

Mineração de Dados em Biologia Molecular

Planejamento e Análise de Experimentos Preditivos

Docente: André C. P. L. F. de Carvalho
PAE: Victor Hugo Barella



Principais tópicos

- Desempenho preditivo
- Partição dos dados
- Reamostragem
- Tipos de erro
- Avaliação do desempenho
- Curvas ROC

© André de Carvalho - ICMC/USP

2

Tarefa preditiva

- Indução de modelo com bom desempenho preditivo para novos dados
 - Escolher algoritmo com viés (bias) adequado para o conjunto de dados utilizado
 - Necessário avaliar o desempenho preditivo de algoritmos
 - Avaliar desempenho preditivo de modelos induzidos pelos algoritmos
 - Modelo induzido pelo algoritmo deve estimar corretamente o rótulo de novos exemplos

© André de Carvalho - ICMC/USP

3

Algoritmos e modelos

- Desempenho preditivo deve ser avaliado para
 - Algoritmos de AM
 - Saída de algoritmos de AM:
 - Modelos, funções, hipóteses preditivos
 - Modelos preditivos
 - Saída de modelos preditivos
 - Classe ou valor para um novo exemplo

© André de Carvalho - ICMC/USP

4

Desempenho de modelos preditivos

- Depende da tarefa:
 - Classificação: considera taxa de exemplos incorretamente classificados
 - Acurácia ou outras similares
 - Regressão: considera diferença entre valor previsto e valor correto
 - MSE ou outras similares
- Média das taxas de erro obtidas em diferentes execuções

© André de Carvalho - ICMC/USP

5

Tarefa de classificação

- Objetivo de um modelo de classificação:
 - Estimar corretamente o rótulo de novos exemplos (errar o mínimo possível)
 - Minimizar taxa de erro de classificação
 - Geralmente não é possível medir com exatidão essa taxa de erro
 - Ela deve ser estimada
 - Amostragem de dados

© André de Carvalho - ICMC/USP

6

Amostragem de dados

- Divide conjunto de dados em subconjuntos
- Subconjunto de treinamento
 - Ajuste de hiper-parâmetros de um algoritmo
 - Comparação de algoritmos ou de hiper-parâmetros de um algoritmo
 - Para dados de treinamento
- Subconjunto de teste
 - Comparação de algoritmos/modelos para novos dados
- Nunca para escolha

© André de Carvalho - ICMC/USP

7

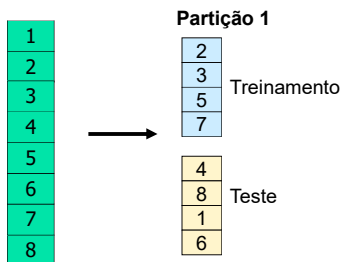
Amostragem de dados

- Permite melhor avaliação do desempenho de um classificador
 - Pode induzir um ou mais modelos
- Alternativas
 - Amostragem única
 - Hold-out
 - Re-amostragem

© André de Carvalho - ICMC/USP

8

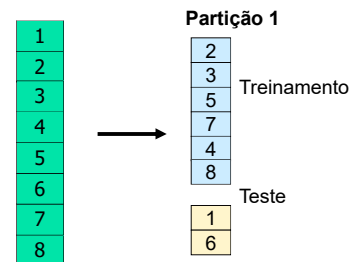
Hold-out



© André de Carvalho - ICMC/USP

9

Hold-out



© André de Carvalho - ICMC/USP

10

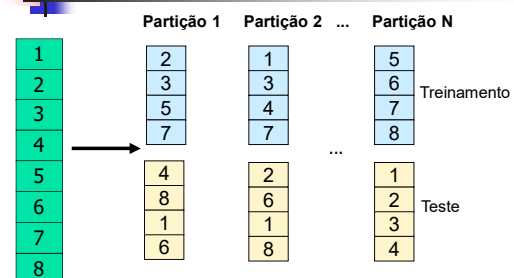
Métodos de re-amostragem

- Utilizam várias partições para os conjuntos de treinamento e teste
 - Random subsampling
 - K-fold Cross-validation
 - Leave-one-out
 - Bootstrap

© André de Carvalho - ICMC/USP

11

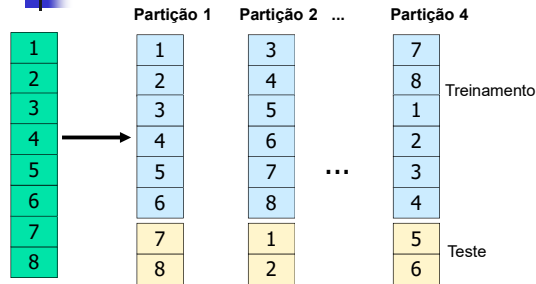
Random subsampling



© André de Carvalho - ICMC/USP

12

K-fold cross-validation



© André de Carvalho - ICMC/USP

13

Leave-one-out

- Estimativa de erro é praticamente não tendenciosa
 - Média das estimativas tende a taxa de erro verdadeira
- Computacionalmente caro
 - Geralmente utilizado para pequenos conjuntos de exemplos
 - 10-fold *cross validation* aproxima *leave-one-out*
- Variância tende a ser elevada

© André de Carvalho - ICMC/USP

14

Bootstrap

- Funciona melhor que *cross-validation* para conjuntos muito pequenos
- Forma mais simples de *bootstrap*:
 - Amostragem com reposição
 - Cada partição é uma amostra aleatória com reposição do conjunto total de exemplos
 - Conjunto de treinamento têm o mesmo número de exemplos do conjunto total
 - Exemplos que restarem são utilizados para teste

© André de Carvalho - ICMC/USP

15

Bootstrap

- Se conjunto original tem N exemplos
 - Amostra de tamanho N tem $\approx 63,2\%$ dos exemplos originais
 - Viés semelhante a 2-fold *cross-validation*
- Processo é repetido k vezes
 - Resultado final = média dos k experimentos
- Existem diversas variações

© André de Carvalho - ICMC/USP

16

Medidas de avaliação

- Acurácia
 - Trata as classes igualmente
 - Pode não ser adequada para dados desbalanceados
 - Pode prejudicar desempenho para classe minoritária
 - Classe rara é geralmente mais interessante que a majoritária
 - Alternativa: acurácia balanceada

© André de Carvalho - ICMC/USP

17

Classificação binária

- Classe de interesse é a classe positiva
- Dois tipos de erro:
 - Classificação de um exemplo N como P
 - Falso positivo (alarme falso)
 - Ex.: Diagnosticado como doente, mas está saudável
 - Classificação de um exemplo P como N
 - Falso negativo
 - Ex.: Diagnosticado como saudável, mas está doente

© André de Carvalho - ICMC/USP

18

Desempenho preditivo

- Matriz de confusão (tabela de contingência) pode ser utilizada para distinguir os erros
 - Base de várias medidas
 - Pode ser utilizada com 2 ou mais classes

Classe verdadeira	Classe predita		
	1	2	3
1	25	0	5
2	10	40	0
3	0	0	20

© André de Carvalho - ICMC/USP

19

Exemplo

- Matriz de confusão para 200 exemplos divididos em 2 classes

Classe verdadeira	Classe predita	
	p	n
P	70	30
N	40	60



Classe verdadeira	Classe predita	
	p	n
P	VP	FN
N	FP	VN

© André de Carvalho - ICMC/USP

20

Medidas de avaliação

$$\text{Taxa de FP (TFP)} = \frac{FP}{FP + VN} \quad \text{(Alarmes falsos)}$$

$$\text{Taxa de FN (TFN)} = \frac{FN}{VP + FN}$$

Erro do tipo I

Classe verdadeira	Classe predita	
	p	n
P	VP	FN
N	FP	VN

Erro do tipo II

Classe verdadeira	Classe predita	
	p	n
P	VP	FN
N	FP	VN

© André de Carvalho - ICMC/USP

21

Medidas de avaliação

$$\text{Taxa de FP (TFP)} = \frac{FP}{FP + VN} \quad \text{(Alarmes falsos)}$$

$$\text{Taxa de VP (TVP)} = \frac{VP}{FN + VP}$$

Custo

Benefício

Classe verdadeira	Classe predita	
	p	n
P	VP	FN
N	FP	VN

Classe verdadeira	Classe predita	
	p	n
P	VP	FN
N	FP	VN

© André de Carvalho - ICMC/USP

22

Exemplo

- Avaliação de 3 classificadores

Classe verdadeira	Classe predita	
	p	n
P	20	30
N	15	35

Classificador 1
TVP =
TFP =

Classe verdadeira	Classe predita	
	p	n
P	70	30
N	50	50

Classificador 2
TVP =
TFP =

Classe verdadeira	Classe predita	
	p	n
P	60	40
N	20	80

Classificador 3
TVP =
TFP =

© André de Carvalho - ICMC/USP

23

Exemplo

- Avaliação de 3 classificadores

Classe verdadeira	Classe predita	
	p	n
P	20	30
N	15	35

Classificador 1
TVP = 0.4
TFP = 0.3

Classe verdadeira	Classe predita	
	p	n
P	70	30
N	50	50

Classificador 2
TVP = 0.7
TFP = 0.5

Classe verdadeira	Classe predita	
	p	n
P	60	40
N	20	80

Classificador 3
TVP = 0.6
TFP = 0.2

© André de Carvalho - ICMC/USP

24

Medidas de avaliação

Medidas frequentemente utilizadas

$$TFP = \frac{FP}{VN + FP}$$

(Erro tipo I)

$$TFN = \frac{FN}{VP + FN}$$

(Erro tipo II)

$$\text{Precisão} = \frac{VP}{VP + FP}$$

$$\text{Especificidade} = \frac{VN}{VN + FP} = 1 - TFP$$

$$TVP = \frac{VP}{VP + FN}$$

Sensibilidade
Revocação (Recall)

$$\text{Acurácia} = \frac{VP + VN}{VP + VN + FP + FN}$$

$$\text{Medida-F1} = \frac{2}{1/\text{prec} + 1/\text{rev}}$$

© André de Carvalho - ICMC/USP

25

Revocação X Precisão

Revocação (*recall*)

- Porcentagem de exemplos positivos classificados como positivos $\frac{VP}{VP + FN}$
 - Nenhum exemplo positivo é deixado de fora

Precisão

- Porcentagem de exemplos classificados como positivos que são realmente positivos $\frac{VP}{VP + FP}$
 - Nenhum exemplo negativo é incluído

© André de Carvalho - ICMC/USP

26

Sensibilidade X Especificidade

Sensibilidade

- Porcentagem de exemplos positivos classificados como positivos $\frac{VP}{VP + FN}$
 - Igual a revocação

Especificidade

- Porcentagem de exemplos negativos classificados como negativos $\frac{VN}{VN + FP}$
 - Nenhum exemplo negativo é deixado de fora

© André de Carvalho - ICMC/USP

27

Medidas de avaliação

Medida-F

- Média harmônica ponderada da precisão e da revocação

$$\frac{(1 + \alpha) \times (\text{prec} \times \text{rev})}{\alpha \times \text{prec} + \text{rev}}$$

Medida-F1

- Precisão e revocação têm o mesmo peso

$$\frac{2 \times (\text{prec} \times \text{rev})}{\text{prec} + \text{rev}} = \frac{2}{1/\text{prec} + 1/\text{rev}}$$

© André de Carvalho - ICMC/USP

28

Exemplo

- Seja um classificador com a seguinte matriz de confusão, definir:

- Acurácia
- Precisão
- Revocação
- Especificidade

		Classe predita	
		p	n
Classe verdadeira	P	70	30
	N	40	60

© André de Carvalho - ICMC/USP

29

Exemplo

$$\text{Acurácia} = \frac{VP + VN}{VP + VN + FP + FN}$$

$$\text{Precisão} = \frac{VP}{VP + FP}$$

$$\text{Revocação} = \frac{VP}{VP + FN}$$

$$\text{Especificidade} = \frac{VN}{VN + FP}$$

		Predito	
		p	n
Verdadeiro	P	VP	FN
	N	FP	VN

© André de Carvalho - ICMC/USP

30

Exemplo

$$\text{Acurácia} = \frac{VP + VN}{VP + VN + FP + FN} = (70 + 60) / (70 + 30 + 40 + 60) = 0.65$$

$$\text{Precisão} = \frac{VP}{VP + FP} = 70 / (70 + 40) = 0.64$$

$$\text{Revocação} = \frac{VP}{VP + FN} = 7 / (70 + 30) = 0.70$$

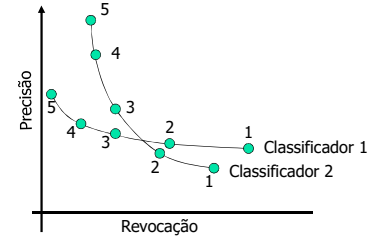
$$\text{Especificidade} = \frac{VN}{VN + FP} = 60 / (40 + 60) = 0.60$$

		Predito	
		p	n
Verdadeiro	P	VP	FN
	N	FP	VN
		p	n
Verdadeiro	P	70	30
	N	40	60

© André de Carvalho - ICMC/USP

31

Observação



© André de Carvalho - ICMC/USP

32

Outras medidas

- Média geométrica de taxas positivas

- G-means

$$\sqrt{\text{precisão} \times \text{revocação}}$$

- Acurácia balanceada
- Kappa

© André de Carvalho - ICMC/USP

33

Gráficos ROC

- Do inglês, *Receiver operating characteristics*
- Medida de desempenho originária da área de processamento de sinais
 - Muito utilizada nas áreas médica e biológica
 - Mostra relação entre custo (TFP) e benefício (TVP)

© André de Carvalho - ICMC/USP

34

Exemplo

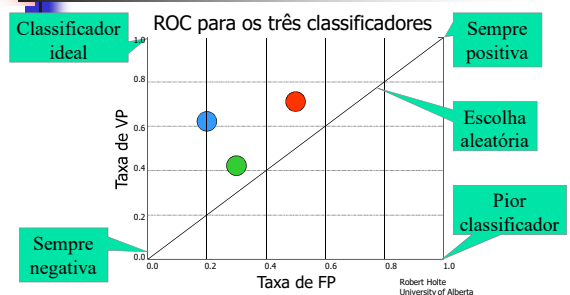
- Colocar no gráfico ROC os 3 classificadores do exemplo anterior

Classificador 1 TFP = 0.3 TVP = 0.4	Classificador 2 TFP = 0.5 TVP = 0.7	Classificador 3 TFP = 0.2 TVP = 0.6

© André de Carvalho - ICMC/USP

35

Gráficos ROC



© André de Carvalho - ICMC/USP

36

Gráficos ROC

- Classificadores discretos produzem um simples ponto no gráfico ROC
 - ADs e conjuntos de regras
- Outros classificadores produzem uma probabilidade ou escore
 - RNAs e NB
 - Permitem gerar vários pontos
- Curvas ROC permitem uma melhor comparação de classificadores
 - Insensíveis a mudanças na distribuição das classes

© André de Carvalho - ICMC/USP

37

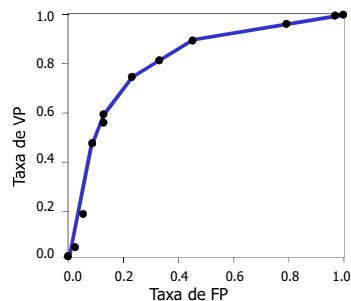
Curvas ROC

- Mostram pontos no gráfico para diferentes variações de um classificador
- Classificadores que geram valores contínuos (*threshold*, probabilidade)
 - Diferentes valores de *threshold* podem ser utilizados para gerar vários pontos
 - Ligação dos pontos gera uma curva ROC
- Classificadores discretos
 - Convertidos internamente ou comitês

© André de Carvalho - ICMC/USP

38

Curvas ROC



© André de Carvalho - ICMC/USP

39

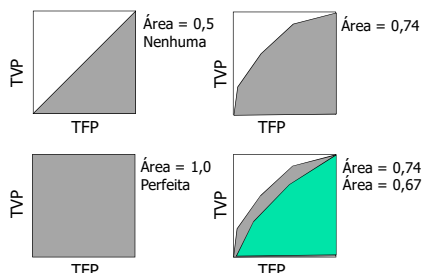
Área sob a curva ROC (AUC)

- Fornecer uma estimativa do desempenho de classificadores
- Gera um valor contínuo no intervalo $[0.0, 1.0]$
 - Quanto maior melhor
 - Adição de áreas de sucessivos trapezóides
- É mais confiável utilizar médias de AUCs

© André de Carvalho - ICMC/USP

40

Área Sob Curvas ROC



© André de Carvalho - ICMC/USP

41

Avaliação de Desempenho

- Teste de Hipóteses
 - Permite afirmar que uma técnica é melhor que outra com X% de confiança
 - Podem assumir que os dados seguem uma dada distribuição de probabilidade
 - Testes paramétricos
 - Testes não paramétricos
 - Número de técnicas comparadas
 - Duas
 - Mais que duas

© André de Carvalho - ICMC/USP


42

Considerações Finais

- Desempenho preditivo
- Avaliação do desempenho
 - Erro
 - Tempo de resposta
 - Memória
 - Representação
- Medidas
- Gráficos e curvas ROC
- Teste de hipóteses

© André de Carvalho - ICMC/USP 43

Perguntas



© André de Carvalho - ICMC/USP 44