

Maximização de Margem

Victor Ruela

Programa de Pós-Graduação em Engenharia Elétrica
Universidade Federal de Minas Gerais

victorspruela@ufmg.br

2 de fevereiro de 2021

Agenda

- 1 Introdução
 - Aprendizado Supervisionado
 - Minimização do Erro
- 2 Maximização de Margem
 - Otimização
 - Padrões Não Separáveis
- 3 Problemas não-linearmente separáveis
 - Exemplos

Problema de Classificação

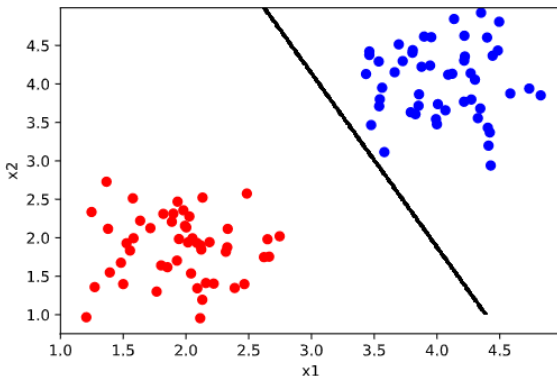
- Dados de treinamento: $\mathcal{T} = \{(\mathbf{x}_1, d_1), \dots, (\mathbf{x}_N, d_N)\}$
- Para um problema de classificação binário: $d_i = \{-1, 1\}$
- Assume-se que os padrões são linearmente separáveis
- Em geral, algoritmos para treinamento de RNAs objetivam minimizar o erro quadrático:

$$\min \sum_{i=1}^N [d_i - \text{sign}(\mathbf{w}^T \mathbf{x}_i + b)]^2 \quad (1)$$

Minimização do Erro

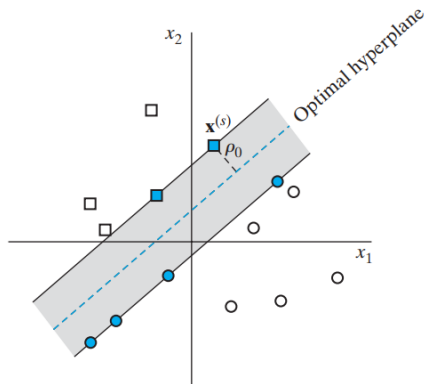
- O resultado da solução do Problema 1 será um conjunto de pesos representando um hiperplano que separa estes padrões:

$$\mathbf{w}^T \mathbf{x} + b = 0 \quad (2)$$



O Hiperplano Ótimo

- Margem de separação: ρ
- Quando a escolha de \mathbf{w} e b maximizam ρ , o hiperplano é dito **ótimo** [Haykin, 2007].



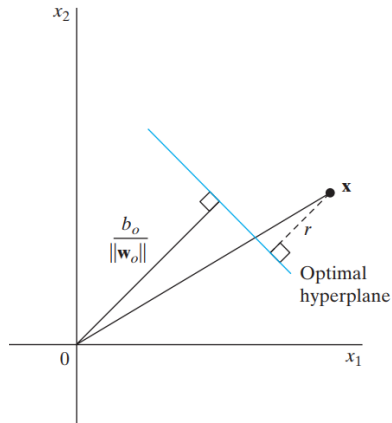
O Hiperplano Ótimo

- A função discriminante ótima é definida como:

$$g(\mathbf{x}) = \mathbf{w}_o^T \mathbf{x} + b_o \quad (3)$$

- E distância deste hiperplano por:

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}_o\|} \quad (4)$$



O Hiperplano Ótimo

- Para otimalidade, o par (\mathbf{w}_o, b_o) deve satisfazer:

$$\begin{cases} \mathbf{w}_o^T \mathbf{x}_i \geq 1, & d_i = +1 \\ \mathbf{w}_o^T \mathbf{x}_i \leq -1, & d_i = -1 \end{cases} \quad (5)$$

- Os pontos (\mathbf{x}_i, d_i) que satisfazem com igualdade estas restrições são chamados de **vetores de suporte**
- Eles são os pontos mais próximos do hiperplano e consequentemente mais difíceis de classificar [Haykin, 2007].

O Hiperplano Ótimo

- Por definição:

$$g(\mathbf{x}^{(s)}) = \mathbf{w}_o^T \mathbf{x}^{(s)} + b_o = \mp 1 \quad (6)$$

- Logo, a distância do vetor de suporte $\mathbf{x}^{(s)}$ é dada por:

$$\begin{cases} \frac{1}{\|\mathbf{w}_o\|} & \text{se } d^{(s)} = +1 \\ \frac{-1}{\|\mathbf{w}_o\|} & \text{se } d^{(s)} = -1 \end{cases} \quad (7)$$

- Finalmente, a margem ótima ρ é:

$$\rho = \frac{2}{\|\mathbf{w}_o\|} \quad (8)$$

Conclusão

Maximizar a margem de separação entre classes binárias é equivalente a minimizar a norma Euclidiana do vetor de pesos \mathbf{w}

Formulação

- Combinando as Equações (6) e (8), podemos formular o problema de encontrar a margem ótima como:

$$\begin{aligned} &\underset{\mathbf{w}, b}{\text{minimizar}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} && (9) \\ &\text{sujeito a} && d_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, && i = 1, 2, \dots, N \end{aligned}$$

- Problema com objetivo quadrático e restrições lineares
- É reformulado e resolvido através da técnica de multiplicadores de Lagrange [Haykin, 2007]

Resultado

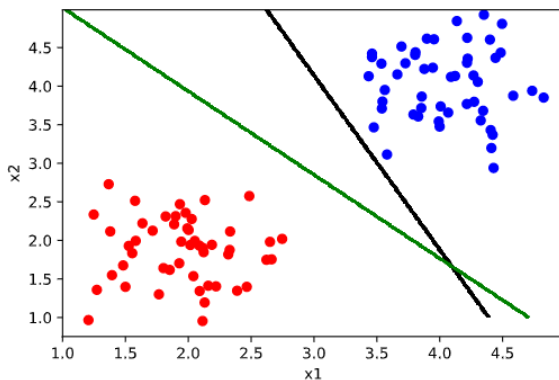


Figura: Superfícies de separação: erro mínimo (preta) e margem máxima (verde)

Padrões Não Separáveis

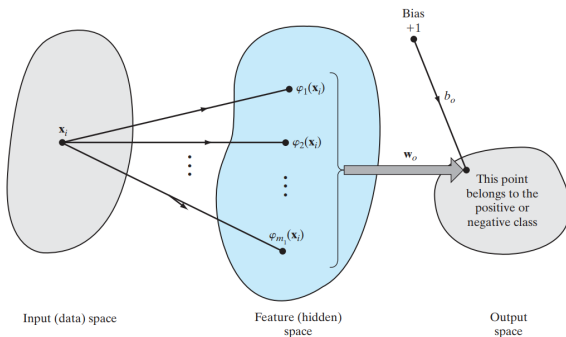
- Em aplicações práticas, não podemos garantir que os dados sejam perfeitamente separáveis
- Logo, é adicionada uma variável de folga ξ para representar estas discrepâncias:

$$\begin{aligned} &\underset{\mathbf{w}, b}{\text{minimizar}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i && (10) \\ &\text{sujeito a} && d_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i && i = 1, 2, \dots, N \\ &&& \xi_i \geq 0 \end{aligned}$$

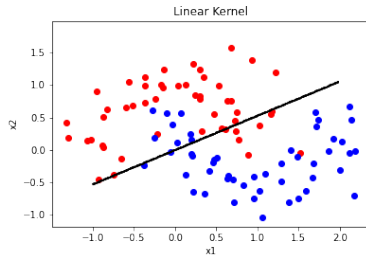
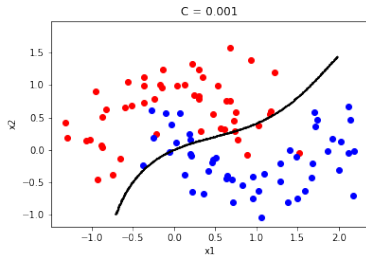
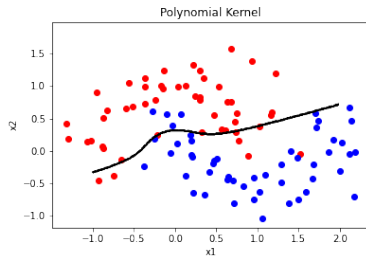
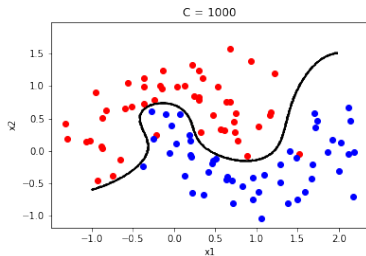
- É reformulado e resolvido de forma similar ao Problema (10)
- A constante C é especificada pelo usuário e controla a generalização do modelo, sendo conhecida como o parâmetro de regularização

Problemas não-linearmente separáveis

- Problemas reais nem sempre podem ser separados linearmente
- Entretanto, podemos mapeá-lo para um novo espaço de alta dimensão onde ele é mais provável de ser linearmente separável: **Teorema de Cover** [Cover, 1965]
- Este mapeamento é feito por funções não-lineares conhecidas como *Kernels*



Exemplo - SVM



Referências



Haykin, S. (2007). Neural Networks: A Comprehensive Foundation (3rd Edition).



Cover, T. M. (1965). Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. IEEE transactions on electronic computers, (3), 326-334.

Obrigado!