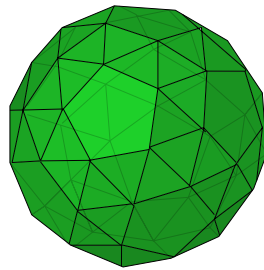


Projects in Mathematics and Applications

SUPPORT VECTOR MACHINE

Phạm Quốc Việt ^{*}, Nguyễn Hà Minh [†]
Nguyễn Đình Bảo Ngọc [‡], Nguyễn Huy Anh [§]

Ngày 27 tháng 8 năm 2018



* Trường THPT Chuyên Amsterdam
† Trường THPT Nguyễn Bình Khiêm
‡ Trường Phổ Thông Năng Khiếu - ĐHQG TP.HCM
§ Campbell County Comprehensive High School

Lời cảm ơn

Chúng em xin chân thành cảm ơn Ban tổ chức PiMA vì đã mang đến một trại hè Toán học vô cùng thú vị và đầy ý nghĩa cho chúng em cũng như là các bạn trại sinh còn lại. Không những chúng em được tìm hiểu sâu và rộng hơn về Toán, chúng em còn được làm quen và học hỏi từ các bạn có cùng chung đam mê dành cho môn Toán từ khắp nơi trên mọi miền đất nước. Tham gia PiMA, chúng em được giới thiệu với các mảng Toán ở Đại học, được học lập trình, nghiên cứu và viết báo cáo, và trên hơn cả là rèn luyện kỹ năng làm việc nhóm. Chúng em rất trân trọng trải nghiệm 10 ngày vừa qua và hy vọng rằng PiMA sẽ tiếp tục truyền cảm hứng và nuôi giữ ngọn lửa đam mê cho các bạn học sinh tham gia như PiMA đã thực hiện thành công trong 3 năm vừa qua.

Chúng em cũng xin gửi lời cảm ơn đến Ban Giám Đốc Đại Học Khoa Học Tự Nhiên TP.HCM và các nhà tài trợ đã tạo điều kiện để trại hè PiMA 2018 được diễn ra thành công.

Chúng em xin cảm ơn đến các thầy diễn giả: GS.TS Hồ Tú Bảo, PGS.TS Đinh Điền, TS. Nguyễn Đức Hoàng Hạ, TS. Nguyễn Minh Triết, PGS.TS Lý Quốc Ngọc. Qua những lời chia sẻ nhiệt huyết từ các thầy, chúng em và các bạn trại sinh còn lại đã được mở mang và tiếp xúc với nhiều khía cạnh khác nhau trong môn Toán, đặc biệt là đối với mảng máy học.

Chúng em xin chân thành gửi lời cảm ơn đến sự hỗ trợ của anh Vũ Lê Thế Anh, là người đã trực tiếp hướng dẫn tụi em trong quá trình tìm hiểu và chuẩn bị bài báo cáo này, cùng sự góp ý của anh Dương Đình Trọng cho bài báo cáo này được hoàn hảo hơn.

Và xin cảm ơn đến các bạn trại sinh còn lại vì đã đồng hành cùng chúng mình trong 10 ngày vừa qua.

Tóm tắt nội dung

Trong khuôn khổ bài báo cáo này, chúng em sẽ trình bày về thuật toán Support Vector Machine và ứng dụng của nó trong các bài toán, đặc biệt là trong hệ thống phân lớp nhị phân. Bài báo cáo được xây dựng trên cơ sở toán học, Soft-margin và Hard-margin SVM, cách tối ưu hệ thống, và mô phỏng thuật toán này trên mô hình dự đoán giao thức tín dụng.

Mục lục

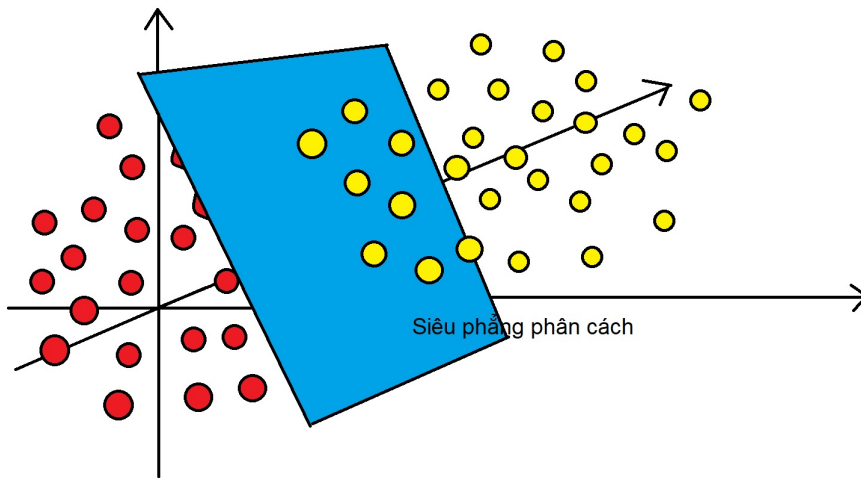
1	Tổng quan	1
2	Một số định nghĩa và tính chất	1
3	Phát biểu bài toán	3
4	Các phương pháp giải bài toán SVM	4
4.1	Tập dữ liệu có thể phân cách tuyến tính: Hard-margin SVM	4
4.2	Trường hợp tập dữ liệu không thể phân cách tuyến tính: Soft-margin SVM và Kernel Trick	6
5	Áp dụng vào mô hình thực tế	12
5.1	Mô tả dữ liệu	12
5.2	Các thuật ngữ quan trọng	12
5.3	Tham số C và kernel trong thuật toán SVM	13
5.4	Đánh giá mô hình thực tế	14
6	Phụ lục	15

1 Tổng quan

Thuật toán Support Vector Machine chính thức được giới thiệu bởi Bernhard E. Boser, Isabelle M. Guyon và Vladimir N. Vapnik vào năm 1992. Thuật toán Support Vector Machine bắt đầu trở nên nổi tiếng khi được áp dụng vào kĩ thuật nhận dạng con số viết bằng tay. Ngày nay, Support Vector Machine có vai trò quan trọng trong máy học và được biết đến như là một ví dụ điển hình của phương pháp kernel.

Support Vector Machine giúp hỗ trợ giải quyết một trong những nhiệm vụ phổ biến của máy học: phân cách dữ liệu. Cụ thể hơn, Support Vector Machine giúp ta phân loại một điểm dữ liệu mới vào một trong hai lớp dữ liệu đã được hình thành từ những điểm dữ liệu cũ bằng cách gắn với mỗi điểm dữ liệu đó một vector p chiều và đi tìm một siêu phẳng $p - 1$ chiều để phân cách các vector trên sao cho khoảng cách từ lề đến siêu phẳng đó là lớn nhất.

Hình sau là một ví dụ của siêu phẳng phân cách trong không gian ba chiều.



2 Một số định nghĩa và tính chất

Định nghĩa 2.1. Một tập hợp A được gọi là lồi nếu với hai phần tử x_i và x_j bất kì thuộc A , phần tử $x_k = \lambda x_i + (1 - \lambda)x_j$ cũng thuộc $A, \forall \lambda \in [0, 1]$.

Tính chất 2.2.

- Một siêu phẳng là một tập lồi.
- Giao của các tập lồi là một tập lồi.

Định nghĩa 2.3. Một hàm số $f : \mathbb{R}^n \rightarrow \mathbb{R}$ là hàm lồi nếu tập xác định (dom) D là một tập lồi và:

$$f(\lambda x_i + (1 - \lambda)x_j) \leq \lambda f(x_i) + (1 - \lambda)f(x_j), \forall x_i, x_j \in D, 0 \leq \lambda \leq 1.$$

Tính chất 2.4. Hàm số $f : \mathbb{R}^n \rightarrow \mathbb{R}$ là một hàm lõm nếu $-f$ là một hàm lồi.

Tính chất 2.5. Hàm norm là một hàm lồi.

Định nghĩa 2.6. Một hàm số $f : \mathbb{R}^n \rightarrow \mathbb{R}$ là hàm lồi chặt nếu f là hàm lồi và dấu bằng không xảy ra $\forall x_i \neq x_j$.

Định nghĩa 2.7. Một bài toán tối ưu lồi là một bài toán tối ưu có dạng:

$$x = \underset{x}{\operatorname{argmin}} f_0, \tag{1}$$

với

$$\begin{aligned} f_i(x) &\leq 0 && \forall i = 1, 2, \dots, n, \\ h_j(x) &= a_j^T x - b_j = 0 && \forall j = 1, 2, \dots, m, \\ x &\in \mathcal{D} = \left(\underset{i=0}{\bigcap} \operatorname{dom} f_i \right) \cap \left(\underset{j=1}{\bigcap} \operatorname{dom} h_j \right). \end{aligned}$$

Trong đó, f_i là các hàm lồi $\forall i = 0, 1, \dots, n$.

Định nghĩa 2.8. Gọi x_0 là nghiệm toàn cục của bài toán tối ưu lồi (1) nếu:

$$x_0 = \underset{x}{\operatorname{argmin}} f_0,$$

$$\text{và } f_i(x_0) \leq 0 \quad \forall i = 1, 2, \dots, n, \quad h_j(x_0) = 0, \quad \forall j = 1, 2, \dots, m.$$

Định nghĩa 2.9. Gọi x_0 là nghiệm cục bộ lân cận " của bài toán tối ưu lồi (1) nếu:

$$x_0 = \underset{x}{\operatorname{argmin}} f_0,$$

$$\text{và } f_i(x_0) \leq 0 \quad \forall i = 1, 2, \dots, n, \quad h_j(x_0) = 0 \quad \forall j = 1, \dots, m, \quad \|x - x_0\|_2 \leq \epsilon.$$

Tính chất 2.10. Nghiệm cục bộ lân cận của một bài toán tối ưu lồi cũng chính là nghiệm toàn cục của bài toán đó.

Chứng minh bằng phản chứng: Giả sử x_0 là một nghiệm cục bộ lân cận " nhưng không là cực trị toàn cục. Gọi y_0 là cực trị toàn cục, tức là $f(y_0) \leq f(x_0)$, $\|y_0 - x_0\|_2 > \epsilon$.

Mặt khác, luôn tồn tại z_0 sao cho:

$$\begin{aligned} \exists z_0 &= y_0 + (1 - \epsilon) x_0 \text{ với } \epsilon \in [0, 1] \text{ đủ nhỏ,} \\ \|z_0 - x_0\|_2 &< \epsilon, \\ z_0 &\text{ thỏa mãn các hàm điều kiện ràng buộc.} \end{aligned}$$

Khi đó:

$$\begin{aligned} f_0(z_0) &= f_0(y_0 + (1 - \epsilon)x_0) \\ &\leq f_0(y_0) + (1 - \epsilon)f_0(x_0) \\ &< f_0(x_0) + (1 - \epsilon)f_0(x_0) \\ &= f_0(x_0). \end{aligned}$$

Điều này mâu thuẫn do x_0 là nghiệm cục bộ lân cận ".

Định nghĩa 2.11. Hàm Lagrangian của bài toán tối ưu lồi (1) là hàm:

$$\mathcal{L}(x, \lambda, \mu) = f_0(x) + \sum_{i=1}^n \lambda_i f_i(x) + \sum_{j=1}^m \mu_j h_j(x),$$

với,

$$\begin{aligned} x &\in \mathcal{D}, \\ \lambda &= [\lambda_1, \lambda_2, \dots, \lambda_n]^T, \quad \lambda_i \geq 0, \\ \mu &= [\mu_1, \mu_2, \dots, \mu_m]^T. \end{aligned}$$

Định nghĩa 2.12. Hàm đối ngẫu của bài toán tối ưu lồi (1) với (λ, μ) là hàm:

$$g(\lambda, \mu) = \inf_{x \in D} \mathcal{L}(x, \lambda, \mu).$$

Nhận xét 2.13. Hàm đối ngẫu trên là một hàm lồi.

Tính chất 2.14. $g(\lambda, \mu) \leq f_0(x_0)$ với x_0 là nghiệm toàn cục của bài toán tối ưu lồi (1).

Định nghĩa 2.15. Bài toán đối ngẫu Lagrange của bài toán (1) là:

$$(\lambda_0, \mu_0) = \operatorname{argmax} g(\lambda, \mu).$$

Định nghĩa 2.16. Đối ngẫu chặt xảy ra khi $g(\lambda_0, \mu_0) = f_0(x_0)$.

Nhận xét 2.17. Khi đối ngẫu chặt xảy ra, nghiệm của bài toán đối ngẫu chính là nghiệm của bài toán gốc.

Định lý 2.18 (Tiêu chuẩn Slater). Khi bài toán gốc là lồi, và tồn tại $x \in \mathcal{D}$ sao cho: $f_i(x) < 0 \forall i = 1, 2, \dots, n$, $h_j(x) = 0 \forall j = 1, 2, \dots, m$ thì đối ngẫu chặt xảy ra.

Định lý 2.19 (Hệ điều kiện KKT). Điều kiện cần và đủ KKT cho nghiệm (x_0, λ_0, μ_0) bài toán đối ngẫu là:

$$\begin{aligned} f_i(x_0) &\leq 0, \forall i = 1, 2, \dots, n, \\ h_j(x_0) &= 0, \forall j = 1, 2, \dots, m, \\ \lambda_i^0 &\geq 0, \forall i = 1, 2, \dots, n, \\ \lambda_i^0 f_i(x_0) &= 0, \forall i = 1, 2, \dots, n, \\ \frac{\partial f_0}{\partial x_0} + \sum_{i=1}^n \lambda_i^0 \frac{\partial f_i}{\partial x_0} + \sum_{j=1}^m \mu_j^0 \frac{\partial h_j}{\partial x_0} &= 0. \end{aligned}$$

3 Phát biểu bài toán

Xét tập dữ liệu $D = \{(x_i, y_i); x_i \in \mathbb{R}^m; y_i \in \{1, -1\} \mid i = 1, 2, \dots, n\}$ với n là số cặp dữ liệu và m là số chiều của dữ liệu (số lượng đặc trưng). Trong đó:

$$y_i = \begin{cases} 1 & \text{nếu cặp dữ liệu thứ } i \text{ thuộc nhóm 1,} \\ -1 & \text{nếu cặp dữ liệu thứ } i \text{ thuộc nhóm 2.} \end{cases}$$

Ta cần xác định siêu phẳng $w^T x + b = 0$ phân chia 2 nhóm sao cho lề là lớn nhất, với định nghĩa lề trong bài toán SVM là giá trị của:

$$\min_{i=1,2,\dots,n} \frac{|w^T x_i + b|}{\|w\|_2}$$

Bài toán trên tương đương với bài toán tối ưu sau:

$$(w, b) = \operatorname{argmax}_{w, b} \min_i \frac{|w^T x_i + b|}{\|w\|_2} \quad (2)$$

Nhận xét 3.1. Ta luôn chọn $\|w\|_2 = 1$ sao cho y_i và $(w^T x_i + b)$ cùng dấu vì nếu không chọn $\|w\|_2 = 1$ thì $y_i(w^T x_i + b) < 0, \forall i = 1, 2, \dots, n$.

Khi đó, (2) tương đương với tìm:

$$(w, b) = \operatorname{argmax}_{w, b} \min_i \frac{y_i(w^T x_i + b)}{\|w\|_2} \quad (3)$$

4 Cực phân hoạch giải bài toán SVM

4.1 Tập dữ liệu cá thể phân cách tuyến tính: Hard-margin SVM

4.1.1 Bài toán tối ưu

Ta cần: tìm cặp (w, b) sao cho $w^T x + b = 0$ phân chia hai tập dữ liệu S_1, S_2 không giao nhau. Gọi (x_0, y_0) là:

$$(x_0, y_0) = \underset{(x_i, y_i)}{\operatorname{argmin}} \frac{y_i (w^T x_i + b)}{\|w\|_2}$$

Nhận xét 4.1. Nếu thay $(w, b) = (\|w\|_2, kb)$ thì $\frac{y_i (w^T x_i + b)}{\|w\|_2}$ không đổi, $i = 1, 2, \dots, n$. Để đơn giản, ta đặt $\|w\|_2 = 1$, thì $y_i (w^T x_i + b) = 1$.

Khi đó (3) trở thành:

$$(w, b) = \underset{w, b}{\operatorname{argmax}} \frac{1}{\|w\|_2}, \quad (4)$$

$$\text{với } y_i (w^T x_i) - y_0 (w^T x_0 + b) = 1, \quad i = 1, 2, \dots, n.$$

(4) tương đương với:

$$(w, b) = \underset{w, b}{\operatorname{argmax}} \frac{1}{2} \|w\|_2^2, \quad (5)$$

$$\text{với } 1 - y_i (w^T x_i) \leq 0, \quad i = 1, 2, \dots, n.$$

Như vậy, sau khi xác định các điều kiện phân hoạch $w^T x + b = 0$, cần tìm $w \in \mathbb{R}^m$ sao cho tồn tại w sao cho $w^T x_i + b = 0$.

4.1.2 Bài toán tối ưu Lagrange

Kiểm tra điều kiện Slater: Với hai tập dữ liệu S_1, S_2 phân cách tuyến tính, tồn tại bài toán tối ưu cá thể, nên tìm cặp (w_0, b_0) sao cho:

$$\begin{cases} 1 - y_i (w_0^T x_i + b_0) < 0, & i = 1, 2, \dots, n, \\ 1 - y_i (2w_0^T x_i + 2b_0) < 0, & i = 1, 2, \dots, n. \end{cases}$$

Chọn $w_1 = 2w_0$ và $b_1 = 2b_0$, suy ra:

$$1 - y_i (w_1^T x_i + b_1) < 0, \quad i = 1, 2, \dots, n$$

Vậy điều kiện Slater thỏa mãn.

Hàm Lagrangian của bài toán (5) là:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i (w^T x_i + b))$$

$$\text{với } \alpha = [\alpha_1, \alpha_2, \dots, \alpha_n]^T \text{ và } \alpha_i \geq 0, \quad i = 1, 2, \dots, n.$$

Hàm Lagrange

$$g(\lambda) = \min_{w, b} L(w, b, \lambda) \quad \forall \lambda \geq 0.$$

Xét $(w_0, b_0, \lambda_0) = \operatorname{argmin}_{w, b} L(w, b, \lambda)$. Khi đó:

$$\frac{\partial L}{\partial w_0} = w_0 - \sum_{i=1}^n \lambda_i y_i x_i = 0 \quad w_0 = \sum_{i=1}^n \lambda_i y_i x_i.$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^n \lambda_i y_i = 0.$$

Thay vào biểu thức của $L(w, b, \lambda)$, ta được:

$$g(\lambda) = \sum_{n=1}^N \lambda_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \lambda_n \lambda_m y_n y_m x_n^T x_m \quad (6)$$

Tới đây, kết hợp hàm Lagrange và các điều kiện ràng buộc của bài toán Lagrange của bài toán (26.3) có dạng

$$= \operatorname{argmax}_{\lambda} g(\lambda)$$

với:

$$\lambda_n \geq 0$$

$$\sum_{n=1}^N \lambda_n y_n = 0$$

Với mọi bài toán lồi và điều kiện chặt yếu, điều kiện KKT là điều kiện cần và đủ cho bộ (w_0, b_0, λ_0) :

$$\lambda_i (1 - y_i (w_0^T x_i + b_0)) = 0, \quad \forall i = 1, 2, \dots, n \quad (7)$$

$$w_0 = \sum_{i=1}^n \lambda_i y_i x_i$$

$$\sum_{n=1}^N \lambda_n y_n = 0 \quad (8)$$

Từ điều kiện (7) ta thấy hoặc $\lambda_i = 0$ hoặc $y_i (w_0^T x_i + b_0) = 1$. Khi đó, hai đẳng thức $w_0^T x_i + b = 1$ chia tách họ n thành hai lớp dữ liệu.

Tới đây ta cần hình ảnh sau:

Hình ảnh 4.2. Hai đẳng thức của bài toán SVM là hai đẳng thức $w_0^T x_i + b = 1$.

Nhận xét 4.3. Số lượng tham số của hai siêu phẳng là hai hình ảnh của số biến đầu vào với tổng số tham số của tập huấn luyện.

ót $S = f_n: n \in \{0, 1, \dots, N_S\}$ là phần tử của tập S
 Theo (8), ta có:

$$w = \sum_{m \in S} a_m y_m X_m,$$

$$b = \frac{1}{N_S} \sum_{n \in S} y_n \sum_{m \in S} a_m y_m X_m^T X_n$$

Vì vậy để tính w và b , ta có thể dùng bảng ngôn ngữ lập trình Python.

4.2 Trùng hợp đặc biệt của hàm mất mát: Soft-margin SVM và Kernel Trick

4.2.1 Bài toán tối ưu: Soft-margin SVM

Trong thực tế, đặc biệt là trong huấn luyện ta cần tìm hàm mất mát không chỉ phụ thuộc vào khoảng cách mà còn phụ thuộc vào độ rộng của khe hở. Khi đó, chúng ta cần sử dụng hàm mất mát Soft-margin SVM.

Với suy nghĩ như vậy, chúng ta cần tìm hàm mất mát Soft-margin SVM. Lúc này, độ rộng của khe hở là biến số cần tối ưu. Chúng ta cần tìm hàm mất mát sao cho độ rộng của khe hở là lớn nhất. Điều này có thể thực hiện bằng cách tối ưu hàm mất mát Soft-margin SVM.

Để tìm hàm mất mát Soft-margin SVM, chúng ta cần tìm hàm mất mát sao cho độ rộng của khe hở là lớn nhất. Điều này có thể thực hiện bằng cách tối ưu hàm mất mát Soft-margin SVM.

Hình ảnh 4.4. Các biến sai của bài toán SVM:

$$\xi_i = \begin{cases} 0 & \text{nếu } y_i \text{ nằm đúng trong biên của } \hat{S} \\ |w^T x_i + b - y_i| & \text{nếu } y_i \text{ không nằm đúng trong biên của } \hat{S} \end{cases}$$

Như hình 4.5. Ta cần tìm hàm mất mát sao cho tổng các biến sai là nhỏ nhất.

Với Soft-margin SVM, hàm mất mát là tổng của các biến sai. Khi đó, hàm mất mát là:

$$\frac{1}{2} \|w\|_2^2 + C \sum_{n=1}^N \xi_n,$$

với C là một hằng số dương.

Như hình 4.5, với $y_i = 1, 2, \dots, n$, ta có:

$$\begin{aligned} \xi_i &= w^T x_i + b - y_i \\ &= w^T x_i + b - y_i |y_i| \\ &= y_i (w^T x_i + b) - 1 \\ &= 1 - y_i (w^T x_i + b) \end{aligned}$$

Để tìm hàm mất mát, ta cần tìm các biến sai sau, với $i = 1, 2, \dots, n$:

$$\xi_i = \max(0, 1 - y_i (w^T x_i + b))$$

Như vậy, bài toán tối ưu gốc cho Soft-margin SVM như sau, $\delta_i = 1, 2, \dots, n$:

$$(w_0, b_0, \alpha_0) = \underset{w, b, \alpha}{\operatorname{argmin}} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \alpha_i, \quad (9)$$

$$\text{với: } 1 - \alpha_i \leq y_i(w^T x_i + b) \leq 1 + \alpha_i, \quad \alpha_i \geq 0.$$

4.2.2 Hàm giá trị quy hoạch thủ công: Bài toán tối ưu Lagrange của Soft-margin SVM

Tương tự bài toán Hard-margin SVM, ta quan tâm đến bài toán tối ưu của nó.

Kiểm tra điều kiện Slater: $\delta_i = 1, 2, \dots, n$ và $\delta(w, b)$, luôn tồn tại các α_i đủ lớn sao cho:

$$y_i(w^T x_i + b) + \alpha_i > 1.$$

Vậy, bài toán tối ưu này thỏa mãn điều kiện Slater.

Hàm Lagrangian của bài toán Soft-margin SVM là:

$$L(w, b, \alpha, \beta) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \alpha_i + \sum_{i=1}^n \alpha_i (1 - \alpha_i - y_i(w^T x_i + b)) + \sum_{i=1}^n \beta_i \alpha_i, \quad (10)$$

$$\text{với } \alpha = [\alpha_1, \alpha_2, \dots, \alpha_n]^T \geq 0 \text{ và } \beta = [\beta_1, \beta_2, \dots, \beta_n]^T \geq 0.$$

Bài toán tối ưu của Soft-margin SVM là:

$$g(\alpha, \beta) = \min_{w, b} L(w, b, \alpha, \beta)$$

Giá trị bài toán tối ưu, với mọi cặp (α, β) , xấp xỉ (w_0, b_0, α_0) thỏa mãn:

$$\frac{\partial L}{\partial w_0} = 0, \quad w_0 = \sum_{i=1}^n \alpha_i y_i x_i \quad (11)$$

$$\frac{\partial L}{\partial b} = 0, \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad (12)$$

$$\frac{\partial L}{\partial \alpha_0} = 0, \quad \alpha_i = C - \beta_i, \quad \delta_i = 1, 2, \dots, n \quad (13)$$

Từ (13) suy ra các quan hệ liên quan tới các cặp (α, β) thỏa mãn $\alpha_i = C - \beta_i$.

Tuyệt vời ta công suy ra $0 \leq \beta_i \leq C, \delta_i = 1, 2, \dots, n$.

Thay vào biểu thức Lagrangian, ta thu được hàm mục tiêu của bài toán tối ưu:

$$g(\alpha, \beta) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j. \quad (14)$$

Nhận xét 4.6. Hàm này không phụ thuộc vào w, b từ (13) suy ra $0 \leq \beta_i \leq C, \delta_i = 1, 2, \dots, n$

Khi α, b bài toán tối ưu trở thành:

$$\alpha_0 = \operatorname{argmax} g(\alpha)$$

$$\text{với: } \sum_{i=1}^n \alpha_i y_i = 0, \quad (15)$$

$$0 \leq \beta_i \leq C, \delta_i = 1, 2, \dots, n. \quad (16)$$

Hãy tìm các điều kiện KKT của bài toán tối ưu Soft-margin SVM, $\delta_i = 1, 2, \dots, n$:

$$1 - \lambda_i^0 (y_i(w_0^T x_i + b_0) - 1) = 0, \quad (17)$$

$$\lambda_i^0 = 0, \quad (18)$$

$$\lambda_i^0 = 0, \quad (19)$$

$$\lambda_i^0 = 0, \quad (20)$$

$$\lambda_i^0 (1 - \lambda_i^0 (y_i(w_0^T x_i + b_0) - 1)) = 0, \quad (21)$$

$$\lambda_i^0 \lambda_i^0 = 0, \quad (22)$$

Ngoài ra các điều kiện (11), (12), (13).

Khi $\lambda_i^0 > 0$, theo (21):

$$y_i(w_0^T x_i + b_0) = 1 - \lambda_i^0. \quad (23)$$

Nếu $0 < \lambda_i^0 < C$, theo (13), $\lambda_i^0 = C - \lambda_i^0 > 0$; ngoài ra theo (22) ta thu được $\lambda_i^0 = 0$. Tiếp theo (23), suy ra $y_i(w_0^T x_i + b_0) = 1$, hay:

$$w_0^T x_i + b_0 = y_i, \delta_i : 0 < \lambda_i^0 < C.$$

Tương tự, nếu $0 < \lambda_i^0 < C$, các điểm x_i nằm trên biên của S . Khi đó, (w_0, b_0) được tìm thấy theo:

$$w_0 = \sum_{i=1}^n \lambda_i^0 y_i x_i, \quad (24)$$

$$b_0 = \frac{1}{N_S} \sum_{i \in S} \lambda_i^0 (y_i - w_0^T x_i)$$

với $S = \{m : 0 < \lambda_m < C\}$ và N_S là phần tử của S .

Khi đó, hàm của x được tìm thấy:

$$w_0^T x + b_0 = \sum_{i \in M} \lambda_i^0 y_i x_i^T x + \frac{1}{N_S} \sum_{j \in S} \lambda_j^0 y_j \sum_{i \in M} \lambda_i^0 y_i x_j^T x \quad (25)$$

Với $M = \{n : 0 < \lambda_n < C\}$ và N_M là phần tử của M .

4.2.3 Hàm giải quy, thuật toán 2: Bài toán tối ưu không ràng buộc cho soft-margin SVM

Bài toán tối ưu không ràng buộc tương ứng: Tìm các điều kiện ràng buộc của:

$$1 - \lambda_i (y_i(w^T x_i + b) - 1) = 0, \quad \lambda_i (1 - y_i(w^T x_i + b)) = 0$$

Một khi, $\lambda_i = 0$ nên ta có bài toán không ràng buộc tương ứng với bài toán (9) như sau:

$$(w_0, b_0, \lambda_0) = \arg \min_{w, b, \lambda} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \lambda_i \quad (26)$$

$$\text{thỏa mãn: } \lambda_i = \max(0, 1 - y_i(w^T x_i + b)), \delta_i = 1, 2, \dots, n$$

Mệnh đề 4.7. Nếu (w_0, b_0, λ_0) là nghiệm của bài toán tối ưu (26), thì:

$$\lambda_i^0 = \max(0, 1 - y_i(w_0^T x_i + b_0)), \delta_i = 1, 2, \dots, n \quad (27)$$

Chứng minh: Giải sử tồn tại sao cho

$$i > \max(0, 1 - y_i(w_0^T x_i + b_0)),$$

Khi đó, hàm $J_i = \max(0, 1 - y_i(w_0^T x_i + b_0))$, ta được một giá trị nhỏ hơn của hàm mục tiêu, trong khi tất cả các ràng buộc vẫn được thỏa mãn (không vi phạm).

Vậy mệnh đề được chứng minh.

Để tìm giá trị tối ưu ta trở về với:

$$(w_0, b_0) = \arg \min_{w, b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b)) \quad (28)$$

thỏa mãn: $J_i = \max(0, 1 - y_i(w^T x_i + b)), \forall i = 1, 2, \dots, n,$

hay:

$$(w, b) = \arg \min_{w, b} \left(\frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b)) \right) \quad (29)$$

Nhận xét 4.8. Bài toán này có thể giải được bằng các phương pháp gradient descent.

Hình ảnh 4.9. Hàm Hinge loss cá nhân:

$$J_i(w, b) = \max(0, 1 - y_i(w^T x_i + b))$$

Ứng dụng hàm mất mát

Hình ảnh 4.10. Hàm Hinge loss trung bình cá nhân:

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n J_i = \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b))$$

Nhận xét 4.11. Xét (w_0, b_0) là nghiệm bài toán. Khi tập dữ liệu cá nhân tùy ý, hàm mất mát của $J(w, b) = 0$. Khi đó, nếu $C > 1$ thì hàm sẽ, (w_0, b_0) cũng là nghiệm bài toán đối với hàm mất mát tổng cộng.

Vậy, ta có thể hình ảnh sau:

Hình ảnh 4.12. Hàm mất mát tổng cộng cá nhân:

$$J(w, b) = L(w, b) + R(w, b),$$

$$\forall C > 0, R(w, b) = \frac{1}{2} \|w\|_2^2.$$

Nhận xét 4.13. Hàm mất mát tổng cộng là một hàm lồi theo (w, b) .

Để tìm hàm mất mát: Với việc tìm hàm mất mát dựa trên gradient descent, việc chọn hàm mục tiêu là một hàm mất mát theo w và b .

Nhận xét 4.14. Hàm mục tiêu tổng cộng không phụ thuộc.

Vậy, ta dùng Gradient Descent theo w và b để tìm hàm mất mát.

Phần tính toán cơ sở, sẽ được trình bày trong code Python. Để tìm các phần hinge loss phụ thuộc vào:

$$\frac{\partial \max(0, 1 - y_i(w^T x_i + b))}{\partial w} = \begin{cases} -y_i x_i & \text{nếu } 1 - y_i(w^T x_i + b) > 0 \\ 0 & \text{nếu } 1 - y_i(w^T x_i + b) < 0 \end{cases}$$

$$\frac{\partial \max(0, 1 - y_i(w^T x_i + b))}{\partial b} = \begin{cases} -y_i & \text{nếu } 1 - y_i(w^T x_i + b) > 0 \\ 0 & \text{nếu } 1 - y_i(w^T x_i + b) < 0 \end{cases}$$

Phần regularization công thức để tìm giá trị:

$$\frac{\partial \frac{1}{2} \|w\|_2^2}{\partial w} = w$$

$$\frac{\partial \frac{1}{2} \|w\|_2^2}{\partial b} = 0$$

Nếu cập nhật bằng gradient descent thông qua hàm mất mát $\ell(w, b)$ (stochastic gradient descent). Nếu $1 - y_i(w^T x_i + b) < 0$, ta không cập nhật w và chuyển sang hàm mất mát theo. Ngược lại khi $1 - y_i(w^T x_i + b) > 0$, ta cập nhật w và b như sau:

$$w \leftarrow w + \eta y_i x_i; \quad b \leftarrow b + \eta y_i$$

$$w \leftarrow w; \quad b \leftarrow b - \eta$$

với η là learning rate

4.2.4 Kernel trick

Để tiếp tục thu được các Support Vector Machine trên các dữ liệu không phân biệt tuyến tính, ngược lại ta cần biến các dữ liệu này lên không gian nhiều chiều sao cho phân biệt tuyến tính hoặc gần như phân biệt tuyến tính trên không gian đó. Để không gian này, biến đổi các dữ liệu thành giá trị quy chuẩn bằng hard/soft-margin SVM. Dựa trên suy nghĩ này, ta sẽ tìm một hàm $\phi(x)$ nhúng biến đổi dữ liệu vào không gian chi phối. Như vậy, việc tìm kiếm các hàm Kernel để mô tả mối quan hệ giữa hai dữ liệu bất kỳ trong không gian mới, thay vì tìm kiếm trực tiếp biến đổi của chúng.

Vì vậy:

Để tổng cho hàm Kernel $X^T h m$ sẽ () sao cho sau khi biến đổi sang không gian mới, mỗi x trở thành (x) , và khi áp dụng dữ liệu mới trở nên gần như bất tuy nhiên. Khi áp dụng bài toán điều kiện của Soft-margin SVM trở thành:

$$= \operatorname{argmax}_{\alpha} \sum_{i=1}^n \alpha_i \sum_{j=1}^n \alpha_j y_i y_j (x_i)^T (x_j),$$

$$\text{thỏa mãn: } \sum_{i=1}^n \alpha_i y_i = 0, \tag{30}$$

$$0 \leq \alpha_i \leq C, \forall i = 1, 2, \dots, n,$$

và hàm của một điểm dữ liệu mới x chính là độ của bài toán:

$$w^T (x) + b = \sum_{j \in S} \alpha_j y_j (x_j)^T (x) + \frac{1}{N_M} \sum_{i \in M} \alpha_i y_i \sum_{j \in S} \alpha_j y_j (x_j)^T (x_i). \tag{31}$$

Nhận xét 4.15. Trong (30) và (31), ta chỉ cần tính $(x_i)^T (x_j)$ với hai điểm dữ liệu x_i, x_j bất kỳ. Vì thế, ta chỉ cần xác định một hàm $k(x_i, x_j) = (x_i)^T (x_j)$.

Khi đó, (30) và (31) trở thành:

$$= \operatorname{argmax}_{\alpha} \sum_{i=1}^n \alpha_i \sum_{j=1}^n \alpha_j y_i y_j k(x_i, x_j)$$

$$\text{thỏa mãn: } \sum_{i=1}^n \alpha_i y_i = 0, \tag{32}$$

$$0 \leq \alpha_i \leq C, \forall i = 1, 2, \dots, n,$$

và hàm của một điểm dữ liệu mới x chính là độ của bài toán:

$$w^T (x) + b = \sum_{j \in S} \alpha_j y_j k(x_i, x_j) + \frac{1}{N_M} \sum_{i \in M} \alpha_i y_i \sum_{j \in S} \alpha_j y_j k(x_i, x_j). \tag{33}$$

- Tính chất của hàm Kernel:
- + Tính chất đối xứng.
- + Thỏa mãn điều kiện Mercer:

$$\sum_{i=1}^n \sum_{j=1}^n k(x_i, x_j) c_i c_j \geq 0, \forall c_k \in \mathbb{R}, k = 1, 2, \dots, n. \tag{34}$$

Nhận xét 4.16. Điều kiện Mercer giúp tìm kiếm bài toán điều kiện (33) dễ dàng.

Một số hàm Kernel thông dụng

Hàm linear (tuy nhiên): $k(x, z) = x^T z$.

Hàm RBF (radial basic function/Gaussian kernel) $k(x, z) = \exp(-\gamma \|x - z\|_2^2), \gamma > 0$

Hàm sigmoid: $k(x, z) = \tanh(x^T z + r)$

Hàm polynomial (đa thức): $k(x, z) = (r + x^T z)^d$

5 Áp dụng vào mô hình thực tế

Bài toán Support Vector Machine trên có thể áp dụng vào mô hình dự đoán lừa đảo qua thẻ tín dụng. Ta sẽ áp dụng thuật toán phân lớp nhị phân vào mô hình này. Mục đích của thuật toán là cho biết giao dịch là lừa đảo (**positive class**) hay tin cậy (**negative class**).

Phần code chi tiết thuật toán bằng Python được đề cập ở mục Phụ lục.

5.1 Mô tả dữ liệu

Dữ liệu đầu vào gồm các giao dịch qua thẻ tín dụng thu thập vào 09/2013 tại Châu Âu bao gồm:

- Có 492 giao dịch lừa đảo trong tổng số 284,807 (chiếm 0.172%).
- 28 đặc trưng ($V = v_1, v_2, \dots, v_{28}$) đã được giảm chiều dữ liệu qua phương thức phân tích thành phần chính (Principal Component Analysis).
- 2 đặc trưng mô tả thời gian và số tiền chuyển giữa hai giao dịch kế nhau.
- Lớp của dữ liệu: là lừa đảo hay tin cậy.

Do dữ liệu về giao dịch lừa đảo quá ít so với tổng số dữ liệu, nên việc dự đoán các giao dịch lừa đảo không được chú trọng. Để giải quyết vấn đề này, ta chỉ giữ lại một số dữ liệu của các giao thức tin cậy sao cho tập dữ liệu cân bằng. Kỹ thuật này gọi là **Undersampling**.

Việc mô phỏng thuật toán SVM trên tập dữ liệu này được thực hiện trên môi trường Python với sự hỗ trợ của thư viện **scikit-learn**.

5.2 Các thuật ngữ quan trọng

Để đánh giá một hệ thống phân lớp, người ta thường phân tập dữ liệu thành 3 phần chính:

- Tập **training**: dùng để huấn luyện mô hình.
- Tập **validation**: tối ưu các tham số khác cho mô hình, sử dụng w được huấn luyện từ tập **training**.
- Tập **test**: dùng để đánh giá thuật toán sau khi đã chọn ra tham số mong muốn.

Việc này nhằm mục đích chọn các tham số tối ưu và đánh giá hệ thống trên một tập dữ liệu chưa thấy (**unseen data**).

Định nghĩa 5.1. Độ chính xác là tỉ lệ giữa số điểm dữ liệu đoán đúng trên tổng số dữ liệu.

Nhận xét 5.2. *Ta thấy rằng vì các biến dữ liệu phân lớp không cân bằng (dữ liệu phân loại có chênh lệch, xác suất phân loại không bằng nhau), chính xác không phải là thang đo phù hợp để đánh giá hiệu suất (do sự thiên lệch dữ liệu dẫn đến kết quả sai).*

Để khắc phục điều này, người ta đưa ra hai thang đo sau.

Định nghĩa 5.3. Precision là tỉ lệ số điểm được dự đoán đúng **positive class** trên tổng số các điểm được dự đoán **positive class**.

Định nghĩa 5.4. Recall là tỉ lệ số điểm được dự đoán đúng **positive class** trên tổng số các điểm thật sự thuộc **positive class**.

Nhận xét 5.5. Precision cao nghĩa vì ví c chính xác c a các i m c d oán positive class là cao. Recall cao nghĩa vì ví c t l b sót các i m th c s thu c positive class là th p.

5.3 Tham số C và kernel trong thuật toán SVM

Tham số C là tham số quyết định việc tổng quát hóa của thuật toán với dữ liệu. Ở góc nhìn khác, tham số C xác định việc xem xét các dữ liệu nhiễu gần siêu phẳng phân cách của thuật toán có đáng kể hay không.

Nhận xét 5.6. Nếu C có giá trị lớn, thuật toán Soft-margin SVM trở thành Hard-margin SVM.

Thư viện **scikit-learn** cho phép ta cung cấp cụ thể các kernel khác nhau nhằm mục đích chọn ra kernel thích hợp để tối ưu hóa (ở bài toán này là ưu tiên tối đa giá trị recall để dự đoán không bỏ sót các giao dịch lừa đảo).

Quá trình chọn hai tham số trên được thực hiện lần lượt như sau:

- Thực hiện thuật toán trên tập **training** để chọn ra tham số w .
- Lần lượt thử w với các tham số C và kernel khác nhau trên tập **validation** và chọn bộ tham số tối ưu nhất.

Ở mô hình này ta chọn $C \in \{0.01, 1, 100\}$ và $\text{kernel} \in \{\text{'linear'}, \text{'rbf'}, \text{'sigmoid'}\}$. Các kết quả đánh giá với các tham số khác nhau:

Kernel = 'linear'	0.01	1	100
A.Precision	0.983	0.985	0.981
A.Accuracy	0.898	0.898	0.901
A.Recall	0.81	0.808	0.819
Standard deviation of recall	0.0475	0.0413	0.049

Kernel = 'rbf'	0.01	1	100
A.Precision	0.326	0.61	0.607
A.Accuracy	0.483	0.598	0.606
A.Recall	0.666	0.577	0.625
Standard deviation of recall	0.471	0.0623	0.0327

Kernel = 'sigmoid'	0.01	1	100
A.Precision	0.326	0.493	0.462
A.Accuracy	0.483	0.483	0.488
A.Recall	0.666	0.666	0.67
Standard deviation of recall	0.471	0.468	0.465

Sau quá trình thử ta chọn được $C = 100$ với kernel là hàm tuyến tính (linear).

Nhận xét 5.7. Ta thấy rằng ví c dùng Hard-margin SVM là hợp lí, vì ví c d oán sót các giao d ch l a o nghiêm tr ng h n ví c d oán nh m các giao d ch tin c y. Lí gi i cho ví c hàm a th c và hàm sigmoid không th c hi n hi u qu b ng hàm tuy n tính, có th là do tình hu ng **Overfitting** trên t p training.

5.4 Đánh giá mô hình thực tế

Thuật toán dự đoán các giao thức lừa đảo đạt được các kết quả sau trên tập **test**:

- Precision đạt được 0.88.
- Recall đạt được 1.00.

	Precision	Recall
Giao dịch lừa đảo	0.88	1.00
Giao dịch tin cậy	1.00	0.85
Độ chính xác	0.929	

