

Aprendizado de Máquina

Máquinas de vetores de suporte

André C. P. L. F. de Carvalho
ICMC-USP



© André de Carvalho - ICMC/USP

4

TAE

- Sejam
 - h : classificador (hipótese, modelo, função)
 - H : conjunto de todos os classificadores que um algoritmo de AM pode induzir
- Algoritmo de AM utiliza conjunto de dados de treinamento para:
 - Induzir um classificador $\hat{h} \in H$
- Assume que dados são gerados de forma i.i.d. de acordo com $P(x,y)$

© André de Carvalho - ICMC/USP

4

Principais tópicos

- Introdução
- Risco empírico e risco estrutural
- Margens
- Margens suaves
- SVMs
- Kernels
- Multiclasses

© André de Carvalho - ICMC/USP

2

TAE

- TAE define condições matemáticas para auxiliar na escolha de uma boa função \hat{h}
 - A partir de um conjunto de dados de treinamento
 - Permite escolher \hat{h} com menor risco esperado
 - Para manter bom desempenho com novos dados, avalia:
 - Desempenho preditivo de h para dados do conjunto de treinamento
 - Complexidade de \hat{h}

© André de Carvalho - ICMC/USP

5

Teoria de Aprendizado Estatístico

- Algoritmos de AM
 - Estimam um função (modelo) a partir de um conjunto finito de exemplos
 - Função (classificador ou regressor)
- TAE estabelece princípios para induzir função com boa generalização
 - Vapnik e Chervonenkis em 1968
 - Base das máquinas de vetores de suporte

© André de Carvalho - ICMC/USP

3

Limites no risco esperado

- Isso é feito pelas máquinas de vetores de suporte (SVMs)
- Estratégia básica
 - Encontrar um hiperplano que maximize margem de separação (margem larga)
 - Distância da fronteira de decisão a um conjunto de "vetores de suporte"
 - Com erro marginal baixo
 - Número mínimo de objetos entre as margens

© André de Carvalho - ICMC/USP

6

Máquinas de Vetores de Suporte (SVMs)

Rede Neural SVMs

© André de Carvalho - ICMC/USP 7

Linearmente separáveis

- SVMs apresentam bons desempenhos para problemas linearmente separáveis
- Não conseguem lidar com problemas não linearmente separáveis
- Alguns conjuntos de dados exigem fronteiras mais complexas que lineares
 - Para isso foram propostas alterações baseadas no teorema de Cover

© André de Carvalho - ICMC/USP 10

SVMs

$\bar{w} \cdot \bar{x} + b = 0$
 $H_2: \bar{w} \cdot \bar{x} + b = -1$ $H_1: \bar{w} \cdot \bar{x} + b = +1$
 Margem máxima $\text{Vetores de suporte (pontos críticos)}$
 $\text{Hiperplano separador ótimo}$ $\text{Margem} = \frac{2}{\|\bar{w}\|^2}$

$\text{signal}(h(x)) = \begin{cases} +1 & \text{if } w \cdot x + b \geq 1 \\ -1 & \text{if } w \cdot x + b \leq -1 \end{cases}$
 $y_i \times (w \cdot x_i + b) \geq 1$

© André de Carvalho - ICMC/USP 8

Teorema de Cover

Conjunto de dados não linearmente separáveis em um espaço podem ser transformados para outro espaço em que, com alta probabilidade, se tornam linearmente separáveis

- Condições:
 - Transformação seja não linear
 - Dimensão do novo espaço seja suficientemente alta

© André de Carvalho - ICMC/USP 11

Variáveis de folga

- Slack variables

© André de Carvalho - ICMC/USP 9

Problemas não lineares

- Generalização para problemas não lineares
 - Mapeamento de dados de entrada para um espaço de maior dimensão

$f(x) = w \cdot \Phi(x) + b$
 Espaço de entradas Espaço de características

© André de Carvalho - ICMC/USP 12

Exemplo

- Suor conjunto de dados com 2 atributos preditivos
- Definir 3 pontos de localização no conjunto original
- Usar esses pontos para transformar 2 atributos originais em 3 outros atributos
 - Ex. Distância entre cada exemplo x e cada um dos 3 pontos de localização

© André de Carvalho - ICMC/USP

13

Funções Kernel

- Em geral, K é menos complexa que Φ
 - É comum definir-se a função K sem conhecer-se explicitamente Φ

Tipos de Kernel	Função $K(x_i, x_j)$ correspondente
Polinomial	$(x_i^T \cdot x_j + 1)^p$ ($p = 1$, linear)
Gaussiano	$\exp(-1/(2\sigma^2) \ x_i - x_j\ ^2)$
Sigmoidal	$\tanh(\beta_0 x_i \cdot x_j + \beta_1)$

© André de Carvalho - ICMC/USP

16

Fronteiras mais complexas

- Computação da função Φ pode ter custo computacional elevado
 - Informação necessária: cálculo do produto escalar entre objetos
 - Pode simplificar usando funções kernel (K)
 - Rcebem 2 pontos no espaço de entradas e calculam produto escalar deles no espaço de características
 - $K(x_i, x_j) \leftrightarrow \Phi(x_i) \cdot \Phi(x_j)$

© André de Carvalho - ICMC/USP

14

Classificação multiclases

- SVMs podem induzir apenas classificadores binários
 - Outros algoritmos de AM têm a mesma limitação
- Existe um grande número de problemas reais com mais que 2 classes
 - Necessidade de estratégias multiclases

© André de Carvalho - ICMC/USP

17

Funções Kernel

- Diversas
 - Gaussiana
 - Polinomial
 - Linear
 - Sigmoidal
 - Para aplicações específicas
- Seguem condições estabelecidas pelo teorema de Mercer
- Híper-parâmetros ajustáveis

© André de Carvalho - ICMC/USP

15

Estratégias multiclases

- Duas abordagens têm sido utilizadas:
 - Algoritmo de classificação é internamente adaptado
 - Modificação de parte de suas operações internas
 - Decomposição do problema multiclases em vários problemas binários
 - Estratégias decomposicionais

© André de Carvalho - ICMC/USP


18

Estratégias decomposicionais

- Etapas
 - Decomposição da tarefa
 - Reconstrução
- Decomposição
 - Geralmente reduz a complexidade da tarefa
 - Permite processamento paralelo
 - Alternativas:
 - Matrizes de códigos (MC)
 - Hierarquias de classificadores

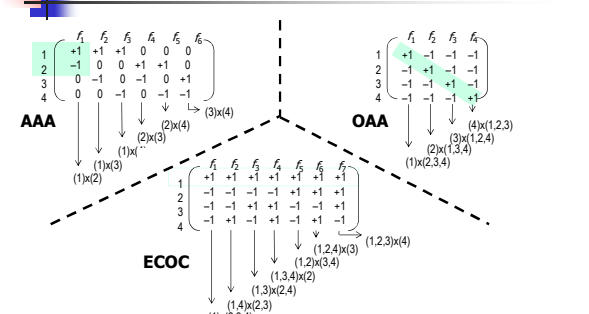
© André de Carvalho - ICMC/USP 19

Perguntas



© André de Carvalho - ICMC/USP 22

Matrizes de códigos



© André de Carvalho - ICMC/USP 20

Exercício

- Utilizando do repositório UCI as bases de dados IRIS e GLASS
 - Investigar Perceptron e MLP
 - Investigar SVMs
 - Três kernels diferentes
 - AAA e OAA
 - Particionar os dados com k-fold crossvalidation
 - Ajustar parâmetros por tentativa e erro
 - De 2 a 4 páginas

© André de Carvalho - ICMC/USP 23

Conclusão

- Teoria de Aprendizado Estatístico
- SVMs
- Problemas não linearmente separáveis
- Classificação binária e multiclases
- Regressão
- Redes profundas

© André de Carvalho - ICMC/USP 21