

Mini-teste 2

Mineração de Dados

Mestrado em Engenharia Informática

Vitor Peixoto - A79175

10/11/2018

Universidade do Minho

Introdução

Neste mini-teste foi atribuído um conjunto de quatro questões baseados na análise de conjuntos de dados. Essa análise é feita pela comparação da *performance* dos diversos classificadores utilizados e quando sujeitos a alterações dos parâmetros e algoritmos que os compõem.

Exercício 1

Este primeiro exercício pede para observar o *dataset labor*, com o algoritmo de classificação *J48* recorrendo a validação cruzada de 10 *folds*.

Este método consiste em 10 divisões das instâncias do *dataset*, sendo que a cada iteração, uma parte será usada como dados de treino e as restantes como dados de teste. Este procedimento é repetido para cada parte dividida, sendo que cada porção do *dataset* são dados de treino uma vez e dados de teste 9 (k-1) vezes.

Os resultados obtidos com estes parâmetros revelaram uma percentagem de cerca de 73.68% de instâncias corretamente classificadas e de 0.695 de AUC (*Area under curve*).

Correctly Classified Instances	42	73.6842 %
Incorrectly Classified Instances	15	26.3158 %
Kappa statistic	0.4415	
Mean absolute error	0.3192	
Root mean squared error	0.4669	
Relative absolute error	69.7715 %	
Root relative squared error	97.7888 %	
Total Number of Instances	57	

=== Detailed Accuracy By Class ===									
	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.700	0.243	0.609	0.700	0.651	0.444	0.695	0.559	bad
	0.757	0.300	0.824	0.757	0.789	0.444	0.695	0.738	good
Weighted Avg.	0.737	0.280	0.748	0.737	0.740	0.444	0.695	0.675	

Figura 1: Resultados obtidos através de *J48* com validação cruzada.

Este exercício pedia também para correr o mesmo algoritmo melhorado com o método de *Bagging*. *Bagging* é uma das técnicas de composição de modelos onde um conjunto de classificadores são combinados para obter um classificador mais forte e com melhores resultados.

Os resultados obtidos com a composição de *Bagging* sobre *J48* foram muito superiores aos obtidos apenas com *J48*. Obtivemos uma percentagem perto dos 86% de instâncias corretamente classificadas e um AUC de 0.884.

Correctly Classified Instances	49	85.9649 %
Incorrectly Classified Instances	8	14.0351 %
Kappa statistic	0.6771	
Mean absolute error	0.2588	
Root mean squared error	0.3533	
Relative absolute error	56.5714 %	
Root relative squared error	73.985 %	
Total Number of Instances	57	

=== Detailed Accuracy By Class ===									
	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.700	0.054	0.875	0.700	0.778	0.686	0.884	0.746	bad
	0.946	0.300	0.854	0.946	0.897	0.686	0.884	0.940	good
Weighted Avg.	0.860	0.214	0.861	0.860	0.855	0.686	0.884	0.872	

Figura 2: Resultados obtidos através de *Bagging* sobre *J48* com validação cruzada.

No exemplo do *Netflix Prize* apresentado na aula, a união entre as equipas e os seus classificadores permitiu um *score* superior. Esse caso verifica-se também na composição de modelos para este *dataset* onde a sobreposição entre *Bagging* e *J48* foi benéfico na classificação das instâncias e no valor da AUC. Esta melhoria deve-se ao facto de o *Bagging* provocar uma melhoria em algoritmos instáveis, como é o caso do *J48*, uma vez que enriquece o conjunto de dados.

Exercício 2

Este exercício pede para analisar e especificar os benefícios que a aplicação de *Bagging* traria num *dataset* onde foi aplicado o classificador de *Naive Bayes* com uma taxa de erro de 0.035.

A técnica de *Bagging* reduz a variância do erro do classificador. Isto funciona para algoritmos instáveis, como é o caso de algoritmos que envolvem árvores de decisão (i.e. *J48*). No entanto, para algoritmos mais estáveis, como é o caso de *Naive Bayes*, esta técnica pode mesmo aumentar a variância, uma vez que este algoritmo apresenta uma variância por si só já baixa.

Tendo então em conta a já insignificativa taxa de erro apresentada pelo algoritmo *Naive Bayes unbagged*, pode-se assumir que a aplicação de *Bagging* neste caso poderá não ser benéfica e aumentar a taxa de erro nas instâncias corretamente classificadas.

Exercício 3

Neste exercício era pedido um estudo sobre o custo de erros no *dataset vote* usando os classificadores *J48* e *Naive Bayes*.

Antes de iniciarmos o estudo, temos de entender que há duas maneiras de tornar um classificador sensível ao custo: por *Classification* ou por *Learning*. O método ***Classification*** usa um classificador *standard*, ajustando o *output* à matriz de custo. O método ***Learning*** aprende um novo classificador ótimo, reajustando os exemplos mal classificados por duplicação de exemplos ou por reajustamento dos pesos.

Para conseguir estudar o custo dos erros, foram definidas 3 matrizes de custo:

$$M1 = \begin{bmatrix} 0 & 2 \\ 1 & 0 \end{bmatrix}$$

$$M2 = \begin{bmatrix} 0 & 1 \\ 2 & 0 \end{bmatrix}$$

$$M3 = \begin{bmatrix} 0 & 2 \\ 2 & 0 \end{bmatrix}$$

Vamos agora passar ao estudo dos custos neste *dataset*:

Usamos o classificador ***J48*** aplicado num modelo sensível ao custo por ***Classification***. Obtivemos os seguintes resultados para as respetivas matrizes de custo:

M1	<p>Correctly Classified Instances 419 96.3218 %</p> <p>Incorrectly Classified Instances 16 3.6782 %</p> <p>Kappa statistic 0.9224</p> <p>Total Cost 24</p> <p>Average Cost 0.0552</p> <p>Mean absolute error 0.0368</p> <p>Root mean squared error 0.1918</p> <p>Relative absolute error 7.7558 %</p> <p>Root relative squared error 39.3895 %</p> <p>Total Number of Instances 435</p>	<pre> === Confusion Matrix === a b <-- classified as 259 8 a = democrat 8 160 b = republican </pre>
M2	<p>Correctly Classified Instances 420 96.5517 %</p> <p>Incorrectly Classified Instances 15 3.4483 %</p> <p>Kappa statistic 0.9277</p> <p>Total Cost 20</p> <p>Average Cost 0.046</p> <p>Mean absolute error 0.0345</p> <p>Root mean squared error 0.1857</p> <p>Relative absolute error 7.271 %</p> <p>Root relative squared error 38.1387 %</p> <p>Total Number of Instances 435</p>	<pre> === Confusion Matrix === a b <-- classified as 257 10 a = democrat 5 163 b = republican </pre>
M3	<p>Correctly Classified Instances 419 96.3218 %</p> <p>Incorrectly Classified Instances 16 3.6782 %</p> <p>Kappa statistic 0.9224</p> <p>Total Cost 32</p> <p>Average Cost 0.0736</p> <p>Mean absolute error 0.0368</p> <p>Root mean squared error 0.1918</p> <p>Relative absolute error 7.7558 %</p> <p>Root relative squared error 39.3895 %</p> <p>Total Number of Instances 435</p>	<pre> === Confusion Matrix === a b <-- classified as 259 8 a = democrat 8 160 b = republican </pre>

Figura 3: Resultados obtidos através de *J48* com *Cost-Sensitive Classification* aplicado nas 3 matrizes.

Apesar de o algoritmo aplicado sobre as três matrizes ser exatamente o mesmo, o número de instâncias corretamente classificadas em M2 é diferente, pois a própria matriz de custo influencia a aprendizagem do classificador, podendo afetar a sua performance. Assim sendo, neste caso, para as matrizes M1, M2 e M3, obtivemos um custo de 24, 20 e 32, respectivamente.

Usamos agora o mesmo classificador *J48* mas com um modelo *Cost-Sensitive* por *Learning*. Os resultados obtidos foram os seguintes:

M1	<p>Correctly Classified Instances 419 96.3218 %</p> <p>Incorrectly Classified Instances 16 3.6782 %</p> <p>Kappa statistic 0.9224</p> <p>Total Cost 24</p> <p>Average Cost 0.0552</p> <p>Mean absolute error 0.0683</p> <p>Root mean squared error 0.1838</p> <p>Relative absolute error 14.3944 %</p> <p>Root relative squared error 37.7491 %</p> <p>Total Number of Instances 435</p>	<pre> === Confusion Matrix === a b <-- classified as 259 8 a = democrat 8 160 b = republican </pre>
M2	<p>Correctly Classified Instances 419 96.3218 %</p> <p>Incorrectly Classified Instances 16 3.6782 %</p> <p>Kappa statistic 0.9229</p> <p>Total Cost 21</p> <p>Average Cost 0.0483</p> <p>Mean absolute error 0.063</p> <p>Root mean squared error 0.1775</p> <p>Relative absolute error 13.2887 %</p> <p>Root relative squared error 36.4589 %</p> <p>Total Number of Instances 435</p>	<pre> === Confusion Matrix === a b <-- classified as 256 11 a = democrat 5 163 b = republican </pre>
M3	<p>Correctly Classified Instances 419 96.3218 %</p> <p>Incorrectly Classified Instances 16 3.6782 %</p> <p>Kappa statistic 0.9224</p> <p>Total Cost 32</p> <p>Average Cost 0.0736</p> <p>Mean absolute error 0.0611</p> <p>Root mean squared error 0.1748</p> <p>Relative absolute error 12.887 %</p> <p>Root relative squared error 35.9085 %</p> <p>Total Number of Instances 435</p>	<pre> === Confusion Matrix === a b <-- classified as 259 8 a = democrat 8 160 b = republican </pre>

Figura 4: Resultados obtidos através de *J48* com *Cost-Sensitive Learning* aplicado nas 3 matrizes.

Neste caso, obtivemos um custo de 24, 21 e 32 para cada uma das matrizes de custo. Estes resultados obtidos com *Learning* foram muito semelhantes aos obtidos com *Classification*.

Passamos agora para *Naive Bayes* aplicado num modelo sensível ao custo por *Classification*. Os resultados obtidos foram os seguintes:

M1	<p>Correctly Classified Instances 390 89.6552 %</p> <p>Incorrectly Classified Instances 45 10.3448 %</p> <p>Kappa statistic 0.7849</p> <p>Total Cost 74</p> <p>Average Cost 0.1701</p> <p>Mean absolute error 0.1034</p> <p>Root mean squared error 0.3216</p> <p>Relative absolute error 21.8131 %</p> <p>Root relative squared error 66.0582 %</p> <p>Total Number of Instances 435</p>	<p>=== Confusion Matrix ===</p> <table> <tr> <td>a</td><td>b</td><td><-- classified as</td></tr> <tr> <td>238</td><td>29</td><td>a = democrat</td></tr> <tr> <td>16</td><td>152</td><td>b = republican</td></tr> </table>	a	b	<-- classified as	238	29	a = democrat	16	152	b = republican
a	b	<-- classified as									
238	29	a = democrat									
16	152	b = republican									
M2	<p>Correctly Classified Instances 396 91.0345 %</p> <p>Incorrectly Classified Instances 39 8.9655 %</p> <p>Kappa statistic 0.8152</p> <p>Total Cost 48</p> <p>Average Cost 0.1103</p> <p>Mean absolute error 0.0897</p> <p>Root mean squared error 0.2994</p> <p>Relative absolute error 18.9047 %</p> <p>Root relative squared error 61.4968 %</p> <p>Total Number of Instances 435</p>	<p>=== Confusion Matrix ===</p> <table> <tr> <td>a</td><td>b</td><td><-- classified as</td></tr> <tr> <td>237</td><td>30</td><td>a = democrat</td></tr> <tr> <td>9</td><td>159</td><td>b = republican</td></tr> </table>	a	b	<-- classified as	237	30	a = democrat	9	159	b = republican
a	b	<-- classified as									
237	30	a = democrat									
9	159	b = republican									
M3	<p>Correctly Classified Instances 392 90.1149 %</p> <p>Incorrectly Classified Instances 43 9.8851 %</p> <p>Kappa statistic 0.7949</p> <p>Total Cost 86</p> <p>Average Cost 0.1977</p> <p>Mean absolute error 0.0989</p> <p>Root mean squared error 0.3144</p> <p>Relative absolute error 20.8437 %</p> <p>Root relative squared error 64.5736 %</p> <p>Total Number of Instances 435</p>	<p>=== Confusion Matrix ===</p> <table> <tr> <td>a</td><td>b</td><td><-- classified as</td></tr> <tr> <td>238</td><td>29</td><td>a = democrat</td></tr> <tr> <td>14</td><td>154</td><td>b = republican</td></tr> </table>	a	b	<-- classified as	238	29	a = democrat	14	154	b = republican
a	b	<-- classified as									
238	29	a = democrat									
14	154	b = republican									

Figura 5: Resultados obtidos através de *Naive Bayes* com *Cost-Sensitive Classification* aplicado nas 3 matrizes.

Neste caso, obtivemos um custo de 74, 48 e 86 para cada uma das matrizes de custo M1, M2 e M3, respetivamente.

Usamos novamente o classificador *Naive Bayes* aplicado agora em num modelo sensível ao custo por *Learning*. Os resultados obtidos foram os seguintes:

M1	<p>Correctly Classified Instances 391 89.8851 %</p> <p>Incorrectly Classified Instances 44 10.1149 %</p> <p>Kappa statistic 0.7899</p> <p>Total Cost 73</p> <p>Average Cost 0.1678</p> <p>Mean absolute error 0.1014</p> <p>Root mean squared error 0.3006</p> <p>Relative absolute error 21.3742 %</p> <p>Root relative squared error 61.728 %</p> <p>Total Number of Instances 435</p>	<p>=== Confusion Matrix ===</p> <table> <tr> <td>a</td><td>b</td><td><-- classified as</td></tr> <tr> <td>238</td><td>29</td><td>a = democrat</td></tr> <tr> <td>15</td><td>153</td><td>b = republican</td></tr> </table>	a	b	<-- classified as	238	29	a = democrat	15	153	b = republican
a	b	<-- classified as									
238	29	a = democrat									
15	153	b = republican									
M2	<p>Correctly Classified Instances 395 90.8046 %</p> <p>Incorrectly Classified Instances 40 9.1954 %</p> <p>Kappa statistic 0.8102</p> <p>Total Cost 50</p> <p>Average Cost 0.1149</p> <p>Mean absolute error 0.098</p> <p>Root mean squared error 0.2954</p> <p>Relative absolute error 20.6591 %</p> <p>Root relative squared error 60.6777 %</p> <p>Total Number of Instances 435</p>	<p>=== Confusion Matrix ===</p> <table> <tr> <td>a</td><td>b</td><td><-- classified as</td></tr> <tr> <td>237</td><td>30</td><td>a = democrat</td></tr> <tr> <td>10</td><td>158</td><td>b = republican</td></tr> </table>	a	b	<-- classified as	237	30	a = democrat	10	158	b = republican
a	b	<-- classified as									
237	30	a = democrat									
10	158	b = republican									
M3	<p>Correctly Classified Instances 392 90.1149 %</p> <p>Incorrectly Classified Instances 43 9.8851 %</p> <p>Kappa statistic 0.7949</p> <p>Total Cost 86</p> <p>Average Cost 0.1977</p> <p>Mean absolute error 0.0995</p> <p>Root mean squared error 0.2977</p> <p>Relative absolute error 20.9815 %</p> <p>Root relative squared error 61.1406 %</p> <p>Total Number of Instances 435</p>	<p>=== Confusion Matrix ===</p> <table> <tr> <td>a</td><td>b</td><td><-- classified as</td></tr> <tr> <td>238</td><td>29</td><td>a = democrat</td></tr> <tr> <td>14</td><td>154</td><td>b = republican</td></tr> </table>	a	b	<-- classified as	238	29	a = democrat	14	154	b = republican
a	b	<-- classified as									
238	29	a = democrat									
14	154	b = republican									

Figura 6: Resultados obtidos através de *Naive Bayes* com *Cost-Sensitive Learning* aplicado nas 3 matrizes.

Neste caso, obtivemos um custo de 73, 50 e 86 para cada uma das matrizes de custo M1, M2 e M3, respetivamente. Os resultados obtidos neste caso, foram novamente semelhantes entre *Learning* e *Classification*.

No cômputo geral, os custos obtidos com *J48* foram melhores aos obtidos com *Naive Bayes*. Isto deve-se às matrizes de confusão obtidas com NB terem mais instâncias como FP e FN. Multiplicando esses valores pelo respetivo valor na matriz de custo, obtém-se o custo total dos

erros, tendo neste caso o custo sido maior com o algoritmo NB.

Exercício 4

Este exercício pedia para considerar o *dataset* **unbalanced** e comparar o desempenho dos algoritmos de *clustering k-means* e *EM*.

Inicialmente corremos o *dataset* com o algoritmo *EM*, ignorando a classe 'outcome' visto ser a classe do resultado. Quando o número de *clusters* como -1 para este algoritmo, ele irá determinar o número de *clusters*.

Correndo então este algoritmo, ele indica 18 como o número de *clusters* através de validação cruzada. Estes 18 *clusters* devem-se ao facto de o algoritmo dividir as instâncias 'active' da classe 'outcome' em vários *clusters* em vez de as por num só, ao passo que as instâncias de 'inactive' estão bem distribuídas. A distribuição das instâncias do 'outcome' pode ser vista no seguinte gráfico:

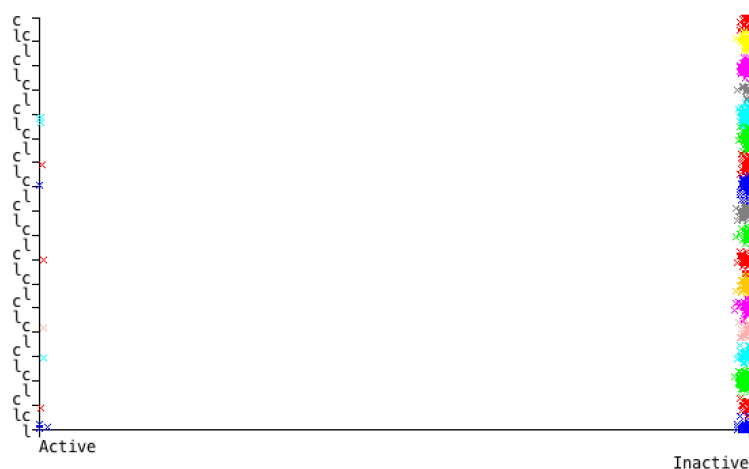


Figura 7: Distribuição de 'outcome' pelos clusters.

A execução do algoritmo *EM* é extremamente lenta (cerca de 48 segundos) dada a necessidade de cálculo do número de *clusters*, verificando-se uma descida significativa (para 0.4 segundos) quando esse número é fornecido.

O algoritmo *EM* foi executado com 2 iterações e apresentou a seguinte colocação das instâncias do *dataset* nos respetivos *clusters*:

Clustered Instances		
0	45	(5%)
1	25	(3%)
2	133	(16%)
3	42	(5%)
4	26	(3%)
5	56	(7%)
6	46	(5%)
7	38	(4%)
8	22	(3%)
9	54	(6%)
10	69	(8%)
11	37	(4%)
12	55	(6%)
13	58	(7%)
14	13	(2%)
15	61	(7%)
16	48	(6%)
17	28	(3%)

Figura 8: Distribuição das instâncias pelo algoritmo *EM*.

Tendo agora um número de *clusters* definido, podemos executar também o algoritmo *k-means*. *K-means* é um algoritmo simples, iterativo e baseado nas distâncias entre instâncias (distância euclidiana, neste caso). Difere do algoritmo *EM* na medida em que o *EM* se baseia na probabilidade (expetativa) de um ponto pertencer a um *cluster* específico.

Este algoritmo correu em apenas 0.05 segundos, com 20 iterações, com uma soma dos erros quadrados de cerca de 321.32 e apresentou a seguinte colocação das instâncias do *dataset* nos respetivos *clusters*:

Clustered Instances		
0	22	(3%)
1	19	(2%)
2	75	(9%)
3	27	(3%)
4	52	(6%)
5	122	(14%)
6	39	(5%)
7	102	(12%)
8	52	(6%)
9	21	(2%)
10	10	(1%)
11	72	(8%)
12	45	(5%)
13	37	(4%)
14	38	(4%)
15	30	(4%)
16	57	(7%)
17	36	(4%)

Figura 9: Distribuição das instâncias pelo algoritmo *k-means*.

Conclusão

Com este mini-teste, foi possível tirar algumas conclusões ao longo do trabalho desenvolvido:

- A técnica de *Bagging* é muito eficiente sobre o algoritmo *J48*, conseguindo melhorar significativamente a *performance* deste algoritmo.
- Nem sempre esta técnica é benéfica na classificação de instâncias de um *dataset*. Em algoritmos estáveis, como o *Naive Bayes*, a técnica de *Bagging* pode até diminuir a *performance* do algoritmo.
- O estudo sobre o custo dos erros é essencial para diminuir os casos de FN e FP numa matriz de confusão.
- Há dois métodos de modelos *Cost-sensitive*, podendo ser aplicados sob um algoritmo (*J48*, *Naive Bayes*, etc.) e a sua performance pode variar de acordo com o *dataset*, sendo importante testar todas as possibilidades e escolher a que minimiza o custo dos erros.
- As técnicas de *clustering* são importantes no agrupamento de instâncias similares entre si.

Resumindo, este trabalho assumiu um papel importante e positivo no acompanhamento e aprofundamento das matérias lecionadas nas aulas da Unidade Curricular de Mineração de Dados.