



David Steinbruch

**Um estudo de algoritmos para classificação automática de
textos utilizando naive-Bayes**

Dissertação de Mestrado

Dissertação apresentada como requisito parcial para
obtenção do título de Mestre pelo Programa de Pós-
Graduação em Informática da PUC-Rio.

Orientador: Daniel Schwabe
Co-orientador: Ruy Luiz Milidiú

Rio de Janeiro, setembro de 2006



David Steinbruch

Um estudo de algoritmos para classificação automática de textos utilizando naive-Bayes

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

Daniel Schwabe

Orientador
PUC-Rio

Ruy Luiz Milidiú

Co-orientador
PUC-Rio

Marcus Vinicius Soledade Poggi de Aragão

PUC-Rio

Eduardo Sany Laber

PUC-Rio

Prof. José Eugenio Leal

Coordenador Setorial do Centro Técnico Científico - PUC-Rio

Rio de Janeiro

05 de setembro de 2006

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

David Steinbruch

Graduou-se em Engenharia de Computação na PUC-Rio em 2003. Atuou como programador no desenvolvimento de soluções web e aplicações de mineração de textos. Possui interesse acadêmico e profissional nas áreas de Inteligência Artificial, Hipertexto e Multimídia.

Ficha Catalográfica

Steinbruch, David

Um estudo de algoritmos para classificação automática de textos utilizando naive-Bayes / David Steinbruch ; orientador: Daniel Schwabe. – 2006.

78 f. : il. ; 30 cm

Dissertação (Mestrado em Informática)– Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2006.

Inclui bibliografia

1. Informática – Teses. 2. Aprendizado de máquina. 3. Categorização de textos. 4. Classificação de textos. 5. Multirótulo. 6. Naive-Bayes. 7. Internet. I. Schwabe, Daniel. II. Milidiú, Ruy Luiz. III. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. IV. Título.

CDD: 004

Este trabalho é dedicado a minha família,
pela compreensão e pelo apoio,
e ao meu orientador e coorientador,
pela motivação.

Agradecimentos

À PUC-Rio, ao departamento de informática e ao CNPq pela oportunidade.

Ao meu orientador, Daniel Schwabe, e co-orientador, Ruy Luiz Milidiú, pela motivação, paciência, confiança e ajuda.

A todos os professores, funcionários do Departamento de Informática pelo apoio dado quando precisei.

A Deborah e Emanuelle pela paciência e ajuda nos processos burocráticos.

Aos meus amigos que cursaram o mestrado comigo e familiares pelo incentivo e apoio nos momentos difíceis.

Aos meus irmãos Daniel, Rachel e Natan pelos ótimos momentos juntos.

Aos meus pais Luna e Beni e meu padrasto Samy por tudo.

Resumo

Steinbruch, David; Schwabe, Daniel. **Um estudo de algoritmos para classificação automática de textos utilizando naive-Bayes**. Rio de Janeiro, 2006. 78p. Dissertação de Mestrado - Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

A quantidade de informação eletrônica vem crescendo de forma acelerada, motivada principalmente pela facilidade de publicação e divulgação que a Internet proporciona. Desta forma, é necessária a organização da informação de forma a facilitar a sua aquisição. Muitos trabalhos propuseram resolver este problema através da classificação automática de textos associando a eles vários rótulos (classificação multirótulo). No entanto, estes trabalhos transformam este problema em subproblemas de classificação binária, considerando que existe independência entre as categorias. Além disso, utilizam limiares ("thresholds"), que são muito específicos para o conjunto de treinamento utilizado, não possuindo grande capacidade de generalização na aprendizagem. Esta dissertação propõe dois algoritmos de classificação automática de textos baseados no algoritmo multinomial naive Bayes e sua utilização em um ambiente on-line de classificação automática de textos com realimentação de relevância pelo usuário. Para testar a eficiência dos algoritmos propostos, foram realizados experimentos na base de notícias Reuters 21758 e na base de documentos médicos Ohsumed.

Palavras-chave

Aprendizado de Máquina; Categorização de Textos, Classificação de Textos; Multirótulo; Naive-Bayes; Internet

Abstract

Steinbruch, David; Schwabe, Daniel. **A study of multilabel text classification algorithms using naive-Bayes.** Rio de Janeiro, 2006. 78p. MSc. Dissertation - Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

The amount of electronic information has been growing fast, mainly due to the easiness of publication and spreading that Internet provides. Therefore, is necessary the organisation of information to facilitate its retrieval. Many works have solved this problem through the automatic text classification, associating to them several labels (multilabel classification). However, those works have transformed this problem into binary classification subproblems, considering there is not dependence among categories. Moreover, they have used thresholds, which are very sepecific of the classifier document base, and so, does not have great generalization capacity in the learning process. This thesis proposes two text classifiers based on the multinomial algorithm naive Bayes and its usage in an on-line text classification environment with user relevance feedback. In order to test the proposed algorithms efficiency, experiments have been performed on the Reuters 21578 news base, and on the Ohsumed medical document base.

Keywords

Machine Learning; Text Categorization; Text Classification; Multilabel; Naive-Bayes; Internet

Sumário

1	Introdução	15
1.1	Caracterização do problema	15
1.2	Objetivo	16
1.3	Trabalhos relacionados	17
1.4	Organização da dissertação	18
2	Fundamentos Teóricos	20
2.1	Aprendizado de Máquina	20
2.2	O problema da classificação	22
2.3	Classificação automática de textos	24
2.4	Classificação unirótulo e multirótulo	25
2.5	Conjuntos de treinamento, validação e teste	25
2.6	Representação de documentos	26
2.7	Pré-processamento de documentos	27
2.7.1	Stopwords	28
2.7.2	Stemming	28
2.7.2.1.	Método de Stemmer S	28
2.7.2.2.	Método de Porter	29
2.7.2.3.	Método de Lovins	29
2.7.3	Frequência de Documentos (DF)	29
2.7.4	Ganho de Informação (IG)	30
2.7.5	Informação Mútua (MI)	30
2.7.6	Estatística χ^2 (CHI)	31
2.8	Classificadores probabilísticos	31
2.8.1	Naive Bayes	32
2.8.1.1.	Modelo binário	33
2.8.1.2.	Modelo multinomial	34
2.9	Medidas de desempenho	34
2.9.1	Matriz de contingência	34

2.9.2 Precision e Recall	35
2.10 Métodos de avaliação de classificadores	37
2.10.1 Resubstituição	37
2.10.2 Holdout	37
2.10.3 Amostragem aleatória	38
2.10.4 K-Fold Cross Validation	38
2.10.5 Leave-One-Out	38
2.10.6 Bootstrap	39
3 Algoritmos propostos	40
3.1 Classificador pseudo-multirótulo	41
3.2 Classificador multirótulo	43
4 Experimentos	48
4.1 Introdução	48
4.2 Bases de dados	48
4.2.1 Reuters-21578	48
4.2.2 Ohsumed	50
4.3 Experimentos Realizados	51
4.3.1 Experimentos com a base Reuters R(10)	51
4.3.2 Experimentos com a base Reuters R(90)	56
4.3.3 Conclusões dos experimentos com as bases da Reuters	62
4.3.4 Experimentos com a base Ohsumed	62
4.3.5 Conclusões dos experimentos com a base Ohsumed	67
4.3.6 Comparação com outros trabalhos	68
4.3.6.1. Reuters	68
4.3.6.2. Ohsumed	69
5 Conclusão	71
5.1 Contribuições	73
5.2 Trabalhos futuros	73
6 Referências bibliográficas	75

Siglas

Micro Recall	Micro Averaged Recall
Macro Recall	Macro Averaged Recall
Micro Precision	Micro Averaged Precision
Macro Precision	Macro Averaged Precision
Micro F1	Micro Averaged F1
Macro F1	Macro Averaged F1
MESH	Medical Subject Heading
R(10)	Subconjunto de 10 categorias da base de dados Reuters 21758
R(90)	Subconjunto de 90 categorias da base de dados Reuters 21758
CONSTRUE	Categorization of News STories, Rapidly, Uniformly and Extensibly
TF	Term frequency
IDF	Inverse Document Frequencie
Reuters 21578	Base de notícias da Reuters com 21578 documentos
Ohsumed	Subconjunto de 348.566 documentos da base MEDLINE
MEDLINE	Medical Literature Analysis and Retrieval System Online
ModApté	Subonjunto de documentos da base Reuters 21578, composto de 9.603 documentos de treinamento e 3.299 documentos de teste
DF	Document Frequency
IG	Information Gain
MI	Mutual Information
CHI	CHI-squared statistic

Lista de tabelas

Tabela 1 – Matriz de contingência de um classificador.	35
Tabela 2 – Matriz de contingência de um classificador binário.	35
Tabela 3 – Resultados do algoritmo pseudo-multirótulo na base R(10) para as 10 categorias que compõem a base.	52
Tabela 4 – Resultados globais do algoritmo pseudo-multirótulo na base R(10).	52
Tabela 5 – Tempo de execução da fase de treinamento e da fase de classificação do algoritmo pseudo-multirótulo na base R(10).	52
Tabela 6 – Resultados do algoritmo multirótulo na base R(10) para as 10 categorias que compõem a base.	52
Tabela 7 – Resultados globais do algoritmo multirótulo na base R(10).	53
Tabela 8 – Tempo de execução da fase de treinamento e da fase de classificação do algoritmo multirótulo na base R(10).	53
Tabela 9 – Resultados do algoritmo pseudo-multirótulo na base R(90) para as 10 categorias com maior quantidade de documentos de treinamento.	57
Tabela 10 – Resultados globais do algoritmo pseudo-multirótulo na base R(90).	57
Tabela 11 – Tempo de execução da fase de treinamento e da fase de classificação do algoritmo pseudo-multirótulo na base R(90).	57
Tabela 12 – Resultados do algoritmo multirótulo na base R(90) para as 10 categorias com maior quantidade de documentos de treinamento.	58
Tabela 13 – Resultados globais do algoritmo multirótulo na base R(90).	58
Tabela 14 – Tempo de execução da fase de treinamento e da fase de classificação do algoritmo multirótulo na base R(90).	58
Tabela 15 – Resultados do algoritmo pseudo-multirótulo na base Ohsumed para 5 categorias.	63
Tabela 16 – Resultados globais do algoritmo pseudo-multirótulo na base Ohsumed.	63
Tabela 17 – Tempo de execução da fase de treinamento e da fase de classificação do algoritmo pseudo-multirótulo na base Ohsumed.	63
Tabela 18 – Resultados do algoritmo multirótulo na base Ohsumed para 5 categorias.	63

Tabela 19 – Resultados globais do algoritmo multirótulo na base Ohsumed.	63
Tabela 20 – Tempo de execução da fase de treinamento e da fase de classificação do algoritmo multirótulo na base Ohsumed	63
Tabela 21 – Resultados de Bennett et al. [2002] na base R(10)	68
Tabela 22 – Resultados de Sebastiani & Debole [2004] na base de dados R(10)	69
Tabela 23 – Resultados de Sebastiani & Debole [2004] na base de dados R(90)	69
Tabela 24 – Resultados de Yang & Liu [1999] na base de dados R(90)	69
Tabela 25 – Resultados de Moschitti [2003b] em 5 categorias da base de dados Ohsumed	70
Tabela 26 – Resultados globais de Moschitti [2003b] na base de dados Ohsumed	70

Lista de figuras

Figura 1 – Hierarquia do Aprendizado	24
Figura 2 - Exemplo da regra de decisão do algoritmo multirótulo proposto.	45
Figura 3 - Resultados Micro Recall para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(10).	54
Figura 4 - Resultados Macro Recall para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(10).	54
Figura 5 - Resultados Micro Precision para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(10).	55
Figura 6 - Resultados Macro Precision para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(10).	55
Figura 7 - Resultados Micro F1 para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(10).	56
Figura 9 – Resultados Micro Recall para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(90).	59
Figura 10 – Resultados Macro Recall para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(90).	59
Figura 11 – Resultados Micro Precision para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(90).	60
Figura 12 – Resultados Macro Precision para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(90).	60
Figura 13 – Resultados Micro F1 para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(90).	61
Figura 14 – Resultados Macro F1 para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base R(90).	61
Figura 15 – Resultados Micro Recall para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base Ohsumed.	64
Figura 16 – Resultados Macro Recall para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base Ohsumed.	65
Figura 17 – Resultados Micro Precision para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base Ohsumed.	65

Figura 18 – Resultados Macro Precision para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base Ohsumed.	66
Figura 19 – Resultados Micro F1 para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base Ohsumed	66
Figura 20 – Resultados Macro F1 para o algoritmo pseudo-multirótulo e para o algoritmo multirótulo na base Ohsumed	67