

Mineração de Dados Biológicos

PROJETO – Segundo Semestre – 2010

Prof. Ricardo J. G. B. Campello

O presente trabalho consiste da aplicação e avaliação de algoritmos de seleção de atributos e classificação no contexto de dados de expressão gênica. As bases de dados que serão utilizadas são provenientes de diversos experimentos associados à detecção de câncer. Nestas bases de dados, os objetos são amostras de tecidos pertencentes a diferentes classes enquanto os atributos correspondem a genes para os quais os níveis de expressão foram mensurados. Para tal avaliação, as implementações dos algoritmos vistos em aula e disponíveis no software de mineração de dados WEKA devem ser utilizadas. Dois diferentes cenários de aplicação e avaliação devem ser considerados, como descrito a seguir¹.

1. Classificadores

Neste cenário os algoritmos de classificação de dados devem ser aplicados e avaliados **sem** a utilização prévia de métodos de seleção de atributos. Os seguintes algoritmos devem ser considerados para avaliação:

- kNN
- Naïve Bayes
- 1Rule
- J48 (implementação do C4.5 disponível no WEKA)

Para o algoritmo kNN, varie o número de vizinhos de 1 à 10. Plote um gráfico na forma *número de vizinhos x acurácia de classificação* e avalie a influencia que o número de vizinhos possui na acurácia do método. Escolha o melhor dentre os 10 classificadores kNN avaliados para posterior comparação com os demais métodos. Para tanto, plote gráfico(s) que permitam a comparação entre os métodos: melhor kNN, Naïve Bayes, 1Rule e J48. Além da acurácia, considere também ao menos a medida F-measure. Por fim comente e apresente sua interpretação para os resultados obtidos.

2. Seleção de Atributos + Classificadores

Neste cenário, algoritmos de seleção de atributos devem ser aplicados previamente aos classificadores². Os algoritmos de seleção de atributos que devem ser avaliados são:

- ReliefF
- CFS
- LVF

Após a aplicação dos filtros de seleção de atributos, cada um dos classificadores anteriormente mencionados devem ser avaliados. Quanto ao algoritmo ReliefF, este gera uma ordenação de atributos, de forma que a escolha de um número de atributos é necessária após sua aplicação. Considerando ainda o algoritmo ReliefF, varie o número de atributos selecionados por este algoritmo em um intervalo determinado (cabe a você escolher o intervalo de valores) e plote gráficos para os classificadores na forma *número de atributos x acurácia*, de forma a mostrar o impacto do número de atributos selecionados nos métodos de classificação. Observe as correlações entre os atributos selecionados e verifique se há presença de atributos redundantes nos ranks de atributos gerados (utilize a parte de visualização do WEKA – aba *visualize*). Compare os resultados de acurácia obtidos entre os pares algoritmo de seleção – algoritmo de

¹ Para ambos os cenários de avaliação utilize o procedimento de validação cruzada de 10 pastas estratificado.

² Aplique os métodos de seleção somente nas pastas de treinamento e não utilizando todo o conjunto de dados.

classificação e comente os resultados obtidos. Por fim, compare os resultados obtidos após a seleção de atributos com aqueles obtidos somente com a aplicação dos classificadores (sem seleção de atributos). Comente sobre os impactos e efeitos da seleção de atributos em cada um dos classificadores utilizados. Além disso, verifique se todos os classificadores obtiveram as mesmas melhoras com a aplicação dos métodos de seleção de atributos, justificando os fatos observados.

Além dos gráficos solicitados, podem ser gerados gráficos e tabelas adicionais a fim de facilitar a comparação e obtenção de conclusões. Sua análise e conclusões devem ser apresentados em forma de um relatório com no máximo 5 páginas. Relatórios contendo somente gráficos e/ou tabelas sem análises profundas e fundamentadas receberão nota **zero**.

- **Bases de Dados**

A cada aluno foi atribuída uma base de dados única. A distribuição das bases de dados é apresentada na tabela abaixo. Note que uma porção de genes irrelevantes já foram removidos das bases. Caso ache necessário aplique técnicas de pré-processamento aos dados (normalização, discretização, ...) e comente sobre seu impacto nos resultados obtidos.

Nome	Base de Dados	Nome	Base de Dados
Alessandra Figueiredo	<i>alizadeh-2000-v1</i>	Gilson Campani Junior	<i>lapointe-2004-v1</i>
Amanda Souza Camara	<i>alizadeh-2000-v2</i>	Gustavo M. Alvares de Lima	<i>lapointe-2004-v2</i>
Americo Tavares Ranzani	<i>alizadeh-2000-v3</i>	Heloisa dos Santos Muniz	<i>liang-2005</i>
Andre Zamith Selvaggio	<i>armstrong-2002-v1</i>	Jaqueline Yu Ting Wang	<i>nutt-2003-v1</i>
Andrea Rodrigues Cardoso	<i>armstrong-2002-v2</i>	Jessica Baleiro Okado	<i>nutt-2003-v2</i>
Anielle Coelho Ranulfi	<i>bhattacharjee-2001</i>	Leticia Marchese	<i>nutt-2003-v3</i>
Bruno L. S. Paula de Mello	<i>bittner-2000</i>	Luis Felipe Santos Mendes	<i>pomeroy-2002-v1</i>
Bruno Yasui Matsuyama	<i>bredel-2005</i>	Luma Godoy Magalhaes	<i>pomeroy-2002-v2</i>
Caio M. Ramos de Oliveira	<i>dyrskjot-2003</i>	Mariana R. Bunoro Batista	<i>ramaswamy-2001</i>
Cesar Maschio Fioravanti	<i>garber-2001</i>	Mira Melke de Oliveira	<i>risinger-2003</i>
Eduardo Cocca Padovani	<i>golub-1999-v1</i>	Pedro Michelao Neuber	<i>shipp-2002-v1</i>
Evandro Jose Mulinari	<i>golub-1999-v2</i>	Ricardo Simionato Boffa	<i>singh-2002</i>
Fernando Takeshi Tanouye	<i>gordon-2002</i>	Thais Panhan Merlo	<i>su-2001</i>
Flavio A. C. da Silva Andrade	<i>khan-2001</i>	Wagner Rafael Correr	<i>tomlins-2006-v1</i>
Francesco B. Teixeira	<i>laiho-2007</i>	Waldomiro Thiago Corsi	<i>tomlins-2006-v2</i>

As bases de dados em formato ARFF encontram-se disponíveis no site da disciplina no arquivo bases.tar. Informações e referências adicionais para cada uma das bases de dados podem ser encontradas no site: <http://algorithmics.molgen.mpg.de/Static/Supplements/CompCancer/datasets.htm>.

Informações Importantes

Os trabalhos podem ser realizados individualmente ou em dupla. Caso seja realizado em dupla, somente uma das bases de dados atribuídas aos alunos da dupla deve ser utilizada durante a realização do trabalho. O relatório deverá ser entregue **pessoalmente** de forma impressa ao estagiário PAE da disciplina. O estagiário PAE pode ser consultado nos horários de atendimento para sanar dúvidas quanto à realização do trabalho. A **data limite de entrega** do relatório (**pessoalmente**) ao estagiário PAE será 29/11/2010. No momento da entrega o aluno (ou a dupla) deverá escolher a data e horário da arguição em grade de disponibilidade pré-estabelecida. A preferência das datas e horários na grade será por ordem de entrega dos trabalhos.