

Segmentação Textual Automática de Atas de Reunião

Resumo. A tarefa de segmentação textual consiste em dividir um texto em porções com significado relativamente independente, de maneira que cada segmento esteja relacionado a um assunto. Muitos pesquisadores avaliam técnicas de segmentação textual usando textos longos como a concatenação de documentos de vários assuntos ou com a transcrição de discursos e reuniões com múltiplos participantes. Este artigo concentra-se na avaliação dos principais algoritmos da literatura aplicando-os ao contexto das atas de reunião em português, as quais além do idioma menos estudado com essas técnicas, possuem estilo de redação mais formal e sucinto. Ao final, chegou-se a uma método capaz de segmentar as atas com performance similar a outros trabalhos.

Abstract. The text segmentation task consists of split a text into parts with relatively independent meaning, so that each segment contains a subject. Many researches evaluate they text segmentation methods using long texts such as concatenation of documents of several subjects or with transcriptions of discourse and multipart meetings. This article concentrates on the main algorithms of the literature employing them to the context of the meeting minutes in Portuguese, which besides the language less studied with this method, has a more formal and succinct writing style. Finally, we presents a method to segment the meeting minutes with similar performance to others works.

1. Introdução

As atas são documentos textuais que em geral descrevem dados não estruturados. Assim, um sistema que responde a consultas do usuário ao conteúdo das atas, retornando trechos de textos relevantes à sua intenção, é um desafio que envolve a compreensão de seu conteúdo [4].

Devido a fatores como a não estruturação e volume dos textos, a localização de assunto em uma ata é uma tarefa custosa. Usualmente, o que se faz são buscas manuais guiadas pela memória ou com uso de ferramentas computacionais baseadas em localização de palavras-chave. Normalmente esse tipo de busca exige a inserção de termos exatos e apresentam ao usuário um documento com as palavras buscadas em destaque, mantendo o texto longo, o que dificulta por exemplo o ranqueamento por relevância.

Uma vez que a ata registra a sucessão de assuntos discutidos na reunião, há interesse em um sistema que aponte trechos de uma ata que tratam de um assunto específico. Tal sistema tem duas principais tarefas: 1) Descobrir quando há uma mudança de assunto. 2) Descobrir quais são esses assuntos. Este trabalho tem como foco principal a primeira tarefa, a detecção de mudança de assuntos, que pode ser atendida pela segmentação automática de textos [8, 17, 6].

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

O interesse por segmentação textual tem crescido em aplicações voltadas a recuperação de informação e sumarização de textos [15]. Essa técnica pode ser usada para melhorar o acesso a informação quando essa é solicitada por meio de uma consulta, onde é possível oferecer porções menores de texto mais relevantes ao invés de exibir um documento grande que pode conter informações menos pertinentes. Além disso, encontrar pontos onde o texto muda de assunto, pode ser útil como etapa de pré-processamento em aplicações voltadas ao entendimento do texto, principalmente em textos longos. A navegação pelo documento pode ser aprimorada, em especial na utilização por usuários com deficiência visual, os quais utilizam sintetizadores de texto como ferramenta de acessibilidade [9]. A sumarização de texto pode ser melhorada ao processar segmentos separados por tópicos ao invés de documentos inteiros [3, 15, 5].

As atas de reunião diferem dos textos comumente estudados em outros trabalhos em alguns pontos. O estilo de escrita mais sucinto, com poucos detalhes dificulta o processo de segmentação [10]. Há também um maior foco no idioma inglês, presente na maioria dos artigos publicados. Diferenças de performance podem ser vistas no mesmo algoritmo quando aplicado em documentos de diferentes idiomas, onde a aplicação em textos em inglês apresenta um taxa de erro significativamente menor que o alemão e o espanhol [14, 20]. Assim, esse trabalho trata da aplicação e avaliação de algoritmos tradicionais ao contexto de documentos em português do Brasil, com ênfase especial nas atas de reuniões.

O restante deste artigo está dividido da seguinte forma. Na Seção 2 o processo de segmentação textual é descrito, bem como os principais algoritmos voltados para esta atividade. Na Seção 3 são apresentados trabalhos recentes que desenvolveram as técnicas existentes em contextos específicos. Na Seção 4 é apresentada a proposta desse trabalho voltada às atas de reunião bem como é detalhada a avaliação dos experimentos e seus resultados. Por fim, na Seção 5 são apresentadas as considerações finais e trabalhos futuros.

2. Referencial Teórico

Um documento textual, sobretudo quando longo, é frequentemente uma sucessão de assuntos. A segmentação textual ou segmentação topical é a tarefa de dividir um texto mantendo em cada parte um assunto com seu significado completo.

O texto deve ser previamente dividido em unidades de informação, que podem ser palavras, parágrafos e mais usualmente em sentenças. Essas unidades são processadas pelos algoritmos como a menor parte de um segmento, ou seja uma unidade não pode ser segmentada. Dessa forma, um segmento pode ser visto como uma sucessão de unidades de informação que compartilham o mesmo assunto e cada ponto entre duas unidades é considerado um candidato a limite entre segmentos.

Trabalhos anteriores se apoiam na ideia de que a mudança de tópicos em um texto é acompanhada de uma proporcional mudança de vocabulário. Essa ideia, chamada de coesão léxica, sugere que a distribuição das palavras é um forte indicador da estrutura do texto. A partir disso, vários algoritmos foram propostos baseados na ideia de que um segmento pode ser identificado pela análise das palavras que o compõe [8, 11, 19].

Uma vez que a coesão léxica é pressuposto básico da maioria dos algoritmos, o cálculo da similaridade entre unidades de informação é fundamental. Uma medida de

similidade frequentemente utilizada é o cosseno, apresentada na Equação 1, onde $f_{x,j}$ é a frequência da palavra j na sentença x e $f_{y,j}$ é a frequência da palavra j na sentença y .

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

2.1. Principais algoritmos

Entre os principais trabalhos da literatura podemos citar o *TextTiling* [12] e o *C99* [9]. O *TextTiling* é um algoritmo baseado em janelas deslizantes, em que, para cada candidato a limite, analisa-se o texto circundante. Um limite ou quebra de segmento é identificado sempre que a similaridade cai abaixo de um limiar.

O *TextTiling* 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 contém as frequências de suas palavras. Então, usa-se cosseno (Equação 1) para calcular a similaridade entre os blocos adjacentes a cada candidato e identifica-se uma transição entre tópicos pelos vales na curva de dissimilaridade.

O *TextTiling* possui baixa complexidade computacional. Por outro lado, algoritmos mais complexos, como os baseados em matrizes de similaridade, apresentam acurácia relativamente superior como apresentado em [9, 14, 16].

O *C99* é um algoritmo baseado em ranking. Embora muitos trabalhos utilizem matrizes de similaridades para pequenos segmentos, o cálculo de suas similaridades não é confiável, pois uma ocorrência adicional de uma palavra causa um impacto que pode alterar significativamente o cálculo da similaridade [9]. Além disso, o estilo da escrita pode não ser constante em todo o texto. Por exemplo, textos iniciais dedicados a introdução costumam apresentar menor coesão do que trechos dedicados a um tópico específico. Portanto, comparar a similaridade entre trechos de diferentes regiões não é apropriado. Devido a isso, as similaridades não podem ser comparadas em valores absolutos. Então, contorna-se esse problema fazendo uso de *rankings* de similaridade para encontrar os segmentos de texto.

Inicialmente é construída uma matriz que contém as similaridades de todas as unidades de texto. Em seguida, cada valor na matriz de similaridade é substituído por seu ranking local. Para cada elemento da matriz, seu *ranking* é o número de elementos vizinhos com valor de similaridade menor que o seu. Então, cada elemento é comparado com seus vizinhos dentro de uma região denominada máscara.

Na Figura 1(a) é destacado um quadro 3 x 3 de uma matriz em que cada elemento é a similaridade entre duas unidades de informação. Tomando como exemplo o elemento com valor 0,5, a mesma posição na matriz de *rankings* terá o valor 4, pois esse é o número de vizinhos com valores inferiores a 0,5 dentro do quadro analisado na matriz de similaridades. Da mesma forma, na Figura 1(b) para o valor 0,2 a matriz de *rankings* conterá o valor 1 na mesma posição.

Finalmente, com base na matriz de *rank*, o *C99* utiliza um método de *clustering*

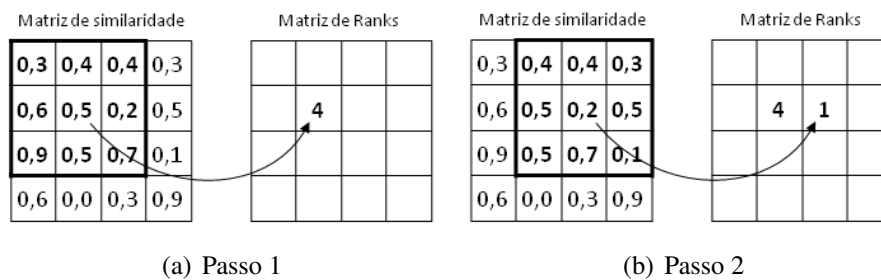


Figure 1. Exemplo de construção de uma matriz de rankings [9].

baseado no algoritmo de maximização de Reynar [18] para identificar os limites entre os segmentos.

2.2. Segmentação de Referência

Para que se possa avaliar um segmentador automático de textos é preciso uma referência, isto é, um texto com os limites entre os segmento conhecidos. Essa referência, deve ser confiável, sendo uma segmentação legítima que é capaz de dividir o texto em porções relativamente independentes, ou seja, uma segmentação ideal.

Entre as abordagens mais comuns para se conseguir essas referências, encontramos: 1) A concatenação aleatória de documentos distintos, onde o ponto entre o final de um texto e o início do seguinte é um limite entre eles. 2) A segmentação manual dos documentos, em que, pessoas capacitadas, também chamadas de juízes, ou mesmo o autor do texto, são consultadas e indicam manualmente onde há uma quebra de segmento. 3) Em transcrição de conversas faladas em reuniões com múltiplos participantes, um mediador é responsável por encerrar um assunto e iniciar um novo, nesse caso o mediador anota manualmente o tempo onde há uma transição de tópico.

Em aplicações em que a segmentação é uma tarefa secundária e quando essas abordagens são custosas ou não se aplicam, é possível, ao invés de avaliar o segmentador, analisar seu impacto na aplicação final.

2.3. Medidas de Avaliação

As medidas de avaliação tradicionais como precisão e revocação computam os erros do algoritmo, isto é, falsos positivos e falsos negativos, a fim de calcular seu desempenho. Além dessas medidas, que consideram apenas se um segmento foi corretamente definido, pode-se também considerar a distância entre o segmento extraído automaticamente e o segmento de referência [14]. Chama-se *near misses* o caso em que um limite identificado automaticamente não coincide exatamente com a referência, mas é necessário considerar a proximidade entre eles.

Na Figura 2 é apresentado um exemplo com duas segmentações extraídas automaticamente e uma referência. Em ambos os casos não há nenhum verdadeiro positivo, o que implica em zero para os valores de precisão, acurácia, e revocação, embora a segunda hipótese possa ser considerada superior à primeira se levado em conta a proximidade dos limites.

Considerando o conceito de *near misses*, algumas soluções foram propostas. As medidas de avaliação mais utilizadas são a P_k e *WindowDiff*. Proposta por [2], P_k , atribui

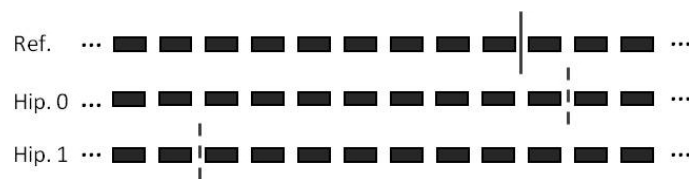


Figure 2. Exemplos de *near missing* e falso positivo puro. Os blocos indicam uma unidade de informação e as linha verticais representam uma transição de assunto.

valores parciais a *near misses*, ou seja, limites sempre receberão um peso proporcional à sua proximidade, desde que dentro de uma janela de tamanho k . Para isso, esse método move uma janela de tamanho k ao longo do texto. A cada passo verifica, na referência e na hipótese, se o início e o final da janela estão ou não dentro do mesmo segmento, então, penaliza o algoritmo caso não concorde com a referência. Ou seja, dado duas palavras de distância k , o algoritmo é penalizado quando não concordar com a segmentação de referência se as palavras estão ou não no mesmo segmento. O valor de k é calculado como a metade da média dos comprimentos dos segmentos reais. Como resultado, é retornado a contagem de discrepâncias dividida pela quantidade de segmentações analisadas. P_k é uma medida de dissimilaridade entre as segmentações e pode ser interpretada como a probabilidade de duas sentenças extraídas aleatoriamente pertencerem ao mesmo segmento.

WindowDiff é uma medida alternativa à P_k . De maneira semelhante, move uma janela pelo texto e penaliza o algoritmo sempre que o número de limites proposto pelo algoritmo não coincidir com o número de limites esperados para aquela janela. Ou seja, o algoritmo é penalizado quando não concordar com a segmentação de referência quanto ao número de segmentos na janela. Assim, consegue manter a sensibilidade a *near misses* e além disso, considerar o tamanho das janelas. A fim de melhor equilibrar o peso dos falsos positivos em relação a *near misses*, dobra-se a penalidade para falsos positivos, evitando-se a supervalorização dessa medida.

As medidas de avaliação tradicionais como precisão e revocação são usadas em recuperação de informação e classificação automática para medir o desempenho de modelos de classificação e predição. São baseadas na comparação dos valores produzidos por uma hipótese com os valores reais. Esses valores são apresentados em uma tabela que permite a visualização do desempenho de um algoritmo, a qual é chamada de matriz de confusão. Na Tabela 1 é apresentada a matriz de confusão para duas classes (Positivo e Negativo). Uma matriz de confusão é uma tabela que permite a visualização do desempenho de um algoritmo. No contexto de segmentação textual, um falso positivo é um

	Predição Positiva	Predição Negativa
Positivo real	VP (Verdadeiro Positivo)	FN (Falso Negativo)
Negativo real	FP (Falso Positivo)	VN (Verdadeiro Negativo)

Table 1. Matriz de confusão.

limite identificado pelo algoritmo que não corresponde a nenhum limite na segmentação de referência, ou seja, o algoritmo indicou que em determinado ponto há uma quebra de segmento, mas na segmentação de referência, no mesmo ponto, não há. De maneira semelhante, um falso negativo é quando o algoritmo não identifica um limite existente

na segmentação de referência, ou seja, em determinado ponto há, na segmentação de referência, um limite entre segmentos, contudo, o algoritmo não o identificou. Um verdadeiro positivo é um ponto no texto indicado pelo algoritmo e pela segmentação de referência como uma quebra de segmentos, ou seja, o algoritmo e a referência concordam que em determinado ponto há uma transição de assunto. Na avaliação de segmentadores, não há o conceito de verdadeiro negativo. Este seria um ponto no texto indicado pelo algoritmo e pela segmentação de referência onde não há uma quebra de segmentos. Uma vez que os algoritmos apenas indicam onde há um limite, essa medida não é necessária.

Nesse sentido, a precisão, é a proporção de limites corretamente identificados pelo algoritmo. É calculada dividindo-se o número de limites identificados automaticamente pelo número de candidatos a limite (Equação 2).

$$Precisão = \frac{VP}{VP + FP} \quad (2)$$

Essa medida varia entre 0,0 e 1,0, que indica a proporção de limites identificados pelo algoritmo que são corretos, ou seja, correspondem a um limite real na segmentação de referência. Porém não diz nada sobre quantos limites reais existem.

A revocação, é a proporção de limites verdadeiros que foram identificados pelo algoritmo. É calculada dividindo-se o número de limites identificados automaticamente pelo número limites verdadeiros.

$$Revocação = \frac{VP}{VP + FN} \quad (3)$$

Pode variar entre 0,0 e 1,0, onde indica que a proporção de limites corretos que foram identificados. Porém não diz nada sobre quantos limites foram indentificados incorretamente. Existe uma relação inversa entre precisão e revocação. Conforme o algoritmo aponta mais segmentos no texto, este tende a melhorar a revocação e ao mesmo tempo, reduzir a precisão. Esse problema de avaliação pode ser contornado utilizado a medida F^1 que é a média harmônica entre precisão e revocação onde ambas tem o mesmo peso.

3. Trabalhos Relacionados

Semelhante a este trabalho, outras abordagens foram propostas no sentido de propor segmentadores para outros idiomas além do inglês, bem como alvaliá-los em contextos específicos como discursos e conversas em reuniões.

A fim de aplicar os algoritmos *TextTiling* e *C99* ao idioma árabe, foi utilizado como *corpus* a concatenação de textos extraídos de notícias de diferentes países como Tunísia, Egito e Algeria, que tratam de assuntos distintos como política, esporte, cultura, história, tecnologia e artes, e que trazem particularidades nos estilos de escrita e até diferenças entre dialetos locais [7]. Verificou-se que, devido as diferenças entre os dialetos de cada país, melhores resultados dependem da escolha de um stemmer adequado.

As técnicas de segmentação também foram aplicadas em conversas em reuniões com múltiplos participantes onde se estuda os textos extraídos dos discursos, ou seja, o texto a ser segmentado é uma transcrição das falas dos participantes durante a reunião. Nesse sentido, como parte do projeto *CALO-MA*¹, foi apresentado um segmentador

¹<http://www.ai.sri.com/project/CALO>

baseado no *TextTiling*, que foi aplicado em reuniões com múltiplos participantes. Os autores apoiaram-se em elementos da fala como pausas, trocas de falantes e entonação para encontrar melhores segmentos [12]. O *corpus* utilizado contém a transcrição da fala dos participantes durante as reuniões que foram conduzidas por um mediador que propunha os tópicos e anotava o tempo em que os participantes mudavam de assunto [1, 21].

A presença de *pistas* pode ser um indicativo que ajuda a aprimorar a detecção de limites entre segmentos, uma vez que alguns elementos são frequentes na transição de tópicos. Essas *pistas* podem ser palavras como “Ok”, “*continuando*” e frases como “Boa noite”, “*Dando prosseguimento*” ou pausas prolongadas que ajudam a indicar uma mudança de tópicos. Essas *pistas* podem ser detectadas por meio de algoritmos de aprendizado de máquina ou anotadas manualmente [22, 5].

4. Proposta: Segmentação Linear Automática de Atas de Reunião

Os algoritmos *TextTiling* e *C99* foram propostos para o inglês, ou seja, a proposta inicial dos autores é trabalhar em qualquer texto nessa língua. Neste trabalho, esses algoritmos foram aplicados em atas de reunião escritas em português, ou seja, em uma língua diferente e dentro de um contexto específico. As implementações dos algoritmos, bem como ferramentas utilizadas e os resultados completos estão disponíveis para utilização e consulta em lasid.sor.ufscar.br/meetingminer/. As subseções seguintes tratam da aplicação desses algoritmos para esse nicho mais específico.

4.1. Coleção de documentos

A fim de obter um conjunto de atas a serem segmentadas automaticamente, coletou-se 6 atas de reunião do Conselho de Pós-Graduação da UFSCar Sorocaba². Selecionou-se atas que registram cinco reuniões ordinárias e uma extraordinária (Ata 6). Tomou-se esse cuidado para manter a representatividade e diversidade do conjunto.

Os documentos apresentam um texto com poucas quebras de parágrafo e sem marcações de estrutura, como capítulos, seções ou quaisquer indicações sobre o assunto do texto. É comum a presença de cabeçalhos, rodapés e numeração de páginas e linhas o que pode prejudicar tanto similaridade entre sentenças como a apresentação dos segmentos ao usuário.

4.2. Pré-processamento

As atas a serem segmentadas são extraídas de documentos do tipo *pdf*, *doc*, *docx* ou *odt* que normalmente possuem formato binário. Aplicou-se um processo para transformar esses arquivos em texto plano. A fim de preparar o texto e selecionar as palavras mais significativas, as atas passaram por processos de transformação os quais serão apresentados a seguir.

1. Remoção de cabeçalhos e rodapés: as atas contém trechos que podem ser considerados pouco informativos e descartados durante o pré-processamento, como cabeçalhos e rodapés que se misturam aos tópicos tratados na reunião, podendo ser inseridos no meio de um tópico prejudicando tanto o algoritmo de segmentação, quanto a leitura do texto pelo usuário.

²http://www.ppgccs.net/?page_id=1150

2. Identificação de finais sentenças: devido ao estilo de pontuação desses documentos, como encerrar sentenças usando um ";" e inserção de linhas extras, foram usadas as regras especiais para identificação de finais de sentença. Os detalhes sobre essas regras estão disponíveis para consulta em lasid.sor.ufscar.br/meetingminer/.
3. Redução de termos: eliminou-se a acentuação, sinais de pontuação, numerais e todos os *tokens* menores que três caracteres. Palavras de uso muito frequente como artigos, preposições e pronomes, chamadas de *stop words*, foram removidas utilizando-se uma lista de 438 palavras.
4. *Stemming*: extraiu-se o radical de cada palavra. Para isso, as letras foram convertidas em caixa baixa e aplicou-se o algoritmo *Orengo*³ para remoção de sufixos.

Na Tabela 2 é mostrado, para cada ata, a quantidade de *tokens* após a extração dos documentos e durante o pré-processamento.

Processo	Ata 1	Ata 2	Ata 3	Ata 4	Ata 5	Ata 6
Extração do texto	809	851	1043	1407	834	496
Remoção de Cabeçalhos e Rodapés	665	704	840	1247	708	392
Redução de termos	461	489	566	872	485	286

Table 2. Quantidade de *tokens* por ata.

Algoritmo 1: Identificação de finais de sentença

Entrada: Texto

Saída : Texto com identificações de finais de sentença

```

1 para todo token, marcá-lo como final de sentença se:
2   Terminar com um !
3   Terminar com um . e não for uma abreviação
4   Terminar em . ? ; e:
5     For seguido de uma quebra de parágrafo ou tabulação
6     O próximo token iniciar com ( { [ " '
7     O próximo token iniciar com letra maiúscula
8     O penúltimo caracter for ) } ] " '
9 fim
```

4.3. Segmentação de Referência

A fim de obter um conjunto de documentos segmentados que possam servir como referência na avaliação, os documentos coletados foram segmentados manualmente por profissionais que participam de reuniões. Para isso, utilizou-se um *software*, desenvolvido com esse objetivo específico, que permitiu aos voluntários visualizar um documento, e indicar livremente as divisões entre segmentos. Com o uso desse *software* foram coletados os dados de seis atas segmentadas por dois participantes das reuniões, os quais serviram como referência para a avaliação dos algoritmos. O *software* desenvolvido para segmentação manual está disponível para utilização e consulta em lasid.sor.ufscar.br/meetingminer/

³<http://www.inf.ufrgs.br/viviane/rs1p/>

A Tabela 3 contém, para cada ata, a quantidade de sentenças e a quantidade de segmentos identificadas pelos participantes.

Ata	Sentenças	Participante 1	Participante 2
Ata 1	18	7	15
Ata 2	26	9	20
Ata 3	24	7	15
Ata 4	32	9	17
Ata 5	25	11	17
Ata 6	10	4	9

Table 3. Quantidade de sentenças e segmentos de referência por ata.

4.4. Configuração experimental

O *TextTiling* permite ajustarmos dois parâmetros, sendo o tamanho da janela e o passo. Por meio de testes empíricos escolheu-se os valores 20, 40 e 60 para o tamanho da janela e 3, 6, 9 e 12 para o passo. Gerando ao final 20 configurações.

O *C99* permite o ajuste de três parâmetros, sendo, o primeiro a quantidade segmentos desejados, uma vez que, não se conhece o número ideal de segmentos e os documentos não apresentam muitos candidatos, calculou-se uma proporção dos candidatos a limite. Para isso atribuiu-se os valores 0,2; 0,4; 0,6; 0,8. O segundo parâmetro, o tamanho do quadro utilizado para gerar a matriz de ranking, atribuiu-se os valores 9 e 11, sendo 11 o valor padrão da apresentado pelo autor. O algoritmo permite ainda indicar se as sentenças serão representados por vetores contendo a frequência ou o peso de cada termo. Ambas as representações foram utilizadas. Considerando todos os parâmetros, foram geradas 16 configurações para o algoritmo *C99*.

4.5. Critérios de avaliação

Para fins de avaliação desse trabalho, um bom método de segmentação é aquele cujo resultado melhor se aproxima de uma segmentação manual, sem a obrigatoriedade de estar perfeitamente alinhado com tal. Ou seja, visto o contexto das atas de reunião, e a subjetividade da tarefa, não é necessário que os limites entre os segmentos (real e hipótese) sejam idênticos, mas que se assemelhem em localização e quantidade.

Os algoritmos foram comparados com a segmentação fornecida pelos participantes das reuniões. Calculou-se as medidas mais aplicadas à segmentação textual, P_k e *WindowDiff*. Além dessas, computou-se também as medidas tradicionais acurácia, precisão, revocação e F^1 para comparação com outros trabalhos que as utilizam.

Inicialmente, calculou-se as medidas configurando cada algoritmo conforme mostrado na Subseção 4.4, sem aplicar o pré-processamento. O teste de Friedman com pós-teste de Nemenyi foi utilizado para gerar um ranking das melhores configurações para cada medida calculada. Com isso, foi possível descobrir quais valores otimizam um algoritmo para uma medida, desconsiderando o pré-processamento.

A fim de conhecer o impacto do pré-processamento, repetiu-se os testes com o texto pré-processado. Com isso, descobriu-se quais valores otimizam os algoritmos para cada medida, considerando essa etapa.

Com os testes anteriores obteve-se, para cada medida, 4 configurações, levando em conta ambos os algoritmos e a presença ou ausência do pré-processamento. Novamente utilizou-se o teste de Friedman e Nemenyi e descobriu-se, para cada medida, qual configuração a otimiza. Os resultados completos estão disponíveis para consulta em lasid.sor.ufscar.br/meetingminer/.

4.6. Resultados

Obteve-se, por meio dos testes estatísticos apresentados, as melhores configurações para as principais medidas de avaliação de segmentadores. Com essas configurações calculou-se a média de cada medida considerando o conjunto de documentos. Na Tabela 4 são apresentadas, as médias obtidas com o *TextTiling* bem como as configurações utilizadas, onde **J** é o tamanho da janela e **P** é o passo.

Medida	Sem Pré-processamento			Com Pré-processamento		
	J	P	Média	J	P	Média
P_k	50	9	0,142	50	9	0,144
<i>WindowDiff</i>	50	6	0,387	40	9	0,396
Acurácia	50	6	0,612	40	9	0,603
Precisão	40	9	0,611	50	12	0,613
Revocação	20	3	0,886	20	3	0,917
F^1	30	6	0,605	40	3	0,648

Table 4. Resultados obtidos com o *TextTiling*

Na Tabela 5 são apresentadas, as médias obtidas com o *C99* bem como as configurações utilizadas, onde **S** é a proporção de segmentos em relação a quantidade de candidatos, **M** é o tamanho do quadro utilizado para criar a matriz de *rankings* e **W** indica se os segmentos são representados por vetores contendo a frequência ou um peso das palavras.

Medida	Sem Pré-processamento				Com Pré-processamento			
	S	M	W	Média	S	M	W	Média
P_k	20	9	Sim	0,134	20	11	False	0,116
<i>WindowDiff</i>	60	9	Sim	0,411	60	9	Sim	0,390
Acurácia	60	9	Sim	0,588	60	9	Sim	0,609
Precisão	40	9	Sim	0,645	20	11	False	0,720
Revocação	80	9	Sim	0,869	80	11	Sim	0,897
F^1	80	9	Sim	0,638	80	11	Sim	0,655

Table 5. Resultados obtidos com o *C99*

De acordo com os últimos testes, o algoritmo *C99* obteve melhor desempenho em acurácia, precisão, F^1 , P_k e *WindowDiff*, enquanto o *TextTiling* obteve o melhor desempenho em revocação como pode ser visto na Tabela 6.

De maneira geral, o algoritmo *C99* apresentou melhores resultados em relação ao *TextTiling*, sobre tudo quando aplicado o pré-processamento, contudo, testes estatísticos realizados indicaram que não houve diferença significativa entre os métodos. A etapa de

Algoritmo	Medida	Parâmetros	Pré-processamento	Média
C99	P_k	S=20 M=11 W=Não	Sim	0,116
C99	WindowDiff	S=60 M=09 W=Sim	Sim	0,390
C99	Acurácia	S=60 M=09 W=Sim	Sim	0,609
C99	Precisão	S=20 M=11 W=Não	Sim	0,720
C99	F^1	S=80 M=11 W=Sim	Sim	0,655
TextTiling	Revocação	J=20 P=3	Sim	0,917

Table 6. Melhores resultados obtidos.

pré-processamento proporciona melhora de desempenho quando aplicada, porém o seu maior benefício é a diminuição do custo computacional, uma vez que não prejudica a qualidade dos resultados.

As medidas de avaliação tradicionais, podem não ser confiáveis, por não considerarem a distância entre os limites, mas penalizam o algoritmo sempre que um limite que não coincide perfeitamente com a referência. Essas medidas podem ser mais adequadas quando necessita-se de segmentações com maior exatidão.

As medidas *WindowDiff* e P_k , consideram a quantidade e proximidade entre os limites, sendo mais tolerantes a pequenas imprecisões. Essa é uma característica desejável, visto que as segmentações de referência possuem diferenças consideráveis. *WindowDiff* equilibra melhor os falsos positivos em relação a *near misses*, ao passo que P_k os penaliza com peso maior. Isso significa que segmentadores melhores avaliados em P_k ajudam a selecionar as configurações que erram menos ao separar trechos de texto com o mesmo assunto, enquanto *WindowDiff* é mais tolerante nesse aspecto.

Observa-se melhores resultados de *WindowDiff* quando os algoritmos aproximam a quantidade de segmentos automáticos da quantidade de segmentos da referência. Por outro lado, observa-se que P_K avalia melhor as configurações que retornam menos segmentos. A configuração do tamanho do passo (P) e da proporção de segmentos em relação ao número de candidatos (S), influenciam os algoritmos na quantidade de segmentos extraídos. Contudo, não é possível definir um valor adequado, uma vez que os segmentadores humanos frequentemente apontam valores diferentes.

5. Conclusão

As atas de reunião, objeto de estudo desse artigo, apresentam características peculiares em relação à discursos e textos em geral. Características como ausência de parágrafos, segmentos curtos e estilo que evita repetição de palavras em benefício da leitura por humanos, dificultam o processamento por computadores.

Os algoritmos *TextTiling* e *C99* foram testados em um conjunto de atas coletadas do Departamento de Computação da UFSCar-Sorocaba. Os segmentos gerados automaticamente foram comparados à segmentação manual fornecida por participantes das reuniões. Por meio da análise dos chegou-se as configurações cujos segmentos melhor se aproximaram das amostras fornecidas. Os resultados obtidos indicam que os algoritmos originalmente propostos para o idioma inglês e usualmente avaliados em outros contextos podem ser aplicados as atas de reunião.

A segmentação de atas de reunião pode ajudar na organização, busca e com-

preensão dos conteúdos nelas contidos. Também outros domínios e aplicações diferentes podem se beneficiar dos resultados apresentados, como aplicações voltadas a resgate de informação, sumarização e acessibilidade. A metodologia apresentada e as ferramentas disponibilizadas podem ser úteis na avaliação e configuração de algoritmos de segmentação em outros contextos.

Em trabalhos futuros, serão investigadas técnicas de classificação e extração de tópicos para descrever os segmentos e com isso aprimorar o acesso ao conteúdo das atas de reunião.

References

- [1] S. Banerjee and A. Rudnicky. A texttiling based approach to topic boundary detection in meetings. *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, 1:57–60, 2006.
- [2] D. Beeferman, A. Berger, and J. Lafferty. Statistical models for text segmentation. *Machine Learning*, 34(1):177–210, 1999.
- [3] N. Bhatia and A. Jaiswal. Automatic text summarization and it's methods - a review. In *2016 6th International Conference - Cloud System and Big Data Engineering (Confluence)*, pages 65–72, Jan 2016.
- [4] M. H. Bokaei, H. Sameti, and Y. Liu. Linear discourse segmentation of multi-party meetings based on local and global information. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 23(11):1879–1891, Nov. 2015.
- [5] M. H. BOKAEI, H. SAMETI, and Y. LIU. Extractive summarization of multiparty meetings through discourse segmentation. *Natural Language Engineering*, 22(1):41–72, 2016.
- [6] P. Cardoso, T. Pardo, and M. Taboada. Subtopic annotation and automatic segmentation for news texts in brazilian portuguese. *Corpora*, 12(1):23–54, 2017.
- [7] A. H. Chaibi, M. Naili, and S. Sammoud. Topic segmentation for textual document written in arabic language. *Procedia Computer Science*, 35:437 – 446, 2014.
- [8] H. Chen, L. Xie, C. C. Leung, X. Lu, B. Ma, and H. Li. Modeling latent topics and temporal distance for story segmentation of broadcast news. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 25(1):112–123, Jan 2017.
- [9] 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.
- [10] F. Y. Y. Choi, P. Wiemer-Hastings, and J. Moore. Latent semantic analysis for text segmentation. In *In Proceedings of EMNLP*, pages 109–117, 2001.
- [11] O. Ferret. Improving text segmentation by combining endogenous and exogenous methods. In *International Conference Recent Advances in Natural Language Processing, RANLP*, pages 88–93, 2009.
- [12] M. Galley, K. McKeown, E. Fosler-Lussier, and H. Jing. Discourse segmentation of multi-party conversation. In *Proceedings of the 41st Annual Meeting on Association for*

Computational Linguistics - Volume 1, ACL '03, pages 562–569, Stroudsburg, PA, USA, 2003. Association for Computational Linguistics.

- [13] 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.
- [14] R. Kern and M. Granitzer. Efficient linear text segmentation based on information retrieval techniques. *Proceedings of the International Conference on Management of Emergent Digital EcoSystems, MEDES '09*, pages 167–171, 2009.
- [15] E. Maziero, G. Hirst, and T. Pardo. Adaptation of discourse parsing models for the portuguese language. *Proceedings - 2015 Brazilian Conference on Intelligent Systems, BRACIS 2015*, pages 140–145, 2016.
- [16] H. Misra, F. Yvon, J. M. Jose, and O. Cappe. Text segmentation via topic modeling: An analytical study. In *Proceedings of the 18th ACM Conference on Information and Knowledge Management, CIKM '09*, pages 1553–1556, New York, NY, USA, 2009. ACM.
- [17] M. Naili, A. H. Chaibi, and H. H. B. Ghezala. Exogenous approach to improve topic segmentation. *International Journal of Intelligent Computing and Cybernetics*, 9(2):165–178, 2016.
- [18] J. C. Reynar. *Topic Segmentation: Algorithms and Applications*. PhD thesis, Philadelphia, PA, USA, 1998. AAI9829978.
- [19] M. Sakahara, S. Okada, and K. Nitta. Domain-independent unsupervised text segmentation for data management. In *2014 IEEE International Conference on Data Mining Workshop*, pages 481–487, Dec 2014.
- [20] L. Sitbon and P. Bellot. Adapting and comparing linear segmentation methods for french. In *Coupling Approaches, Coupling Media and Coupling Languages for Information Retrieval*, RIAO '04, pages 623–637, Paris, France, France, 2004. LE CENTRE DE HAUTES ETUDES INTERNATIONALES D'INFORMATIQUE DOCUMENTAIRE.
- [21] G. Tur, A. Stolcke, L. Voss, S. Peters, D. Hakkani-Tür, J. Dowding, B. Favre, R. Fernández, M. Frampton, M. Frandsen, C. Frederickson, M. Graciarena, D. Kintzing, K. Leveque, S. Mason, J. Niekrasz, M. Purver, K. Riedhammer, E. Shriberg, J. Tien, D. Vergyri, and F. Yang. The calo meeting assistant system. *Trans. Audio, Speech and Lang. Proc.*, 18(6):1601–1611, Aug. 2010.
- [22] P. yun Hsueh, J. D. Moore, and S. Renals. Automatic segmentation of multiparty dialogue. In *EACL*, 2006.