

# Avaliação de Técnicas de Recuperação de Informação para Organização e Extração de Conhecimento de Documentos de Reunião

Ovídio José Francisco

Orientadora: Prof.<sup>a</sup> Dr. Katti Faceli

Coorientador: Prof. Dr. Rafael Geraldeli Rossi



August 20, 2018

- 1 Contextualização
- 2 Objetivos
- 3 Proposta
- 4 Avaliação Experimental
- 5 Conclusão
- 6 Trabalhos Futuros

- 1 Contextualização
- 2 Objetivos
- 3 Proposta
- 4 Avaliação Experimental
- 5 Conclusão
- 6 Trabalhos Futuros

- As atas registram assuntos discutidos em reuniões;
  - Uma ata pode conter diversos assuntos discutidos.
- 
- Um assunto pode ser discutido em diversas reuniões ao longo do tempo;
  - Recuperar o histórico de um assunto pode apoiar tomadas de decisão.
- 
- Informações contidas em grandes quantidades de texto;
  - Inerentemente não estruturados;
  - Documentos com múltiplos assuntos;
  - Necessidade de ferramentas automáticas.

Desafios na recuperação de informações em atas:

- Apresentar uma ata integral pode dificultar o usuário encontrar determinado assunto;
- É desejável agrupar trechos de atas com assuntos semelhantes.

## Algoritmos de Segmentação Textual:

- Dividem um texto em trechos com um único assunto completo;
- Úteis em aplicações com textos sem indicações de quebras de assunto, como transcrições de áudio, e diálogos em chats;
- Podem ser uma etapa de pré-processamento para outras aplicações;
- Não dão indicações sobre o conteúdo dos segmentos.

## Modelos de Extração de Tópicos:

- Estimam o assunto de cada documento de uma coleção;
- Agrupam documentos por tópico;
- Identificam palavras para descrever os tópicos;
- Incorporam conhecimento de domínio aos dados.

---

(ZAMIR; ETZIONI, 1998)	Web document clustering (1998)
(MASAO; KôITI, 2000)	Multi-topic multi-document summarization (2000)
(JEONG; TITOV, 2010)	Multi-document topic segmentation (2010)
(NGUYEN, 2011)	A Study on Statistical Generation of a Hierarchical Structure of Topic-information for Multi-documents (2011)
(TAGARELLI; KARYPIS, 2013)	A segment-based approach to clustering multi-topic documents (2013)

---

Trabalho	Divisão	Método de inferência	Representação	Idioma
(ZAMIR; ETZIONI, 1998)	Sentenças	Agrupamento		Inglês
(MASAO; KôITI, 2000)	Sentenças	Agrupamento	Rotulação manual	Inglês
(JEONG; TITOV, 2010)	Segmentação	Modelo Bayesiano		Inglês
(NGUYEN, 2011)	Segmentação	Agrupamento hierárquico	Sumarização e palavras-chave	Inglês
(TAGARELLI; KARYPIS, 2013)	Parágrafos	Agrupamento e classificação		Inglês



- 1 Contextualização
- 2 **Objetivos**
- 3 Proposta
- 4 Avaliação Experimental
- 5 Conclusão
- 6 Trabalhos Futuros

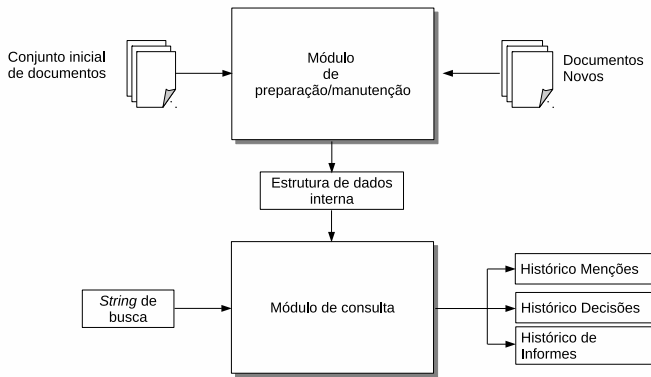
Propor uma solução para identificar, organizar e consultar assuntos registrados em atas de reunião.

Utilizar técnicas de Segmentação Textual em conjunto com modelos de Extração de Tópicos para:

- Gerar uma estrutura mais organizada que a coleção original;
- Utilizar a estrutura latente dos segmentos para Recuperação de Informação.

- 1 Contextualização
- 2 Objetivos
- 3 Proposta**
- 4 Avaliação Experimental
- 5 Conclusão
- 6 Trabalhos Futuros

# Visão Geral do Sistema



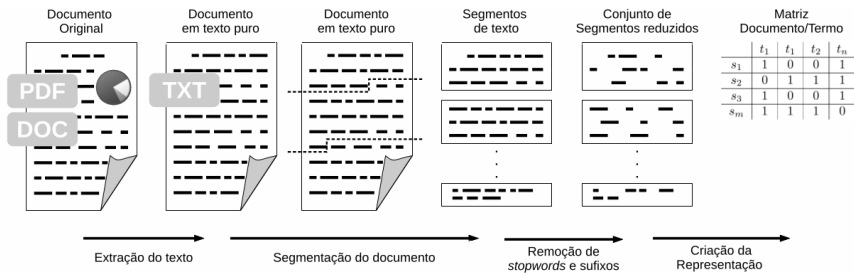
## Preparação

- Extração de texto plano;
- Segmentação;
- Remoção de termos;
- Representação dos Segmentos;
- Extração de Conhecimento.
  - Extração de tópicos;
  - Classificação;

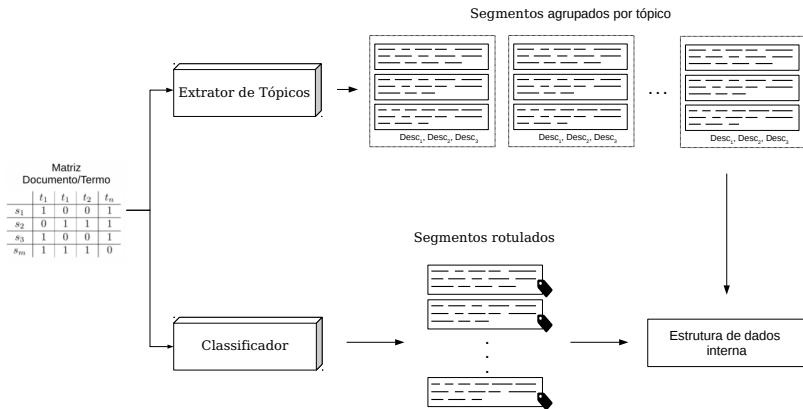
## Manutenção

Realimentação do sistema

# Preparação



# Extração de Conhecimento



# Estrutura de dados interna

Coleção de documentos (D)



Coleção de segmentos (S)



Segmentos agrupados por tópicos



Documento-tópico  
(W)

Termo-tópico  
(Z)



Obtém-se uma estrutura:

- Mais organizada que a coleção original;
- Assuntos concentrados em grupos;
- Acrescida de novos atributos;
- Distribuição dos tópicos conhecida.

# Distribuição de tópicos em uma ata real



UNIVERSIDADE FEDERAL DE SÃO CARLOS – Campus Sorocaba

BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO

Rodovia João Leme dos Santos, Km 110 (SP-264)

Bairro do Itinga – Sorocaba-SP – CEP 13052-780

Telefone: (15) 3202-2022 / [www.ufscar.br](http://www.ufscar.br)

**Ata da 17ª Reunião Extraordinária do Conselho do Curso de Bacharelado em Ciência da Computação, UFScar – Sorocaba**

**Local, Dia e Hora:** No laboratório de Pedagogia, situada nas dependências da Universidade Federal de São Carlos – Campus Sorocaba, à Rodovia João Leme dos Santos, quilômetro 110, Bairro do Itinga, na cidade de Sorocaba, Estado de São Paulo e realizada aos vinte e cinco dias do mês de agosto de dois mil e dez, às 14h00. **Convocação e Presença:** A convocação foi realizada sete dias antes da data de sua realização, estando presentes os membros do Conselho do Curso de Bacharelado em Ciência da Computação – Sorocaba, sendo eles os professores

representantes discentes Sr.

(2010) **Comunicação da**

**Presidência:** A prof. informou que na última reunião do ConCam foram indicados os nomes do prof. titular e suplente respectivamente, como representantes do campus Sorocaba no ConsUni (Conselho Universitário). A prof.<sup>a</sup> esclareceu também, que o campus já possui um assento ocupado pelo prof. e que agora teremos dois assentos no respectivo conselho. A prof.<sup>a</sup> também informou que na mesma reunião do ConCam foram indicados os nomes do prof. titular e suplente, para o assento do campus Sorocaba no CoAd (Conselho de Administração da UFScar), destacando que agora o campus Sorocaba não está mais no plano de implementação do campus, e irá concorrer com os outros campi da universidade pela distribuição de verbas. A prof.<sup>a</sup> também informou que no próximo dia trinta de agosto haverá uma reunião com a ProGрад em Sorocaba para retorno sobre o Projeto Político Pedagógico do curso no qual a coordenação do curso irá

**participar. Comunicação dos membros:** O prof. comunicou que houveram diversas problemáticas com o lançamento do edital para o concurso público da vaga docente para Banco de Dados. O prof. comunicou que já solicitou uma retificação do edital, mas que até o momento tal solicitação está em trâmite no setor de recursos humanos. O Sr. comunicou que foi criado uma lista de e-mails com os e-mails de todos os alunos do curso e docentes do curso, sendo que esta foi uma solicitação da coordenação do curso para que houvesse maior facilidade na comunicação entre alunos e docentes. Foi destacado que os docentes, por padrão, não receberão os e-mails da lista, estando autorizados apenas ao envio de mensagens. Caso desejem receber e-mails da lista o docente deve se comunicar com o Sr. Rubens.

**Ordem do Dia: (I) APROVAÇÃO DAS FICHAS DE CARACTERIZAÇÃO. (II) A** prof.<sup>a</sup> apresentou as fichas de caracterização que seriam analisadas, esclarecendo que todas elas foram criadas considerando fielmente o Projeto Político Pedagógico do curso. As fichas contemplam as disciplinas obrigatórias do curso que ainda não foram oferecidas e as disciplinas optativas que poderão ser oferecidas no primeiro semestre de 2011, sendo respectivamente as obrigatórias: Física para Computação, Algoritmos e Complexidade, Trabalho de Graduação 1, Estágio Supervisionado 1, Trabalho de Graduação 2, Estágio Supervisionado 2 e Seminários de Computação; e as optativas: Tópicos Avançados em Ciência da Computação, Tópicos Avançados em Redes de Computadores e Sistemas Distribuídos, Segurança e Auditoria de Sistemas e Mineração de Dados. Além das fichas foram apresentados os checklists das respectivas disciplinas, que foi uma solicitação da Coordenação Acadêmica para este semestre. DISCIPLINAS OPTATIVAS. (II.I) A prof.<sup>a</sup> destacou que haviam sido elaboradas fichas de quatro disciplinas optativas, porém somente três disciplinas serão oferecidas. Desta maneira haverá uma maior flexibilidade para decidir dentre as quatro disciplinas quais seriam as três oferecidas no próximo semestre. DISCIPLINAS DE FÍSICA PARA COMPUTAÇÃO E ALGORITMOS E COMPLEXIDADE. (II.III) A prof.<sup>a</sup> também colocou pra análise as fichas de caracterização das disciplinas de Física para Computação e Algoritmos e

dia; realizada; chamada; estado; conselho;

cursadas; conselho; coordenação; computação; presidente;

docentes; técnica; administrativo; presidente; dia;



UNIVERSIDADE FEDERAL DE SÃO CARLOS – Campus Sorocaba

BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO

Rodovia João Leme dos Santos, Km 110 (SP-264)

Bairro do Itinga – Sorocaba-SP – CEP 13052-780

Telefone: (15) 3202-2022 / [www.ufscar.br](http://www.ufscar.br)

Complexidade, explicando que estas disciplinas foram readequadas segundo discussões anteriores sobre a omissão e omissão das mesmas. Juntamente com as fichas de caracterização foi colocada um ofício justificando a alteração das disciplinas. **AVALIAÇÃO DAS FICHAS. (II.IV)** A prof.<sup>a</sup> colocou as fichas para que fosse avaliadas pelos membros do conselho. O prof.<sup>a</sup> questionou se quando o aluno vai se matricular na disciplina de Estágio ou de Trabalho de Graduação o sistema verifica se o pré-requisito de créditos mínimos cursados é checado. O prof.<sup>a</sup> respondeu que acreditava que esta verificação era realizada. A prof.<sup>a</sup> colocou que no checklist da disciplina de Algoritmos e Complexidade estava que a disciplina era pré-requisito para outras disciplinas do curso, perguntando se realmente havia alguma disciplina que utilizasse essa como pré-requisito. A prof.<sup>a</sup> consultou o Projeto Político Pedagógico e verificou que existe uma disciplina optativa, Tópicos Avançados em Teoria da Computação que tem como pré-requisito a disciplina de Algoritmos e Complexidade. (II.V) As fichas foram aprovadas pelo Conselho. **Encerramento:** Estando todos de acordo e nada mais havendo a deliberar, lavra-se, lida-se, aprova-se e assina-se esta Ata por todos os membros do Conselho do Curso de Bacharelado em Ciência da Computação, Campus Sorocaba, participantes desta reunião que, em 02 (duas) vias, será levada a registro e arquivamento junto à Coordenação do Conselho do Curso de Bacharelado em Ciência da Computação – Sorocaba, ficando ali à disposição para consulta restrita aos professores da UFScar – Sorocaba. Nada mais.

Prof.<sup>a</sup> Dr.<sup>a</sup>

Presidente desta Reunião

Prof. Dr.

Professor Associado

Prof. Dr.

Professor Adjunto

Profa. Dra.

Professora Adjunta

Profa. Dra.

Professora Adjunta

Prof. Dr.

Professor Adjunto

Prof. Dr.

Professor Adjunto

Prof. Dr.

Professor Adjunto

Representante Discente - 2008

Prof. Dr.

Professor Adjunto

Representante Discente - 2009

Representante Discente - 2010

disciplinas; cursadas; libras; conselho; aprovado;

computação; conselho; aprovado; acordado; ficou;

representante; discente; presidente; secretária; turma;

# Distribuição de tópicos em uma ata real



UNIVERSIDADE FEDERAL DE SAO CARLOS – Campus Sorocaba

BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO

Rodovia João Leme dos Santos, Km 110 (SP-264)

Bairro do Itinga - Sorocaba-SP - CEP 18052-780

Telefone: (15) 3202-2022 / [www.ufscar.br](http://www.ufscar.br)

## Ata da 17ª Reunião Extraordinária do Conselho do Curso de Bacharelado em Ciência da Computação, UFSCar – Sorocaba

**Local, Dia e Hora:** No laboratório de Pedagogia, situada nas dependências da Universidade Federal de São Carlos - campus Sorocaba, à Rodovia João Leme dos Santos, quilômetro 110, Bairro do Itinga, na cidade de Sorocaba, Estado de São Paulo e realizada aos vinte e cinco dias do mês de agosto de dois mil e dez, às 14h00. **Convocação e Presença:** A convocação foi realizada sete dias antes da data de sua realização, estando presentes os membros do Conselho do Curso de Bacharelado em Ciência da Computação - Sorocaba, sendo eles os professores

representantes discentes Sr.

(2010). **Comunicação da**

**Presidência:** A prof.<sup>1</sup> informou que na última reunião do ConCam foram indicados os nomes do prof. \_\_\_\_\_, titular e suplente respectivamente, como representantes do campus Sorocaba no ConsUni (Conselho Universitário). A prof.<sup>2</sup> esclareceu também, que o campus já possui um assento

dia; realizada; chamada; estado; conselho;

# Distribuição de tópicos em uma ata real

data de sua realização, estando presentes os membros do Conselho do Curso de Bacharelado em Ciência da Computação - Sorocaba, sendo eles os professores

representantes discentes Sr.

(2010). **Comunicação da**

**Presidência:** A prof.<sup>a</sup> informou que na última reunião do ConCam foram indicados os nomes do prof. , titular e suplente respectivamente, como representantes do campus Sorocaba no ConsUni (Conselho Universitário). A prof.<sup>a</sup> esclareceu também, que o campus já possui um assento ocupado pelo prof. ; e que agora teremos dois assentos no respectivo conselho. A prof.<sup>a</sup> também informou que na mesma reunião do ConCam foram indicados os nomes do prof. ,

titular e suplente, para o assento do campus Sorocaba no CoAd (Conselho de Administração da UFSCar), destacando que agora o campus Sorocaba não está mais no plano de implantação do campus, e irá concorrer com os outros campi da universidade pela distribuição de verbas. A prof.<sup>a</sup> também informou que no próximo dia trinta de agosto haverá uma reunião com a ProGrad em Sorocaba para retorno sobre o Projeto Político Pedagógico do curso no qual a coordenação do curso irá participar.

**Comunicação dos membros:** O prof. comunicou que houveram diversos problemas com o lançamento do edital para o concurso público da vaga docente para Banco de Dados. O prof. comunicou que já solicitou uma retificação do edital, mas que até o momento tal solicitação está em trâmite no setor de recursos humanos. O Sr. comunicou que foi criado uma lista de e-mails com os e-mails de todos os alunos do curso e docentes do curso, sendo que este foi

cursadas; conselho; coordenação; computação; presidente;

## Módulo de Consulta:

- Utiliza a Estrutura de Dados Interna como base;
- Os tópicos são representados por seus descritores;
- Os segmentos de atas são agrupados em tópicos;
- Usa o Modelo de Espaço Vetorial para ranquear os tópicos;
- Exibe os segmentos atribuídos ao primeiro tópico do ranking;

## Interface do Sistema após uma consulta

The screenshot displays the 'Meeting Explorer' application window. At the top, there's a search bar with the text 'Digite um assunto: concurso professor' and an 'Explorar' button. Below the search bar, a tree view on the left lists various topics under the selected category 'área: concurso; problemas: enviar; cursadas:'. The main area on the right shows three document excerpts related to the search. The first excerpt is titled '21ª Reunião Extraordinária CoC-CCS 16-03-11.doc' and discusses a professor's situation. The second excerpt is titled '21ª Reunião Ordinária CoC-CCS 03-08-11.doc' and mentions a communication from the Presidency. The third excerpt is titled '27ª Reunião Ordinária CoC-CCS 05-12-12.doc' and discusses a vestibular for refugees. The fourth excerpt is titled '30ª Reunião Ordinária CoC-CCS 22-05-13.doc' and discusses a correction of the salary of a collaborator. At the bottom, a status bar shows '175 documentos na base de dados', '70 topics extracted', and 'PLSA Parametric'.

Meeting Explorer

Manutenção Configurações

Digite um assunto: concurso professor Explorar

175 documentos na base de dados 70 topics extracted PLSA Parametric

12 trechos relacionados

(1.1) A prof. colocou que temos a seguinte situação: a prof. prestou o concurso de Banco de Dados, e explicou que o prof. Márcio tinha um perfil que navegava em outras disciplinas, além daquelas das áreas de Banco de Dados, e que portanto a pessoa que vier a ocupar a vaga dele também deverá ter um perfil que atenda as áreas básicas, esclareceu que o concurso de Banco de Dados de Dadas tem uma comissão avaliadora. Desde essa data, não houve mais nenhuma reunião.

21ª Reunião Extraordinária CoC-CCS 16-03-11.doc

Comunicação da Presidência: A prof. agradeceu a presença de todos e informou sobre o ingresso do Prof. como professor efetivo no quadro do curso, referente ao concurso para a vaga na área de Estrutura de Dados, Algoritmos e Complexidade. Informou também sobre os problemas ocorridos no período de matrícula causados pela greve os técnicos administrativos, descrevendo que a coordenação de curso foi quem cuidou da ausência do material. Comunicação dos membros do Banco de Dados. Dia 11/08/2011 de 14h00 às 16h00.

21ª Reunião Ordinária CoC-CCS 03-08-11.doc

Vestibular para refugiados: O processo seletivo para refugiados será realizado em São Carlos e existe um candidato para o curso de Ciência da Computação. O candidato é refugiado de Angola, mas já faz o Ensino Médio no Brasil.

27ª Reunião Ordinária CoC-CCS 05-12-12.doc

Correção do saldo da verba empenhada: No dia 02 de maio de 2013 a Coordenação do Curso de Ciência da Computação enviou ao Departamento Financeiro (IfeFin) e ao Departamento Contábil (DeCont) o ofício 14/2013 CCCS-Sor (ANEXO 01), para averiguação e correção do valor real do saldo disponível para utilização no empenho para Pagamentos à Colaboradores. A coordenação informou ainda que por esse motivo, os pagamentos dos colaboradores Menezes, Lucílio, Domingos, e Tavares, Basso, Bruno, Alberto, estão em andamento.

30ª Reunião Ordinária CoC-CCS 22-05-13.doc

A estagiária foi contratada em 11 de junho e desde então tem auxiliado nas tarefas da secretária dos cursos Ciência da Computação e Ciências Econômicas. 1.2.Verba: Os empenhos em Compra de materiais e pagamento estudantil estão em andamento, conforme definido na 30ª Reunião Ordinária do CoCCCS.

O sistema permite:

- Receber uma base de dados não estruturada;
- Identificar os assuntos tratados em cada ata;
- Agrupar segmentos por tópico;
- Adicionar novos atributos (descritores) aos segmentos;
- Expandir o espaço de busca;
- Retornar trechos relevantes à consulta.

- 1 Contextualização
- 2 Objetivos
- 3 Proposta
- 4 Avaliação Experimental**
- 5 Conclusão
- 6 Trabalhos Futuros



A Avaliação Experimental foi dividida em duas etapas:

- Segmentação Textual;
- Extração de Tópicos.

Os segmentadores foram avaliados objetivamente.

- Processo de anotação manual em segmentos;
  - Criação de segmentações de referência;
  - Configuração dos algoritmos;
  - Medidas de desempenho;
    - Acurácia e  $F^1$ ;
    - WindowDiff e  $P_k$ .
- 
- *Testes de significância estatística (Friedman e Nemenyi)*

A tarefa dos anotadores consistiu em:

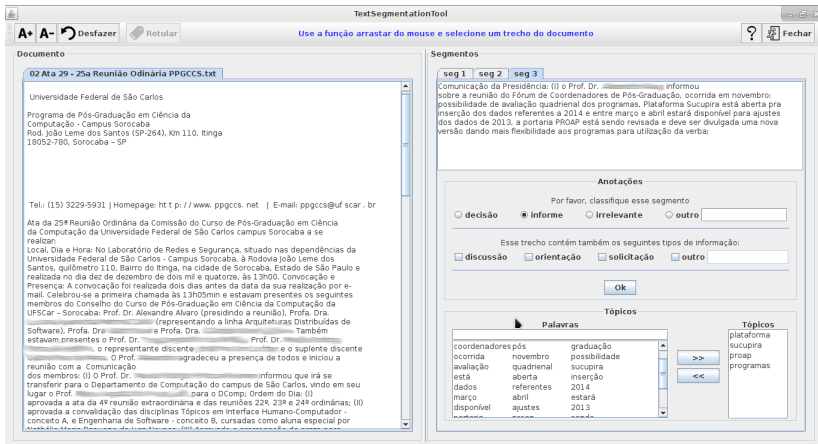
- Selecionar trechos com um único assunto;
- Rotular os trechos selecionados;
  - Tipo comunicação;
  - Contexto onde se gerou o assunto;
  - Descrição do assunto.

Utilizou-se:

- 12 atas da UFSCar;
- 09 anotadores;
- Ferramenta para anotações em segmentos.

# Processo de anotação manual em segmentos

## Ferramenta desenvolvida para anotação em segmentos



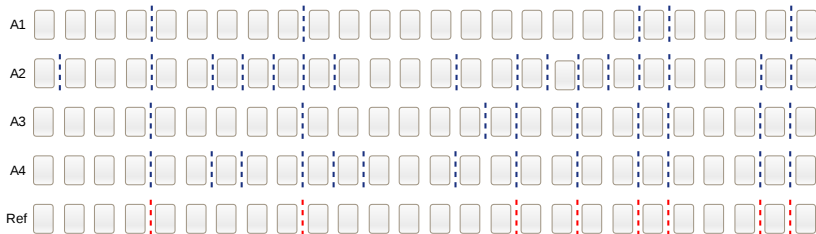
# Descrição dos resultados obtidos com anotadores

## Descrição dos resultados obtidos com anotadores

<b>Ata</b>	<b>#Sent.</b>	<b>A1</b>	<b>A2</b>	<b>A3</b>	<b>A4</b>	<b>A5</b>	<b>A6</b>	<b>A7</b>	<b>A8</b>	<b>A9</b>
Ata 1	25	7	4	11	6	16	8	8	15	16
Ata 2	17	4	4	8	6	11	6	6	15	14
Ata 3	26	6	6	8	4	15	9	10	18	14
Ata 4	26	5	5	10	6	14	17	7	11	12
Ata 5	33	4	4	6	5	17	22	9	18	16
Ata 6	11	3	4	6	4	9	9	4	7	5
Ata 7	20	3	7	5	4	11	14	5	5	4
Ata 8	35	4	8	3	8	12	17	5	11	9
Ata 9	24	3	5	3	6	11	11	3	9	9
Ata 10	50	4	5	4	7	31	29	5	9	8
Ata 11	43	4	7	5	7	29	19	5	9	12
Ata 12	56	3	10	4	16	33	25	4	13	11
<b>Total</b>	<b>366</b>	<b>50</b>	<b>69</b>	<b>73</b>	<b>79</b>	<b>209</b>	<b>186</b>	<b>71</b>	<b>140</b>	<b>130</b>

# Exemplo de criação de uma segmentação de referência

Referência criada a partir da concordância entre segmentações manuais.



# Configuração Experimental - Segmentadores

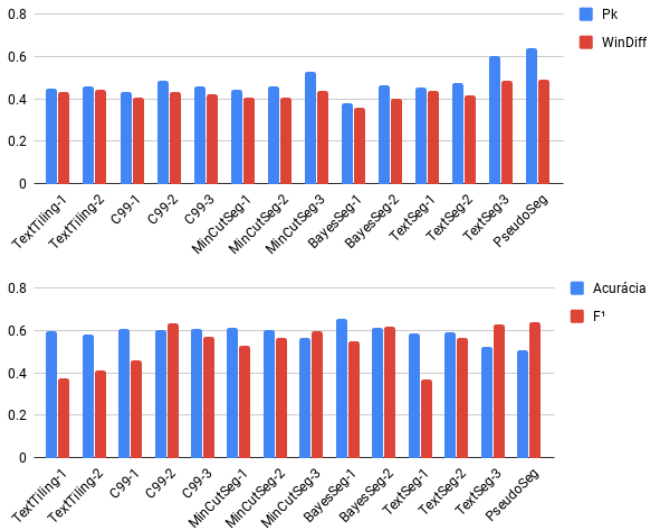
<b>Algoritmo</b>	<b>Parâmetros (Configuração)</b>					
<i>Text Tiling</i>	<b>Win</b>	(20-60)	<b>Step</b>	(30-55)		
<i>C99</i>	<b>SR</b>	(.2-.7)	<b>W</b>	(sim/não)	<b>RS</b>	(3-7)
<i>BayesSeg</i>	<b>SR</b>	(auto, .3-.9)	<b>Prior</b>	(.08-.11)	<b>Disp</b>	(.1-.7)
<i>MinCut</i>	<b>SR</b>	(.2-.7)	<b>LenCut</b>	(5-15)		
<i>TextSeg</i>	<b>SR</b>	(auto, .1-.9)				
<i>PseudoSeg</i>						

## Resumo dos melhores resultados obtidos por cada configuração

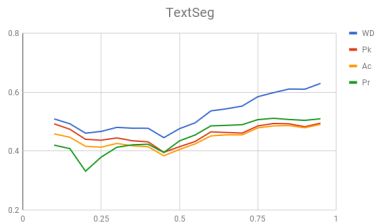
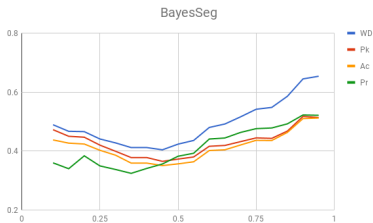
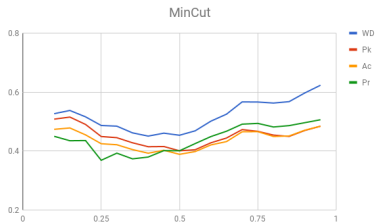
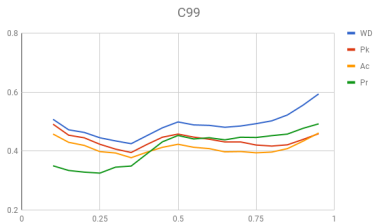
Algoritmo		Step	Win	$P_k$	WD	Ac	$F^1$	#Segs
TextTiling-1		20	30	0.461	0.444	0.581	<b>0.411</b>	8.833
TextTiling-2		30	45	<b>0.450</b>	<b>0.435</b>	<b>0.596</b>	0.373	6.417
Algoritmo	RS	W	SRate	$P_k$	WD	Ac	$F^1$	#Segs
C99-1	3	true	0.300	<b>0.434</b>	<b>0.407</b>	0.607	0.457	9.250
C99-2	3	true	0.700	0.485	0.431	0.602	<b>0.633</b>	21.417
C99-3	5	true	0.500	0.460	0.421	<b>0.609</b>	0.571	15.500
Algoritmo		Cut	SRate	$P_k$	WD	Ac	$F^1$	#Segs
MinCutSeg-1		9	0.400	<b>0.444</b>	0.408	<b>0.614</b>	0.526	11.917
MinCutSeg-2		11	0.500	0.459	<b>0.407</b>	0.603	0.563	15.000
MinCutSeg-3		5	0.700	0.528	0.438	0.567	<b>0.599</b>	21.000
Algoritmo	Prior	Disp.	SRate	$P_k$	WD	Ac	$F^1$	#Segs
BayesSeg-1	0.080	0.500	Auto	<b>0.380</b>	<b>0.361</b>	<b>0.655</b>	0.551	10.000
BayesSeg-2	0.110	0.100	0.600	0.462	0.399	0.615	<b>0.619</b>	18.417
Algoritmo			SRate	$P_k$	WD	Ac	$F^1$	#Segs
TextSeg-1			Auto	<b>0.455</b>	0.439	0.585	0.368	6.417
TextSeg-2			0.500	0.475	<b>0.417</b>	<b>0.594</b>	0.566	15.500
TextSeg-3			0.900	0.604	0.484	0.524	<b>0.627</b>	27.500
Algoritmo			SRate	$P_k$	WD	Ac	$F^1$	#Segs
PseudoSeg			1.000	<b>0.640</b>	<b>0.490</b>	<b>0.506</b>	<b>0.638</b>	30.500



## Resumo dos melhores resultados obtidos por cada configuração



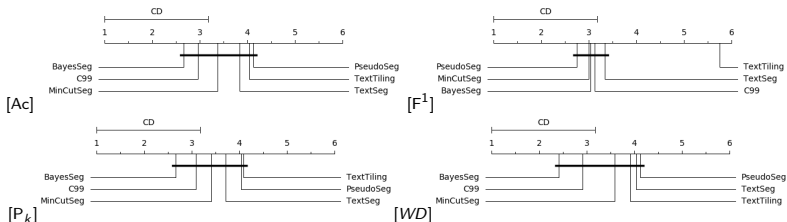
## Influência da taxa de segmentos na eficiência dos algoritmos



Testes de significância estatística de Friedman com pós teste de Nemenyi

- Inicialmente com as configurações de cada algoritmo.
- Novamente com as melhores configurações de cada algoritmo.

Diagramas de Diferença Crítica sobre *ranking* considerando valores de Acurácia,  $F^1$ , *WindowDiff*, e  $P_k$ .



Não há diferença significativa entre os métodos.

Os modelos de Extração de Tópicos foram avaliados junto aos usuários.

- Resultados de 2 consultas ao Sistema usando
- 3 Extratores (K-Means, LDA, PLSA);
- Impressões dos usuários coletadas via questionários.

Entrada:

- 1 “defesa de dissertação”;
- 2 “compra de equipamentos”.

*Corpus*

- Formado por 175 atas;
- Segmentadas com o *BayesSeg*;
- Conjunto de 1276 segmentos.

Resultados utilizando 3 modelos para Extração de Tópicos:

- 1 K-Means;
- 2 LDA;
- 3 PLSA.

Os resultados foram apresentados à um grupo de avaliadores:

- 24 profissionais da UFSCar;
- 13 profissionais de escolas técnicas;
- 03 profissionais de escolas do Ensino Fundamental.

Perfil:

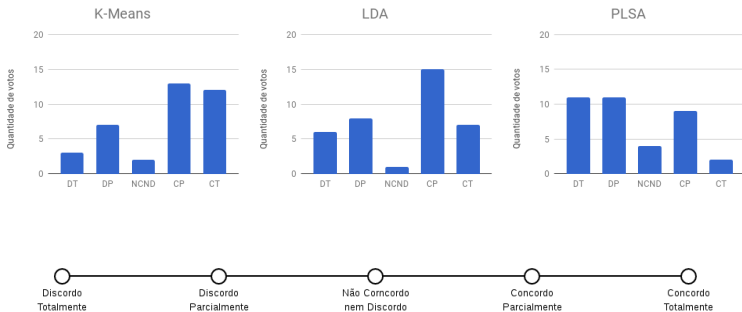
- 17 membros de conselhos;
- 12 gestores;
- 05 administrativos;
- 03 professores;
- 03 sem afinidade com atas (descartados).

Coletar respostas referentes à:

- 1 Coesão dos tópicos;
- 2 Representatividade dos descritores;
- 3 Coesão dos segmentos;
- 4 Completude dos segmentos.

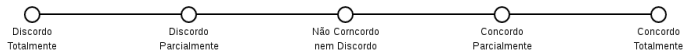
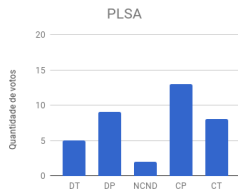
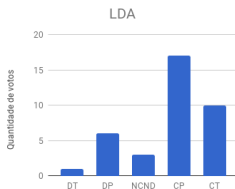
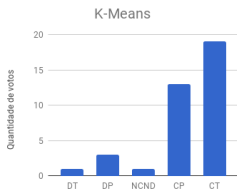


Primeira questão: *“Todos os trechos apresentados compartilham um mesmo assunto.”*.

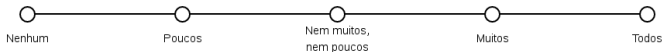
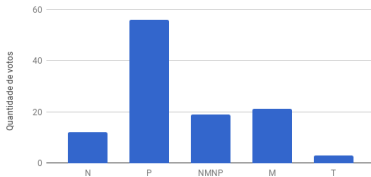


# Representatividade dos descritores

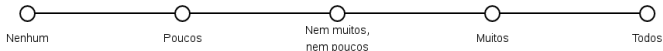
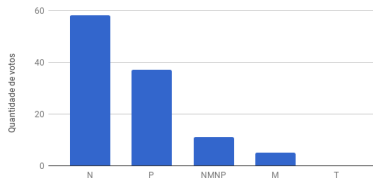
Segunda questão: “As palavras <descritores> resumem bem o assunto tratado nos trechos.”.



Terceira questão: *“Existem trechos que não tratam de um único assunto?”*.



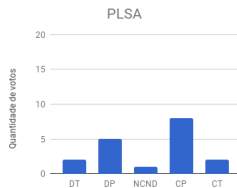
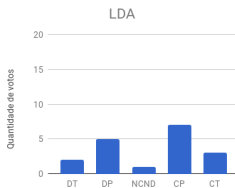
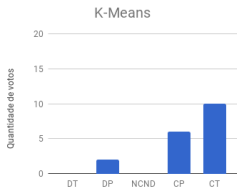
Quarta questão: *“Existem trechos incompletos e insuficientes para compreensão do assunto do trecho?”*.



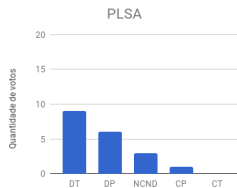
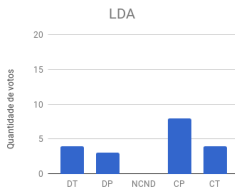
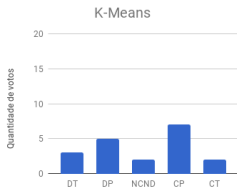
## Coesão dos tópicos

Primeira questão: *"Todos os trechos apresentados compartilham um mesmo assunto."*

### Primeira Consulta: *"defesa de dissertação"*



### Segunda Consulta: *"compra de equipamentos"*

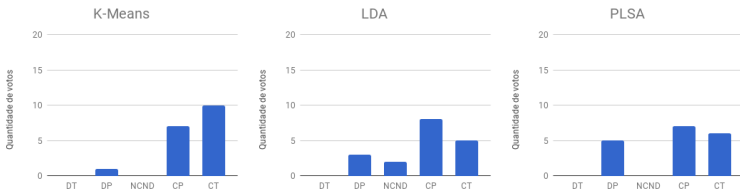


# Comportamento dos extratores em consultas diferentes

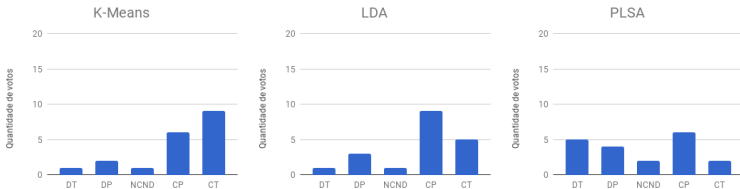
## Representatividade dos descritores

Segunda questão: "As palavras <descritores> resumem bem o assunto tratado nos trechos."

### Primeira Consulta: "defesa de dissertação"



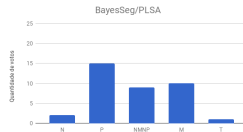
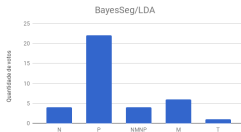
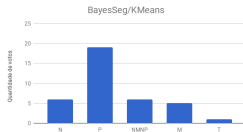
### Segunda Consulta: "compra de equipamentos"



# Qualidade dos segmentos apresentados em diferentes técnicas

## Coesão dos segmentos

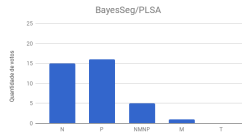
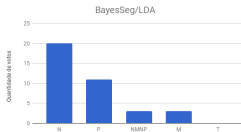
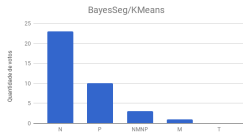
Terceira questão: “Existem trechos que não tratam de um único assunto?”.



# Qualidade dos segmentos apresentados em diferentes técnicas

## Compleitude dos Segmentos

Quarta questão: “Existem trechos incompletos e insuficientes para compreensão do assunto do trecho?”.





- 1 Contextualização
- 2 Objetivos
- 3 Proposta
- 4 Avaliação Experimental
- 5 Conclusão**
- 6 Trabalhos Futuros

A metodologia utilizada nesse trabalho:

- Conecta as técnicas de segmentação textual aos modelos de Extração de Tópicos;
- Gera um estrutura derivada de um *corpus* não estruturado;
- Utiliza variáveis latentes em conjunto com técnicas de Recuperação de Informação.

## Segmentação

### Resultados:

- Medidas abaixo do esperado;
- Impressões satisfatórias dos usuários;
  - Completude;
  - Coesão;

### Possíveis melhorias:

- Segmentação de referência com mais anotadores;
- Treinamento dos anotadores;
- Maior concordância entre anotadores;
- Segmentação de referência mais confiável e representativa.

## Extração de Tópicos

- Melhores resultados com o K-Means;
  - Coesão dos grupos;
  - Capacidade representativa do descritores;

## Contribuições:

- O método para extração de conhecimento em documentos multi-temáticos na língua portuguesa;
- O corpus de atas anotadas;
- A ferramenta para segmentação e anotação manual;
- O sistema proposto e sua implementação;
- As avaliações dos Segmentadores e Extratores de Tópicos.

- 1 Contextualização
- 2 Objetivos
- 3 Proposta
- 4 Avaliação Experimental
- 5 Conclusão
- 6 Trabalhos Futuros**

Adição de novas técnicas ao sistema:

- Algoritmos de agrupamento incremental
- Categorização dos segmentos (decisão, informe, irrelevante);

Melhorias:

- Inclusão de novos corpora (transcrições de conversas, diálogos em chats, discursos e atas de outras organizações)
- Fontes externas para melhorar os métodos de segmentação textual (*thesaurus* e *clue words*);
- Testes voltados a experiência do usuário.

- Submissão de artigo.

Obrigado!