



**Ian Monteiro Nunes**

**Agrupamento de Registros Textuais  
Baseado em Similaridade Entre Textos**

**Dissertação de Mestrado**

Dissertação apresentada ao Programa de Pós-graduação em Informática da PUC-Rio como requisito parcial para obtenção do título de Mestre em Informática. Aprovada pela Comissão Examinadora abaixo assinada.

Orientador: Prof. Ruy Luiz Milidiú

Rio de Janeiro

Setembro de 2008



**Ian Monteiro Nunes**

## **Agrupamento de Registros Textuais**

### **Baseado em Similaridade Entre Textos**

Dissertação apresentada como requisito parcial para a obtenção do grau de Mestre pelo Programa de Pós-graduação em Informática do Departamento de Informática do Centro Técnico e Científico da PUC-Rio. Aprovada pela Comissão Examinadora abaixo assinada.

**Prof. Ruy Luiz Milidiú**

Orientador

Departamento de Informática – PUC-Rio

**Prof. Marco Antonio Casanova**

Departamento de Informática – PUC-Rio

**Prof. Rubens Nascimento Melo**

Departamento de Informática – PUC-Rio

**Prof. José Eugenio Leal**

Coordenador Setorial do Centro

Técnico Científico – PUC-Rio

Rio de janeiro, 04 de abril de 2008

Todos os direitos reservados. É proibida a reprodução total ou parcial do trabalho sem autorização da universidade, do autor e do orientador.

### Ian Monteiro Nunes

Graduou-se em Engenharia de Computação na PUC-Rio em 2003. Desenvolveu diversos sistemas para o meio corporativo, sendo hoje usados por empresas de abrangência nacional. Responsável pela área de pesquisa e desenvolvimento da empresa Dínamo DM desde 2004.

#### Ficha catalográfica

Nunes, Ian Monteiro

Agrupamento de registros textuais baseado em similaridade entre textos / Ian Monteiro Nunes ; orientador: Ruy Luiz Milidiú. – 2008.

69 f.; 30 cm

Dissertação (Mestrado em Informática)–Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2008.

Inclui bibliografia

1. Informática – Teses. 2. Aprendizado de máquina. 3. Mineração de textos. 4. Deduplicação. 5. Recuperação de informação. I. Milidiú, Ruy Luiz. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. III Título.

CDD: 004

## Agradecimentos

Aos meus pais, Ivan e Isabel, por todo carinho, atenção, amor, apoio, paciência e tudo mais ao longo de toda a minha vida.

Ao meu orientador professor Ruy Luiz Milidiú pela parceria, disposição e ensinamentos em todas as etapas deste trabalho.

Um agradecimento especial ao meu amigo Julio Brafman por toda sua colaboração no desenvolvimento deste trabalho.

Aos meus amigos Christian Nunes, Lucas Sigaud e Eduardo Gouveia por todas as discussões e ajuda dispensada a este trabalho.

Aos professores que prontamente aceitaram participar da comissão examinadora.

À PUC-Rio pelo auxílio e estrutura fornecidos, sem os quais este trabalho não poderia ter sido realizado.

Ao professor Marcus Poggi por toda a orientação e ensinamentos ao longo de toda a minha vida acadêmica.

Aos professores e funcionários do Departamento de Informática.

A todos os amigos que direta ou indiretamente participaram deste trabalho, mesmo que não tenham sido citados nominalmente, sou igualmente grato.

## Resumo

Nunes, Ian Monteiro; Milidiú, Ruy Luiz. **Agrupamento de Registros Textuais Baseado em Similaridade Entre Textos**. Rio de Janeiro, 2008, 69p. Dissertação de Mestrado. Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro

O presente trabalho apresenta os resultados que obtivemos com a aplicação de grande número de modelos e algoritmos em um determinado conjunto de experimentos de agrupamento de texto. O objetivo de tais testes é determinar quais são as melhores abordagens para processar as grandes massas de informação geradas pelas crescentes demandas de *data quality* em diversos setores da economia. O processo de deduplicação foi acelerado pela divisão dos conjuntos de dados em subconjuntos de itens similares. No melhor cenário possível, cada subconjunto tem em si todas as ocorrências duplicadas de cada registro, o que leva o nível de erro na formação de cada grupo a zero. Todavia, foi determinada uma taxa de tolerância intrínseca de 5% após o agrupamento. Os experimentos mostram que o tempo de processamento é significativamente menor e a taxa de acerto é de até 98,92%. A melhor relação entre acurácia e desempenho é obtida pela aplicação do algoritmo K-Means com um modelo baseado em trigramas.

## Palavras-chave

Aprendizado de máquina; mineração de textos; deduplicação; recuperação de informação.

## Abstract

Nunes, Ian Monteiro; Milidiú, Ruy Luiz. **Clustering Text Structured Data Based on Text Similarity**. Rio de Janeiro, 2008, 69p. MSc. Dissertation. Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

This document reports our findings on a set of text clustering experiments, where a wide variety of models and algorithms were applied. The objective of these experiments is to investigate which are the most feasible strategies to process large amounts of information in face of the growing demands on data quality in many fields. The process of deduplication was accelerated through the division of the data set into individual subsets of similar items. In the best case scenario, each subset must contain all duplicates of each produced register, mitigating to zero the cluster's errors. It is established, although, a tolerance of 5% after the clustering process. The experiments show that the processing time is significantly lower, showing a 98.92% precision. The best accuracy/performance relation is achieved with the K-Means Algorithm using a trigram based model.

## Keywords

Machine learning; text mining; deduplicate; information retrieval.

## Sumário

<b>1</b>	<b>INTRODUÇÃO.....</b>	<b>13</b>
<b>2</b>	<b>FRAMEWORK PARA APRENDIZADO DE MÁQUINA .....</b>	<b>17</b>
<b>2.1.</b>	<b>PROCESSO, CARACTERÍSTICAS E BENEFÍCIOS.....</b>	<b>17</b>
<b>2.2.</b>	<b>DESCRIÇÃO DAS ETAPAS DO PROCESSO .....</b>	<b>19</b>
2.2.1.	LIMPEZA.....	19
2.2.2.	VALIDAÇÃO.....	20
2.2.3.	PADRONIZAÇÃO.....	22
2.2.4.	QUALIFICAÇÃO.....	25
2.2.5.	SEGMENTAÇÃO.....	26
2.2.6.	DEDUPLICAÇÃO .....	28
<b>3</b>	<b>APRENDIZADO DE MÁQUINA - CONCEITOS BÁSICOS.....</b>	<b>30</b>
<b>4</b>	<b>SOLUÇÕES.....</b>	<b>32</b>
<b>4.1.</b>	<b>ESTRATÉGIAS.....</b>	<b>32</b>
4.1.1.	SEM PARTICIONAMENTO .....	32
4.1.2.	COM PARTICIONAMENTO .....	32
<b>4.2.</b>	<b>BASE VERDADE .....</b>	<b>33</b>
<b>4.3.</b>	<b>BASE DE TESTE .....</b>	<b>33</b>
<b>4.4.</b>	<b>ERRO.....</b>	<b>33</b>
<b>4.5.</b>	<b>ESTRATÉGIA VETORIAL .....</b>	<b>34</b>
4.5.1.	AGRUPAMENTOS.....	34
4.5.1.1.	Algoritmos .....	34
4.5.1.1.1.	K-Means .....	34
4.5.1.1.1.1.	Descrição .....	34
4.5.1.1.1.2.	Pseudocódigo.....	35
4.5.1.1.2.	K-Medoid .....	35
4.5.1.1.2.1.	Descrição .....	35
4.5.1.1.2.2.	Pseudocódigo.....	36
4.5.1.1.3.	GNG .....	36
4.5.1.1.3.1.	Descrição .....	36
4.5.1.1.3.2.	Pseudocódigo.....	38
4.5.2.	REPRESENTAÇÃO .....	39
4.5.2.1.	Completa.....	39
4.5.2.2.	Tipos de <i>token</i> .....	39

4.5.2.2.1. Caractere.....	39
4.5.2.2.2. Digrama.....	40
4.5.2.2.3. Trigrama.....	40
4.5.2.3. Pesos .....	40
4.5.2.4. Tamanhos .....	41
4.5.2.5. Heurísticas Lingüísticas.....	41
4.5.3. EXPERIMENTOS .....	43
4.5.3.1. <i>Corpus</i> .....	43
4.5.3.2. Metodologia.....	44
4.5.3.2.1. Primeira etapa .....	44
4.5.3.2.2. Segunda etapa .....	47
4.5.3.3. Resultados .....	48
4.5.3.4. K-Means.....	49
4.5.3.5. GNG.....	52
4.5.3.6. K-Medóide.....	54
<b>4.6. ESTRATÉGIA NÃO-VETORIAL.....</b>	<b>54</b>
4.6.1. AGRUPAMENTO .....	55
4.6.1.1. Algoritmos .....	55
4.6.1.1.1. Modóide.....	55
4.6.1.1.2. Medóide.....	55
4.6.1.1.2.1. Pseudocódigo.....	56
4.6.2. REPRESENTAÇÃO .....	56
4.6.3. CENTRÓIDE .....	56
4.6.3.1. Seleção estatística .....	56
4.6.3.1.1. Moda posicional.....	56
4.6.3.1.1.1. Moda do registro inteiro .....	59
4.6.3.1.1.2. Moda de cada palavra e média de palavras.....	60
4.6.3.1.2. Medóide.....	61
4.6.3.1.2.1. Aleatório .....	61
4.6.3.1.2.2. Múltiplos .....	61
4.6.3.1.3. Combinando Medóide e Modóide .....	61
4.6.3.2. Medidas de confiança.....	62
4.6.3.3. Heurísticas Lingüísticas.....	62
4.6.4. EXPERIMENTOS .....	62
4.6.4.1. Metodologia.....	62
4.6.4.2. Resultados .....	63
4.6.4.2.1. K-Modóide .....	63



4.6.4.2.2. K-Modóide híbrido com K-Medóide.....	64
<b>5 RESULTADO FINAL .....</b>	<b>65</b>
<b>6 CONCLUSÕES.....</b>	<b>66</b>
<b>7 OS PRÓXIMOS PASSOS.....</b>	<b>67</b>
<b>8 REFERÊNCIAS .....</b>	<b>68</b>

## Índice de tabelas

Tabela 1: Tempos para comparação dois-a-dois .....	15
Tabela 2: Identificação de termos impróprios e/ou obscenos.....	19
Tabela 3: Evolução da validação dos dados .....	21
Tabela 4: Endereços validados após tratamento .....	22
Tabela 5: Endereços sem validação .....	22
Tabela 6: Padronização de unidades monetárias e de medidas .....	23
Tabela 7: Padronização de veículos .....	24
Tabela 8: Qualificação .....	26
Tabela 9: Deduplicação .....	29
Tabela 10: Exemplo de representação completa. ....	39
Tabela 11: Tokens por caractere. ....	39
Tabela 12: Tokens por digrama. ....	40
Tabela 13: Tokens por trígama. ....	40
Tabela 14: Pesos para caracteres inicial e final .....	41
Tabela 15: Diferentes pesos até o centro.....	41
Tabela 16: Fonemas do português brasileiro .....	43
Tabela 17: Taxas de erro da primeira rodada da primeira etapa.....	45
Tabela 18: Taxas de erro da segunda rodada da primeira etapa .....	45
Tabela 19: Resultados para K-Means com caracteres como tokens.....	49
Tabela 20: Resultados para K-Means com digramas como tokens.....	50
Tabela 21: Resultados para K-Means com trigramas como tokens.....	51
Tabela 22: Resultados para GNG com caracteres como tokens.....	52
Tabela 23: Resultados para GNG com digramas como tokens.....	52
Tabela 24: Resultados para GNG com trigramas como tokens.....	53
Tabela 25: Resultados para K-Medóide .....	54
Tabela 26: Dados para exemplo .....	57

Tabela 27: Matriz de frequência total .....	57
Tabela 28: Matriz de frequência para a primeira palavra .....	58
Tabela 29: Matriz de frequência para segunda palavra.....	58
Tabela 30: Moda de todo o registro – dados originais.....	59
Tabela 31: Moda de todo o registro – novos dados .....	59
Tabela 32: Dados para exemplo .....	60
Tabela 33: Matriz de frequências .....	60
Tabela 34: Resultados obtidos pelo K-Modóide .....	63
Tabela 35: K-Modóide híbrido com K-Medóide .....	64
Tabela 36: Resultado final .....	65

## Índice de figuras

Figura 1: CRM no mercado brasileiro .....	14
Figura 2: Data quality.....	18
Figura 3: Padronização de nomes de universidades.....	24
Figura 4: Resultado da Padronização de cargos.....	25
Figura 5: Ação segmentada de venda .....	28
Figura 6: Erro.....	33
Figura 7: Pseudocódigo para o <i>K-Means</i> .....	35
Figura 8: Pseudocódigo para o K-Medóide .....	36
Figura 9: Pseudocódigo para o GNG original pelo autor. ....	38
Figura 10: Pseudocódigo para o K-Modóide .....	56