

# Máquinas de suporte vetorial e sua aplicação na detecção de spam

Antonio Carlos dos Santos  
Orientador: Paulo J. S. Silva (IME-USP)

Universidade de São Paulo  
Instituto de Matemática e Estatística  
Departamento de Ciência da Computação

MAC499 Trabalho de Formatura Supervisionado 2005

# Conteúdo

- 1 Motivação
  - Importância da detecção de spam
- 2 Aprendizado Computacional
  - Conceitos
- 3 Máquinas de suporte vetorial
  - Dados linearmente separáveis
  - Dados não-linearmente separáveis
- 4 Aplicação na detecção de spam
- 5 BCC

## Informações sobre spams

Spams (***Stupid Pointless Annoying Messages***): e-mails indesejados que as pessoas recebem.

- 31 milhões de e-mails enviados diariamente em 2004 [Sfr05]
- 12,4 milhões eram spams [Sfr05]
- Perdas mundiais chegarão a **US\$50 bilhões** em 2005, principalmente devido à queda de produtividade dos funcionários [Fer05]

## Informações sobre spams

Spams (***Stupid Pointless Annoying Messages***): e-mails indesejados que as pessoas recebem.

- 31 milhões de e-mails enviados diariamente em 2004 [Sfr05]
- 12,4 milhões eram spams [Sfr05]
- Perdas mundiais chegarão a **US\$50 bilhões** em 2005, principalmente devido à queda de produtividade dos funcionários [Fer05]

## Informações sobre spams

Spams (***Stupid Pointless Annoying Messages***): e-mails indesejados que as pessoas recebem.

- 31 milhões de e-mails enviados diariamente em 2004 [Sfr05]
- 12,4 milhões eram spams [Sfr05]
- Perdas mundiais chegarão a **US\$50 bilhões** em 2005, principalmente devido à queda de produtividade dos funcionários [Fer05]

## Informações sobre spams

Spams (***Stupid Pointless Annoying Messages***): e-mails indesejados que as pessoas recebem.

- 31 milhões de e-mails enviados diariamente em 2004 [Sfr05]
- 12,4 milhões eram spams [Sfr05]
- Perdas mundiais chegarão a **US\$50 bilhões** em 2005, principalmente devido à queda de produtividade dos funcionários [Fer05]

# Problemas com spams

- Tempo gasto analisando mensagens para decidir se é spam ou não
- Risco de apagar mensagens que não são spams
- Prejuízo por vírus e outros programas maliciosos

## Pergunta:

- Por que não fazemos um programa para classificar automaticamente os spams?
- Dificuldades:
  - Caracterizar um spam
  - Acompanhar o surgimento de novos formatos de spam

## Pergunta:

- Por que não fazemos um programa para classificar automaticamente os spams?
- Dificuldades:
  - Caracterizar um spam
  - Acompanhar o surgimento de novos formatos de spam

# Idéia

- Criarmos um programa que classifique as mensagens como spam ou não através da análise prévia de algumas mensagens (exemplos) que já tenham sido analisadas, semelhante a como aprendemos a fazer algumas tarefas

# Conceitos

- **Mundo real:** conjunto no qual estamos interessados
- **Dado ou exemplo:** elemento do mundo real
- **Máquina de aprendizado:** dispositivo (um programa) capaz de analisar dados do mundo real e classificá-los (dar rótulos a eles).
- **Conceito:** a classificação correta dos dados que a máquina tenta aprender
- **Hipótese:** cada uma das possíveis classificações que a máquina faz sobre os dados

# Objetivo

- Fazer com que a máquina de aprendizado gere uma hipótese que melhor aproxima o conceito desejado
- Para isso, devemos treinar a máquina de aprendizado com exemplos já classificados corretamente

# Fases

- 1 Treinamento: definir o padrão de classificação
- 2 Teste: classificar novos dados

# Máquinas de suporte vetorial

## Apresentação

- Desenvolvidas principalmente por V. Vapnik, usam idéias de aprendizado estatístico
- Têm algumas vantagens sobre outros métodos de classificação
  - Sem mínimos locais
  - Fase de testes é rápida

# Máquinas de suporte vetorial

Sejam  $I$  dados de treinamento (amostras), cada um formado por um vetor  $\mathbf{x}_i \in \mathbb{R}^d$  e um rótulo  $y_i \in \{-1, 1\}$

- Os dados possuem uma distribuição de probabilidade  $P(\mathbf{X}, Y)$  desconhecida e são independentes e identicamente distribuídos
- Mas iremos supor que os rótulos são fixos

# Máquinas de suporte vetorial

Sejam  $I$  dados de treinamento (amostras), cada um formado por um vetor  $\mathbf{x}_i \in \mathbb{R}^d$  e um rótulo  $y_i \in \{-1, 1\}$

- Os dados possuem uma distribuição de probabilidade  $P(\mathbf{X}, Y)$  desconhecida e são independentes e identicamente distribuídos
- Mas iremos supor que os rótulos são fixos

# Máquinas de suporte vetorial

Sejam  $I$  dados de treinamento (amostras), cada um formado por um vetor  $\mathbf{x}_i \in \mathbb{R}^d$  e um rótulo  $y_i \in \{-1, 1\}$

- Os dados possuem uma distribuição de probabilidade  $P(\mathbf{X}, Y)$  desconhecida e são independentes e identicamente distribuídos
- Mas iremos supor que os rótulos são fixos

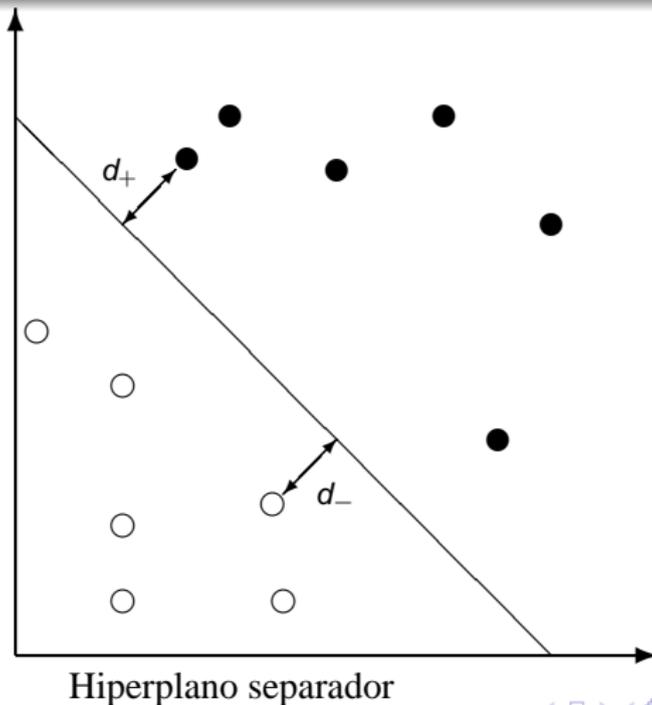
# Modelos

- Dados linearmente separáveis
- Dados não-linearmente separáveis

# Hiperplano separador

- Um hiperplano em  $\mathbb{R}^d$  que separa as amostras positivas ( $y_i = +1$ ) das negativas ( $y_i = -1$ )
- $\langle \mathbf{w}, \mathbf{x} \rangle + b = 0$  é uma equação para este hiperplano
- $\frac{|b|}{\|\mathbf{w}\|}$  é a distância perpendicular do hiperplano da origem
- $d_+$  será a menor distância dos pontos classificados como positivos ao hiperplano e  $d_-$  será a menor distância dos pontos classificados como negativos ao hiperplano
- Definiremos a **margem** por  $(d_+ + d_-)$

## Exemplo



- Mudando a escala entre  $|b|$  e  $\|\mathbf{w}\|$ , teremos:

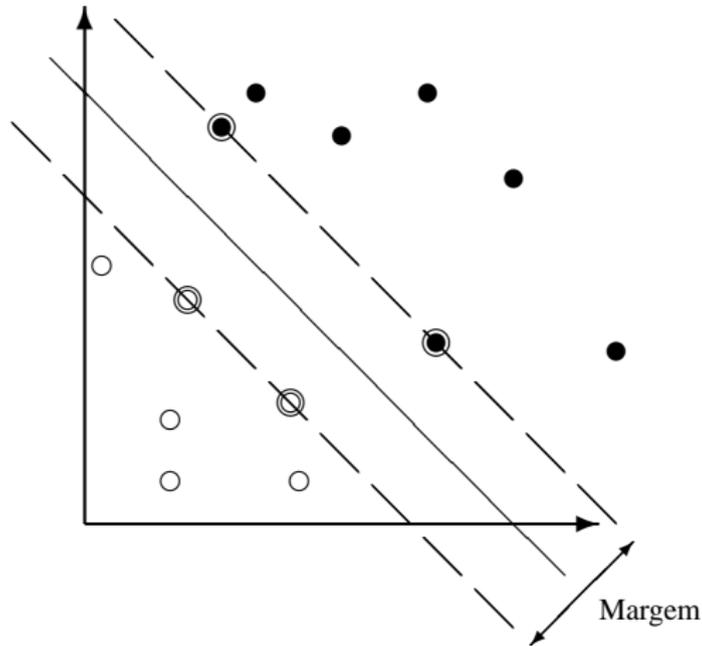
$$\left. \begin{array}{l} \langle \mathbf{x}_i, \mathbf{w} \rangle + b \geq +1 \quad \text{se } y_i = +1 \\ \langle \mathbf{x}_i, \mathbf{w} \rangle + b \leq -1 \quad \text{se } y_i = -1 \end{array} \right\} \quad i = 1, \dots, l.$$

- Podemos combinar as duas restrições em uma só:

$$y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1 \geq 0 \quad i = 1, \dots, l.$$

- **Vetores suporte:** dados de treinamento para os quais é válida a igualdade nas restrições acima
- A margem será  $\frac{2}{\|\mathbf{w}\|}$

# Exemplo



## Objetivo

- Máquinas de suporte vetorial procuram maximizar a margem do hiperplano separador
- Portanto, temos o seguinte problema:

$$\begin{aligned} \max \quad & 2/\|\mathbf{w}\| \\ \text{s.a} \quad & y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1 \geq 0 \quad i = 1, \dots, l. \end{aligned}$$

## Objetivo

- Máquinas de suporte vetorial procuram maximizar a margem do hiperplano separador
- Portanto, temos o seguinte problema:

$$\begin{aligned} \max \quad & 2/\|\mathbf{w}\| \\ \text{s.a} \quad & y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1 \geq 0 \quad i = 1, \dots, l. \end{aligned}$$

## Problema

- Sendo  $\|\mathbf{w}\| \geq 0$ , o problema é equivalente a:

$$\begin{aligned} \min \quad & \|\mathbf{w}\|^2/2 \\ \text{s.a} \quad & y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1 \geq 0 \quad i = 1, \dots, l. \end{aligned}$$

- A função Lagrangeana associada é:

$$\begin{aligned} L_P(\mathbf{w}, b, \alpha) &= \frac{\|\mathbf{w}\|^2}{2} - \sum_{i=1}^l \alpha_i y_i (\langle \mathbf{x}_i, \mathbf{w} \rangle + b) + \sum_{i=1}^l \alpha_i \\ \alpha_i &\geq 0, \quad i = 1, \dots, l. \end{aligned}$$

em que  $\alpha$  é o vetor dos multiplicadores de Lagrange

## Problema dual

- A partir das condições KKT, temos o seguinte problema dual equivalente:

$$\begin{aligned} \max \quad & L_D(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{s.a} \quad & \alpha_i \geq 0 \quad i = 1, \dots, l. \end{aligned}$$

no qual os dados de treinamento só aparecem como produtos internos.

## Classificação de novos dados

- Após termos treinado a máquina, como sabemos a classe de um novo dado?
- Basta ver em qual lado do hiperplano o dado está:

$$\text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b) = \text{sgn} \left( \sum_{i=1}^l \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle \right)$$

## Dados não-linearmente separáveis

- Para dados não linearmente separáveis, nós introduzimos uma variável de folga  $\xi_i \geq 0$  na função objetivo para cada dado de treinamento, de forma a penalizar a violação das restrições do problema original
- Novo problema:

$$\begin{aligned} \min \quad & \frac{\|\mathbf{w}\|^2}{2} + C \sum_{i=1}^l \xi_i \\ \text{s.a} \quad & y_i(\langle \mathbf{x}_i, \mathbf{w} \rangle + b) - 1 + \xi_i \geq 0 \quad i = 1, \dots, l \\ & \xi_i \geq 0 \quad i = 1, \dots, l, \end{aligned}$$

em que  $C$  e  $k$  são parâmetros escolhidos pelo usuário e determinam a importância dos erros de classificação feitos pelo algoritmo.

## Aplicação na detecção de spam

- O estudo foi baseado principalmente no artigo *Support Vector Machines for Spam Categorization* [HWV99]
- **Característica**: uma palavra de um e-mail
- A cada mensagem, teremos um vetor de características  $\mathbf{x}$  associado, formado por palavras de um dicionário gerado pela análise de todas as mensagens.
- Apenas palavras que aparecem em pelo menos 3 e-mails diferentes

# Abordagens

- Frequência da palavra: número que cada palavra do dicionário ocorre no texto
- Representação binária: se cada palavra ocorre ou não no texto

Foram usadas todas as palavras ao invés de selecionarmos apenas algumas

## Disciplinas importantes:

- MAC0122 Princípios de Desenvolvimento de Algoritmos
- MAC0338 Análise de Algoritmos
- MAC0315 Programação Linear
- MAC0300 Métodos Numéricos de Álgebra Linear
- MAC5732 Aprendizado Computacional
- MAT0121 Cálculo Diferencial e Integral II
- MAT0139 Álgebra Linear para Computação
- MAE0212 Introdução à Probabilidade e à Estatística II

## Para quem quiser saber mais: I

-  C. J. C. Burges. A Tutorial on Support Vector Machines for Pattern Recognition. *Data Mining and Knowledge Discovery*, 2:121-167, 1998.
-  N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines and other kernel-based methods*. Cambridge University Press, 2002.
-  Ferris Research. <http://www.ferris.com>, Outubro de 2005.
-  H. Druker, D. Wu, V. N. Vapnik. Support Vector Machines for Spam Categorization. *IEEE Transactions on Neural Networks*, 10(5):1048-1054, 1999.

## Para quem quiser saber mais: II



Spam Filter Review.

<http://spam-filter-review.toptenreviews.com>, Outubro de 2005.



V. N. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, 1998.