

Segmentação topical automática de atas de reunião

Ovídio José Francisco
ovidiojf@gmail.com

RESUMO

Keywords

1. INTRODUÇÃO

Frequentemente atas de reunião tem a característica de apresentar um texto com poucas quebras de parágrafo e sem marcações de estrutura, como capítulos, seções ou quaisquer indicações sobre o tema do texto.

A tarefa de segmentação textual consiste dividir um texto em partes que contenham um significado relativamente independente. Em outras palavras, é identificar as posições onde há uma mudança significativa de tópicos.

É útil em aplicações que trabalham com textos sem quebras de assunto, ou seja, não apresentam parágrafos, seções ou capítulos, como transcrições automáticas de áudio e grandes documentos que contêm assuntos não idênticos como atas de reunião e notícias.

O interesse por segmentação textual tem crescido em em aplicações voltadas a recuperação de informação e sumarização de textos. Essa técnica pode ser usada para aprimorar o acesso a informação quando essa é solicitada por um usuário por meio de uma consulta, onde é possível oferecer porções menores de texto mais relevante ao invés de exibir um documento maior que pode conter informações menos pertinente. A sumarização de texto também pode ser aprimorada ao processar segmentos separados por tópicos ao invés de documentos inteiros.

Assim, esse trabalho trata da adaptação e avaliação de algoritmos tradicionais ao contexto de documentos em português do Brasil, com ênfase especial nas atas de reuniões.

2. TRABALHOS RELACIONADOS

Os principais algoritmos de segmentação textual baseiam-se na ideia de coesão léxica entre assuntos. Isto é, a mudança de tópicos é acompanhada de uma proporcional mudança de vocabulário. A partir disso, vários algoritmos foram propostos. Dessa forma, assumem o pressuposto que um segmento pode ser identificado e delimitado pela análise das palavras

que o compõe.

Uma vez que coesão léxica é pressuposto básico da maioria dos algoritmos, o cálculo da similaridade entre textos é fundamental. Uma medida de similaridade frequentemente utilizada é a *cosine*, a qual pode ser vista na equação 1, sendo $f_{x,j}$ a frequência da palavra j na sentença x e $f_{y,j}$ sendo a frequência da palavra j na sentença y .

$$Sim(x, y) = \frac{\sum_j f_{x,j} \times f_{y,j}}{\sqrt{\sum_j f_{x,j}^2 \times \sum_j f_{y,j}^2}} \quad (1)$$

Entre os trabalhos mais influentes podemos citar o *Text-Tiling* [3] proposto por Hearst. Ela propõe um algoritmo baseado em janelas deslizantes, onde para cada candidato a limite, analisa-se o texto circundante. Um limite ou quebra se segmento é identificado quando a similaridade entres os blocos apresenta uma queda considerável.

Ela propõe um algoritmo baseado em janelas deslizante, para analisar blocos de texto adjacentes e identificar os limites com base nas similaridades dos blocos.

O algoritmo recebe uma lista de candidatos a limite, usualmente finais de parágrafo ou finais de sentenças. Para cada posição candidata são construídos 2 blocos, um contendo sentenças que a precedem e outro com as que a sucedem. O tamanho desses blocos é um parâmetro a ser fornecido ao algoritmo e determina o tamanho mínimo de um segmento. Em seguida, os blocos de texto são representados por vetores que contêm as frequências de suas palavras. Então, usa-se *cosini* (equação 1) para calcular a similaridade entre os blocos.

Finalmente, os limites são identificados sempre que a similaridade entre blocos adjacentes entre cada candidato ultrapassa um determinado *threshold*

Apresenta baixa complexidade computacional, devido a simplicidade do algoritmo e baixa eficiência quando comparado a outros métodos mais sofisticados como mostrando em [2, 4].

Choi [2] apresenta um trabalho que usa *cosine*, a qual é exibida na equação 1, como medida similaridade e apresenta um esquema de ranking em seu algoritmo, o *C99*. Embora muitos dos melhores trabalho utilizarem matrizes de similaridades, o autor traz observações. Ele aponta que para pequenos segmentos, o cálculo de suas similaridades não é confiável. Pois uma ocorrência adicional de uma palavra causa um impacto desproporcional no cálculo. Além disso, o estilo da escrita pode não ser constante em todo o texto. Choi sugere que, por exemplo, textos iniciais dedicados a introdução costumam apresentar menor coesão do que trechos

dedicados a um t pico espec fico.

Portanto comparar a similaridade entre trechos de diferentes regi es, n o   apropriado. Devido a isso, as similaridades n o podem ser comparadas em valores absolutos, ent o, o autor apresenta um esquema de ranking para contornar esse problema.

Cada valor na matriz de similaridade   substituído por seu ranking local. Onde ranking   o n mero de elementos vizinhos com similaridade menor, o qual   calculado com a equa  o ???. Um exemplo   mostrado na Figura ?? abaixo.

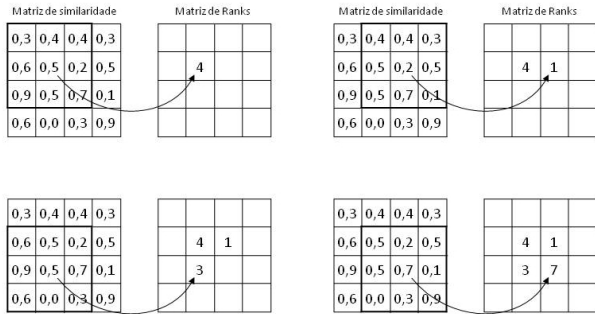


Figure 1: Exemplo de constru  o de uma matriz de rank

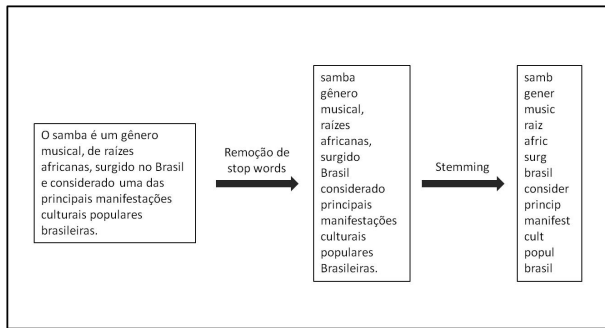


Figure 2: Exemplo de pr processamento

$$r(x, y) = \frac{\text{Numero de elementos com similaridade menor}}{\text{Numero de elementos examinados}} \quad (2)$$

Finalmente, na etapa de *clustering*, Choi utiliza um m todo baseado no algoritmo de maximiza  o de Reynar [?] para identificar os limites entre os segmentos.

Semelhante a esse trabalho, outras abordagens foram propostas como ...

[1] faz uma adapta  o do *TextTiling* ao contexto das conversas em reuni es com m ltiplos participantes.

3. ADAPTA  O  S ATAS DE REUNI O

3.1 Pr processamento

3.2 Identifica  o de senten as

4. AN LISE DOS RESULTADOS

5. AVALIA  O

Definir o que   um bom algoritmo de segmenta  o avalia  o todos precisam de um gold text

1 - Concatena  o 2 - Ju zes concordam ou n o 3 - Mediador na reuni o 4 - N o avaliar o segmentador e sim o resultado da aplica  o final.

De acordo com [?] h  duas principais dificuldades na avalia  o de segmentadores autom ticos. A primeira   conseguir um refer ncia confi vel de texto segmentado, ou seja, uma segmenta  o ideal, j  que ju zes humanos costumam n o concordar entre si, sobre onde os limites est o. A segunda   que tipos diferentes de erros devem ter pesos diferentes de acordo com a aplica  o. H  casos onde certa imprecis o   toler vel e outras como a segmenta  o de not cias, onde a precis o   mais importante.

Para contornar essas dificuldades, algumas abordagens podem ser utilizadas. Algumas autores preferem detectar a segmenta  o em textos formados pela concatena  o de documentos distintos, para que n o haja diferen as subjetivas [?]. H  ainda outros que n o avaliam o algoritmo diretamente, mas seu impacto na aplica  o final[?, ?, ?]. Outras abordagens apenas atribuem um segmento cada quebra de par grafo [?]

O vocabul rio das reuni es, ainda que em t picos diferentes, compartilham certo vocabul rio pertencente ao ambiente onde as se deram as reuni es. Isso   um fator que diminui a o princ pio da coes o l xica entre os segmentos.

5.1 Medidas de Avalia  o

5.1.1 Pk

5.1.2 WindowDiff

No trabalho de [?], os autores apontam problemas na avalia  o mais tradicional Pk, como a demasiada penaliza  o dos falsos negativos e a desconsidera  o de *near misses*, quando um limite entre t picos n o casa exatamente com esperado mas fica pr ximo a ele.

A ideia   mover uma janela pelo texto e penalizar o algoritmo sempre que o n mero de limites (proposto pelo algoritmo) n o coincidir com o n mero de limites (reais) para aquela janela de texto.

6. CONCLUS O

7. REFERENCES

- [1] S. Banerjee and A. Rudnicky. A texttiling based approach to topic boundary detection in meetings. volume 1, pages 57–60, 2006. cited By 3.
- [2] F. Y. Y. Choi. Advances in domain independent linear text segmentation. In *Proceedings of the 1st North American Chapter of the Association for Computational Linguistics Conference*, NAACL 2000, pages 26–33, Stroudsburg, PA, USA, 2000. Association for Computational Linguistics.
- [3] M. A. Hearst. Multi-paragraph segmentation of expository text. In *Proceedings of the 32Nd Annual Meeting on Association for Computational Linguistics*, ACL ’94, pages 9–16, Stroudsburg, PA, USA, 1994. Association for Computational Linguistics.

- [4] R. Kern and M. Granitzer. Efficient linear text segmentation based on information retrieval techniques. pages 167–171, 2009. cited By 10.