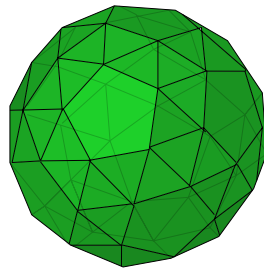


Projects in Mathematics and Applications

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

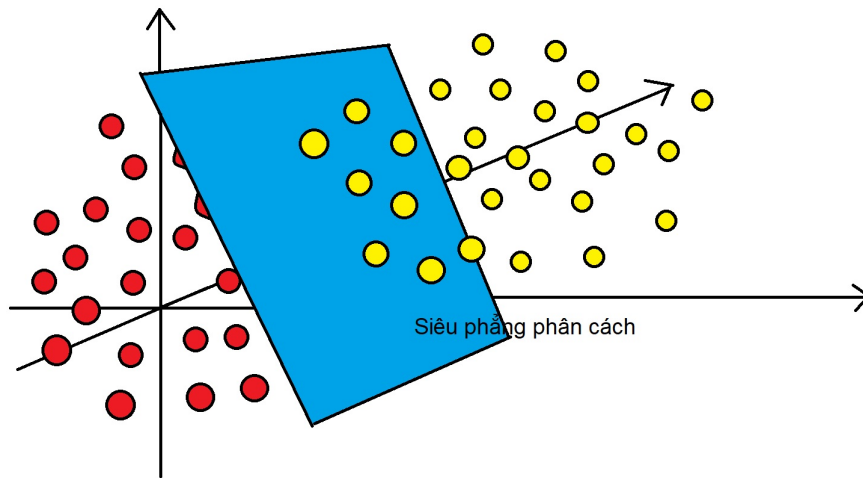
1	Tổng quan	1
2	Một số định nghĩa và tính chất	1
3	Phát biểu bài toán	3
4	Các phương pháp giải bài toán SVM	4
4.1	Tập dữ liệu có thể phân cách tuyến tính: Hard-margin SVM	4
4.2	Trường hợp tập dữ liệu không thể phân cách tuyến tính: Soft-margin SVM và Kernel Trick	6
5	Áp dụng vào mô hình thực tế	12
5.1	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 $x_k = \epsilon x_i + (1 - \epsilon)x_j$ cũng thuộc A , $\forall \epsilon \in [0, 1]$.

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:

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

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:

$$\mathbf{x} = \underset{\mathbf{x}}{\operatorname{argmin}} f_0, \quad (1)$$

với

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

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

Định nghĩa 2.8. Gọi \mathbf{x}_0 là nghiệm toàn cục của bài toán tối ưu lồi (1) nếu:

$$\mathbf{x}_0 = \underset{\mathbf{x}}{\operatorname{argmin}} f_0,$$

$$\text{và } f_i(\mathbf{x}_0) \leq 0 \quad \forall i = 1, 2, \dots, n, \quad h_j(\mathbf{x}_0) = 0, \quad \forall j = 1, 2, \dots, m.$$

Định nghĩa 2.9. Gọi \mathbf{x}_0 là nghiệm cục bộ lân cận ε của bài toán tối ưu lồi (1) nếu:

$$\mathbf{x}_0 = \underset{\mathbf{x}}{\operatorname{argmin}} f_0,$$

$$\text{và } f_i(\mathbf{x}_0) \leq 0 \quad \forall i = 1, 2, \dots, n, \quad h_j(\mathbf{x}_0) = 0 \quad \forall j = 1, \dots, m, \quad \|\mathbf{x} - \mathbf{x}_0\|_2 \leq \varepsilon.$$

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 \mathbf{x}_0 là một nghiệm cục bộ lân cận ε nhưng không là cực trị toàn cục. Gọi \mathbf{y}_0 là cực trị toàn cục, tức là $f(\mathbf{y}_0) \leq f(\mathbf{x}_0)$, $\|\mathbf{y}_0 - \mathbf{x}_0\|_2 > \varepsilon$. Mặt khác, luôn tồn tại \mathbf{z}_0 sao cho:

$$\begin{cases} \mathbf{z}_0 = \theta \mathbf{y}_0 + (1 - \theta) \mathbf{x}_0 \text{ với } \theta \in [0, 1] \text{ đủ nhỏ,} \\ \|\mathbf{z}_0 - \mathbf{x}_0\|_2 < \varepsilon, \\ \mathbf{z}_0 \text{ thỏa mãn các hàm điều kiện ràng buộc.} \end{cases}$$

Khi đó:

$$\begin{aligned} f_0(\mathbf{z}_0) &= f_0(\theta \mathbf{y}_0 + (1 - \theta) \mathbf{x}_0) \\ &\leq \theta f_0(\mathbf{y}_0) + (1 - \theta) f_0(\mathbf{x}_0) \\ &< \theta f_0(\mathbf{x}_0) + (1 - \theta) f_0(\mathbf{x}_0) \\ &= f_0(\mathbf{x}_0). \end{aligned}$$

Điều này mâu thuẫn do \mathbf{x}_0 là nghiệm cục bộ lân cận ε .

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

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

với,

$$\begin{aligned} \mathbf{x} &\in \mathcal{D}, \\ \boldsymbol{\lambda} &= [\lambda_1, \lambda_2, \dots, \lambda_n]^T, \quad \boldsymbol{\lambda} \geq 0, \\ \boldsymbol{\nu} &= [\nu_1, \nu_2, \dots, \nu_m]^T. \end{aligned}$$

Định nghĩa 2.12. Hàm đối ngẫu của bài toán tối ưu lồi (1) với $(\boldsymbol{\lambda}, \boldsymbol{\nu})$ là hàm:

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

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

Tính chất 2.14. $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq f_0(\mathbf{x}_0)$ với \mathbf{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:

$$(\boldsymbol{\lambda}_0, \boldsymbol{\nu}_0) = \operatorname{argmax}_{\boldsymbol{\lambda}, \boldsymbol{\nu}} g(\boldsymbol{\lambda}, \boldsymbol{\nu}).$$

Định nghĩa 2.16. Đối ngẫu chặt xảy ra khi $g(\boldsymbol{\lambda}_0, \boldsymbol{\nu}_0) = f_0(\mathbf{x}_0)$.

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

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

Định lý 2.19 (Hệ điều kiện KKT). Điều kiện cần và đủ KKT cho nghiệm $(\mathbf{x}_0, \boldsymbol{\lambda}_0, \boldsymbol{\nu}_0)$ bài toán đối ngẫu là:

$$\begin{aligned} f_i(\mathbf{x}_0) &\leq 0, \forall i = 1, 2, \dots, n, \\ h_j(\mathbf{x}_0) &= 0, \forall j = 1, 2, \dots, m, \\ \lambda_i^0 &\geq 0, \forall i = 1, 2, \dots, n, \\ \lambda_i^0 f_i(\mathbf{x}_0) &= 0, \forall i = 1, 2, \dots, n, \\ \frac{\partial f_0}{\partial \mathbf{x}_0} + \sum_{i=1}^n \lambda_i^0 \frac{\partial f_i}{\partial \mathbf{x}_0} + \sum_{j=1}^m \nu_j^0 \frac{\partial h_j}{\partial \mathbf{x}_0} &= 0. \end{aligned}$$

3 Phát biểu bài toán

Xét tập dữ liệu $D = \{(\mathbf{x}_i, y_i); \mathbf{x}_i \in \mathbb{R}^m; y_i \in \{1, -1\} \mid i = 1, 2, \dots, n\}$ 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 $\mathbf{w}^T \mathbf{x} + b = 0$ phân chia 2 nhóm sao cho lề là lớn nhất, với định nghĩa lề trong bài toán SVM là giá trị của:

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

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

$$(\mathbf{w}, b) = \operatorname{argmax}_{\mathbf{w}, b} \left\{ \min_i \frac{|\mathbf{w}^T \mathbf{x}_i + b|}{\|\mathbf{w}\|_2} \right\} \quad (2)$$

Nhận xét 3.1. Ta luôn chọn được (\mathbf{w}, b) sao cho y_i và $(\mathbf{w}^T \mathbf{x}_i + b)$ cùng dấu vì nếu không chỉ cần đổi dấu của $y_i, \forall i = 1, 2, \dots, n$.

Khi đó, (2) tương đương với tìm:

$$(\mathbf{w}, b) = \operatorname{argmax}_{\mathbf{w}, b} \left\{ \min_i \frac{y_i (\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|_2} \right\} \quad (3)$$

4 Các phương pháp giải bài toán SVM

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

4.1.1 Bài toán tương đương

Ta có: luôn tồn tại siêu phẳng $\mathbf{w}^T \mathbf{x} + b = 0$ phân chia hoàn toàn hai nhóm. Gọi (\mathbf{x}_0, y_0) là:

$$(\mathbf{x}_0, y_0) = \underset{(\mathbf{x}_i, y_i)}{\operatorname{argmin}} \frac{y_i(\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|_2}$$

Nhận xét 4.1. Nếu thay $(\mathbf{w}, b) = (k\mathbf{w}, kb)$ thì $\frac{y_i(\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|_2}$ không đổi, $\forall i = 1, 2, \dots, n$. Từ đó, không mất tính tổng quát, giả sử $y_0(\mathbf{w}^T \mathbf{x}_0 + b) = 1$.

Khi đó (3) trở thành tìm:

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

$$\text{với } y_i(\mathbf{w}^T \mathbf{x}_i) \geq y_0(\mathbf{w}^T \mathbf{x}_0 + b) = 1, \quad \forall i = 1, 2, \dots, n.$$

(4) tương đương với tìm:

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

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

Như vậy, sau khi xác định được siêu phẳng $\mathbf{w}^T \mathbf{x} + b = 0$, nhóm của $\mathbf{x}_i \in \mathbb{R}^m$ bất kì trùng với dấu của $\mathbf{w}^T \mathbf{x}_i + b = 0$.

4.1.2 Bài toán đối ngẫu Lagrange

Kiểm tra tiêu chuẩn Slater: Vì hai lớp dữ liệu là có thể phân cách tuyến tính, tức bài toán có nghiệm, nên luôn tồn tại cặp (\mathbf{w}_o, b_o) sao cho:

$$\begin{aligned} 1 - y_i(\mathbf{w}_o^T \mathbf{x}_i + b_o) &\leq 0, \quad \forall i = 1, 2, \dots, n, \\ \Leftrightarrow 2 - y_i(2\mathbf{w}_o^T \mathbf{x}_i + 2b_o) &\leq 0, \quad \forall i = 1, 2, \dots, n. \end{aligned}$$

Chọn $\mathbf{w}_1 = 2\mathbf{w}_o$ và $b_1 = 2b_o$, suy ra:

$$1 - y_i(\mathbf{w}_1^T \mathbf{x}_i + b_1) \leq -1 < 0, \quad \forall i = 1, 2, \dots, n$$

Vậy tiêu chuẩn Slater thỏa mãn.

Hàm Lagrangian của bài toán (5) là:

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \lambda_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

với $\boldsymbol{\lambda} = [\lambda_1, \lambda_2, \dots, \lambda_n]^T$ và $\lambda_i \geq 0, \forall i = 1, 2, \dots, n$.

Hàm đối ngẫu Lagrange là

$$g(\boldsymbol{\lambda}) = \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) \text{ với } \boldsymbol{\lambda} \succeq 0.$$

Xét $(\mathbf{w}_0, b_0, \boldsymbol{\lambda}_0) = \operatorname{argmin}_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda})$. Khi đó:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}_0} &= \mathbf{w}_0 - \sum_{i=1}^n \lambda_i y_i \mathbf{x}_i = 0 \Rightarrow \mathbf{w}_0 = \sum_{i=1}^n \lambda_i y_i \mathbf{x}_i. \\ \frac{\partial \mathcal{L}}{\partial b_0} &= \sum_{i=1}^n \lambda_i y_i = 0. \end{aligned}$$

Thay vào biểu thức của $\mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda})$, ta được:

$$g(\boldsymbol{\lambda}) = \sum_{n=1}^N \lambda_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \lambda_n \lambda_m y_n y_m \mathbf{x}_n^T \mathbf{x}_m \quad (6)$$

Từ đó, kết hợp hàm đối ngẫu Lagrange và các điều kiện ràng buộc của $\boldsymbol{\lambda}$, ta sẽ thu được **bài toán đối ngẫu** Lagrange của bài toán (26.3) có dạng

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

với:

$$\begin{aligned} \boldsymbol{\lambda} &\succeq 0 \\ \sum_{n=1}^N \lambda_n y_n &= 0 \end{aligned}$$

Vì đây là một bài toán lồi và đối ngẫu chặt xảy ra, **hệ điều kiện KKT** là điều kiện cần và đủ cho bộ $(\mathbf{w}_0, b_0, \boldsymbol{\lambda}_0)$:

$$\begin{aligned} 1 - y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0) &\leq 0, \forall i = 1, 2, \dots, n \\ \lambda_i &\geq 0, \forall i = 1, 2, \dots, n \\ \lambda_i(1 - y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0)) &= 0 \forall i = 1, 2, \dots, n \end{aligned} \quad (7)$$

$$\mathbf{w}_0 = \sum_{i=1}^n \lambda_i y_i \mathbf{x}_i$$

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

Từ điều kiện (7) ta thấy hoặc $\lambda_i = 0$ hoặc $y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0) = 1$. Khi đó, hai đường thẳng $\mathbf{w}_0^T \mathbf{x}_i + b = \pm 1$ chia tách hoàn toàn hai lớp dữ liệu.

Từ đây ta có định nghĩa sau:

Định nghĩa 4.2. Hai đường biên của bài toán SVM là hai đường thẳng $\mathbf{w}_0^T \mathbf{x}_i + b = \pm 1$.

Nhận xét 4.3. Số điểm nằm trên cả hai siêu phẳng là nhỏ hơn rất nhiều so với tổng số điểm của tập huấn luyện.

Đặt $S = \{n : \lambda_n \neq 0\}$ và N_S là số phần tử của tập S
 Theo (8), ta tính \mathbf{w} :

$$\mathbf{w} = \sum_{m \in S} \lambda_m y_m \mathbf{x}_m,$$

$$b = \frac{1}{N_S} \sum_{n \in S} \left(y_n - \sum_{m \in S} \lambda_m y_m \mathbf{x}_m^T \mathbf{x}_n \right)$$

Về việc tính λ_0 thế nào, mô hình sẽ được dựng bằng ngôn ngữ lập trình Python.

4.2 Trường hợp tập dữ liệu không thể phân cách tuyến tính: Soft-margin SVM và Kernel Trick

4.2.1 Bài toán tương đương: Soft-margin SVM

Trong thực tế, tập dữ liệu huấn luyện ta có gần như không thể *phân cách tuyến tính*. Khi đó, nếu dùng thuật toán Hard-margin SVM sẽ dẫn đến việc kết quả *lề* rất nhỏ, thậm chí không tìm được siêu phẳng phân cách.

Với suy nghĩ như vậy, thuật toán *Soft-margin SVM* ra đời. Lúc này, để thu được *lề* lớn hơn, ta chấp nhận một vài điểm dữ liệu bị *sai* rơi vào bên trong hoặc sang hẳn phía bên kia của *lề*. Để thấy sự chấp nhận *sai* cần được hạn chế, nếu không, ta có thể tạo ra một *lề* rất lớn bằng cách chấp nhận hầu hết các điểm bị *sai*.

Từ đó, hàm mục tiêu nên là một sự kết hợp để tối đa *lề* cũng như tối thiểu sự *sai*.

Như đã nói trong *hard-margin SVM*, việc tối đa *lề* có thể đưa về việc tối thiểu $\|\mathbf{w}\|_2^2$. Để tính sự *sai*, định nghĩa thêm một biên đo độ *sai* với mỗi điểm \mathbf{x}_i trong tập dữ liệu huấn luyện.

Định nghĩa 4.4. α_i là biên *sai* của biên x_i nếu như:

$$\alpha_i = \begin{cases} 0 & \text{nếu điểm dữ liệu } x_i \text{ nằm trong biên của nó,} \\ |\mathbf{w}^T \mathbf{x}_i + b - y_i| & \text{nếu điểm dữ liệu } x_i \text{ không nằm trong biên của nó.} \end{cases}$$

Nhận xét 4.5. Ta có thể định nghĩa biên *sai* như trên do α_i chỉ cần tỉ lệ với khoảng cách từ x_i đến siêu phẳng

Với *soft-margin SVM*, hàm mục tiêu sẽ có thêm một số hạng nữa giúp tối thiểu tổng độ *sai*. Khi đó, hàm mục tiêu mới là:

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

với C là một hằng số dương.

Nhận thấy $|y_i| = 1, \forall i = 1, 2, \dots, n$, ta có:

$$\begin{aligned} \alpha_i &= \left| \mathbf{w}^T \mathbf{x}_i + b - y_i \right| \\ &= \left| \mathbf{w}^T \mathbf{x}_i + b - y_i \right| |y_i| \\ &= \left| y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 \right| \\ &\geq 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) \end{aligned}$$

Từ đó ta thu được điều kiện ràng buộc mới sau, $\forall i = 1, 2, \dots, n$:

$$1 - \alpha_i - y_i (\mathbf{w}^T \mathbf{x}_i + b) \leq 0,$$

và $\alpha_i \geq 0$.

Như vậy, bài toán tối ưu gốc cho *Soft-margin SVM* phát biểu như sau, $\forall i = 1, 2, \dots, n$:

$$\begin{aligned} (\mathbf{w}_0, b_0, \alpha_0) &= \underset{\mathbf{w}, b, \alpha}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \alpha_i, \\ \text{với: } 1 - \alpha_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) &\leq 0, \\ -\alpha_n &\leq 0. \end{aligned} \quad (9)$$

4.2.2 Hướng giải quyết thứ nhất: Bài toán đối ngẫu Lagrange của *Soft-margin SVM*

Tương tự bài toán *Hard-margin SVM*, ta quan tâm đến bài toán đối ngẫu của nó.

Kiểm tra tiêu chuẩn Slater: $\forall i = 1, 2, \dots, n$ và $\forall (\mathbf{w}, b)$, luôn tồn tại các số dương α_i , đủ lớn sao cho:

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) + \alpha_i > 1.$$

Vì vậy, bài toán tối ưu này thỏa mãn tiêu chuẩn Slater.

Hàm Lagrangian của bài toán *Soft-margin SVM* là:

$$\begin{aligned} \mathcal{L}(\mathbf{w}, b, \alpha, \lambda, \nu) &= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \alpha_i + \sum_{i=1}^n \lambda_i (1 - \alpha_i - y_i(\mathbf{w}^T \mathbf{x}_i + b)) - \sum_{i=1}^n \nu_i \alpha_i, \\ \text{với } \lambda &= [\lambda_1, \lambda_2, \dots, \lambda_n]^T \succeq 0 \text{ và } \nu = [\nu_1, \nu_2, \dots, \nu_n]^T \succeq 0. \end{aligned} \quad (10)$$

Bài toán đối ngẫu của *Soft-margin SVM* là:

$$g(\lambda, \nu) = \min_{\mathbf{w}, b, \alpha} \mathcal{L}(\mathbf{w}, b, \alpha, \lambda, \nu)$$

Để giải bài toán trên, với mỗi cặp (λ, ν) , xét bộ $(\mathbf{w}_0, b, \alpha_0)$ thỏa mãn:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_0} = 0 \Leftrightarrow \mathbf{w}_0 = \sum_{i=1}^n \lambda_i y_i \mathbf{x}_i \quad (11)$$

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

$$\frac{\partial \mathcal{L}}{\partial \alpha_0} = 0 \Leftrightarrow \lambda_i = C - \nu_i \quad \forall i = 1, 2, \dots, n \quad (13)$$

Từ (13) suy ra chỉ cần quan tâm tới các cặp (λ, ν) thỏa mãn $\lambda_i = C - \nu_i$.

Từ đây ta cũng suy ra $0 \leq \lambda_i, \nu_i \leq C, \forall i = 1, 2, \dots, n$.

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

$$g(\lambda, \nu) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j. \quad (14)$$

Nhận xét 4.6. Hàm này không phụ thuộc vào ν và từ (13) suy ra $0 \leq \lambda_i \leq C, \forall i = 1, 2, \dots, n$

Khi đó, bài toán đối ngẫu trở thành tìm:

$$\lambda_0 = \underset{\lambda}{\operatorname{argmax}} g(\lambda)$$

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

$$0 \leq \lambda_i \leq C, \forall i = 1, 2, \dots, n. \quad (16)$$

Hệ điều kiện KKT của bài toán tối ưu *Soft-margin SVM* là, $\forall i = 1, 2, \dots, n$:

$$1 - \alpha_i^0 - y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0) \leq 0, \quad (17)$$

$$-\alpha_i^0 \leq 0, \quad (18)$$

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

$$\nu_i^0 \geq 0, \quad (20)$$

$$\lambda_i^0(1 - \alpha_i^0 - y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0)) = 0, \quad (21)$$

$$\nu_i^0 \alpha_i^0 = 0, \quad (22)$$

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

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

$$y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0) = 1 - \alpha_i^0. \quad (23)$$

Nếu $0 < \lambda_i^0 < C$, theo (13), $\nu_i^0 = C - \lambda_i^0 > 0$; ngoài ra theo (22) ta thu được $\alpha_i^0 = 0$. Tiếp tục theo (23), suy ra $y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0) = 1$, hay:

$$\mathbf{w}_0^T \mathbf{x}_i + b_0 = y_i, \forall i : 0 < \lambda_i^0 < C.$$

Tóm lại, nếu $0 < \lambda_i^0 < C$, các điểm \mathbf{x}_i nằm chính xác trên các lề. Khi đó, (\mathbf{w}_0, b_0) được tính theo:

$$\begin{aligned} \mathbf{w}_0 &= \sum_{i=1}^n \lambda_i^0 y_i \mathbf{x}_i, \\ b_0 &= \frac{1}{N_S} \sum_{i \in \mathcal{S}} (y_i - \mathbf{w}_0^T \mathbf{x}_i) \end{aligned} \quad (24)$$

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

Khi đó, nhóm của \mathbf{x}^* chính là đầu của:

$$\mathbf{w}_0^T \mathbf{x}^* + b_0 = \sum_{i \in \mathcal{M}} \lambda_i^0 y_i \mathbf{x}_i^T \mathbf{x}^* + \frac{1}{N_S} \sum_{j \in \mathcal{S}} \left(y_j - \sum_{i \in \mathcal{M}} \lambda_j^0 y_j \mathbf{x}_j^T \mathbf{x}^* \right) \quad (25)$$

Với $\mathcal{M} = \{n : 0 < \lambda_n\}$ và $N_{\mathcal{M}}$ là số phần tử của \mathcal{M} .

4.2.3 Hướng giải quyết thứ 2: Bài toán tối ưu không ràng buộc cho soft-margin SVM

Bài toán tối ưu không ràng buộc tương đương: Xét điều kiện ràng buộc đầu tiên:

$$1 - \alpha_i - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0 \Leftrightarrow \alpha_i \geq 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)$$

Mặt khác, $\alpha_i \geq 0$ nên ta có bài toán ràng buộc tương đương với bài toán (9) như sau:

$$(\mathbf{w}_0, b_0, \alpha_0) = \arg \min_{\mathbf{w}, b, \alpha} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \alpha_i \quad (26)$$

$$\text{thỏa mãn: } \alpha_i \geq \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)), \forall i = 1, 2, \dots, n$$

Mệnh đề 4.7. Nếu $(\mathbf{w}_0, b_0, \alpha_0)$ là nghiệm của bài toán tối ưu (26), thì:

$$\alpha_i^0 = \max(0, 1 - y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0)), \forall i = 1, 2, \dots, n \quad (27)$$

Chứng minh: Giả sử tồn tại i sao cho

$$\alpha_i > \max(0, 1 - y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0)),$$

Khi đó, chọn $\alpha'_i = \max(0, 1 - y_i(\mathbf{w}_0^T \mathbf{x}_i + b_0))$, ta được một giá trị nhỏ hơn của hàm mục tiêu, trong khi tất cả các ràng buộc vẫn được thỏa mãn (mâu thuẫn).

Vậy mệnh đề được chứng minh.

Từ đó, bài toán tối ưu trở thành:

$$(\mathbf{w}_0, b_0, \alpha_0) = \arg \min_{\mathbf{w}, b, \alpha} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)) \quad (28)$$

thỏa mãn: $\alpha_i = \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)), \forall i = 1, 2, \dots, n,$

hay:

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

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

Định nghĩa 4.9. Hàm Hinge loss có dạng:

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

Xây dựng hàm mất mát

Định nghĩa 4.10. Hàm Hinge loss trung bình có công thức:

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

Nhận xét 4.11. Xét (\mathbf{w}_0, b_0) là nghiệm bài toán Khi tập dữ liệu có thể phân cách tuyến tính, giá trị tối ưu của $J(\mathbf{w}, b) = 0$. Khi đó dễ thấy $\forall c > 1$ là hằng số, $(c\mathbf{w}_0, cb_0)$ cũng là nghiệm bài toán dẫn tới nghiệm có thể lớn tùy ý.

Vì vậy, ta có định nghĩa sau:

Định nghĩa 4.12. Hàm mất mát tổng cộng có công thức là:

$$J(\mathbf{w}, b) = L(\mathbf{w}, b) + \beta R(\mathbf{w}, b),$$

với $\beta > 0$, $R(\mathbf{w}, b) = \frac{1}{n} \frac{\beta}{2} \|\mathbf{w}\|_2^2$.

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

Tối ưu hàm mất mát: Vì việc tối ưu hàm mất mát dựa trên gradient descent, việc chính của mục này là tính đạo hàm của hàm mất mát theo \mathbf{w} và b .

Nhận xét 4.14. Đạo hàm của hàm mất mát tổng cộng không quá phức tạp.

Vì thế ta dùng Gradient Descent theo \mathbf{w} và b để tính giá trị tối ưu hàm mất mát.

Phần tính toán cụ thể, sẽ đề cập trong code Python. Đạo hàm của phần *hinge loss* không quá phức tạp:

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

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

Phần *regularization* cũng có đạo hàm tương đối đơn giản:

$$\frac{\partial \left(\frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right)}{\partial \mathbf{w}} = \lambda \mathbf{w}$$

$$\frac{\partial \left(\frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right)}{\partial b} = 0$$

Nếu cập nhật bằng gradient descent thông qua chỉ một điểm dữ liệu (\mathbf{x}_i, y_i) (*stochastic gradient descent*). Nếu $1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$, ta không cập nhật gì và chuyển sang điểm tiếp theo. Ngược lại biểu thức cập nhật cho \mathbf{w} , b được cho bởi

$$\mathbf{w} \leftarrow \mathbf{w} - \eta(-y_i \mathbf{x}_i + \lambda \mathbf{w}); \quad b \leftarrow b + \eta y_i \quad \text{nếu } 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 0$$

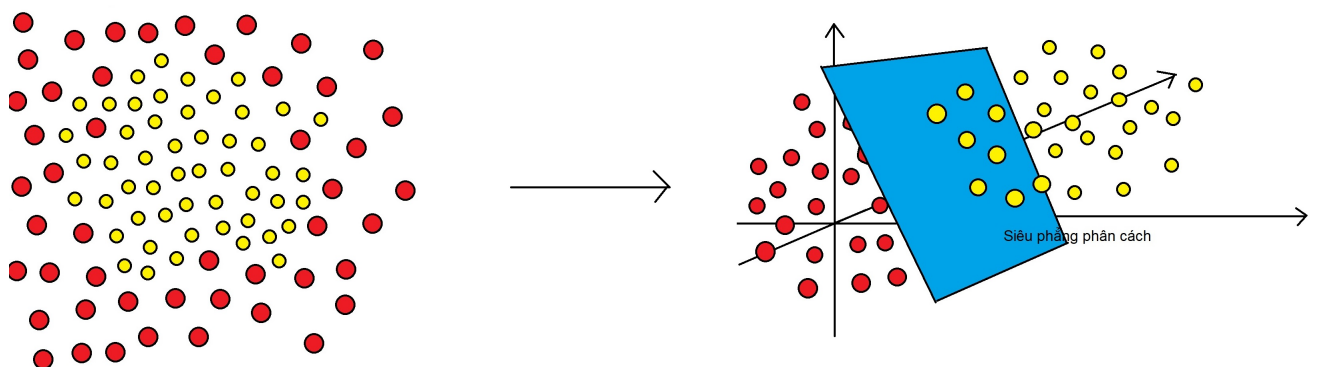
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \lambda \mathbf{w}; \quad b \leftarrow b \quad \text{nếu } 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$$

với η là *learning rate*.

4.2.4 Kernel trick

Để áp dụng thuật toán Support Vector Machine lên các dữ liệu không phân biệt tuyến tính, người ta cần biến các điểm dữ liệu này lên không gian nhiều chiều sao cho nó phân biệt tuyến tính hoặc gần như phân biệt tuyến tính trên không gian đó. Ở không gian này, bài toán có thể được giải quyết bằng hard/soft-margin SVM. Dựa trên suy nghĩ đó, ta đi tìm một hàm số $\phi(x)$ nhằm biến đổi dữ liệu x để thỏa các tính chất trên. Nhằm tăng hiệu năng tính toán và tiết kiệm bộ nhớ, các hàm Kernel được sử dụng để mô tả mối quan hệ giữa hai điểm dữ liệu bất kì trong không gian mới, thay vì tính toán trực tiếp biến đổi của chúng.

Ví dụ:



Ý tưởng cho hàm Kernel Xét hàm số $\phi(\cdot)$ sao cho sau khi được biến đổi sang không gian mới, mỗi \mathbf{x} trở thành $\phi(\mathbf{x})$, và khi đó tập dữ liệu mới trở nên gần phân biệt tuyến tính. Khi đó bài toán đối ngẫu của Soft-margin SVM trở thành:

$$\begin{aligned} \boldsymbol{\lambda} &= \arg \max_{\boldsymbol{\lambda}} \sum_{i=1}^n \lambda_i - \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j), \\ \text{thỏa mãn: } \sum_{i=1}^n \lambda_i y_i &= 0, \\ 0 \leq \lambda_i &\leq C, \forall i = 1, 2, \dots, n, \end{aligned} \quad (30)$$

và nhóm của một điểm dữ liệu mới \mathbf{x}^* chính là dấu của biểu thức:

$$\mathbf{w}^T \phi(\mathbf{x}^*) + b = \sum_{j \in S} \lambda_j y_j \phi(\mathbf{x}_j)^T \phi(\mathbf{x}^*) + \frac{1}{N_{\mathcal{M}}} \sum_{i \in \mathcal{M}} \left(y_i - \sum_{j \in S} \lambda_j y_j \phi(\mathbf{x}_j)^T \phi(\mathbf{x}_i) \right). \quad (31)$$

Nhận xét 4.15. Trong (30) và (31), ta chỉ cần tính $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ với hai điểm dữ liệu $\mathbf{x}_i, \mathbf{x}_j$ bất kì. Vì thế ta chỉ cần xác định một hàm $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$.

Khi đó, (30) và (31) trở thành:

$$\begin{aligned} \boldsymbol{\lambda} &= \arg \max_{\boldsymbol{\lambda}} \sum_{i=1}^n \lambda_i - \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{thỏa mãn: } \sum_{i=1}^n \lambda_i y_i &= 0, \\ 0 \leq \lambda_i &\leq C, \forall i = 1, 2, \dots, n, \end{aligned} \quad (32)$$

và nhóm của một điểm dữ liệu mới \mathbf{x}^* chính là dấu của biểu thức:

$$\mathbf{w}^T \phi(\mathbf{x}^*) + b = \sum_{j \in S} \lambda_j y_j k(\mathbf{x}_j, \mathbf{x}^*) + \frac{1}{N_{\mathcal{M}}} \sum_{i \in \mathcal{M}} \left(y_i - \sum_{j \in S} \lambda_j y_j k(\mathbf{x}_i, \mathbf{x}_j) \right). \quad (33)$$

Tính chất của hàm Kernel:

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

$$\sum_{i=1}^n \sum_{j=1}^n k(\mathbf{x}_i, \mathbf{x}_j) c_i c_j \geq 0, \forall c_k \in \mathbb{R}, k = 1, 2, \dots, n. \quad (34)$$

Nhận xét 4.16. Điều kiện Mercer giúp đảm bảo bài toán đối ngẫu (33) là lồi.

Một số hàm Kernel thông dụng

- Hàm linear (tuyến tính): $k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z}$.
- Hàm RBF (radial basic function/Gaussian kernel): $k(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|_2^2), \gamma > 0$
- Hàm sigmoid: $k(\mathbf{x}, \mathbf{z}) = \tanh(\gamma \mathbf{x}^T \mathbf{z} + r)$
- Hàm polynomial (đa thức): $k(\mathbf{x}, \mathbf{z}) = (r + \gamma \mathbf{x}^T \mathbf{z})^d$

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 ($\vec{v} = \vec{v}_1, \vec{v}_2, \dots, \vec{v}_{28}$) 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 với các bộ dữ liệu phân lớp không cân bằng (dữ liệu chẩn đoán có bệnh, xác định phần tử khủng bố), độ chính xác không là thang đo phù hợp để đánh giá hệ thống (do sự thiếu hụt 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 đồng nghĩa với việc độ chính xác của các điểm được dự đoán **positive class** là cao. Recall cao đồng nghĩa với việc tỉ lệ bỏ sót các điểm thực sự thuộc **positive class** là thấp.

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 đủ 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 việc dùng Hard-margin SVM là hợp lí, vì việc dự đoán sót các giao dịch lừa đảo nghiêm trọng hơn việc dự đoán nhầm các giao dịch tin cậy. Lí giải cho việc hàm đa thức và hàm sigmoid không thực hiện hiệu quả bằng hàm tuyến tính, có thể là do tình huống **Overfitting** trên tập **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	

6 Phụ lục

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import os
5
6 from sklearn.model_selection import train_test_split, KFold, cross_validate
7 from sklearn import svm
8 from sklearn.metrics import recall_score, accuracy_score, classification_report
9
10 main_file_path = "../input/creditcard.csv"
11 data = pd.read_csv(main_file_path)
12
13 #class = 1: fraudulent, class = 0: genuine
14 fraud_count = len(data[data.Class == 1])
15 fraud_indices = np.array(data[data.Class == 1].index)
16 genui_indices = data[data.Class == 0].index
17 #Random sample
18 random_sample = np.random.choice(genui_indices, fraud_count, replace = False)
19 genui_indices = np.array(random_sample)
20 #indices of chosen sample
21 indices = np.concatenate([fraud_indices, genui_indices])
22 filtered_data = data.iloc[indices, :]
23 #undersampling, ratio of positive:negative = 1:1
24 X = filtered_data.loc[:, filtered_data.columns != 'Class']
25 y = filtered_data.loc[:, filtered_data.columns == 'Class']
26 y = np.array(y).T[0]
27
28 #train, test splitting
29 train_X, test_X, train_y, test_y = train_test_split(X, y, test_size=0.1, random_state=0)
30 #train, validation splitting
31
32 #cross-validation
33 for i in kernels:
34     for j in C:
35         SVM = svm.SVC(C = j, kernel = i, gamma = 'scale')
36         Kfold = KFold(n_splits = 3, shuffle = True, random_state = 0)
37         results = cross_validate(estimator=SVM, X=train_X, y=train_y, scoring=['accuracy', 'recall', 'precision'], cv=kfold)
38         average_recall = np.mean(results['test_recall'])
39         average_accuracy = np.mean(results['test_accuracy'])
40         average_precision = np.mean(results['test_precision'])
41         standard_deviation_recall = np.std(results['test_recall'])
42         print('Parameter C: ', j, '; Kernel: ', i)
43         print('Average recall: ', average_recall, '; average accuracy: ', average_accuracy, '; average precision: ', average_precision)
44         print('Standard deviation of recall: ', standard_deviation_recall)
45
46 C = 100
47 kernel = 'linear'
48 SVM = svm.SVC(C, kernel, gamma = 'scale')
49 SVM.fit(train_X, train_y)
50 y_pred = SVM.predict(test_X)
51 print(classification_report(test_y, y_pred, target_names = ['Fraud', 'Genuine']))
52 print('Accuracy score: ', accuracy_score(test_y, y_pred))

```

Tài liệu

- [1] Vũ Hữu Tiệp. *Machine Learning Cơ Bản*. 2017.
- [2] Machine Learning Group ULB. Credit card fraud detection dataset, 2013.