 [cs.CL].
- Voorhees, Ellen M. 2008. Evaluating question answering system performance. Em *Advances in Open Domain Question Answering*, 409–430. doi 10.1007/978-1-4020-4746-6_13.
- Wen, Tsung-Hsien, David Vandyke, Nikola Mrksic, Milica Gasic, Lina M. Rojas-Barahona, Pei-Hao Su, Stefan Ultes & Steve Young. 2017. A network-based end-to-end trainable task-oriented dialogue system. Em *15th Conference of the European Chapter of the Association for Computational Linguistics*, 438–449.
- Zhang, Tianyi, Varsha Kishore, Felix Wu, Kilian Q Weinberger & Yoav Artzi. 2020. Bertscore: Evaluating text generation with BERT. Em *8th International Conference on Learning Representations (ICLR)*, arXiv:1904.09675 [cs.CL].

<http://www.linguamatica.com/>

linguamatica

Artigos de Investigação

Uso de tecnologias linguísticas para estudar a evolução dos sufixos -ÇOM e -VEL no galego-português medieval a partir de *corpora* históricos

Pablo Gamallo, José Ramon Pichel, José Martinho Montero Santalha & Marco Neves

AIA-BDE: um Corpo de Perguntas, Variações e outras Anotações

Hugo Gonçalo Oliveira & Ana Alves