

## ENCADEAR

## ENCADEAMENTO AUTOMÁTICO DE NOTÍCIAS

CARLA ABREU, JORGE TEIXEIRA E EUGÊNIO OLIVEIRA

## ABSTRACT

This work aims at defining and evaluating different techniques to automatically build temporal news sequences. The approach proposed is composed by three steps: (i) near duplicate documents detection; (ii) keywords extraction; (iii) news sequences creation. This approach is based on: Natural Language Processing, Information Extraction, Name Entity Recognition and supervised learning algorithms. The proposed methodology got a precision of 93.1% for news chains sequences creation.

## [1] INTRODUÇÃO

Diariamente são publicadas grandes quantidades de notícias *online*, o que pode conduzir a uma sobrecarga de informação para o leitor. Para estar informado e atualizado de um determinado acontecimento, o leitor depara-se com um vasto conjunto de artigos noticiosos, artigos esses que, em muitos casos, descrevem um mesmo evento, podendo apresentar apenas pequenas variações textuais. A situação agrava-se quando o leitor pretende saber mais detalhes sobre uma dada história ou sequência de eventos. Um exemplo concreto é o desaparecimento do avião da *Malaysia Airlines* a 8 de março de 2014. Considerando o dia 6 de outubro de 2014 a pergunta (*query*) “avião *Malaysia*” pesquisada no *Google News* (*news.google.pt*) retorna uma lista com mais de 50 notícias relacionadas. Dessas notícias retiramos a informação de que as buscas pelo avião foram retomadas. Como é possível observar pelos seguintes títulos: *Retomadas buscas pelo avião da Malaysia Airlines* (*Renascença*, 06/10/2014) e *Recomeçam as buscas pelo avião desaparecido da Malaysia Airlines* (*Jornal de Notícias*, 06/10/2014) o evento noticiado é o mesmo, mas pelo facto das notícias serem provenientes de fontes noticiosas diferentes apresentam variações textuais.

Este problema da sobrecarga de informação agrava-se quando o leitor quer perceber a história do desaparecimento do avião como um todo, e informar-se sobre todos os eventos que ocorreram relativamente a este acontecimento. A pergunta (*query*) “*desaparecimento Malaysia Airlines*” sem delimitações temporais ao *Google News* apresenta mais de 4.500 resultados. Neste conjunto de resultados torna-se complicado ou até mesmo humanamente impossível não só a deteção de todos os eventos como apenas os mais relevantes para a história. Por conseguinte,

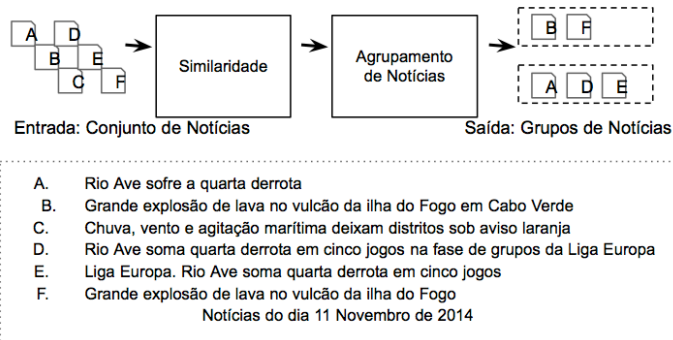


FIGURA 1: Detecção e agrupamento de notícias similares

o leitor não consegue ter a perceção de toda a história, descrita em mais de 4.500 notícias diferentes.

O objetivo deste trabalho é colmatar este problema: automaticamente detetar e agrupar notícias similares e automaticamente criar histórias a partir de notícias relacionadas temporalmente. Proporciona-se deste modo ao leitor uma nova forma de navegação entre eventos relativos a um mesmo acontecimento.

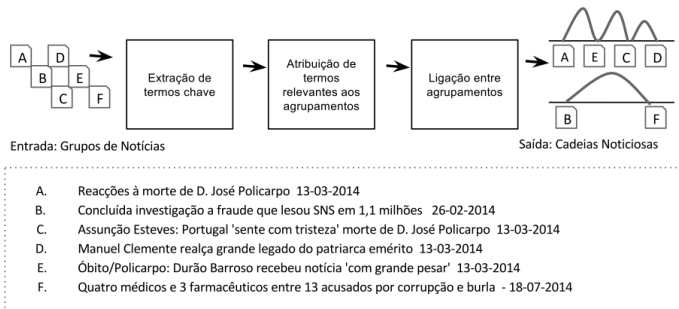


FIGURA 2: Construção de cadeias noticiosas

Pretendemos, numa primeira fase, detetar e agrupar notícias duplicadas (ver Figura 1). Utilizamos métodos de processamento de linguagem natural, algoritmos de medição de distância entre *strings*<sup>1</sup> (para o cálculo da proximidade entre notícias) e algoritmos supervisionados de aprendizagem automática (para a determinação da similaridade entre notícias). Numa segunda fase (ver Figura 2), com vista à formação automática de cadeias noticiosas, extraímos termos relevantes das notícias como por exemplo o tópico principal da notícia, as entidades, os locais e os nomes das personalidades; e ligamos os grupos de notícias pela medição da distância entre os mesmos. Utilizamos algoritmos de aprendizagem supervisi-

[1] Sequência de caracteres.

onada para ligar notícias de forma sequencial para criar uma história temporalmente lógica e contextualizada.

Este artigo encontra-se organizado da seguinte forma: na secção [2] apresentaremos o essencial sobre trabalhos relacionados. Na secção [3] vamos expor detalhadamente todos os passos da metodologia aplicada. Na secção [4] vamos enunciar os recursos linguísticos utilizados. Seguem-se a descrição das experiências realizadas (secção [5]) e a apresentação e discussão dos resultados na secção [6]. Na secção [7] é apresentada a interface gráfica desenvolvida como prova de conceito. Por fim são apresentadas as conclusões e o trabalho futuro na secção [8].

## [2] TRABALHOS RELACIONADOS

### [2.1] *Detetar Notícias Duplicadas*

Notícias quase duplicadas são notícias publicadas por fontes distintas mas cujo conteúdo e data de publicação são muito semelhantes. A publicação deste tipo de notícias é bastante comum mas não traz nenhuma mais valia ao leitor. Adicionalmente, o seu armazenamento tem elevados custos computacionais. Devido a estes constrangimentos torna-se necessária a deteção deste tipo de notícias (Kumar & Govindarajulu 2009).

São várias as abordagens propostas para a resolução do problema de deteção de notícias quase duplicadas, entre elas encontram-se: a abordagem baseada no léxico, a abordagem baseada no URL e a abordagem baseada na semântica. A abordagem baseada no léxico não requer nenhum conhecimento linguístico. O objetivo é perceber a existência de termos em comum entre documentos. A abordagem baseada no URL visa detetar notícias duplicadas pela comparação do endereço URL. Porém esta abordagem continua a não ser suficiente. Isto porque, não existe um padrão estabelecido pelas diversas fontes noticiosas de como criar um URL e, portanto, podendo este conter ou não informação útil. A abordagem semântica é uma abordagem mais completa, esta inclui a necessidade de pré-processamento implicando: *tokenization*, *stemming* e remoção das *stop-words*. Após o pré-processamento do texto, as notícias são comparadas através de uma função de similaridade. Esta função tem como objetivo medir o grau de semelhança entre pares de notícias. O valor retornado por esta função varia entre [0,1], e é tanto maior quanto maior for a semelhança existente entre as notícias.

No trabalho intitulado *Duplicate Record Detection: A Survey*, Elmagarmid et al. (2007) explicam todo o fluxo necessário à deteção de documentos duplicados. Este trabalho refere-se à abordagem semântica. As notícias são inicialmente processadas, seguindo-se a determinação dos campos a comparar; é, depois, medido o grau de semelhança entre pares de notícias; e por fim, com base no resultado obtido é determinado se os documentos são ou não similares. Os autores ilustram quatro métricas: similaridade de *strings* baseada em caracteres; similaridade baseada em *tokens*; similaridade fonética e similaridade numérica.

A similaridade baseada em caracteres foi desenvolvida para detetar erros tipográficos. Alguns exemplos dessas métricas são: algoritmos de edição de distância (Hamming (He et al. 2004) e Levenshtein (Levenshtein 1965)) que visam calcular o número de adições, substituições e remoções necessárias para converter uma *string* numa outra, como por exemplo “futebol” e “futbol”; distância Affine Gap (Waterman et al. 1976) que consiste em abrir ou estender um espaço, para transformar uma *string* noutra, como: “C Ronaldo” e “Cristiano Ronaldo”; a métrica de distância Jaro (Bilenko et al. 2003) que mede a semelhança entre duas *strings* tendo em conta o comprimento das mesmas, o número de caracteres em comum e o número de transposições necessárias; e a métrica Q-grams (Ullmann 1977) que consiste na divisão das *strings* iniciais em *substrings* de tamanho  $q$ , a medição de similaridade entre documentos consiste na medição de *substrings* em comum entre as duas notícias.

Para o cálculo da similaridade entre pares de notícias utilizamos uma abordagem baseada em algoritmos de aprendizagem automática.

Infelizmente, existem poucos estudos desenvolvidos no sentido de verificar a eficiência da utilização de métricas de distância (Elmagarmid et al. 2007). Existem, por exemplo, alguns estudos que mencionam a eficiência da métrica de distância Jaro (Bilenko et al. 2003; Yancey 2005) na comparação de nomes.

Para a deteção e agrupamento de notícias similares é também recorrente a utilização de abordagens de *clustering* (Banerjee et al. 2007; Vadrevu et al. 2011). Nesta abordagem o documento é caracterizado por um conjunto de palavras, usualmente representado por um vetor de frequência da ocorrência dos termos. A determinação da similaridade entre agrupamentos e respetivo agrupamento efetua-se após a aplicação de um algoritmo de *clustering* sobre a coleção. Existem duas abordagens de *clustering* que podem ser aplicadas: a supervisionada, onde os tópicos são conhecidos, e a não supervisionada, onde não existe conhecimento inicial. Existem dois grandes problemas associados à aplicação de técnicas de *clustering* supervisionado, estes são: definição de categorias, tornam o sistema rudimentar, pois ao longo do tempo há uma tendência para o aparecimento de novas categorias; uma categoria abrange não só notícias duplicadas, como abrange também notícias que se referem ao mesmo tema. O problema relacionado com o *clustering* não supervisionado é o de não conhecermos os elementos responsáveis pela elaboração dos agrupamentos.

O nosso contributo, na componente da deteção de notícias quase duplicadas, diz respeito ao estudo da eficiência de alguns algoritmos de edição de distância para textos estruturados de dimensão variável, pela utilização de uma abordagem baseada na semântica.

As etapas necessárias para a elaboração deste módulo pode ser observada na Figura 1.

## [2.2] *Geração Automática de Histórias*

Diversos trabalhos tem sido conduzidos com o objetivo de criarem histórias a partir de vários documentos como: notícias (Shahaf & Guestrin 2010; Mei & Zhai 2005), blogs (Lin et al. 2012; Qamra et al. 2006) e resultados de pesquisas (Kumar et al. 2004). Em alguns trabalhos, antes da criação da história noticiosa o leitor tem que indicar o tema de pesquisa (Shahaf & Guestrin 2010; Lin et al. 2012). Outros trabalhos porém, visam ser mais abrangentes, e determinar dentro do seu conjunto de dados todas as histórias existentes (Allan et al. 1998b; McKeown et al. 2002). A primeira abordagem é utilizada em estudos relacionados com o tópico “Geração da História” sendo que a segunda abordagem é mais popular em estudos de “Detecção de Tópicos e Monitorização”. Em relação a estes dois tópicos, é de notar que existem poucos estudos sobre o primeiro, mas, no entanto, o segundo tópico tem vindo a ser extensivamente estudado (Lin & Liang 2008). Segundo Allan et al. (1998b), o conhecimento inicial dado ao sistema para a criação das histórias pode não ser adequado à monitorização das mesmas uma vez que o tema de discussão associado a um evento muda frequentemente.

Outra área que visa organizar e estruturar informação é a classificação hierárquica (Sun & Lim 2001; Lawrie & Croft 2000; Li et al. 2007). A estrutura hierarquia impõe uma estrutura a um conjunto de dados. Porém, não identificamos nenhum estudo realizado de forma a perceber se essa estrutura reflete as relações existentes entre os diversos documentos (Nallapati et al. 2004).

A nossa abordagem para a geração automática de histórias a partir das notícias baseia-se nas etapas utilizadas nos diferentes trabalhos com o mesmo propósito. As diferentes etapas consideradas, bem com o seu fluxo, podem ser observadas na Figura 2.

### *Geração da História*

O trabalho intitulado *Connecting the Dots Between News* (Shahaf & Guestrin 2010) visa encontrar uma história coerente num conjunto de artigos noticiosos a partir de um ou mais tópicos indicados pelo utilizador. O método utilizado neste trabalho é aplicável a outros domínios como: *emails*, artigos científicos e inteligência militar. Neste trabalho os autores introduziram a noção de coerência, e *feedback* do utilizador. A abordagem proposta pelos autores consistiu na identificação de ligações entre notícias, tendo em conta: palavras omissas, palavras que estão relacionadas com as palavras do texto embora não apareçam no mesmo, e a importância das palavras. O problema da formação das cadeias de notícias foi solucionado recorrendo a uma abordagem de programação linear.

Outro trabalho desenvolvido com o propósito de gerar uma linha temporal de uma história é o *A Graph Teoretic Approach to Extract Storylines from Serach Results* (Kumar et al. 2004). Neste trabalho os resultados de pesquisa são representados numa estrutura de grafos, onde, os nós representam a informação associada ao

documento, e as ligações entre os nós, representam o peso de ligação. Para a elaboração das cadeias, os autores recorrem à utilização de um algoritmo de pesquisa local sobre a estrutura definida.

### *Deteção de Tópicos e Monitorização*

Existem três tarefas associadas a deteção de tópicos e monitorização, são elas: monitorização de eventos conhecidos (eventos já detetados pelo sistema), deteção de novos eventos, e segmentação das notícias em histórias. O grande objetivo dos estudos de deteção de tópicos e monitorização é o de identificar todas e quaisquer notícias relacionadas com um dado evento (Allan et al. 1998a).

Para o nosso trabalho, a componente mais interessante deste estudo é a forma como é executado o monitoramento de uma história nas notícias. A abordagem de monitoramento utilizada em “On-line News event detection and tracking” (Allan et al. 1998b) começa por reduzir o conteúdo noticioso a um conjunto de entre 10 a 20 *features*. Os autores acreditam que poucas *features* são necessárias para o monitoramento de notícias uma vez que o essencial de uma história tende a ser descrito por um conjunto pequeno de palavras ou frases. Neste trabalho, as cadeias são obtidas pelo cálculo de semelhança entre as *queries* que caracterizam cada notícia.

## [3] METODOLOGIA

### [3.1] *Similaridade*

Abordamos a similaridade entre artigos noticiosos em quatro passos distintos: (i) normalização do conteúdo noticioso; (ii) identificação dos elementos a comparar; (iii) comparação entre pares de notícias; (iv) tomada de decisão.

### *Normalização*

A normalização de textos é uma etapa tradicional em NLP para simplificar a análise posterior dos mesmos. Realizamos as seguintes tarefas de normalização:

- 1) Remoção de símbolos de pontuação, como: <, >, /, “, ”, (, ), -;
- 2) Remoção de padrões redundantes e que no âmbito deste trabalho, não são informativos, como: “Lusa - Esta notícia foi escrita nos termos do Acordo Ortográfico”;
- 3) Remoção de *stop-words*, através da utilização de uma lista disponibilizada pelo *snowball*<sup>2</sup> (para a língua portuguesa);
- 4) Redução das palavras à sua raiz através da utilização do *Porter Stemmer* para língua portuguesa, disponibilizado pelo *PTStemmer* (Oliveira 2008).

Na Tabela 1 apresentamos um exemplo da normalização, desde a notícia original até à sua versão normalizada.

[2] <https://snowball.tartarus.org>

Operação	Exemplo
<i>Notícia original</i>	Nova Deli, 02 jan (Lusa) - A Índia anunciou que vai permitir a cidadãos estrangeiros investirem no seu mercado de ações.
1- <b>Pontuação</b>	Nova Deli 02 jan Lusa A Índia anunciou que vai permitir a cidadãos estrangeiros investirem no seu mercado de ações.
2- <b>Padrões</b>	A Índia anunciou que vai permitir a cidadãos estrangeiros investirem no seu mercado de ações.
3- <b>Stop-words</b>	Índia anunciou vai permitir cidadãos estrangeiros investirem mercado ações.
4- <b>Stemm</b>	Índi anunc va permit cidadã estrangeir invest merc açõ.

TABELA 1: Exemplo do fluxo da normalização.

### Identificação dos elementos a comparar

Identificamos cinco conteúdos essenciais nos artigos noticiosos publicados em formato digital: título, corpo da notícia, data de publicação, URL e metadados (*tags*).

URL: provenientes de diferentes domínios têm uma composição distinta. A Tabela 2 apresenta três pares <título, URL>. O primeiro URL é composto pelo título da notícia; já o segundo dá-nos a indicação das áreas a que a notícia está associada, não explicitando em concreto o acontecimento presente; o terceiro exemplo não nos consegue transmitir nenhuma informação concreta para além do domínio.

Al Qaeda reivindica atentados em quartel militar do Iêmen

<http://visao.sapo.pt/al-qaeda-revindica-atentados-em-quartel-militar-do-iemen=f803958>

Plantel empenhado na vitória em Barcelos

[http://www.record.xl.pt/Futebol/Nacional/1a\\_liga/academica/interior.aspx?content\\_id=919169](http://www.record.xl.pt/Futebol/Nacional/1a_liga/academica/interior.aspx?content_id=919169)

Cidade chinesa gera energia com queima de notas de banco

[http://diariodigital.sapo.pt/news.asp?id\\_news=750321](http://diariodigital.sapo.pt/news.asp?id_news=750321)

TABELA 2: Exemplos de URL

Corpo da notícia: o título ou corpo da notícia, como componentes isolados, podem não ser suficientes para a determinação da similaridade. Identificamos o cabeçalho da notícia, tipicamente o primeiro parágrafo, como sendo um elemento adicional a considerar para o cálculo da similaridade entre notícias (ver Figura 3). Este cabeçalho corresponde muitas vezes ao resumo da notícia e como tal é muito informativo.



FIGURA 3: Campos da notícia a serem comparados.

Data de publicação: as notícias contêm informação temporal importante para a contextualização do evento. Assumimos que existe um intervalo de tempo restrito dentro do qual há uma maior tendência para o aparecimento de notícias duplicadas. Por exemplo, é mais provável a existência de notícias duplicadas com intervalo de datas de publicação de 24 horas do que numa semana. Deste modo, o fator tempo serve como delimitador do intervalo temporal de notícias comparáveis.

### Comparação de Notícias

Podem ser utilizadas diferentes métricas para o cálculo da similaridade. Neste trabalho, consideramos as seguintes: Hamming (He et al. 2004), Levenshtein (Levenshtein 1965) e Jaro (Bilenko et al. 2003).

De forma a que os resultados destas métricas possam ser comparáveis, é necessário proceder à normalização dos mesmos, aplicamos a seguinte fórmula (Expressão 1) aos resultados retornados pelos métodos de edição de distância.

$$D'(s, t) = 1 - \frac{D(s, t)}{\max(|s|, |t|)}, D \in \mathbb{Q} | D \in [0; 1] \quad (1)$$

Onde:

$D(s, t)$  é a distância obtida pela métrica de edição de distância entre a string  $s$  e  $t$ ;

$\max(|s|, |t|)$  é o comprimento da string de maior dimensão entre  $s$  e  $t$ ;

$D'(s, t)$  é a distância normalizada entre  $s$  e  $t$ .

Para cada par de notícias é calculado o  $D'$ . A decisão sobre a similaridade é decidida no passo posterior.

### *Decisão da similaridade entre notícias*

Usamos diversos métodos de aprendizagem supervisionada para a classificação de notícias duplicadas. Os algoritmos usados foram: Support Vector Classifier (SVC), SVC Linear, Decision Tree e Random Forest. Estes algoritmos estão disponíveis, através de bibliotecas python, no *scikit learn* (Pedregosa et al. 2011).

A partir das distâncias calculadas na secção [3.1.3] tiramos partido de algoritmos de aprendizagem supervisionada para classificar pares de notícias como duplicadas ou não duplicadas.

### [3.2] *Agrupamento de Notícias*

Este módulo é responsável pela criação de grupos de notícias duplicadas usando os resultados dos pares de notícias previamente classificadas (ver secção [3.1.4]).

Um exemplo ilustrativo dos passos necessários desde a receção das notícias até à composição dos agrupamentos pode ser ilustrado pela Figura 1. Neste caso, estamos perante seis notícias (A,B,C,D,E,F) que formam quinze pares distintos (AB, AC, AD, AE, AF, BC, BD, BE, BF, CD, CE, BF, DE, DF, EF). Estes pares de notícias são comparados na secção [3.1] e deste módulo são considerados como duplicados os pares AD, AE, BF e DE. Pela observação do exemplo, constatamos que são formados dois grupos (BF e ADE).

### [3.3] *Extração de Termos Chave*

Para cada grupo de notícias é necessário e essencial, sintetizar a informação contida nesses grupos.

Na nossa abordagem, vamos representar as notícias por um conjunto de termos chave. Os termos chave podem ser considerados termos que transmitem informação relevante do texto, como: o tópico da notícia, nomes de personalidades, locais e outros. Consideramos três tipos de termos chave: (i) palavras isoladas (*uni-grams*) (ii) expressões relevantes (*n-grams*) e (iii) entidades.

#### *Palavras Isoladas*

As palavras isoladas correspondem a palavras compostas por um *token* que aparecem explicitamente no conteúdo noticioso. De forma a obtermos estas palavras executamos três tarefas: POS Tagger, normalização e análise da frequência da palavra.

POS Tagger: visa a identificação das categorias gramaticais das palavras que compõe o texto da notícia. Utilizamos nesta tarefa o *TreeTagger* (Schmid 1994) adaptado para a língua portuguesa, disponibilizado por Garcia & Gamallo (2013).

Normalização: corresponde à remoção de padrões linguísticos e frases recorrentes do corpo da notícia obtidos por inspeção manual, como: expressões de datas

(Porto, 12 Agosto 2014), resultados de futebol (2-1) e padrões jornalísticos (Porto, 12 Agosto 2014 (Lusa)).

Análise da frequência da palavra: pela utilização da métrica estatística *Term Frequency-Inverse Document Frequency* (TF-IDF). No seu cálculo, esta métrica relaciona o aparecimento de um termo na notícia com o aparecimento do mesmo na coleção permitindo assim detetar a existência de termos relevantes.

Da análise da frequência de palavras no texto resulta uma lista de palavras com peso associado. Consideramos como palavras relevantes, aquelas com maior peso e pertencentes à categoria gramatical nome.

#### *Expressões Relevantes*

As expressões relevantes correspondem a *ngrams* que aparecem explicitamente no conteúdo noticioso e que de uma forma simplificada podem transmitir informação relevante contida no texto.

Para a extração deste elemento do texto foi adicionado um passo intermédio à abordagem apresentada na secção [3.3.1]. Para tal, após a normalização foi aplicado um filtro de forma a obter expressões do texto. As expressões são *ngrams*, que obedecem a certos padrões gramaticais, como: sequências de nomes (Domingos Paciência), nome e adjetivo (homens encapuzados) entre outros.

A análise da frequência é neste caso efetuada sobre os padrões. O resultado retornado pela análise de frequência indica-nos quais as expressões relevantes para a notícia em questão. A última etapa consiste na atribuição das expressões relevantes à notícia.

#### *Entidades*

O reconhecimento de entidades mencionadas, nomeadamente o nome de personalidades, é essencial no contexto de extração de termos e expressões chave das notícias.

Existem disponíveis vários recursos para o reconhecimento de entidades mencionadas para a língua portuguesa, como os mencionados pela Linguateca<sup>3</sup>. No entanto e no âmbito deste trabalho, estamos perante um domínio muito dinâmico, as notícias, onde constantemente aparecem novas entidades (Charlie Hebdon, Fukushima). Optamos por implementar um sistema que se adapta a estas características.

Foi implementado um algoritmo com o objetivo de verificar, numa primeira fase, quais as palavras no texto que se iniciam com um carácter maiúsculo. Das palavras encontradas, se a palavra maiúscula estiver posicionada no início da frase é verificado se a palavra é ou não uma *stop-word*, e caso seja, então não é considerada. Para as palavras que passarem a fase anterior é verificado se são precedidas

[3] <http://www.linguateca.pt/LivroSegundoHAREM/>

de outras palavras capitalizadas, sendo permitido uma palavra de ligação entre termos capitalizados inicializada a minúscula. Um exemplo de entidades extraídas pelo algoritmo é dado pelos seguintes termos: “Passos”, “Paulo Portas”.

De forma a enriquecer os termos chave extraídos para o conjunto de expressões e entidades extraídas de cada notícia tentamos identificar quais desses termos relevantes são nomes de personalidades. Para tal comparamos esses termos com um recurso externo, o Verbetes<sup>4</sup>.

### [3.4] *Atribuição de termos relevantes aos agrupamentos*

Depois da junção de notícias similares em agrupamentos (secção [3.2]) e após realizada a extração de termos relevantes de cada notícia (secção [3.3]), é possível fazer a atribuição dos termos chave aos agrupamentos de notícias.

Os termos chave associados a cada agrupamento correspondem aos termos relevantes que estão associados a cada uma das notícias do agrupamento. É de referir que cada termo chave tem um peso ( $w$ ), que está relacionado com a sua frequência ( $f$ ) no agrupamento. A importância de um termo é dado pela relação entre o número de notícias em que o termo aparece e número total de notícias que compõe o agrupamento. Um exemplo de palavras relevantes associadas a um agrupamento e respetiva importância é dado por:

reclusos[f=9;w=1];presos[f=9;w=1];  
cárcere[f=7;w=0.78];sudoeste[f=7;w=0.78];  
representantes[f=6;w=0.67];  
violação[f=6;w=0.67];cadeia[f=5;w=0.56];  
quilómetros[f=4;w=0.44];irmãos[f=4;w=0.44];

Neste agrupamento, o termo *reclusos* é mais representativo do conjunto do que o termo *irmãos*. Isto porque, considerando que o agrupamento em questão tem nove notícias, o primeiro termo aparece associado a todas as notícias do agrupamento ( $f = 9$ ), tendo um peso de  $w = \frac{9}{9}$ , ou seja 1; enquanto o segundo termo só se encontra associado a 4 notícias do conjunto ( $f = 4$ ), tendo um peso de  $w = 0.44$ .

### [3.5] *Ligações entre Agrupamentos*

Este módulo visa identificar as ligações entre os agrupamentos de notícias duplicadas previamente calculadas com os respetivos termos relevantes associados (ver Figura 2).

Partimos do pressuposto que as cadeias noticiosas só podem existir para a mesma categoria de notícias, de forma a simplificar esta tarefa. Para isso, fizemos a atribuição das categorias aos grupos de notícias, através de uma fonte de conhecimento externo que mapeia as *tags* atribuídas pelos jornalistas com a categoria a que a notícia fica associada. As categorias indicam de uma forma geral a

[4] <https://store.servicos.sapo.pt/pt/Catalog/other/free-api-information-retrieval-verbetes>

área a que a notícia pertence como: desporto, sociedade, política, economia, entre outros.

Detalhamos nas subsecções apresentadas de seguida a abordagem utilizada para o processo de ligação de pontos entre os agrupamentos. Este foi realizado em duas etapas: cálculo da distância entre termos relevantes e determinação das ligações entre agrupamentos.

#### *Similaridade de termos relevantes*

Começamos por fazer a normalização dos termos relevantes. Para as palavras isoladas, expressões, entidades e personalidades, o texto é convertido para letra minúscula. Para as palavras isoladas que são constituídas apenas por *uni-grams* também se efetua a redução ao seu radical. Após a normalização do texto, é efetuado o cálculo da similaridade entre os termos de cada agrupamentos através do cálculo da distância entre: palavras isoladas, expressões, entidades e personalidades.

Para o cálculo da similaridade entre palavras isoladas, entidades e personalidades, consideramos o peso de cada palavra individual no agrupamento que é dada pelas Expressões 2 e 3.

$$D_1(a, b) = 0.3 \times \frac{|k_a| \wedge |k_b|}{\max(|k_a|, |k_b|)} + 0.7 \times \frac{\sum_{i=1}^{|k_a|} (\sum_{j=1 \wedge a_j=b_i}^{|k_b|} Wk_{a_j} \times Wk_{b_i})}{|k_a| \wedge |k_b|} \quad (2)$$

$$D_2(a, b) = \frac{|k_a| \wedge |k_b|}{\max(|k_a|, |k_b|)} \times \frac{\sum_{i=1}^{|k_a|} (\sum_{j=1 \wedge a_j=b_i}^{|k_b|} Wk_{a_j} \times Wk_{b_i})}{|k_a| \wedge |k_b|} \quad (3)$$

Onde:

$Wk_{a_j}$  é o peso da palavra-chave  $j$  no agrupamento  $a$ ;

$Wk_{b_i}$  é o peso da palavra-chave  $i$  no agrupamento  $b$ ;

$|k_a|$  e  $|k_b|$  são o número de palavras-chave iguais entre os agrupamentos  $a$  e  $b$ ;

$\max(|k_a|, |k_b|)$  é o número máximo de palavras-chave distintas.

As distâncias  $D_1(a, b)$  e  $D_2(a, b)$  têm em conta a percentagem de termos em comum entre os dois agrupamentos e a relação dos pesos que os termos em comum têm nos seus agrupamentos.  $D_1(a, b)$  estabelece um peso entre as duas parcelas, dando um maior relevo à parcela que mede o relacionamento dos pesos das palavras em comum; em  $D_2(a, b)$  não existem pesos associados às parcelas, mas sim, uma relação entre elas.

Para o cálculo da similaridade entre as expressões relevantes a abordagem utilizada foi distinta. Para este caso, a normalização incluiu um passo adicional,

remoção das *stop-words*. Após esta tarefa foi construída uma *string* com todas as expressões pertencentes a cada agrupamento, não considerando para este tipo de termo relevante o seu peso. O cálculo da similaridade entre as expressões foi baseado num algoritmo de edição de distância o *qgrams* (Ullmann 1977) ( $q = 3$ ).

#### *Determinação das ligações entre agrupamentos*

Esta etapa tem como objetivo determinar a partir dos valores de similaridade calculados anteriormente quais as ligações mais relevantes. É a partir destas ligações que se formam as cadeias noticiosas.

Para a ligação de agrupamentos, utilizamos algoritmos de aprendizagem supervisionada. Estes algoritmos recebem um conjunto de treino manualmente anotado com ligações relevantes entre agrupamentos, sobre o qual vão inferir regras para determinar, a existência de ligações válidas e relevantes. Utilizamos como características (*features*) a distância entre as palavras isoladas, expressões, entidades e personalidades. Os algoritmos utilizados foram: *Support Vector Classifier* (SVC), *SVC Linear*, *Decision Tree* e o *Random Forest*.

Ao longo desta secção apresentamos a metodologia utilizada na deteção de notícias duplicadas e na geração automática de cadeias noticiosas.

#### [4] RECURSOS LINGUÍSTICOS

Nesta secção caracterizamos o conjunto de dados e as fontes de conhecimento externo utilizadas na elaboração deste trabalho.

##### [4.1] *Caracterização do conjunto de dados*

Para a realização deste trabalho foram utilizadas notícias publicadas *online*, escritas na língua portuguesa e provenientes de diversas fontes noticiosas da imprensa portuguesa. O conjunto de dados compreende mais de 4 milhões de notícias publicadas entre 2008 e 2014.

As notícias são provenientes de 73Número de fontes com mais de 100 notícias publicadas. fontes noticiosas distintas e compostas em média<sup>5</sup> por: 9 palavras no título; 204 palavras no conteúdo; 10 frases no conteúdo.

Na imprensa portuguesa são publicadas *online* diariamente aproximadamente 2.500 notícias<sup>6</sup>. A Figura 4 representa a distribuição de notícias durante mês de Março de 2014. Através da observação da mesma é possível constatar que tendencialmente são publicadas menos notícias durante o fim de semana.

Estima-se que aproximadamente 45%<sup>7</sup> das notícias publicadas diariamente se-

[5] Análise de aproximadamente 74000 notícias selecionadas de um mês aleatório de 2014.

[6] Dados relativos às notícias publicadas na imprensa portuguesa, no formato digital, no mês de Março de 2014

[7] Número médio de notícias *online* diárias duplicadas, publicadas na imprensa portuguesa, de 10 a 15 de Março de 2014

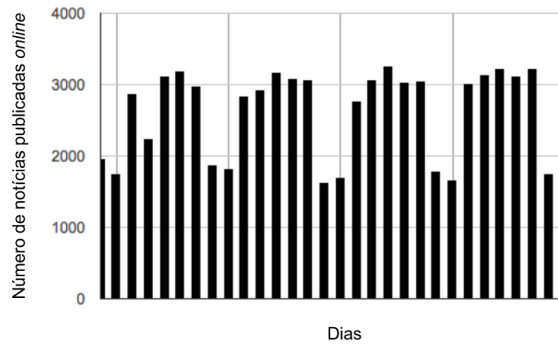


FIGURA 4: Número de notícias publicadas por dia no mês de Março de 2014.

jam duplicadas ou quase duplicadas. A relação entre o número de notícias publicadas mensalmente com o número de notícias utilizadas para a criação dos agrupamentos pode ser visualizada na Figura 5. Para os primeiros oito meses de 2014 o número médio de notícias por grupo é de 3,8, os dados referentes ao número médio de notícias por grupo relativo a cada mês pode ser observado na Figura 6.

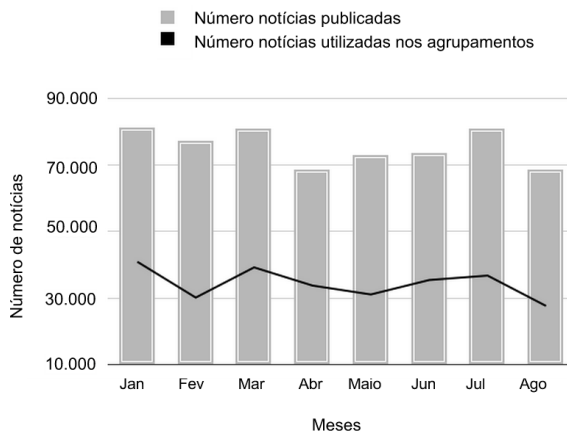


FIGURA 5: Relação entre o número de notícias publicadas por mês com o número de notícias utilizadas na criação dos agrupamentos (Janeiro a Agosto de 2014)

Na Figura 7 podemos constatar que maioritariamente os grupos são constituídos por 2 notícias similares. É possível observar que o número de grupos existentes é inversamente proporcional ao número de notícias que o compõe.

Definimos nove categorias associadas aos agrupamentos que são as categorias tipicamente usadas nos media digitais para organizar as notícias publicadas

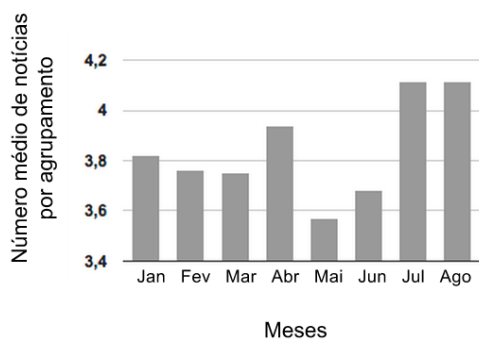


FIGURA 6: Número médio de notícias por agrupamento (Janeiro a Agosto de 2014)

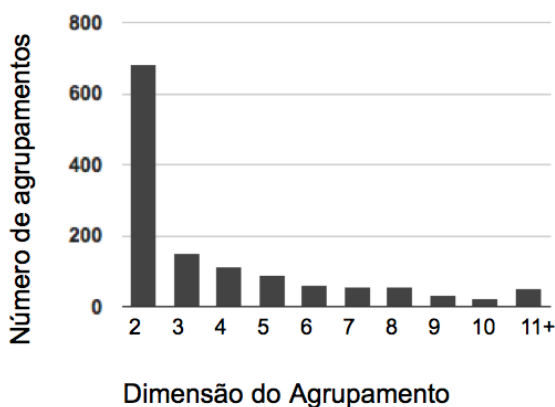


FIGURA 7: Constituição dos agrupamentos (seleção aleatória de 5 dias de 2014)

*online*: política, economia, desporto, saúde, ciências e tecnologias, sociedade, cultura, local e educação. Dos agrupamentos com apenas uma categoria associada a distribuição dos mesmos por áreas pode ser observado na Figura 8. É possível observar que a categoria com maior expressão é a categoria desporto (54.4%) e assim sucessivamente.

#### [4.2] *Enunciação de fontes de conhecimento externo*

No decorrer deste trabalho foram utilizadas as seguintes fontes de conhecimento:

**Lista *stop-words*:** Lista de *stop-words* específica para a língua portuguesa disponibilizada pela snowball.

**Verbetes:** O Verbetes é um sistema de recolha automática de informação a par-

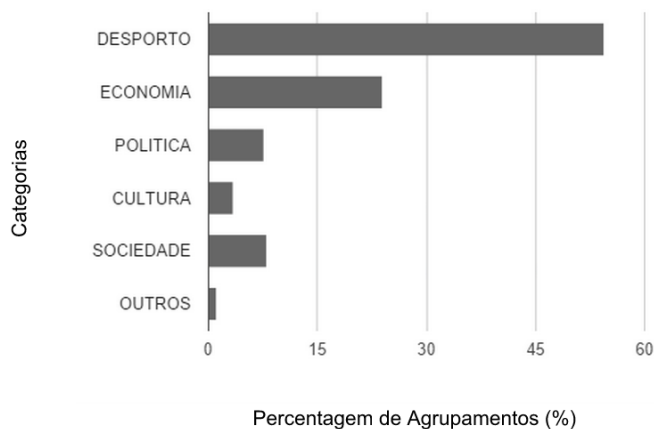


FIGURA 8: Distribuição dos agrupamentos por categoria

tir das notícias. Para este trabalho utilizamos uma lista de personalidades extraída deste sistema.

Lista de *Tags* e Categorias: Lista elaborada manualmente por jornalistas que relaciona a *tag* associada à notícia com a sua categoria principal.

Nesta secção foi caracterizado o conjunto de dados e as fontes de conhecimento externo utilizadas na elaboração deste trabalho.

## [5] EXPERIMENTAÇÃO

Nesta secção são referidas as diferentes métricas de avaliação utilizadas e descrito o conjunto de experiências realizadas.

### [5.1] Métricas de Avaliação

Para avaliar o módulo de similaridade (ver secção [3.1]) e ligações entre agrupamentos (ver secção [3.5.2]), foram utilizadas quatro métricas de avaliação: a precisão (*precision*), a abrangência (*recall*), a *accuracy* e a *F-measure* ( $F_1$ ). No contexto deste trabalho, a precisão indica a taxa de notícias consideradas similares que realmente o são e a taxa de ligações efetuadas entre agrupamentos que realmente existem. A abrangência (*recall*) indica-nos, neste contexto, taxa de notícias duplicadas encontradas face às realmente existentes mas que não conseguimos identificar manualmente. A medida  $F_1$  estabelece uma relação entre a precisão e a abrangência. A *accuracy* indica-nos a avaliação geral do sistema.

A avaliação aos termos relevantes focou-se em avaliar, dos termos extraídos, quais são de facto realmente representativos da notícia. A avaliação foi realizada

usando a Expressão 4. A avaliação geral do sistema é dada pelo somatório percentagem de termos representativos das notícias analisadas, Expressão 5.

$$E(n_i) = \frac{\text{Termos}_{\text{Representativos}}}{\text{Termos}_{\text{Atribuídos}}} \quad (4)$$

$$\text{Avaliação} = \frac{\sum_{i=1}^{\|N\|} (E(n_i))}{\|N\|} \quad (5)$$

Onde:

$\text{Termos}_{\text{Representativos}}$  corresponde ao número de termos relevantes ou entidades atribuídos pelo método, que realmente representam o conteúdo noticioso;

$\text{Termos}_{\text{Atribuídos}}$  corresponde ao número total de termos relevantes ou entidades atribuídas ao documento;

$\|N\|$ : número de notícias da coleção N;

$n_i$ : corresponde à notícia de índice i do conjunto de notícias N.

### [5.2] *Enunciação e definição das experiências*

Nesta secção são apresentadas as cinco experiências realizadas. Começamos por apresentar três experiências relativas à determinação da similaridade entre notícias. Na primeira experiência pretendemos perceber qual o algoritmo mais adequado ao cálculo da similaridade entre notícias. A segunda experiência visa entender qual a influência do fator tempo neste domínio, ou seja, se as notícias duplicadas ou quase duplicadas surgem em intervalos temporais longos ou curtos. Por fim a terceira experiência tem como objetivo perceber qual o método de aprendizagem supervisionado mais apto para a determinação da similaridade entre notícias.

A quarta experiência enunciada está relacionada com os termos chaves extraídos. Por fim a quinta experiência refere-se às ligações entre agrupamentos.

Será usado  $Exp_{i,j}$  para representar a  $j$ -ésima configuração de parâmetros para a experiência  $i$ .

### *Similaridade - Algoritmos de Edição de Distância*

A similaridade entre notícias é obtida através do cálculo da:

- Similaridade do título (ST) que corresponde à percentagem de semelhança entre os títulos;
- Similaridade do 1º parágrafo (SB) que corresponde ao resultado de comparação entre a parte das notícias que foca o evento em si;
- Similaridade de conteúdo noticioso (SC) que corresponde ao resultado da comparação do corpo das respetivas notícias.

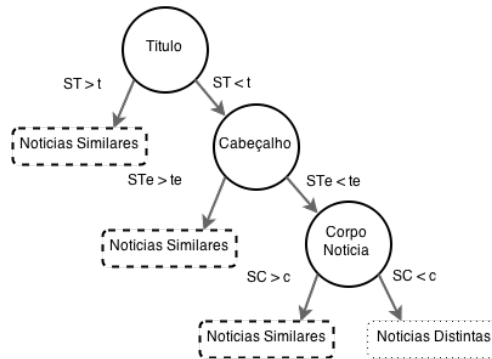


FIGURA 9: Árvore de decisão elaborada para verificar se um par de notícias é ou não similar.

Esta experiência —  $Exp_1$  — visou perceber qual o algoritmo com o melhor desempenho para o cálculo da similaridade entre pares de notícias. Esta experiência foi efetuada sobre uma estrutura em forma de árvore de decisão, representada na Figura 9. Esta foi criada manualmente, onde os valores  $t$ ,  $f$ ,  $c$ , correspondem aos valores de decisão para o título, foco e conteúdo da notícia. L, H, J correspondem respetivamente aos algoritmos Levenshtein, Hamming e Jaro. A parametrização usada nesta experiência encontra-se enunciada na Tabela 3. Por exemplo, a  $Exp_{1,1}$  é efetuada individualmente para os algoritmos Levenshtein, Hamming e Jaro, com um *threshold* de 0.6 para  $t$ ,  $f$  e  $c$ . As diferentes experiências visam perceber a influência que os diferentes *thresholds* têm nos algoritmos.

Exp	Algoritmos	t	f	c
1, 1	L H J	0,60	0,60	0,60
1, 2	L H J	0,70	0,60	0,60
1, 3	L H J	0,70	0,70	0,60
1, 4	L H J	0,70	0,70	0,70
1, 5	L H J	0,80	0,70	0,70
1, 6	L H J	0,80	0,80	0,70
1, 7	L H J	0,80	0,80	0,80

TABELA 3: Parametrização para a experiência do cálculo da similaridade.

Para a realização desta experiência foram comparadas aleatoriamente 124750 notícias, para um dia aleatório de 2014.

#### Similaridade - Fator Tempo

A experiência sobre o fator tempo (intervalo temporal) tem como objetivo verificar a influência do intervalo temporal no que diz respeito à identificação e clas-

sificação de notícias similares. Para tal, foram considerados cinco intervalos de tempo distintos: 3, 6, 12, 24, 48 horas; e foram utilizados quatro métodos de classificação para a determinação da similaridade: SVC, SVC Linear, *Decision Tree* e o *Random Forest*. Esta experiência foi elaborada utilizando uma técnica de avaliação cruzada, o *k-fold cross validation* ( $k = 5$ ). O conjunto de dados utilizado resulta da seleção aleatória de 500 notícias de dois dias distintos e consecutivos, anotadas manualmente.

#### *Similaridade - Decisão da similaridade entre notícias*

Foi efetuada uma experiência com o objetivo de perceber qual o algoritmo de aprendizagem supervisionada com o melhor desempenho na determinação da similaridade entre pares de notícias. A experiência foi efetuada em 500 notícias selecionadas de forma aleatória de um dia aleatório de 2014.

#### *Extração de Termos relevantes*

Esta experiência tem como objetivo testar a abordagem utilizada para a extração de termos chave (palavras isoladas, expressões e entidades). Para a realização desta experiência foi selecionado aleatoriamente um dia de cada mês do ano 2012. De cada dia foi selecionado um intervalo de três horas, e dessas três horas foram selecionadas aleatoriamente dez notícias sobre as quais se efetuou a inspeção manual das palavras-chave atribuídas.

#### *Ligações entre agrupamentos*

Para a determinação das ligações entre agrupamentos de notícias, é realizado o cálculo da distância entre: palavras isoladas, expressões, entidades e personalidades.

Esta experiência — *Exp<sub>2</sub>* — tem como objetivo avaliar qual a abordagem mais adequada para o cálculo da similaridade e qual o método de aprendizagem supervisionado mais eficiente para a determinação das ligações. Todas as experiências consideraram o cálculo distância pelo algoritmo Q-grams, para as expressões. A avaliação resultante das diferentes experiências realizadas entre grupos de notícias ao longo do tempo, para a formação de ligações entre agrupamentos de notícias, encontra-se na Tabela 4. O conjunto de dados é composto por agrupamentos pertencentes aos meses de março e abril de 2014. Desses agrupamentos, foram selecionados aleatoriamente 10 cadeias de notícias com tamanho variável para cada uma das seguintes categorias: desporto, economia, política, cultura e sociedade. O conjunto de dados compreende, em média, 317 comparações por categoria.

Nesta secção foram apresentadas as diferentes métricas de avaliação utilizadas e descrito o conjunto de experiências realizadas.

Exp	Palavras	Entidades	Personalidades
2, 1	$D_1$	$D_2$	$D_1$
2, 2	$D_2$	$D_2$	$D_1$
2, 3	$D_1$	$D_1$	$D_1$
2, 4	$D_1$	$D_2$	$D_2$

TABELA 4: Descrição das experiências para o cálculo das ligações.

## [6] RESULTADOS E ANÁLISE

## [6.1] Experiências

*Similaridade - Algoritmos de Edição de Distância*

Os resultados obtidos nesta experiência —  $Exp_1$  — podem ser observados na Tabela 5. Desta tabela excluímos os resultados obtidos para algoritmo Jaro, devido ao seu desempenho constante.

Exp	Levensthein			Hamming		
	P	R	F	P	R	F
1, 1	0,941	0,761	0,841	0,941	0,289	0,442
1, 2	0,950	0,655	0,775	0,940	0,284	0,436
1, 3	0,951	0,645	0,769	0,940	0,284	0,436
1, 4	<b>0,972</b>	<b>0,637</b>	<b>0,770</b>	0,940	0,284	0,436
1, 5	0,965	0,507	0,665	0,939	0,279	0,430
1, 6	0,964	0,483	0,643	0,939	0,279	0,430
1, 7	0,962	0,463	0,625	0,938	0,279	0,430

TABELA 5: Resultados dos testes aos algoritmos de edição de distância.

Da comparação entre o algoritmo *Levensthein* e o *Hamming* em  $Exp_{1,1}$  podemos verificar que os valores da precisão são semelhantes, o que indica que a percentagem de notícias consideradas similares que realmente o são (*true positive*) é igual. Para o mesmo caso podemos verificar uma melhoria para o algoritmo *Levensthein* para o *recall*.

*Similaridade - Fator Tempo*

O resultado obtido desta análise pode ser observado no gráfico apresentado na Figura 10. Como podemos constatar pela análise do gráfico, o aumento do intervalo de tempo faz com que os valores se tornem constantes. Ao alargar o intervalo de tempo de 24 para 48 horas não há variação nos valores de *precision*, *recall* e da métrica  $F_1$ .

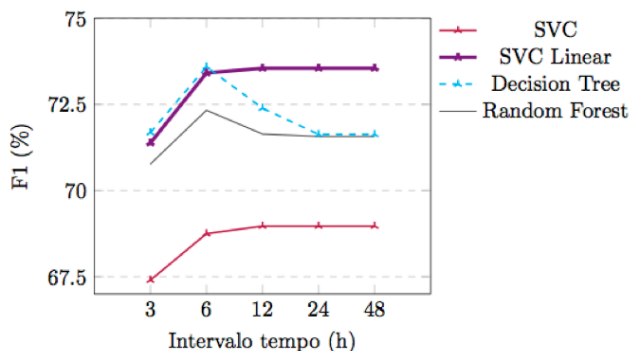


FIGURA 10: Valor da métrica  $F_1$  obtido pelos diferentes algoritmos nos diferentes intervalos de tempo.

#### Similaridade - Determinação Semelhança

Os resultados dos algoritmos de aprendizagem supervisionados na determinação da similaridade podem ser observados na Tabela 6. Pela visualização da tabela é possível constatar que apesar do valor do *recall* ser baixo, o valor obtido pela *precision* é alto, o que garante a elevada qualidade da informação recolhida. O algoritmo que apresenta um melhor desempenho é o SVC Linear.

	P	R	$F_1$	A
<i>Decision Tree</i>	0,863	0,679	0,760	0,998
<i>SVC</i>	0,931	0,508	0,657	0,997
<b><i>SVC Linear</i></b>	<b>0,938</b>	<b>0,561</b>	<b>0,702</b>	<b>0,998</b>
<i>Random Forest</i>	0,803	0,542	0,647	0,998

TABELA 6: Resultado médio das métricas de avaliação obtidas pelo *k fold cross validation*.

#### Extração de Termos Relevantes

Os resultados da extração de termos relevantes pode ser observado na Tabela 7. A representatividade das palavras extraídas face à informação contida nas notícias é de: 73,2% para as palavras isoladas, 76,2% para as expressões e 80.4% para as entidades.

#### Ligações entre agrupamentos

Na Tabela 8 são apresentados os resultados da precisão para as ligações entre agrupamentos. A partir da análise dos resultados podemos verificar que o método

Avaliação	
Palavras	0,732
Expressões	0,762
Entidades	0,804

TABELA 7: Avaliação dos termos chave.

com um melhor desempenho é o *SVC Linear* e que em 93.3% dos casos analisados as ligações entre notícias são verdadeiras.

Exp	SVC	Decision	Random
		Tree	Forest
2, 1	<b>0.931</b>	0.849	0.859
2, 2	0.921	0.821	0.852
2, 3	0.906	0.764	0.824
2, 4	0.931	0.834	0.858

TABELA 8: Valor da precisão na determinação de ligações entre agrupamentos de notícias.

## [6.2] Análise dos resultados obtidos

### *Similaridade - Algoritmos de Edição de Distância*

Dos resultados obtidos nestas experiências, podemos observar na Tabela 5 que o algoritmo *Jaro* é o que apresenta a nível global um pior desempenho. No entanto, segundo estudos realizados, este algoritmo tem um melhor desempenho aquando da comparação de pequenas *strings* (Bilenko et al. 2003), o que não acontece no domínio das notícias. Os valores da precisão entre a utilização do algoritmo *Levenshtein* e o *Hamming* são muito próximos, obtendo o algoritmo *Levenshtein* ao longo das diferentes experiências um melhor desempenho nesta métrica. Comparando as restantes métricas de avaliação, para estes dois algoritmos, é possível observar que o *Levenshtein* obtém uma melhor performance a nível da métrica *recall*, o que significa que consegue detetar mais casos do que o *Hamming*. Uma razão para que isto suceda está relacionado com uma particularidade deste último algoritmo que é a comparação de *strings* do mesmo comprimento; a nível da métrica  $F_1$ , também o *Levenshtein* obtém um melhor resultado. Através da análise efetuada a estes três algoritmos é possível concluir que o *Levenshtein* é o algoritmo mais indicado para o cálculo da similaridade entre pares de notícias.

### *Similaridade - Fator Tempo*

Um fator importante para a comparação das notícias é a sua data de publicação. Dos resultados apresentados, os algoritmos que apresentam uma melhor precisão são o SVC e o SVC Linear. Sendo que destes dois, o SVC Linear tem um desempenho superior a nível do *recall* e da métrica  $F_1$ . Relativamente à questão temporal, podemos perceber, que todos os algoritmos têm um comportamento semelhante à medida que o intervalo temporal aumenta. Pela análise do gráfico é possível verificar que não existem variações dos resultados quando o intervalo de tempo é alargado de 24 para 48 horas. Isto pode indicar que os casos de notícias duplicadas ou quase duplicadas surgem quase sempre num intervalo inferior ou igual a 24 horas. Com base nos resultados obtidos constatou-se que um intervalo de tempo de 24 horas era o mais adequado para a comparação de notícias.

### *Similaridade - Determinação Semelhança*

Para a determinação da similaridade das notícias, os algoritmos que apresentam um melhor desempenho, considerando o  $\Delta T = 24$  horas, são: a nível da precisão o SVC Linear (93.8%) e SVC (93.1%); em relação à métrica *recall* e a métrica  $F_1$  o *Decision Tree* (67.9% e 76.0%) e SVC Linear (56.1% e 70.2%). Comprando o desempenho dos diferentes algoritmos para as diferentes fases de processamento e tendo em conta as opções escolhidas a nível de algoritmo de cálculo da similaridade e intervalo de tempo considerado, podemos constatar que o algoritmo que apresenta um melhor desempenho a nível global é o SVC Linear.

### *Extração de Termos Relevantes*

Foi efetuada uma avaliação manual à relevância das palavras-chave extraídas. A avaliação consistiu em analisar a representatividade dos termos extraídos do texto em relação ao conteúdo da notícia. O resultado da avaliação efetuada a estes elementos pode ser observada na Tabela 7. Os resultados indicam que 73,2% das palavras, 76,2% das expressões e 80,5% das entidades são representativas do conjunto. Através da análise ao teor dos termos extraídos foi possível constatar que as palavras relevantes dizem respeito a palavras que descrevem de uma forma muito genérica o conteúdo da notícia; e, por sua vez, as expressões relevantes já transmitem com mais especificidade o assunto da notícia. Consideremos de novo o exemplo das notícias sobre o desaparecimento do Avião da Malaysia Airlines, temos como palavra relevante *avião* e como expressão *avião Malaysia Airlines*.

### *Ligações entre agrupamentos*

Da análise aos resultados obtidos pela comparação da  $Exp_{2,1}$  com a  $Exp_{2,2}$ , em que o que foi modificada a fórmula de cálculo da distância entre as palavras isoladas, é possível observar que todos os algoritmos conseguem um melhor desempenho considerando a fórmula de cálculo  $D1$ ; face à diferença da precisão entre os

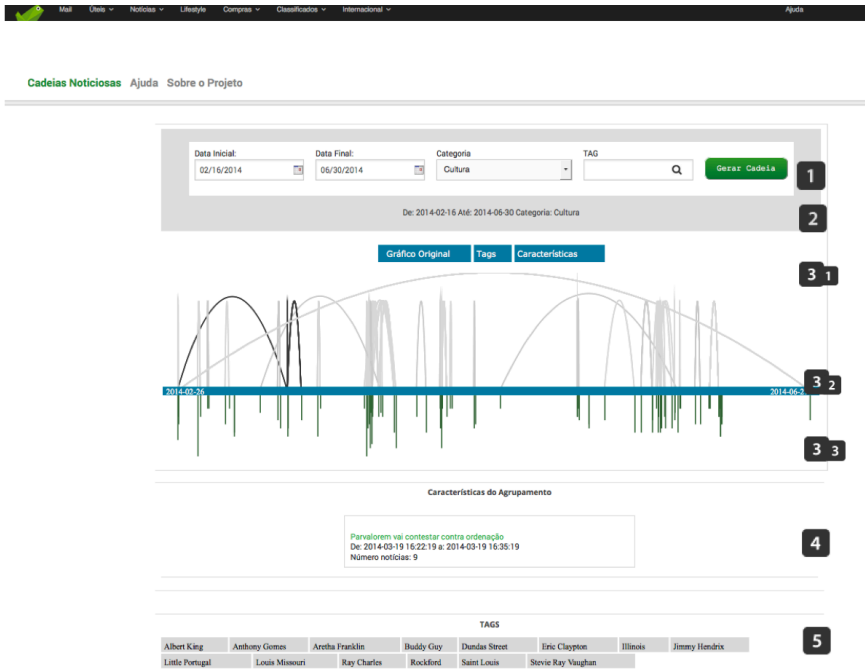


FIGURA 11: Interface do sistema.

algoritmos: 0.010 no SVC Linear; 0.028 no *Decision Tree* e 0.007 no *Random Forest*. Estabelecendo uma comparação entre as experiências  $Exp_{2,1}$  e  $Exp_{2,3}$ , que divergem apenas na fórmula de cálculo da distância entre as entidades, temos que: a utilização da fórmula  $D2$  no cálculo da proximidade de entidades entre dois conjuntos reflete um aumento de desempenho. Confrontando os valores obtidos para a experiência  $Exp_{2,1}$  em relação à experiência  $Exp_{2,3}$  é possível constatar que independentemente do algoritmo de aprendizagem supervisionada os resultados da  $Exp_{2,1}$  são os que apresentam um melhor desempenho. Os valores da precisão obtidos para a experiência  $Exp_{2,1}$  e a  $Exp_{2,4}$  são bastante próximos. Esta experiência difere da primeira na fórmula de cálculo da distância entre personalidades. A partir dos resultados obtidos conclui-se que as personalidades não têm grande impacto na formação das ligações comparativamente com as palavras isoladas e entidades, uma vez que a mudança de cálculo para este elemento não reflete uma variação considerável no resultado. Podemos ainda observar que o melhor desempenho continua a ser o resultante da experiência  $Exp_{2,1}$ . Após o estudo dos resultados obtidos, podemos concluir que a fórmula mais apta para cada tipo de palavra-chave é a seguinte:  $D1$  — personalidades e palavras isolada;  $D2$  — entidades; sendo que esta combinação se refere à experiência  $Exp_{2,1}$ . Comparando os resultados obtidos pelos diferentes métodos de aprendizagem supervisionada

para  $Exp_{2,1}$  podemos observar que o método com um melhor desempenho é o SVC Linear (93.1%).

## [7] INTERFACE

Desenvolvemos uma interface *web* para permitir ao leitor a navegação entre cadeias de notícias. A interface que elaboramos pode ser observada na Figura 11.

A interface é composta por cinco secções distintas. A primeira secção permite que o utilizador defina as características das cadeias de notícias a visualizar. É permitido definir o intervalo temporal, a categoria das notícias e ainda as palavras-chave. A segunda secção, informa o utilizador quais as características das histórias que estão representadas na interface.

As histórias são representadas visualmente na terceira secção. O gráfico com a representação das histórias pode ser repartido em três elementos interconectados. Começando pela parte inferior do gráfico, em 3.3, as linhas representam os agrupamentos de notícias existentes. O comprimento destas barras varia consoante o número de notícias que compõe cada agrupamento. Na parte superior do gráfico, em 3.1, os arcos representam as ligações existentes entre os agrupamentos de notícias (em 3.3). A barra situada em 3.2 posiciona temporalmente a informação apresentada (em 3.1 e 3.3).

A informação presente na quarta secção varia consoante a interação do utilizador com o gráfico. Se o utilizador navegar sobre a parte 3.3 do gráfico a informação que aparece nesta secção informa o utilizador das características do agrupamento. Porém, se o utilizador navegar na parte 3.1 do gráfico, a informação contida na secção quatro informará o utilizador da história noticiosa. A quinta secção apresenta a lista de palavras-chave mais relevantes dentro do intervalo temporal considerado.

A Figura 12 apresenta parte de uma cadeia obtida pelo sistema (para a categoria Cultura de 31 de Janeiro até 17 de Fevereiro de 2014). A interface será brevemente lançada ao público.

## [8] CONCLUSÕES E TRABALHO FUTURO

Este artigo pretende definir e avaliar técnicas para o encadeamento automático de notícias com vista à construção de histórias noticiosas temporais. A abordagem utilizada para a criação das cadeias baseia-se: (i) deteção de notícias (quase) duplicadas e (ii) a criação de ligações entre notícias relacionadas ao longo do tempo.

Para a deteção de notícias duplicadas usamos uma abordagem baseada na semântica para o cálculo da similaridade entre notícias. Foi também utilizado um algoritmo de aprendizagem supervisionado na determinação da semelhança entre as mesmas. Adicionalmente, as notícias incluem informação temporal e, tal como acreditávamos, existe um intervalo onde há uma maior tendência para o apare-

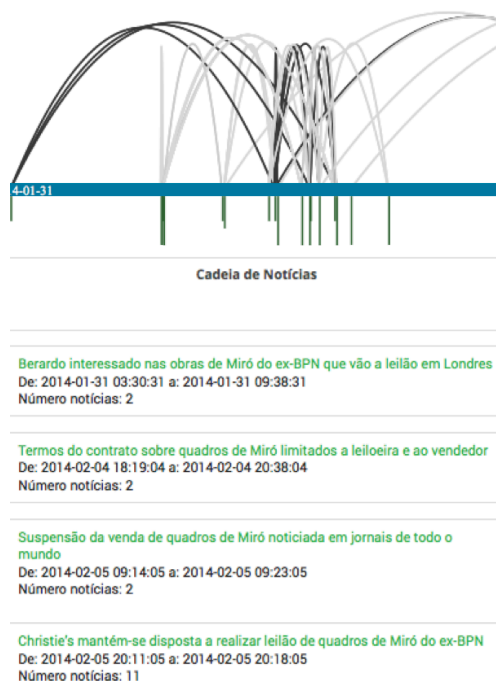


FIGURA 12: Parte de uma cadeia obtida pelo sistema.

cimento de notícias cujo grau de similaridade aponta para a (quase) duplicação. O nosso estudo indicou que tendencialmente as notícias consideradas duplicadas aparecem num intervalo inferior a 24 horas. A nossa abordagem, para a determinação de notícias cujo grau de similaridade as classifica como (quase) duplicadas, num intervalo de tempo de 24 horas, obteve uma precisão de 93.8% quando usado o par Levenshtein, SVC Linear.

Para a criação de ligações entre grupos de notícias similares, a nossa abordagem consistiu na medição do grau de semelhança entre os diferentes grupos. Para esta etapa, sugerimos uma nova forma de medição de distância que tem em conta os termos em comum e a expressão de cada termo nos agrupamentos de notícias similares. Para a determinação das ligações, foram também utilizados algoritmos de aprendizagem supervisionada. A abordagem proposta para a realização desta segunda tarefa apresenta uma precisão de 93.1%. Este resultado, não representa, no entanto a precisão global do sistema, uma vez que há propagação de erro entre as várias etapas.

Como trabalho futuro será importante criar testes mais exaustivos e objetivos para as cadeias de notícias. Tais testes, consistirão, entre outros melhoramentos, na medição da familiaridade do leitor com um tema em específico antes e depois da utilização da plataforma e na medição do erro propagado pelo sistema.

Também pretendemos melhorar o sistema através da:(i) introdução de sumários das notícias, (ii) deteção de novos factos e (iii) hierarquização de notícias.

#### AGRADECIMENTOS

Agradecemos a colaboração do Labs SAPO UP pela disponibilização dos dados utilizados neste trabalho.

#### REFERÊNCIAS

- Allan, James, Jaime G. Carbonell, George Doddington, Jonathan Yamron & Yiming Yang. 1998a. Topic detection and tracking pilot study final report. Em *Proceedings of the DARPA Broadcast News Transcription and Understanding Workshop*, 194–218.
- Allan, James, Ron Papka & Victor Lavrenko. 1998b. On-line new event detection and tracking. Em *Proceedings of the 21st annual international ACM SIGIR conference on research and development in information retrieval*, 37–45. ACM.
- Banerjee, Somnath, Krishnan Ramanathan & Ajay Gupta. 2007. Clustering short texts using Wikipedia. Em *Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '07*, 787–788. ACM.
- Bilenko, Mikhail, Raymond Mooney, William Cohen, Pradeep Ravikumar & Stephen Fienberg. 2003. Adaptive Name Matching in Information Integration. *IEEE Intelligent Systems* 18(5). 16–23.
- Elmagarmid, Ahmed K., Panagiotis G. Ipeirotis & Vassilios S. Verykios. 2007. Duplicate record detection: A survey. *IEEE Transactions on Knowledge and Data Engineering* 19(1). 1–16.
- Garcia, Marcos & Pablo Gamallo. 2013. FreeLing e TreeTagger: um estudo comparativo no âmbito do Português. Relatório técnico. ProLab Technical Report, vol. 01. [http://gramatica.usc.es/~gamallo/artigos-web/PROLNAT\\_Report\\_01.pdf](http://gramatica.usc.es/~gamallo/artigos-web/PROLNAT_Report_01.pdf).
- He, Matthew X., Sergei V. Petoukhov & Paolo E. Ricci. 2004. Genetic code, Hamming distance and stochastic matrices. *Bulletin of mathematical biology* 66(5). 1405–1421.
- Kumar, J. Prasanna & P. Govindarajulu. 2009. Duplicate and Near Duplicate Documents Detection: A Review. *European Journal of Scientific Research* 32. 514–527.
- Kumar, Ravi, Uma Mahadevan & Alan D. Sivakumar. 2004. A Graph-theoretic Approach to Extract Storylines from Search Results. Em *Proceedings of the tenth ACM*

- SIGKDD international conference on knowledge discovery and data mining, KDD'04*, 216–225. ACM.
- Lawrie, Dawn & W Bruce Croft. 2000. Discovering and Comparing Topic Hierarchies. Em *Proceedings of the RIAO 2000 conference*, 314–330.
- Levenshtein, Vladimir. 1965. Binary Codes Capable of Correcting Deletions, Insertions and Reversals. *Doklady Akademii Nauk SSSR* 163. 845–848.
- Li, Tao, Shenghuo Zhu & Mitsunori Ogihara. 2007. Hierarchical document classification using automatically generated hierarchy. *Journal of Intelligent Information Systems* 29(2). 211–230.
- Lin, Chen, Chun Lin, Jingxuan Li, Dingding Wang, Yang Chen & Tao Li. 2012. Generating Event Storylines from Microblogs. Em *Proceedings of the 21st ACM International Conference on Information and Knowledge Management, CIKM'12*, 175–184. ACM.
- Lin, Fu-ren & Chia-Hao Liang. 2008. Storyline-based summarization for news topic retrospection. *Decision Support Systems* 45(3). 473–490.
- McKeown, Kathleen R., Regina Barzilay, David Evans, Vasileios Hatzivassiloglou, Judith L. Klavans, Ani Nenkova, Carl Sable, Barry Schiffman & Sergey Sigelman. 2002. Tracking and summarizing news on a daily basis with Columbia's Newsblaster. Em *Proceedings of the second international conference on Human Language Technology Research*, 280–285.
- Mei, Qiaozhu & ChengXiang Zhai. 2005. Discovering Evolutionary Theme Patterns from Text: An Exploration of Temporal Text Mining. Em *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining, KDD'05*, 198–207. ACM.
- Nallapati, Ramesh, Ao Feng, Fuchun Peng & James Allan. 2004. Event threading within news topics. Em *Proceedings of the thirteenth ACM international conference on Information and knowledge management*, 446–453. ACM.
- Oliveira, Pedro. 2008. Ptstemmer - a stemming toolkit for the portuguese language. Obtido em Maio 2014. <http://code.google.com/p/ptstemmer>.
- Pedregosa, Fabian, Gael Varoquaux, Alexandre Gramfort, Vicent Michel, Bertrand Thirion, Oliver Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vicent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot & Édouard Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12. 2825–2830.

- Qamra, Arun, Belle Tseng & Edward Y Chang. 2006. Mining blog stories using community-based and temporal clustering. Em *Proceedings of the 15th ACM international conference on Information and knowledge management*, 58–67. ACM.
- Schmid, Helmut. 1994. Probabilistic Part-of-Speech Tagging Using Decision Trees. Em *International Conference on New Methods in Language Processing*, 44–49.
- Shahaf, Dafna & Carlos Guestrin. 2010. Connecting the Dots Between News Articles. Em *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '10*, 623–632. New York, NY, USA: ACM.
- Sun, Aixin & Ee-Peng Lim. 2001. Hierarchical text classification and evaluation. Em *Data Mining, 2001. ICDM 2001, Proceedings IEEE International Conference on*, 521–528. IEEE.
- Ullmann, Julian R. 1977. A binary n-gram technique for automatic correction of substitution, deletion, insertion and reversal errors in words. *The Computer Journal* 20(2). 141–147.
- Vadrevu, Srinivas, Choon Hui Teo, Suju Rajan, Kunal Punera, Byron Dom, Alexander J. Smola, Yi Chang & Zhaohui Zheng. 2011. Scalable clustering of news search results. Em *Proceedings of the fourth ACM International Conference on Web Search and Data Mining, wsdm'11*, 675–684. ACM.
- Waterman, Michael S., Temple F. Smith & William A. Beyer. 1976. Some biological sequence metrics. *Advances in Mathematics* 20(3). 367–387.
- Yancey, William E. 2005. Evaluating string comparator performance for record linkage. Relatório técnico. Statistical Research Division. <http://www.census.gov/srd/papers/pdf/rrs2005-05.pdf>.

## CONTACTOS

Carla Abreu

Faculdade de Engenharia da Universidade do Porto  
[cfma@fe.up.pt](mailto:cfma@fe.up.pt)

Jorge Teixeira

Faculdade de Engenharia da Universidade do Porto  
[jft@fe.up.pt](mailto:jft@fe.up.pt)

Eugénio Oliveira

Faculdade de Engenharia da Universidade do Porto  
[eco@fe.up.pt](mailto:eco@fe.up.pt)