

# Inteligência Artificial

Raoni F. S. Teixeira

## Aula 11 - SVM

### 1 Introdução

Support Vector Machine (SVM) é um classificador que separa os dados encontrando o hiperplano com a maior margem possível.

Nesta aula, consideramos um problema de classificação binária com rótulos  $\{+1, -1\}$ . O classificador é definido pela equação:

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b),$$

onde  $\mathbf{w}$  é o vetor de pesos,  $\mathbf{x}$  é o vetor de entrada,  $b$  é o termo bias, e  $\text{sign}$  é uma função que devolve o sinal  $+$  ou  $-$  dos rótulos.

### 2 SVM

Existem infinitos hiperplanos capazes de separar um conjunto de dados linearmente separável. A pergunta central é: *Qual é o melhor hiperplano?*

O SVM responde a essa questão ao escolher o hiperplano que maximiza a distância para os pontos mais próximos de ambas as classes, ou seja, o **hiperplano com margem máxima**.

A Figura 1 ilustra o hiperplano de margem máxima em um conjunto de dados bidimensional. Na Figura 1a, comparamos dois hiperplanos para o mesmo conjunto de dados: o hiperplano de margem máxima (vermelho) e um hiperplano alternativo (azul). A Figura 1b destaca os pontos mais próximos ao limite de decisão, conhecidos como **vetores de suporte**, e a margem  $\gamma$ , que representa a distância entre o hiperplano e esses pontos.

#### 2.1 Margem

A margem  $\gamma$  é a distância entre o hiperplano e os pontos mais próximos das duas classes.

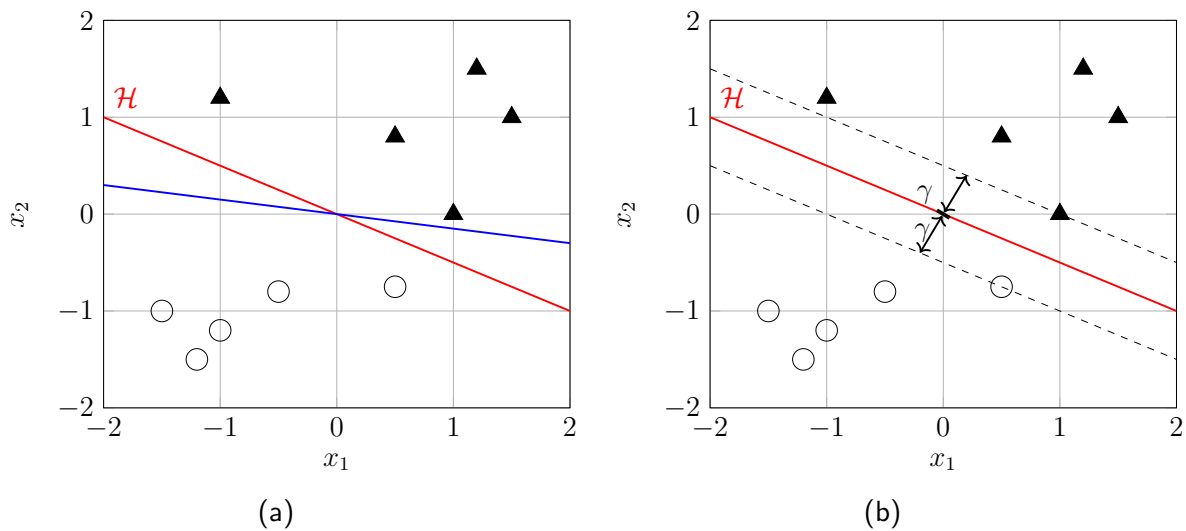


Figura 1: Hiperplano de margem máxima em um conjunto de dados bidimensional: (a) compara dois hiperplanos e (b) destaca os vetores de suporte e a margem.

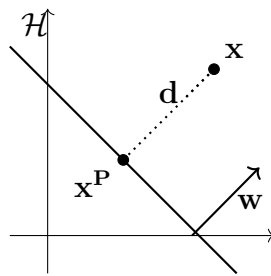


Figura 2: Distância entre um ponto  $x$  e o hiperplano  $w^T x + b$ .

### Como calcular a distância de um ponto ao hiperplano?

Considere um hiperplano definido por  $w$  e  $b$ , descrito como  $\mathcal{H} = \{x \mid w^T x + b = 0\}$ . Considere também o vetor  $d$  conectando o ponto  $x$  ao hiperplano  $\mathcal{H}$ . A distância de um ponto  $x$  até  $\mathcal{H}$  é dada por:

$$\|d\| = \frac{|w^T x + b|}{\|w\|_2}.$$

A Figura 2 ilustra esses conceitos em um exemplo bidimensional. A projeção  $x^P$  de  $x$  no hiperplano é definida por:

$$x^P = x - d.$$

Como  $d$  é paralelo a  $w$ , podemos escrever:

$$d = \alpha w.$$

Sabendo que  $\mathbf{x}^P \in \mathcal{H}$ , temos:

$$\mathbf{w}^T \mathbf{x}^P + b = 0.$$

Substituindo  $\mathbf{x}^P$  em termos de  $\mathbf{d}$ :

$$\mathbf{w}^T (\mathbf{x} - \mathbf{d}) + b = \mathbf{w}^T (\mathbf{x} - \alpha \mathbf{w}) + b = 0.$$

Logo,

$$\alpha = \frac{\mathbf{w}^T \mathbf{x} + b}{\mathbf{w}^T \mathbf{w}}.$$

O comprimento de  $\mathbf{d}$  é então:

$$\|\mathbf{d}\|_2 = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|_2}.$$

E a margem de um hiperplano  $\gamma(\mathbf{w}, b)$  é a menor distância para todo conjunto de pontos:

$$\gamma(\mathbf{w}, b) = \min_{\mathbf{x} \in D} \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|_2}.$$

Como a margem e o hiperplano são invariantes à escala, temos:

$$\gamma(\beta \mathbf{w}, \beta b) = \gamma(\mathbf{w}, b), \quad \forall \beta \neq 0.$$

Se o hiperplano maximiza  $\gamma$ , ele estará centralizado entre as duas classes, garantindo a maior separação possível entre elas.

### 3 Classificador de Margem Máxima

Podemos formular nossa busca pela margem máxima como um problema de otimização restrito. O objetivo é maximizar a margem sob as restrições de que todos os pontos de dados devem estar no lado correto do hiperplano:

$$\underbrace{\max_{\mathbf{w}, b} \gamma(\mathbf{w}, b)}_{\text{maximiza margem}} \quad \text{s.t.} \quad \underbrace{\forall i y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 0}_{\text{garante a separação}}$$

Substituindo a definição de  $\gamma$ , obtemos:

$$\underbrace{\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \min_{\mathbf{x} \in D} |\mathbf{w}^T \mathbf{x} + b|}_{\substack{\gamma(\mathbf{w}, b) \\ \text{maximiza margem}}} \quad \text{s.t.} \quad \underbrace{\forall i y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 0}_{\text{garante separação}}$$

Como o hiperplano é invariante em escala, podemos fixar a escala de  $\mathbf{w}$  e  $b$  como quisermos. Vamos ser estratégicos e escolher de modo que:

$$\min_{\mathbf{x} \in D} |\mathbf{w}^T \mathbf{x} + b| = 1.$$

Podemos adicionar esse redimensionamento como uma restrição de igualdade. Assim, nosso objetivo torna-se:

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2} \cdot 1 = \min_{\mathbf{w}, b} \|\mathbf{w}\|_2 = \min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w}.$$

Aqui, utilizamos o fato de que  $f(z) = z^2$  é uma função monotonicamente crescente para  $z \geq 0$  e  $\|\mathbf{w}\| \geq 0$ . Logo, o  $\mathbf{w}$  que maximiza  $\|\mathbf{w}\|$  também maximiza  $\mathbf{w}^T \mathbf{w}$ .

O novo problema de otimização se torna:

$$\min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} \quad \text{s.t.} \quad \forall i \, y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1.$$

Esta formulação resulta em um problema de otimização quadrática, onde a função objetivo é quadrática e as restrições são lineares. Podemos resolvê-lo de forma eficiente usando métodos apropriados. A solução busca o hiperplano mais simples (com menor  $\mathbf{w}^T \mathbf{w}$ ) que garanta uma margem de no mínimo 1 para todos os pontos.

Alguns pontos de treinamento terão restrições rígidas, o que significa que estarão exatamente na margem de 1 do hiperplano. Isso ocorre porque, se todos os pontos de treinamento tivessem desigualdades estritas, seria possível diminuir os parâmetros  $\mathbf{w}$  e  $b$  até que as restrições se tornassem rígidas, resultando em um valor objetivo ainda menor. Esse é o caso dos vetores de suporte.

Os vetores de suporte definem a margem máxima do hiperplano para o conjunto de dados, determinando assim sua forma. Se um desses pontos fosse movido e o SVM fosse retreinado, o hiperplano resultante mudaria. Em contraste, para pontos que não são vetores de suporte (desde que não sejam movidos excessivamente), o hiperplano permaneceria inalterado.

## 4 SVM com Restrições Suaves

Em dimensões baixas, frequentemente não existe um hiperplano que separe as classes. Para contornar isso, introduzimos variáveis de folga ( $\xi_i$ ), que permitem algumas violações nas margens:

$$\min_{\mathbf{w}, b, \xi} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i, \quad \text{sujeito a} \quad y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \forall i.$$

Essas variáveis penalizam as violações na função objetivo. O parâmetro  $C$  controla o rigor do modelo:

- Valores altos de  $C$  forçam o modelo a minimizar violações, resultando em margens rígidas.
- Valores baixos permitem maior flexibilidade, priorizando uma solução mais simples.

## 4.1 Formulação Irrestrita

Podemos expressar  $\xi_i$  como:

$$\xi_i = \max(1 - y_i(\mathbf{w}^T \mathbf{x}_i + b), 0).$$

Substituindo na função objetivo, obtemos a versão irrestrita:

$$\min_{\mathbf{w}, b} \underbrace{\mathbf{w}^T \mathbf{w}}_{\text{regularização L2}} + C \sum_{i=1}^n \underbrace{\max[1 - y_i(\mathbf{w}^T \mathbf{x}_i + b), 0]}_{\text{hinge-loss}}.$$

Essa é função **hinge loss** com  $\|\mathbf{w}\|^2$  atuando como regularizador. Essa formulação permite otimizar os parâmetros por métodos como descida de gradiente, similar à regressão logística, mas com uma perda distinta.

## 5 Problemas Não Linearmente Separáveis

Para dados não linearmente separáveis, o SVM utiliza *kernels* para mapear os dados para espaços de dimensão mais alta (uma estratégia semelhante ao mapeamento da aula 7), onde a separação linear se torna possível. Funções de kernel comuns incluem:

- **Linear:**  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$ ,
- **Polinomial:**  $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + c)^d$ ,
- **RBF (Radial Basis Function):**  $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$ .

Esses kernels permitem que o SVM lide com dados complexos e não linearmente separáveis, preservando o poder de generalização.

## 6 Exercícios

1. Analise o seguinte conjunto de dados bidimensional linearmente separável:
  - Classe +1: (2, 3) e (4, 5).
  - Classe -1: (1, 1) e (3, 2).

Responda:

- a) Escreva uma equação genérica para o hiperplano que separa as duas classes.
  - b) Calcule a margem do hiperplano.
  - c) Identifique os vetores de suporte.
2. Um SVM com restrições suaves foi treinado para classificar e-mails como "spam" ou "não spam". O parâmetro  $C$  influencia o desempenho do modelo.
- a) Descreva como  $C$  controla o equilíbrio entre margem rígida e margem flexível.
  - b) Explique como ajustar  $C$  nos seguintes cenários:
    - Quando a precisão é prioritária em relação à generalização.
    - Quando a generalização é mais importante que a precisão.