

Mineração de Dados em Biologia Molecular

Planejamento e
Análise de Experimentos

André C. P. L. F. de Carvalho
Monitor: Valéria Carvalho



Principais tópicos

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

27/09/2012

André de Carvalho - ICMC/USP

2

Estimativa de erro

- Depende do problema:
 - Classificação: considera taxa de exemplos incorretamente classificados
 - Acurácia
 - Regressão: considera diferença entre valor produzido e valor esperado
 - Agrupamento: diferentes critérios
- Média dos erros obtidos em diferentes execuções de um experimento

27/09/2012

André de Carvalho - ICMC/USP

3

Estimativa de Erro de Classificação

- Processo de treinamento é utilizado para seleção do modelo
 - Modelo com a complexidade correta (sem overfitting)
- Após construção do modelo, ele pode ser testado com novos exemplos
 - Evitar modelo otimista
 - Conjunto de teste
 - Estimativa não tendenciosa de erro de generalização
 - Comparação de modelos utiliza desempenho em dados de teste
 - Métodos de amostragem

27/09/2012

André de Carvalho - ICMC/USP

4

Métodos de amostragem

- Utilizados para avaliar desempenho de um classificador
 - Hold-out
 - Random subsampling
 - Cross validation
 - Leave-one-out
 - Bootstrap

27/09/2012

André de Carvalho - ICMC/USP

5

Hold-out

- Também conhecido como *split-sample*
- Técnica mais simples para estimativa de erro
- Faz uma única partição da amostra em:
 - Conjunto de treinamento: geralmente 1/2 ou 2/3 dos dados
 - Conjunto para teste: os dados restantes

27/09/2012

André de Carvalho - ICMC/USP

6

Hold-out

- Indicado para grande quantidade de dados (ex.: mais de 1000)
- Pequena quantidade de dados
 - Poucos exemplos são usados no treinamento
 - Modelo pode depender da composição dos conjuntos de treinamento e teste
 - Quanto menor conjunto de treinamento, maior a variância do modelo
 - Quanto menor conjunto de teste, menos confiável a acurácia estimada para ele

27/09/2012

André de Carvalho - ICMC/USP

7

Hold-out

- Conjuntos de treinamento e teste não são independentes
 - Classe sub-representada em um conjunto será super-representada no outro
 - E vice-versa
- Aproximação pessimista
- Resultados obtidos podem ser pouco significativos
- Solução: utilizar reamostragem

27/09/2012

André de Carvalho - ICMC/USP

8

Métodos de reamostragem

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

27/09/2012

André de Carvalho - ICMC/USP

9

Random subsampling

- Diferentes partições treinamento-teste são escolhidas de forma aleatória
 - $D_{\text{Trein}} \cap D_{\text{Teste}} = \emptyset$
 - Taxa de erro é calculada para cada partição
 - Taxa de erro estimada é a média dos erros para as diferentes partições
- Pode obter uma estimativa de erro mais precisa para o desempenho de um modelo

27/09/2012

André de Carvalho - ICMC/USP

10

Cross-validation

- Validação cruzada
- Classe de métodos para estimativa da taxa de erro verdadeira
 - K-fold cross-validation
 - Cada objeto participa o mesmo número de vezes do treinamento
 - E apenas uma vez do teste
 - Estratificado
 - Leave-one-out ($K = N$)

27/09/2012

André de Carvalho - ICMC/USP

11

Leave-one-out

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

27/09/2012

André de Carvalho - ICMC/USP

12

5 x 2 Cross-validation

- Conjuntos de treinamento e teste com mesmo tamanho
- Dietterich, 1998

Seja um conjunto de N exemplos
 Para $i = 1$ até 5
 Dividir N aleatoriamente em duas metades
 Usar metade 1 para treinamento e metade 2 para teste
 Usar metade 2 para treinamento e metade 1 para teste

27/09/2012

André de Carvalho - ICMC/USP

13

Bootstrap

- Funciona melhor que cross-validation para conjuntos muito pequenos
- Forma mais simples de bootstrap:
 - Ao invés de usar sub-conjuntos dos dados, usar sub-amostras
 - Cada sub-amostra é uma amostra aleatória com substituição do conjunto total de exemplos
 - Cada conjunto de treinamento têm o mesmo número de exemplos do conjunto total
 - Os exemplos que restarem são utilizados para teste

27/09/2012

André de Carvalho - ICMC/USP

14

Bootstrap

- Se conjunto original tem N exemplos
 - Amostra de tamanho N tem $\approx 63,2\%$ dos exemplos originais
- Processo é repetido b vezes
 - Resultado final = média dos b experimentos
- Existem diversas variações
 - Como calcular a acurácia do classificador
 - .632 bootstrap

27/09/2012

André de Carvalho - ICMC/USP

15

Erro de classificação

- Principal objetivo de um modelo é classificar corretamente para novos exemplos
 - Erro o mínimo possível
 - Minimizar taxa de erro
 - Geralmente não é possível medir com exatidão essa taxa de erro
 - Ela deve ser estimada

27/09/2012

André de Carvalho - ICMC/USP

16

Estimativa de erro de classificação

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

27/09/2012

André de Carvalho - ICMC/USP

17

Classificação binária

- 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

27/09/2012

André de Carvalho - ICMC/USP

18

Tipos de erro

- Classe positiva é, em geral, a classe de maior interesse
 - Ou com menos exemplos
- Em alguns casos, os erros têm igual importância
 - Em outros, erros diferentes têm consequências diferentes
 - No exemplo anterior, qual é pior (tem mais custo)? falso negativo ou falso positivo?

27/09/2012

André de Carvalho - ICMC/USP

19

Tipos de erro

- 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

27/09/2012

André de Carvalho - ICMC/USP

20

Avaliação de desempenho

- 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

27/09/2012

André de Carvalho - ICMC/USP

21

Avaliação de desempenho

- Medidas de desempenho

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

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

Erro do tipo I

Erro do tipo II

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

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

27/09/2012

André de Carvalho - ICMC/USP

22

Avaliação de desempenho

- Medidas de desempenho

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

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

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

27/09/2012

André de Carvalho - ICMC/USP

23

Exemplo

- Avaliação de 3 classificadores

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

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

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

Classificador 1
TVP =
TFP =

Classificador 2
TVP =
TFP =

Classificador 3
TVP =
TFP =

27/09/2012

André de Carvalho - ICMC/USP

24

Exemplo

Avaliação de 3 classificadores

Classe verdadeira	Classe predita		Classe verdadeira	Classe predita		Classe verdadeira	Classe predita	
	p	n		p	n		p	n
P	20	30	P	70	30	P	60	40
N	15	35	N	50	50	N	20	80

Classificador 1
TVP = 0.4
TFP = 0.3

Classificador 2
TVP = 0.7
TFP = 0.5

Classificador 3
TVP = 0.6
TFP = 0.2

27/09/2012

André de Carvalho - ICMC/USP

25

Avaliação de desempenho

Medidas frequentemente utilizadas

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

(Erro tipo I)

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

(Erro tipo II)

$$Precisão = \frac{VP}{VP + FP}$$

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

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

Sensibilidade
Revocação (Recall)

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

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

27/09/2012

André de Carvalho - ICMC/USP

26

Revocação X Precisão

Revocação (*recall*)

- Taxa com que classifica como positivos todos os exemplos que são positivos
 - Nenhum exemplo positivo é deixado de fora

Precisão

- Taxa com que todos os exemplos classificados como positivos são realmente positivos
 - Nenhum exemplo negativo é incluído

27/09/2012

André de Carvalho - ICMC/USP

27

Sensibilidade X Especificidade

Sensibilidade

- Taxa com que os exemplos positivos são classificados como positivos
 - Igual a revocação

Especificidade

- Taxa com que exemplos negativos são classificados como negativos
 - Nenhum exemplo negativo é deixado de fora

27/09/2012

André de Carvalho - ICMC/USP

28

Avaliação de desempenho

Medida-F

- Média harmônica ponderada da precisão e da revocação

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

Medida-F1

- Precisão e revocação têm o mesmo peso

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

27/09/2012

André de Carvalho - ICMC/USP

29

Exemplo

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

- Acurácia
- Precisão
- Revocação
- Especificidade

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

27/09/2012

André de Carvalho - ICMC/USP

30

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
		p	
		P	N
		70	30
		40	60

27/09/2012

André de Carvalho - ICMC/USP

31

Exemplo

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

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

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

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

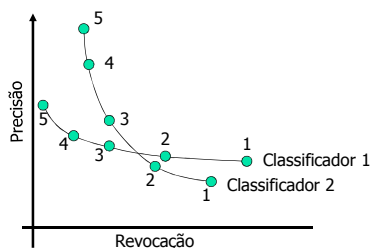
		Predito	
		p	n
Verdadeiro	P	VP	FN
	N	FP	VN
		p	
		P	N
		70	30
		40	60

27/09/2012

André de Carvalho - ICMC/USP

32

Observação



27/09/2012

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)

27/09/2012

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

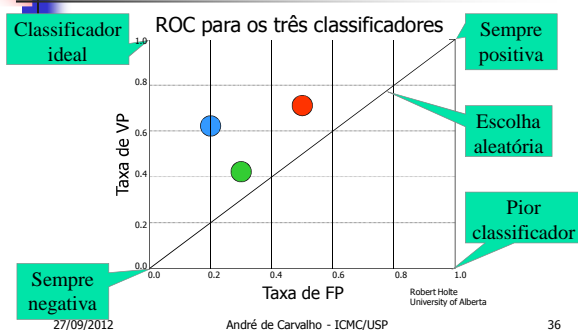
Classificador 3
TFP = 0.2
TVP = 0.6

27/09/2012

André de Carvalho - ICMC/USP

35

Gráficos ROC



27/09/2012

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
- Curvas ROC permitem uma melhor comparação de classificadores
 - São insensíveis a mudanças na distribuição das classes

27/09/2012

André de Carvalho - ICMC/USP

37

Curvas ROC

- Mostram ROC para diferentes variações
- Classificadores que geram escores ou probabilidades
 - Diferentes valores de *threshold* podem ser utilizados para gerar vários pontos
 - Cada valor de *threshold* produz um ponto diferente
 - Ligação dos pontos gera uma curva ROC

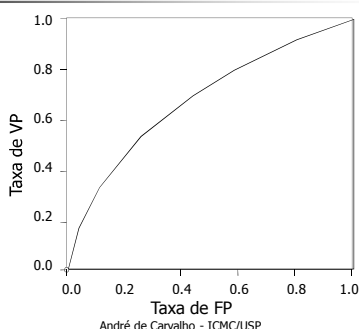
27/09/2012

André de Carvalho - ICMC/USP

38

Curvas ROC

Classificador
Escore/
Probabilístico



27/09/2012

André de Carvalho - ICMC/USP

39

Curvas ROC

- Classificadores que geram valores discretos
 - Podem ser convertidos internamente para gerar escores
 - Para ADs, diferentes thresholds para números de exemplos positivos que tornam a classe positiva
 - Podem ser combinados em comitês
 - Threshold para votos dos classificadores individuais forma escore

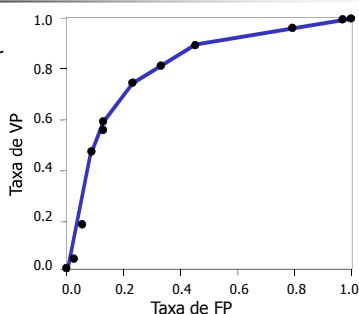
27/09/2012

André de Carvalho - ICMC/USP

40

Curvas ROC

Classificador
Discreto



27/09/2012

André de Carvalho - ICMC/USP

41

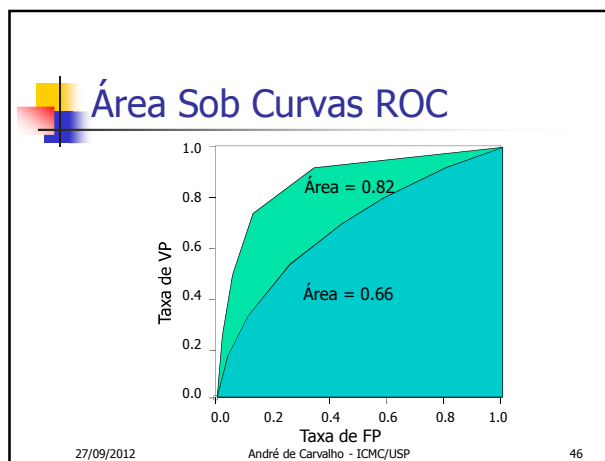
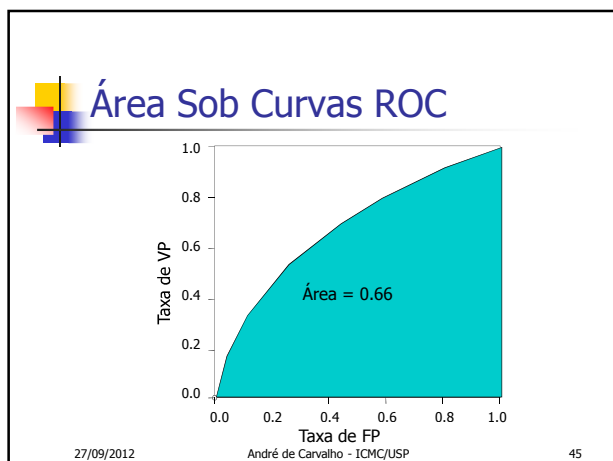
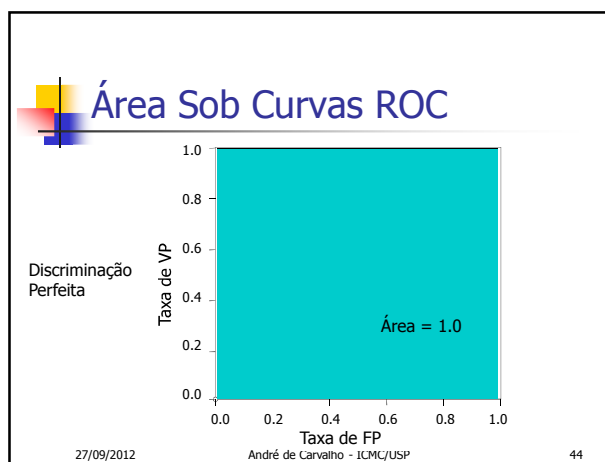
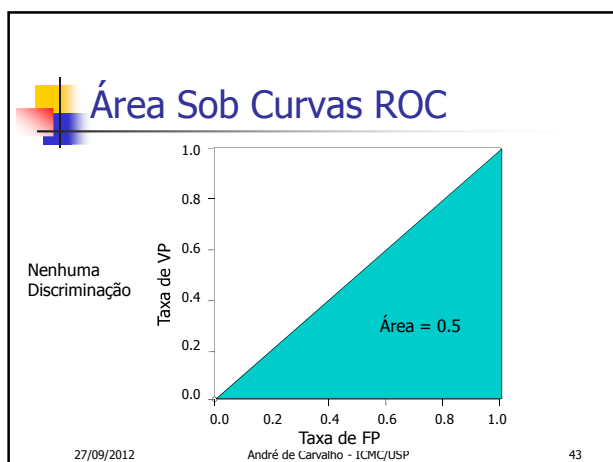
Área sob a curva ROC (AUC)

- Fornece uma estimativa do desempenho de classificadores
- Gera um valor contínuo no intervalo $[0, 1]$
 - Quanto maior melhor
 - Adição de áreas de sucessivos trapezóides
- Um classificador com maior AUC pode apresentar AUC pior em trechos da curva
- É mais confiável utilizar médias de AUCs

27/09/2012

André de Carvalho - ICMC/USP

42



Avaliação de Desempenho


- Teste de Hipóteses
 - Compara dois desempenhos
 - Teste t
 - Teste McNemar
 - Teste t pareado de 5x2 CV
 - Compara vários desempenhos
 - Teste Feelders e Verkooijen
 - Teste de Friedman
 - ANOVA

27/09/2012 André de Carvalho - ICMC/USP 47


Considerações Finais

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

27/09/2012 André de Carvalho - ICMC/USP 48



Perguntas



27/09/2012

André de Carvalho - ICMC/USP

49