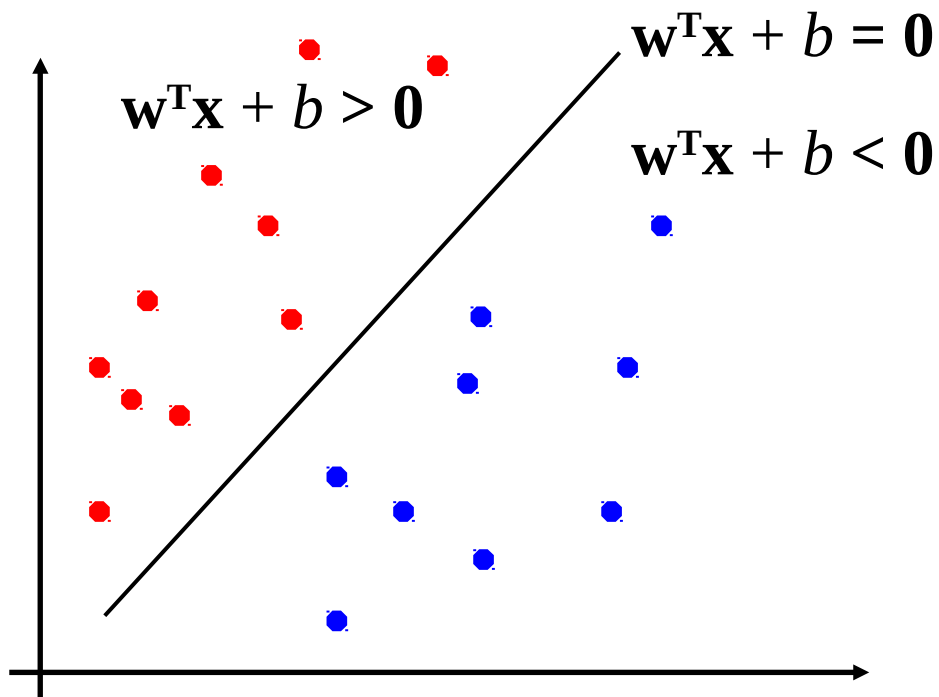


SVM

Máquinas de Vetores de Suporte

Introdução

Classificação > tarefa de separar classes no espaço de entrada



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

Introdução

Dado um conjunto de dados $D_L = \{\mathbf{x}_i, y_i\}_{i=1}^N$,

Onde $y_i \in \{-1, +1\}$

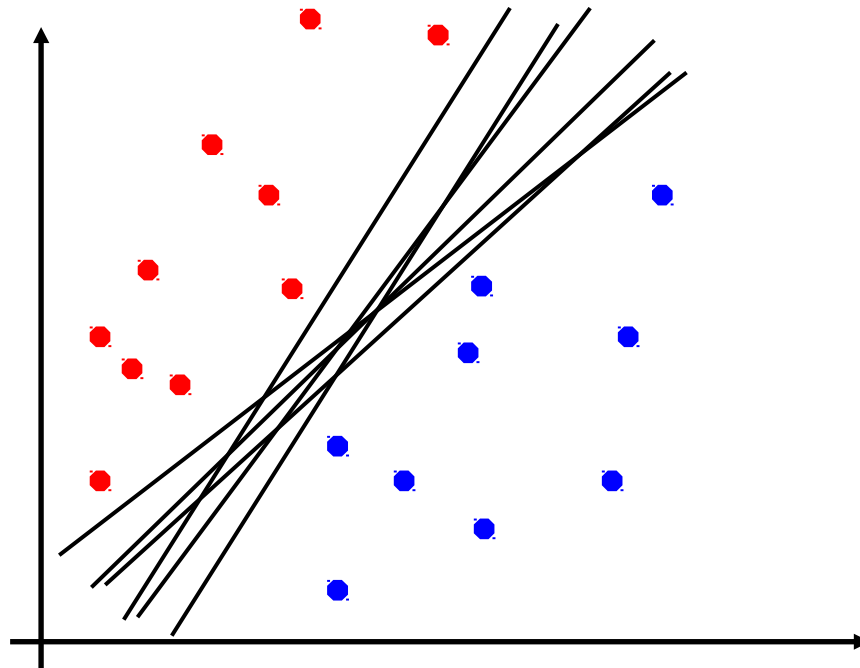
Se $\text{sign}(\mathbf{w}^T \mathbf{x} + b)$ tiver o mesmo sinal de y então a classificação estará correta

Assim, para que todos os vetores \mathbf{x}_i sejam classificados corretamente a seguinte desigualdade deve ser satisfeita para todos os N pares (\mathbf{x}_i, y_i)

$$y(\mathbf{w}^T \mathbf{x} + b) \geq 1$$

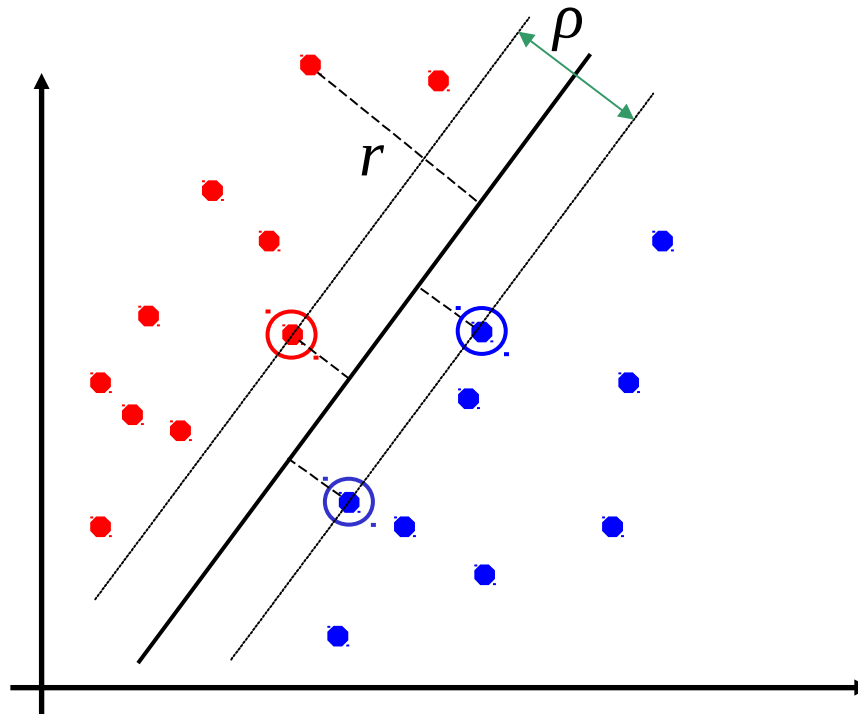
Separador linear

Mas qual separador é o melhor?



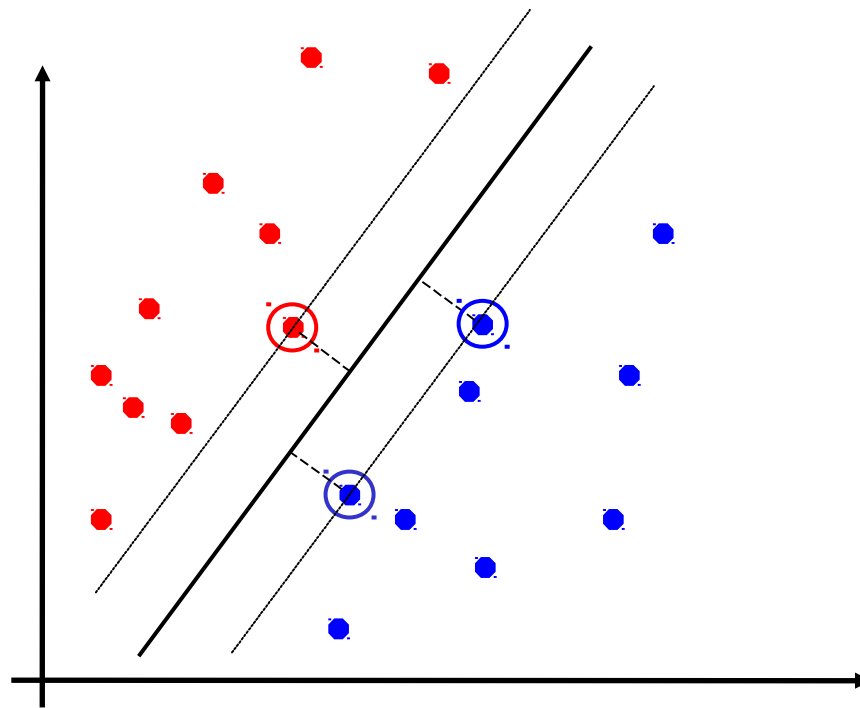
Margem de classificação

- A distância de uma amostra x_i ao separador é $r = \frac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|}$
- Amostras próximas ao hiperplano são chamados Vetores de Suporte
- ρ é a margem máxima entre vetores de suporte.



Margem larga

Apenas vetores de suporte interessam – todo o resto pode ser ignorado



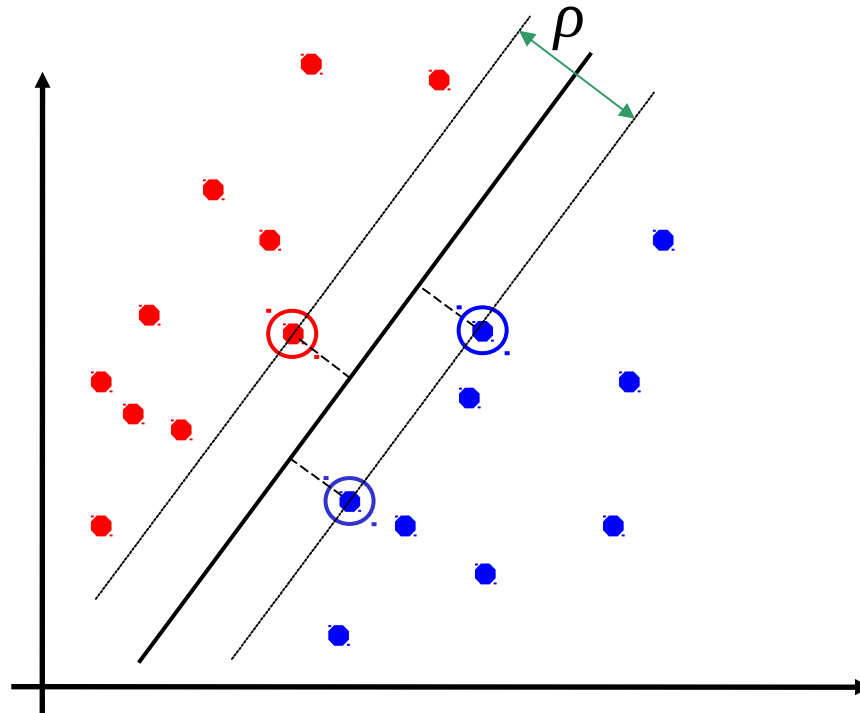
Objetivo – maximizar a margem

SVM linear

Dado um conjunto de dados $D_L = \{\mathbf{x}_i, y_i\}_{i=1}^N$,

Onde $y_i \in \{-1, +1\}$, para cada par (\mathbf{x}_i, y_i) teremos:

$$\begin{aligned} \mathbf{w}^T \mathbf{x}_i + b &\leq -\rho/2 & \text{se } y_i = -1 \\ \mathbf{w}^T \mathbf{x}_i + b &\geq \rho/2 & \text{se } y_i = 1 \end{aligned} \quad \Leftrightarrow \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq \rho/2$$



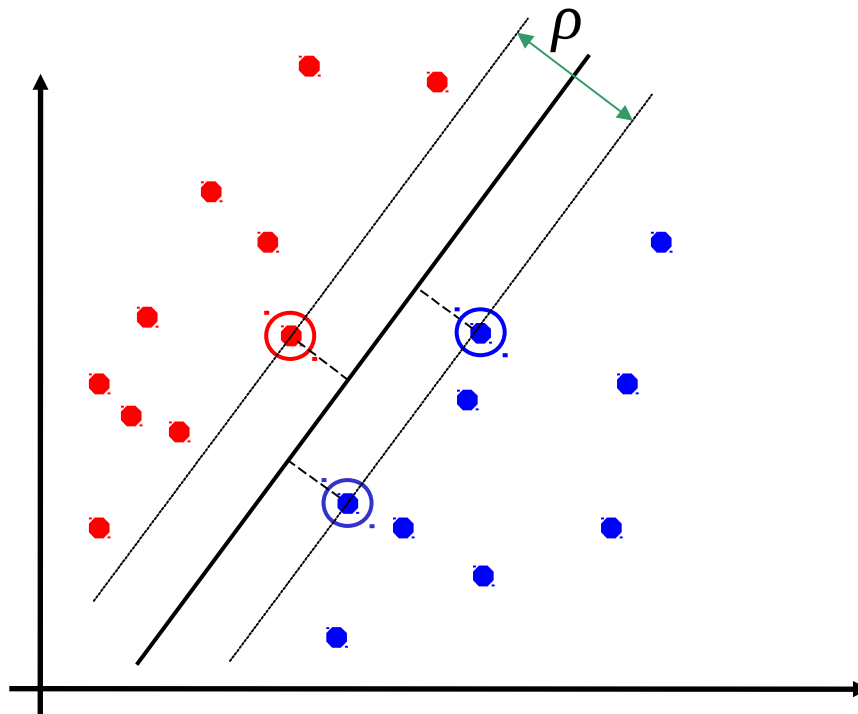
SVM linear

Para cada vetor de suporte \mathbf{x}_s a desigualdade anterior é uma igualdade.

Reescalando \mathbf{w} e b por $\rho/2$ temos que a distância entre o hiperplano e \mathbf{x}_s será:

$$r = \frac{y_s(\mathbf{w}^T \mathbf{x}_s + b)}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

Então a margem será: $\rho = 2r = \frac{2}{\|\mathbf{w}\|}$



SVM linear

Então podemos formular o seguinte problema quadrático de otimização:

Encontrar w e b tal que

$$\rho = \frac{2}{\|\mathbf{w}\|} \quad \text{é maximizado,}$$

$$\text{Para todo } (x_i, y_i), i=1, \dots, n : y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Que por sua vez pode ser reformulado como:

Encontrar w e b tal que

$$\Phi(\mathbf{w}) = \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w} \quad , \quad \text{é minimizado,}$$

$$\text{Para todo } (x_i, y_i), i=1, \dots, n : y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Resolvendo o problema de otimização

Encontrar w e b tal que

$$\Phi(\mathbf{w}) = \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w} \quad , \quad \text{é minimizado,}$$

$$\text{Para todo } (x_i, y_i), i=1, \dots, n : y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

- Necessidade de otimizar uma função quadrática sujeita a restrições lineares.
- Problemas de otimização quadrática são uma classe bem conhecida de problemas de programação matemática para os quais existem vários algoritmos (não triviais).
- A solução envolve a construção de um problema duplo, onde um multiplicador de Lagrange α_i está associado a todas as restrições de desigualdade no problema primal (original):

Encontrar $\alpha_1 \dots \alpha_n$ tal que

$$Q(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad \text{é máximo e}$$

- $\sum \alpha_i y_i = 0$
- $\alpha_i \geq 0$ para todo α_i

Resolvendo o problema de otimização

- Dada uma solução $\alpha_1 \dots \alpha_n$ para o problema duplo, a solução para o primal é:

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i \quad b = y_j - \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_j \quad \text{para todo } \alpha_j > 0$$

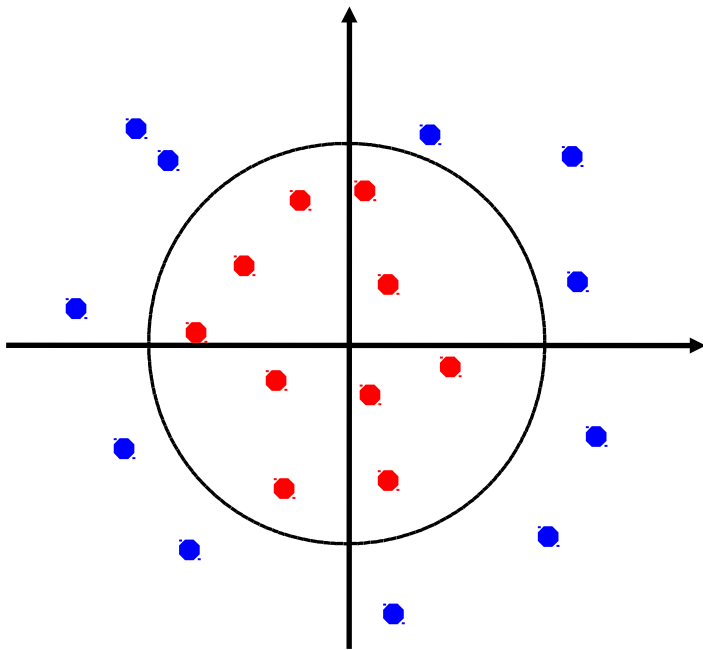
- Cada α_i diferente de zero indica que \mathbf{x}_i correspondente é um vetor de suporte. Em seguida, a função de classificação é (observe que não precisamos de \mathbf{w} explicitamente):

$$\hat{y} = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

- Observe que ele depende de um produto interno entre o ponto de teste \mathbf{x} e os vetores de suporte \mathbf{x}_i . Lembre-se também de que a solução do problema de otimização envolveu a computação dos produtos internos $\mathbf{x}_i^T \mathbf{x}_j$ entre todos os pontos de treinamento.

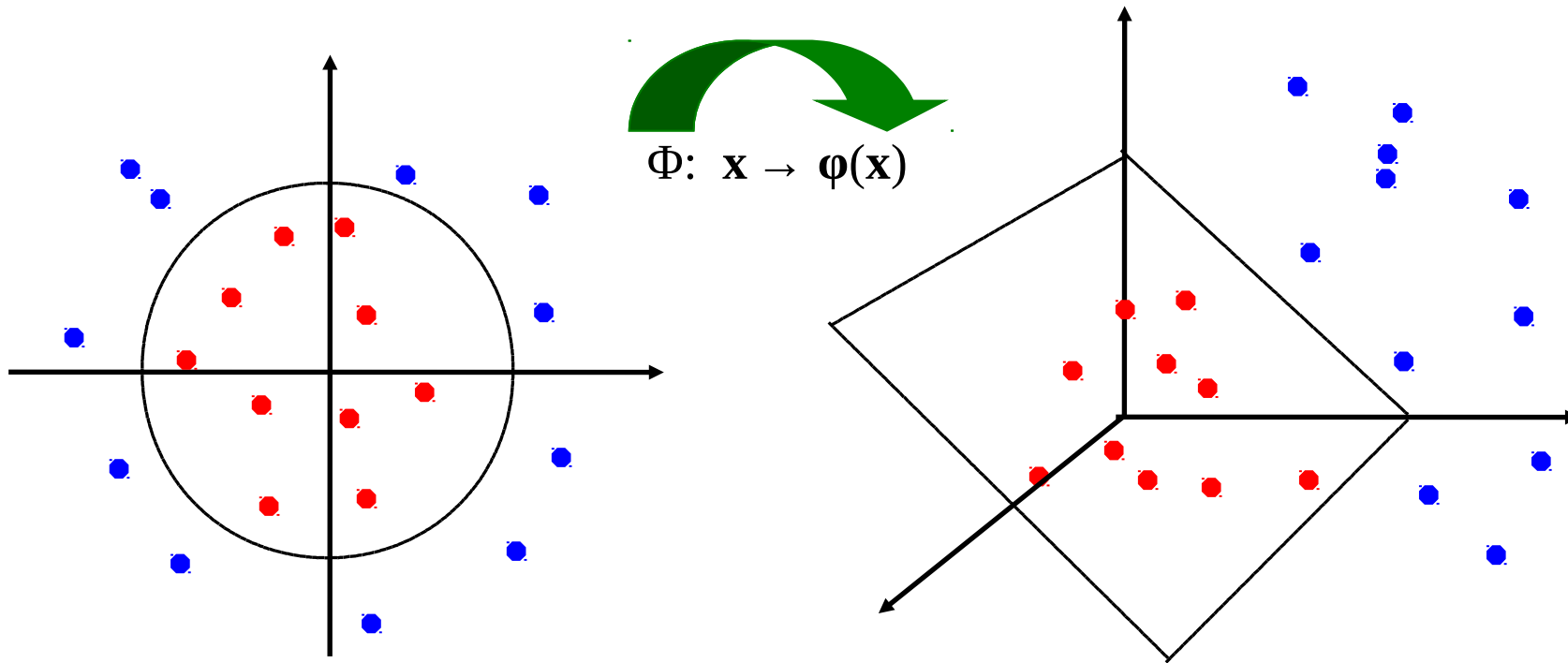
SVM não-linear

E se o problema não for linearmente separável?



SVM não-linear

As amostras X no espaço de entrada podem ser mapeadas para outro espaço onde o problema pode ser linearmente separável.



SVM não-linear

O classificador linear depende do produto interno entre os vetores

$$K(x_i, x_j) = x_i^T x_j$$

Se todo ponto de dados for mapeado no espaço de alta dimensão através de alguma transformação $\Phi: x \rightarrow \phi(x)$, o produto interno se tornará:

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

Uma função do kernel é uma função que é equivalente a um produto interno em algum espaço de variáveis

Assim, uma função do kernel mapeia implicitamente os dados para um espaço de alta dimensão (sem a necessidade de calcular cada $\phi(x)$ explicitamente).

SVM não-linear

Para algumas funções $K(x_i, x_j)$, verificar se $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ pode ser complicado.

Teorema de Mercer:

Toda função simétrica definida semi-positiva é um kernel

As funções simétricas definidas semi-positivas correspondem a uma matriz simétrica definida semi-positiva:

$K =$

$K(\mathbf{x}_1, \mathbf{x}_1)$	$K(\mathbf{x}_1, \mathbf{x}_2)$	$K(\mathbf{x}_1, \mathbf{x}_3)$	\dots	$K(\mathbf{x}_1, \mathbf{x}_n)$
$K(\mathbf{x}_2, \mathbf{x}_1)$	$K(\mathbf{x}_2, \mathbf{x}_2)$	$K(\mathbf{x}_2, \mathbf{x}_3)$		$K(\mathbf{x}_2, \mathbf{x}_n)$
\dots	\dots	\dots	\dots	\dots
$K(\mathbf{x}_n, \mathbf{x}_1)$	$K(\mathbf{x}_n, \mathbf{x}_2)$	$K(\mathbf{x}_n, \mathbf{x}_3)$	\dots	$K(\mathbf{x}_n, \mathbf{x}_n)$

SVM não-linear

Linear: $K(x_i, x_j) = x_i^T x_j$

Mapeamento $\Phi: x \rightarrow \varphi(x)$, onde $\varphi(x)$ é o próprio x

Polinômio de potência p : $K(x_i, x_j) = (1 + x_i^T x_j)^p$

Mapeamento $\Phi: x \rightarrow \varphi(x)$, onde $\varphi(x)$ tem $\binom{d+p}{p}$ dimensões

Gaussiana (função de base radial): $K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$

Mapeamento $\Phi: x \rightarrow \varphi(x)$, onde $\varphi(x)$ é de dimensão infinita: todo ponto é mapeado para uma função (um gaussiano); A combinação de funções para vetores de suporte é o separador.

SVM não-linear

Então o problema dual pode ser formulado como:

Encontrar $\alpha_1 \dots \alpha_n$ tal que

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$ é máximo e

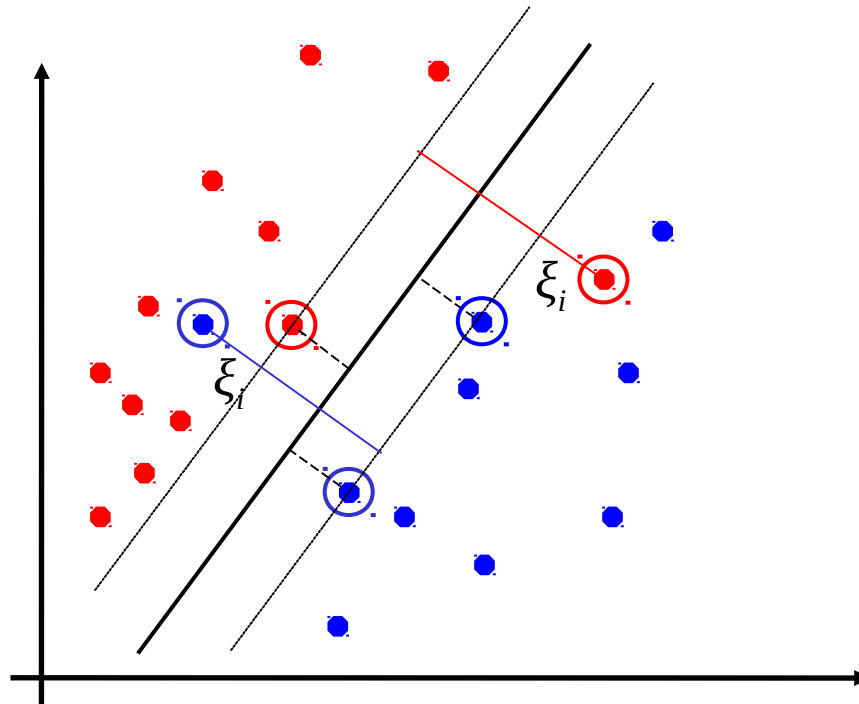
- $\sum \alpha_i y_i = 0$
- $\alpha_i \geq 0$ para todo α_i

E a solução será:

$$f(\mathbf{x}) = \sum \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_j) + b$$

SVM com variável de folga

Variáveis de folga ξ_i podem ser adicionadas para permitir a classificação incorreta de exemplos difíceis ou ruidosos, a margem resultante é denominada suave (soft margin).



SVM linear com variável de folga

O problema pode ser formulado da seguinte forma:

Encontrar \mathbf{w} e b tal que

$\Phi(\mathbf{w}) = \mathbf{w}^T \mathbf{w} + C \sum \xi_i$ é mínimo

e para todo $(\mathbf{x}_i, y_i), i=1..n$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$, $\xi_i \geq 0$

O parâmetro C pode ser visto como uma maneira de controlar o *overfitting*: “trade off” entre a importância de maximizar a margem e ajustar os dados de treinamento.

SVM linear com variável de folga

Resolvendo o problema dual teremos:

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i$$

$$b = y_k(1 - \xi_k) - \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_k \quad \text{para qualquer } k \text{ sujeito a } \alpha_k > 0$$

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

SVM não-linear com variável de folga

Resolvendo o problema dual teremos:

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i$$

$$b = y_k(1 - \xi_k) - \sum \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_j) \quad \text{para qualquer } k \text{ sujeito a } \alpha_k > 0$$

$$f(\mathbf{x}) = \sum \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_j) + b$$