

Conversor XML para TXT: ferramenta para análise de presença social associada a clustering no *software Eureka!*

Yasmin C. Martins, Breno F. T. Azevedo, Hélvia P. P. Bastos

Núcleo de Informática na Educação (NIE)
Instituto Federal Fluminense Campus Campos-Centro (IFF) – Campos dos Goytacazes, RJ – Brasil
{helviabastos, brenoter}@gmail.com, nim_asay@hotmail.com

Abstract: The paper presents a tool able to reconfigure XML files resulting from the automatic analysis of text indicators of social presence performed by the software Presente! (Kambara-Silva, 2011). For clustering experiments in the software Eureka! (Wives, 2004), these XML files must be converted into TXT format. The converter described in the study enables such files to be graphically analyzed, thus helping in decision-making processes.

Resumo: O trabalho apresenta uma ferramenta que realiza a reconfiguração de arquivos XML resultantes do processo de análise automática de marcas textuais de presença social realizado pelo software Presente! (Kambara-Silva, 2011). Para experimentos de agrupamento no software Eureka! (Wives, 2004), os arquivos XML precisam ser formatados em TXT. O conversor descrito neste estudo possibilita que esses arquivos possam ser analisados graficamente, fato que auxilia a tomada de decisão por parte do analista.

1. Introdução

Este trabalho resulta de atividades de pesquisa realizadas para o projeto “Mineração de Dados Textuais Eletrônicos” (Núcleo de Informática na Educação – NIE). A mineração de textos (MT) pode ser definida como a extração de informações úteis de uma série de documentos textuais. Constituindo uma vasta área de estudos, a MT tem aplicação em diversos campos, inclusive a busca de dados textuais em bases eletrônicas de apoio a projetos pedagógicos.



Secretaria de Educação
Profissional e Tecnológica



Ministério
da Educação





O objetivo do artigo é apresentar uma ferramenta desenvolvida para otimizar tarefas de clustering (agrupamento) realizada pelo programa Eureka! [Wives 2004]. O agrupamento de textos (clustering) é uma técnica utilizada para agrupar documentos semelhantes. No agrupamento, os documentos podem aparecer em vários subtópicos nos resultados da pesquisa. Um algoritmo de agrupamento básico cria um vetor de tópicos para cada texto, e calcula os pesos para identificar em qual grupo um documento deve fazer parte [Gupta e Lehal, 2009; Kobayashi e Aono, 2004].

O experimento de conversão descrito neste estudo utilizou arquivos no formato XML resultantes do processo de mineração feito pelo software Presente! [Kambara-Silva 2011]. Esse programa foi desenvolvido para identificar pistas textuais de presença social em ambientes virtuais de aprendizagem conforme modelo de Bastos [2012]. O objetivo desse processo computacional é auxiliar o docente no acompanhamento da participação dos alunos em cursos a distância. As diferentes partes que compõem o programa, seus objetivos e funcionamento são descritas na Seção 2.

A Seção 3 descreve como a ferramenta criada na pesquisa realiza a conversão de arquivos XML para TXT, possibilitando uma verificação e visualização mais correta e objetiva dos dados de agrupamento feitos no Eureka!.

2. Software Presente!

A principal função do software Presente! é fazer análise automática de marcas linguísticas de presença social (PS), segundo as classes designadas no Modelo Presença Plus [Bastos 2012]. Presença social pode ser compreendida como o grau de sentimento, percepção e reação dos indivíduos ao interagirem por meio de recursos de comunicação mediada por computador [Tu e McIsaac 2002].

Para identificar essas pistas textuais nas postagens em fóruns e chats, [Kambara-Silva 2011] desenvolveu softwares auxiliares para gerar alguns arquivos de entrada: o conversor de HTML para XML, o construtor de categorias/gerador de classes.



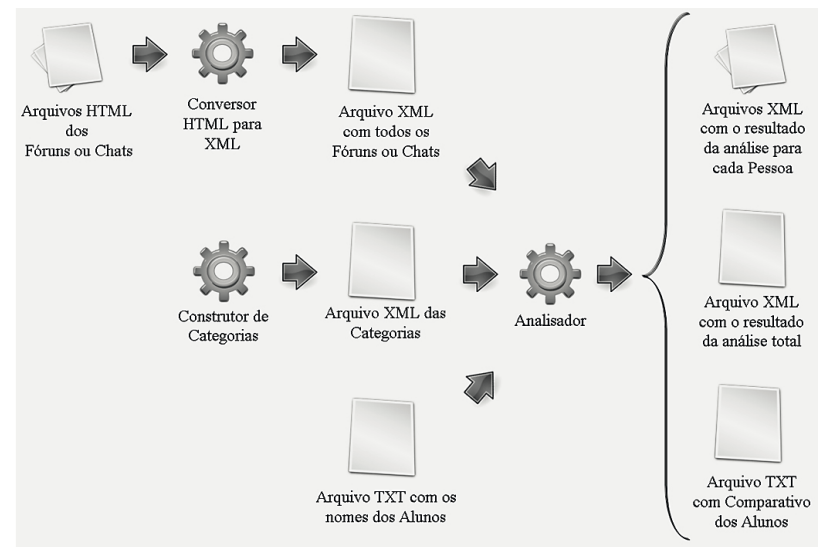


Figura 1: Esquema de funcionamento do Analisador de Presença Social [Kambara-Silva 2011]

A figura acima ilustra a ordem e que tipos de arquivos são analisados pelas ferramentas auxiliares para originar arquivos de resultado. Essas ferramentas são explicadas a seguir.

2.1. Conversor HTML para XML

Este software transforma as páginas web de fóruns e chats, encontradas em determinado ambiente virtual de aprendizagem, em documentos estruturados que possuem somente as tags necessárias para uma análise posterior. Essa operação é feita pela retirada somente de informações como o nome do autor da postagem, o assunto, a data, e o conteúdo da postagem.

O uso da página web completa seria inviável porque ela possui muitas informações que não possuem relação com o que pretende ser analisado, como tags de marcação de estética (como



div e span), tags de imagem, as que indicam ênfase (negrito, cores e tamanhos de fonte diferentes), entre outras. O conversor prepara os arquivos essenciais para a análise correta. Essa etapa é seguida da geração de arquivo XML com todos os posts e informações necessárias contidas neles.

2.2. Gerador de classes

Esse software é utilizado para definir as pistas ou palavras/expressões que podem estar no texto dos posts, classificando-as em relação a uma determinada classe e subclasse de PS. O Modelo PPlus possui quatro classes de indicadores de PS: afetividade, interatividade, coesão social e força. Cada classe apresenta subclasses, entre elas: emoção, humor, termos de sociabilidade e ênfase.

Pode-se retirar essas possibilidades de um dicionário da língua que será utilizada na análise, mas esse uso é opcional. Quanto mais completa e detalhada for a classificação das palavras e expressões nas categorias acima, mais precisa será a análise e, provavelmente, os resultados e as interpretações acerca das intenções e participação do aluno no ambiente virtual.

Nessa etapa também é gerado um arquivo XML com o mapeamento de todas as categorias e subcategorias, com suas respectivas pistas textuais.

Após obter os arquivos necessários nos softwares citados, alimenta-se o software denominado Analisador, gerando, assim, um arquivo XML relativo a cada estudante, com a quantidade de pistas encontradas para cada categoria e subcategoria com base nos posts realizados por esse aluno. A Figura 2 apresenta como esses arquivos são visualizados.



Secretaria de Educação
Profissional e Tecnológica



Ministério
da Educação



```

<?xml version="1.0" encoding="UTF-8" ?>

<classes>
  <quantidade>18</quantidade>
  <afetividade>
    <quantidade>5</quantidade>
    <emocao>
      <quantidade>2</quantidade>
      <interjeicoes_e_expressoes_interjetivas>
        <quantidade>1</quantidade>
        <expr_simples>
          <quantidade>1</quantidade>
        </expr_simples>
        <inicio_da_frase>
          <quantidade>0</quantidade>
        </inicio_da_frase>
      </interjeicoes_e_expressoes_interjetivas>
      <onomatopeias>
        <quantidade>0</quantidade>
      </onomatopeias>
      <emoticons_e_gifs>
        <quantidade>0</quantidade>
      </emoticons_e_gifs>
      <pontuacao_repetida>
        <quantidade>0</quantidade>
      </pontuacao_repetida>
    </emocao>
  </afetividade>
</classes>

```

Figura 2: Ilustração de um dos arquivos XML gerados por aluno

3. Metodologia

Uma das formas de se organizar e analisar o conteúdo dos arquivos XML é fazer sua classificação em grupos que possuem algo em comum. Isso é feito a partir de algoritmos que fazem correspondências métricas e análises qualitativas em relação aos documentos textuais de entrada, como ocorre no software de mineração e agrupamento Eurekha!

Como este software de agrupamento não possui compatibilidade com arquivos XML, a princípio, foi feita apenas a troca da extensão para texto, mantendo as tags com os nomes e os números das quantidades. Entretanto, o software não fornece relevância a números, que no caso seria a principal informação que influenciaria no agrupamento.

Os algoritmos de classificação utilizam principalmente palavras. Para fazer o software interpretar corretamente a informação que se quer usar como parâmetro de classificação – as quantidades numéricas das classes e subclasses, o processo de transformação foi baseado na repetição dos nomes das mesmas de acordo com suas quantidades.



Assim, o algoritmo continuará fazendo a classificação atribuindo peso às palavras. No entanto, como as palavras agora correspondem às quantidades, ao comparar os arquivos, a classificação ocorrerá de forma adequada.

```
<?xml version="1.0" encoding="UTF-8" ?>
<classes>
  <quantidade>18</quantidade>
  <afetividade>
    <quantidade>5</quantidade>
    <emocao>
      <quantidade>2</quantidade>
      <interjeicoes_e_expressoes_interjetivas>
        <quantidade>1</quantidade>
      <expr_simples>
        <quantidade>1</quantidade>
      </expr_simples>
    </emocao>
  </afetividade>
</classes>
```

Figura 3: Trecho do arquivo XML antes de passar pela ferramenta

```
afetividade afetividade afetividade afetividade afetividade
emocao emocao
interjeicoes_e_expressoes_interjetivas
expr_simples
```

Figura 4: O mesmo trecho acima após ser transformado pela ferramenta

As Figuras 3 e 4 apresentam um trecho do documento XML após passar pela ferramenta de preparação. Assim, eles podem ser interpretados e agrupados corretamente pelo **Corekha!**



4. Resultados

Após obter os arquivos no analisador de presença social, foram carregados os arquivos XML resultantes já preparados pela ferramenta. Para fazer o agrupamento, basta escolher a opção “identificar relacionamentos”. Para realizar essa operação foi escolhido o algoritmo Best-star (que foi o que melhor classificou os grupos). Desse modo, logo se visualiza uma matriz de similaridades, conforme a figura 5.

	AGRUPAMENTO_A	AGRUPAMENTO_B	AGRUPAMENTO_C	AGRUPAMENTO_D	AGRUPAMENTO_E	AGRUPAMENTO_F	AGRUPAMENTO_G	AGRUPAMENTO_H	AGRUPAMENTO_I	AGRUPAMENTO_J	AGRUPAMENTO_K	AGRUPAMENTO_L	AGRUPAMENTO_M
AGRUPAMENTO_A	1	0,434230864048004	0,352077752351761	0,518039405345917	0,367406249046326	0,0899342745542526	0,238415330648422	0,36414062976					
AGRUPAMENTO_B		1	0,405088752508163	0,439743727445602	0,356813907623291	0,0453553721308708	0,220185965299606	0,58523499965					
AGRUPAMENTO_C			1	0,360512614250183	0,260577768087387	0,0945191234350204	0,290008187294006	0,39148166775					
AGRUPAMENTO_D				1	0,453053295612335	0,0849220156669617	0,29801806807518	0,404452830553055					
AGRUPAMENTO_E					1	0,0184446685016155	0,260649830102921	0,250761359930038					
AGRUPAMENTO_F						1	0,299796372652054	0,0657588914036751					
AGRUPAMENTO_G							1	0,170526504516602					
AGRUPAMENTO_H								1					
AGRUPAMENTO_I									1				
AGRUPAMENTO_J										1			
AGRUPAMENTO_K											1		
AGRUPAMENTO_L												1	
AGRUPAMENTO_M													1

Figura 5: Representação da matriz de similaridades gerada no experimento

Essa matriz é gerada através da comparação entre os conteúdos de todos os arquivos carregados. Por isso, quando são comparados os mesmos arquivos, o grau de igualdade é máximo. Para os diferentes, podem-se ter resultados entre 0.03% e 52% aproximadamente, ou seja, podem-se ter as mesmas palavras, nas mesmas quantidades ou não. Quanto mais diferentes forem, mais próximos de zero estarão.

Esta operação pode ser utilizada para saber quais alunos possuem um grau de envolvimento e participação semelhante. No entanto, esta informação seria generalizada, pois até então não





foram identificadas as diferenças nos conteúdos e sim as informações quantitativas gerais dos arquivos.

Outro resultado que o software provê é o gráfico de colunas com o percentual relativo de documentos em cada cluster (grupo). Através deste gráfico, pode ser percebida a principal diferença

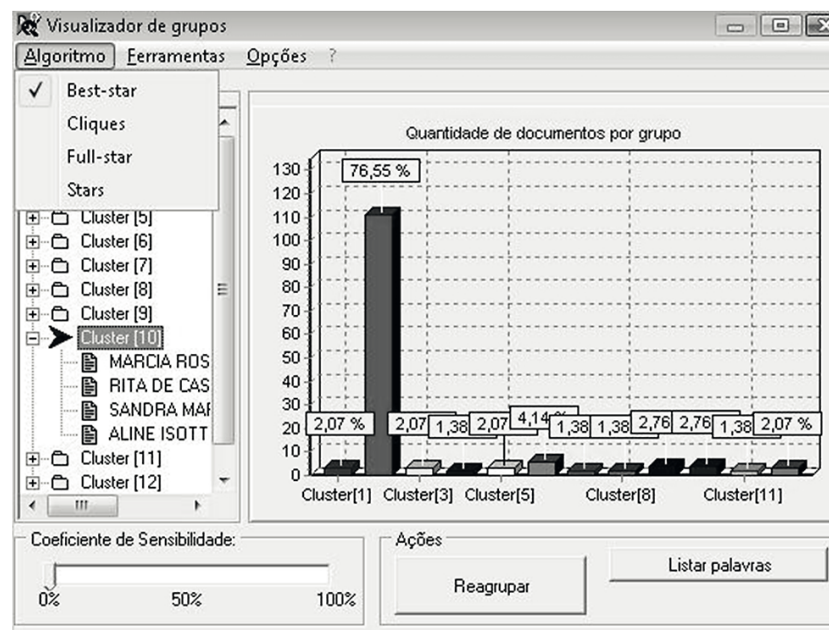


Figura 6: Experimento de clustering sem passar os arquivos na ferramenta

Na figura 6 se percebe que houve um agrupamento não muito compreensível e coerente com os números e as divergências nos conteúdos dos arquivos. No grupo em verde, que corresponde a aproximadamente 77% dos arquivos basicamente, só havia os arquivos que continham zero ocorrência de quantidade, ou seja, praticamente não houve participação desses alunos. Nos



outros casos, não houve um padrão coerente de classificação ao visualizar as características dos arquivos que estavam num mesmo grupo.

Na figura 7, verifica-se que houve diferenças menos contrastantes e expressivas em relação aos percentuais de arquivos nos grupos. De acordo com a nova configuração e distribuição da informação nos arquivos, os grupos apresentaram uma lógica de explicação na sua classificação.

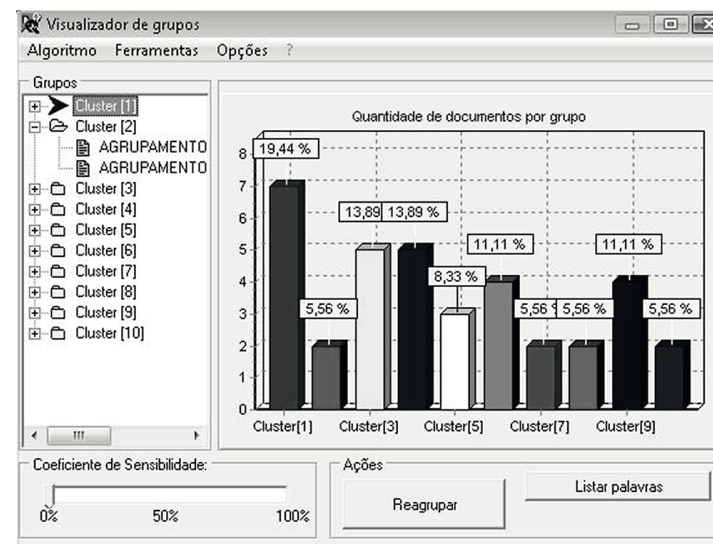


Figura 7: Experimento utilizando a ferramenta nos arquivos de entrada

Para saber o que há de comum nos arquivos que pertenceram a um mesmo grupo, assim como o que há de diferente entre eles e os dos demais grupos, basta executar a ação "Listar palavras". O software mostrará que naquele grupo as palavras mais relevantes possuem suas frequências em igual quantidade.



Dependendo das palavras com mais ocorrências e mais relevância no conjunto, pode-se descobrir a classificação do arquivo do aluno (em qual grupo, classe e subclasse ele mais se caracterizou). Para analisar as diferenças internas dos grupos foram utilizados 3 dos 10 grupos que apresentaram percentuais distintos (os grupos 1, 2 e 4).

4.1. Interpretação da classificação do grupo 1

A Figura 8 apresenta a lista das palavras com ocorrências e relevâncias iguais para os arquivos que a este grupo pertencem. Neste grupo, pode-se verificar que houve ocorrências consideráveis de classes e subclasses diferentes.

Palavra	Ocorrências	Relevância
✓ COESAO	007	0,179334282875061
✓ FORCA	007	0,0876649447849819
✓ NOMES_PROPRI...	006	0,068271164383207
✓ PRONOMES	007	0,0655464913163866
✓ PRONOMES_SIM...	006	0,0501835175922939
✓ INTENSIDADE	007	0,0485853297369821
✓ AFETIVIDADE	007	0,0471889291490827
✓ EMOCAO	007	0,0446810168879373
✓ EXPRESOES_A...	007	0,0435361010687692
✓ ADVERBIOS_E	007	0,0435361010687692
✓ QUANTIDADE	007	0,0302064993551799
✓ PRONOME_INDE...	007	0,0302064993551799
✓ INTERATIVIDADE	006	0,0282120023454939
✓ EXPRESOES_F...	005	0,0231570558888572
✓ SAUDACDES	005	0,0231570558888572
✓ NOMES_GENERI...	006	0,0223595712866102
✓ REALCE	006	0,0199056736060551

Total de palavras.....: 42

Figura 8: Perfil de classificação do grupo/cluster 1



A ordem das palavras está de acordo com a ordem decrescente do seu nível de relevância em relação a todas as palavras encontradas. Por isso, como a subclasse de nomes próprios foi considerada mais relevante nos textos, mesmo tendo uma frequência maior que intensidade, por exemplo, ele ficou acima deste. Destaca-se que neste grupo foram encontradas 42 palavras.

4.2. Interpretação da classificação do grupo 2

A figura 9 mostra a lista de palavras do grupo 2, comuns a todos os arquivos que a ele pertencem. Pode-se perceber que apesar da mistura de classes e subclasses, as ocorrências foram parecidas. No entanto, elas foram muito menores que as do grupo anterior.

Palavra	Ocorrências	Relevância
✓ COESAO	002	0,130819365382195
✓ AFETIVIDADE	002	0,104083530604839
✓ EMOCAO	002	0,101290233433247
✓ FORCA	002	0,0612197406589985
✓ PRONOMES	002	0,0580606535077095
✓ REALCE	001	0,0446927361190319
✓ PONTUACAO_RE...	001	0,0446927361190319
✓ EXPRESSOES_F...	002	0,0413008779287338
✓ PRONOMES_SIM...	002	0,0346169210970402
✓ INTERATIVIDADE	002	0,0321894139051437
✓ PRONOME_INDE...	002	0,029396116733551
✓ QUANTIDADE	002	0,029396116733551
✓ EXPRESSOES_A...	002	0,0290303267538548
✓ INTENSIDADE	002	0,0290303267538548
✓ ADVERBIOS_E	002	0,0290303267538548
✓ SAUDACDES	001	0,0238095242530107
✓ SUJEITO OCULTO	002	0,0234437361359596

Total de palavras.....: 36

Figura 9: Perfil de classificação para o grupo/cluster 2.



Este grupo tem em comum o fato de que as classes de PS mais encontradas para os alunos foi o das categorias coesão, afetividade, emoção etc., com no máximo 2 ocorrências. Obtiveram-se também ocorrências de outros tipos de PS. Elas foram parecidas entre si, porém com valores pequenos. Nesse caso, os arquivos continham 30 palavras em comum nesse grupo.

4.3. Interpretação da classificação do grupo 4

A figura 10 mostra a lista de palavras do grupo 4, também com o mesmo comportamento para todos os arquivos do grupo. Pode-se perceber que houve ocorrências de poucos tipos de classes e subclasses. Algumas destas apresentaram um número considerável, enquanto outras obtiveram uma quantidade bem baixa.

Palavra	Ocorrências	Relevância
✓ COESAO	005	0,330602765083313
✓ EXPRESSOES_FATIC...	005	0,209124231338501
✓ SAUDACOES	005	0,190942394733425
✓ NOMES_PROPRIOS	003	0,106093907356262
✓ MARCADORES_CONV...	001	0,028571429848671
✓ MANUTENCAO_DO	001	0,028571429848671
✓ INTERATIVIDADE	001	0,028571429848671
✓ DIALOGO	001	0,028571429848671
✓ DESPEDIDAS	001	0,018181818723678E
✓ PRONOMES_SIMPLES	001	0,015384615957737
✓ PRONOMES	001	0,015384615957737

Total de palavras.....: 11

Figura 10: Perfil de classificação do grupo/cluster 4.





Nesse grupo houve poucas classes/subclasses diferentes encontradas. As que foram encontradas obtiveram uma frequência considerável, principalmente nas de coesão, expressões fáticas e saudações. As outras classes/subclasses encontradas apresentaram um número bem menor. Nesse caso, existiam somente 11 palavras em comum.

Outro fator importante é que, antes da ferramenta, os documentos que continham zero ocorrência também foram processados. No entanto, após passar pelo programa eles desapareceram, pois não podiam ser repetidos.

A ferramenta pode ser usada como apoio na análise, interpretação e apresentação de resultados mais claros para os utilizadores do Analisador de PS, auxiliando na tomada de decisões a partir da visualização do comportamento dos grupos.

5. Conclusão

Após os experimentos percebe-se que a ferramenta auxiliou de forma significativa na compreensão dos perfis de comportamento dos alunos, ao se comunicarem através de ambientes virtuais de aprendizagem. Isso demonstra que o software é uma possível solução para melhorar a apresentação dos resultados e a tomada de decisões por parte do usuário, quando lidar com arquivos no formato XML no software Eureka!

Referências

Bastos, H. P. P. (2012). Presença Plus: modelo de identificação de presença social em ambientes virtuais de ensino e aprendizagem. Tese de doutorado. Programa de Pós-Graduação em Informática na Educação. Universidade Federal do Rio Grande do Sul. Porto Alegre, RS.

Gupta, V.; Lehal, G. S. (2009) "A Survey of Text Mining Techniques and Applications". Journal of Emerging Technologies in Web Intelligence, v. 1, n. 1.



Kambara-Silva, J.K. (2011) "Automatização do processo de aquisição de presença social em fóruns e chats". Trabalho de Conclusão de Curso. Instituto de Informática. Universidade Federal do Rio Grande do Sul. Porto Alegre, RS

Kobayashi, M.; Aono, M. (2004) "Vector Space Models for Search and Cluster Mining". In: BERRY, M. W. (ed.). Survey of text mining: clustering, classification, and retrieval. New York: Springer-Verlag. p. 103-122.

Tu C.H; Mclsaac (2001) "The relationship of social presence and interaction in online classes". In The American Journal of Distance Education, v.16, n.3, p. 131-50.

Wives, L. K. (2004). Utilizando conceitos como descritores de textos para o processo de identificação de conglomerados (clustering) de documentos. Exame de Qualificação (doutorado em Ciência da Computação) - Instituto de Informática, Universidade Federal do Rio Grande do Sul, Porto Alegre, RS.



Secretaria de Educação
Profissional e Tecnológica



Ministério
da Educação

