

SVM: o Problema de Otimização, as Garantias de Aprendizado e o Truque do Kernel

Alek Fröhlich & Douglas Soares Gonçalves

Universidade Federal de Santa Catarina (UFSC)

alek.frohlich@posgrad.ufsc.br, douglas@mtm.ufsc.br



Resumo

Support Vector Machine (SVM) é um dos modelos mais fascinantes dentro do contexto de Aprendizagem de Máquina, usufruindo de fortes garantias teóricas e de boa performance prática. Neste trabalho, coletamos resultados teóricos instigantes à respeito do treino e do uso de SVMs.

Introdução

O SVM foi introduzido na década de 90 por Vapnik e Cortes [2] para lidar com o problema de classificação. De forma simplificada, o problema consiste em: dada uma amostra $S = (x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$ gerada de maneira i.i.d. a partir de uma distribuição \mathcal{D} sobre $\mathbb{R}^n \times \{\pm 1\}$, determinar uma função de classificação $h : \mathbb{R}^n \rightarrow \{\pm 1\}$ que minimize o erro esperado

$$R(h) = P[h(x) \neq y]. \quad (1)$$

Como não conhecemos \mathcal{D} , na prática utilizamos alguma forma de erro empírico, por exemplo, $\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{h(x_i) \neq y_i}$ e buscamos relacionar $R(h)$ e $\hat{R}_S(h)$. Na prática de aprendizagem de máquina é comum buscar por h dentro de uma classe \mathcal{H} . No caso do SVM, a classe corresponde aos classificadores lineares:

$$\mathcal{H} = \{x \mapsto \text{sign}(w^T x + b) : w \in \mathbb{R}^n, b \in \mathbb{R}\}. \quad (2)$$

Dada uma amostra S , o algoritmo de treinamento do SVM busca encontrar um hiperplano (w, b) que minimize o erro de treino e maximize o tamanho da margem. O segundo objetivo é um diferencial do método e contribui para suas garantias de aprendizado.

Treinando SVMs

O problema de treinamento pode ser expresso como o problema de otimização a seguir:

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} w^T w + C \sum_{i=1}^m \xi_i \\ \text{s.a.} \quad & y_i (w^T x_i + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0, \end{aligned} \quad (3)$$

em que ξ representa as violações da margem e C controla a penalização. O problema é um programa quadrático convexo (CQP) e sempre admite solução, a qual nem sempre é única [1]. As condições de otimalidade para o problema garantem que $w = \sum_{i=1}^m \alpha_i y_i x_i$ e $\alpha^T y = 0$, isso somado à dualidade forte nos permite expressar o problema na sua forma dual:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i^T x_j) - \sum_{i=1}^m \alpha_i \\ \text{s.a.} \quad & 0 \leq \alpha \leq C, \\ & \alpha^T y = 0. \end{aligned} \quad (4)$$

O dual também é um CQP com solução (domínio compacto). Essa é a versão tipicamente resolvida, uma vez que permite a troca de $x_i^T x_j$ por $K(x_i, x_j)$, em que K é um kernel. Os métodos comumente utilizados para resolver (4) são baseados no algoritmo SMO de Platt [6]. O algoritmo resolve uma sequência de subproblemas, nos quais apenas duas variáveis são otimizadas por vez (α_i, α_j) . Essa abordagem escala melhor do que métodos de programação quadrática usuais, uma vez que o número de variáveis cresce linearmente com o tamanho do conjunto de treino. Por fim, a corretude do algoritmo é garantida pelo Teorema de Osuna [5].

Sobre a generalização das SVMs

A noção da margem ρ viabiliza cotas para $R(h)$ que independem da dimensão do espaço de entrada \mathcal{X} , sem as quais o Truque do Kernel não seria possível.

Definição 1. Fixe $\rho > 0$, definimos a função de perda L_ρ como sendo $L_\rho(y, y') = \Phi_\rho(y y')$, em que

$$\Phi_\rho(x) = \begin{cases} 1, & \text{se } x \leq 0 \\ 1 - \frac{x}{\rho}, & \text{se } 0 \leq x \leq \rho \\ 0, & \text{se } \rho \leq x, \end{cases}$$

e denotamos por $\hat{R}_{S, \rho}(h)$ o erro empírico de h sobre uma amostra S com respeito à L_ρ .

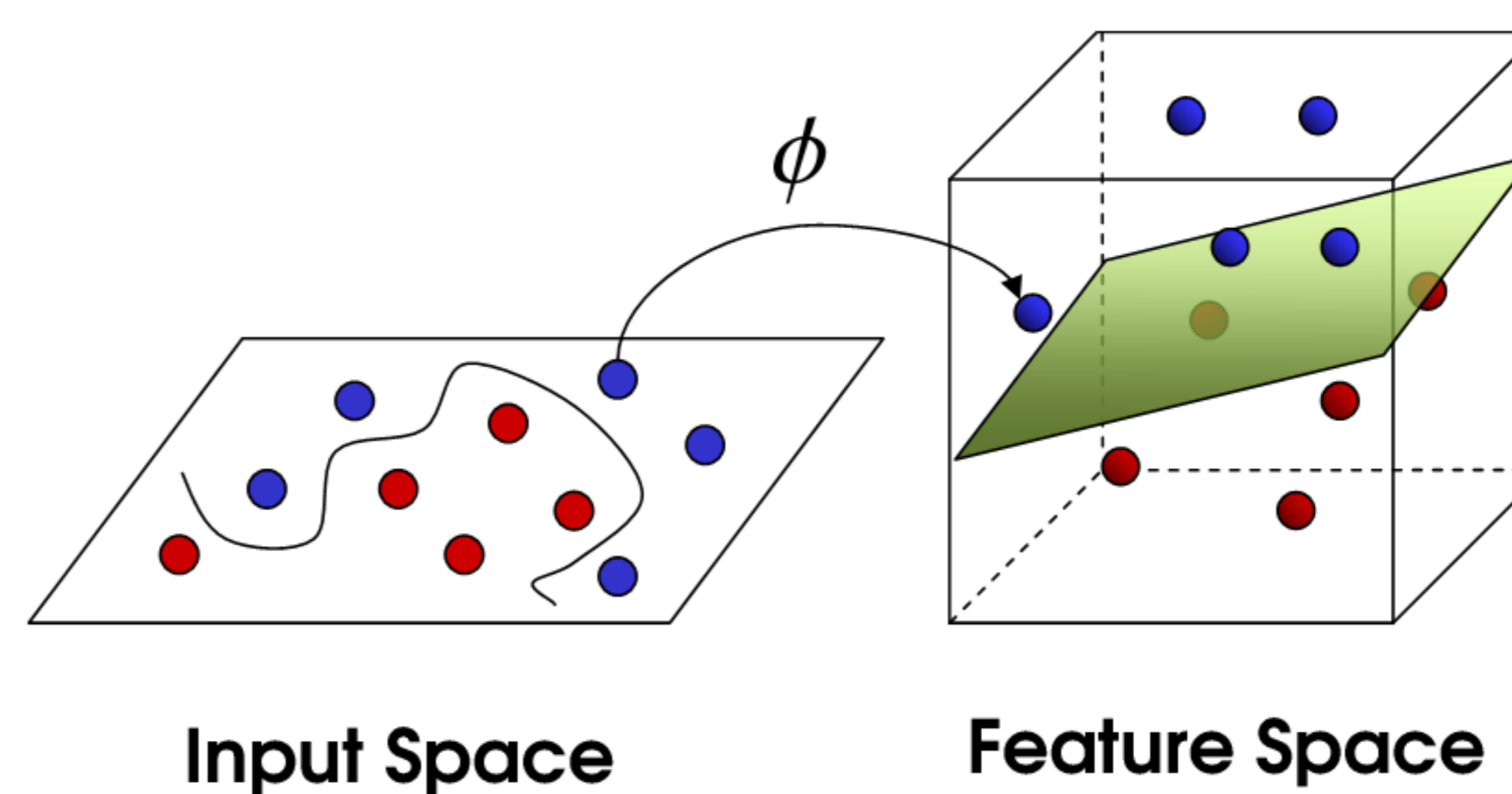
Utilizando cotas de generalização baseadas em Complexidade de Rademacher e uma estimativa para a complexidade da classe $\mathcal{H} = \{x \mapsto w^T x : \|w\| \leq \Lambda\}$, conseguimos estabelecer o seguinte resultado:

Teorema 1. Seja $\mathcal{H} = \{x \mapsto w^T x : \|w\| \leq \Lambda\}$ e assumamos que $\mathcal{X} \subset \{x : \|x\| \leq R\}$. Fixe $\rho > 0$, então, para $\delta > 0$, com probabilidade ao menos $1 - \delta$, temos que o seguinte vale para toda $h \in \mathcal{H}$:

$$R(h) \leq \hat{R}_{S, \rho}(h) + 2\sqrt{\frac{\Lambda^2 R^2}{m \rho^2}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$$

O Truque do Kernel

Embora o SVM clássico seja um modelo linear, é possível adaptá-lo para contextos não lineares através da introdução de núcleos positivo-definidos (kernels). Essa conexão se dá através do Truque do Kernel, que corresponde a implicitamente mapear o conjunto de treino a um espaço de Hilbert de dimensão maior (potencialmente infinita, no caso do kernel RBF) e então aplicar o algoritmo linear sob os dados representados nesse outro espaço.



A principal motivação para o truque vem do fato que é mais fácil separar pontos em \mathbb{R}^N do que em \mathbb{R}^n para $N \gg n$ [3]. O seguinte teorema dá uma caracterização simples para os kernels:

Teorema 2. Seja $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ uma função simétrica e positiva definida, então existe um espaço de Hilbert \mathcal{H} e uma função $\phi : \mathcal{X} \rightarrow \mathcal{H}$ tal que $K(x, y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}}$.

Referências

- [1] C. Burges and D. Crisp. Uniqueness of the svm solution. In *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999.
- [2] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20:273–297, 1995.
- [3] T. Cover. Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. *IEEE Transactions on Electronic Computers*, EC-14(3):326–334, June 1965.
- [4] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of Machine Learning*. MIT Press, 2nd edition, 2018.
- [5] E. Osuna, R. Freund, and F. Girosi. An improved training algorithm for support vector machines. *Neural Networks for Signal Processing VII. Proceedings of the 1997 IEEE Signal Processing Society Workshop*, pages 276–285, 1997.
- [6] J. Platt. Sequential minimal optimization: A fast algorithm for training support vector machines. Technical Report MSR-TR-98-14, Microsoft, April 1998.
- [7] R. Rifkin, M. Pontil, and A. Verri. A note on support vector machine degeneracy. Technical report, MIT Artificial Intelligence Laboratory, 1999.