

SVM

1. Idea of SVM

maximize: margin(w, b)
 w, b

assume correctly separate
 these points

margin(w, b): given a (hyperplane $H: \underline{w^T x + b = 0}$),

the minimum distance between hyperplane and point

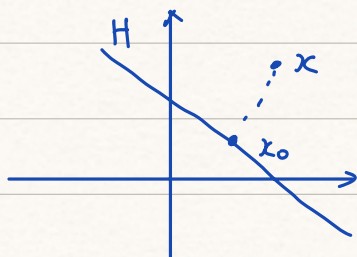
Observation: if hyperplane \tilde{H} is such an optimal hyperplane,
 then it must satisfy:

min distance between \tilde{H} and positive class points
 $=$ min distance between \tilde{H} and negative class points

2. Question: How to calculate Distance between Hyperplane and Point?

a) Self-proposed method

given point x & hyperplane $H: w^T x + b = 0$



$$\begin{aligned} \text{distance} &= | \langle x - x_0, w \rangle | / \|w\| \\ &= \frac{|w^T x + b|}{\|w\|} \end{aligned}$$

b) More sophisticated way

①. consider $\forall x_0 \in H$, the normal cone
$$N_H(x_0) = \{ \lambda \omega : \lambda \in \mathbb{R} \}$$

$$N_H(x) := \{ z : \langle z, y - x \rangle \leq 0, \forall y \in H \}$$

②. then, for $\forall x \in \mathbb{R}^n$, its projection \bar{x} on H satisfy:

$$\bar{x} \in \Pi_H(x)$$

$$\Leftrightarrow x - \bar{x} \in N_H(x)$$

$$\Leftrightarrow x - \bar{x} = \lambda \cdot \omega \quad \text{for some } \lambda$$

③. $\bar{x} \in H \Rightarrow \omega^T \bar{x} + b = 0$

$$\Rightarrow \omega^T (x - \lambda \omega) + b = 0$$

$$\Rightarrow \lambda = \frac{\omega^T x + b}{\omega^T \omega}$$

$$\Rightarrow \text{distance} = \|x - \bar{x}\| = |\lambda \cdot \|\omega\||$$

$$= \frac{|\omega^T x + b|}{\|\omega\|_2}$$

3. Formulation of SVM

$$\rightarrow \begin{cases} \max_{\omega, b} & \min_{i \in [n]} \frac{|\omega^T x_i + b|}{\|\omega\|} \rightarrow \text{margin for dataset} \\ \text{s.t.} & y_i (\omega^T x_i + b) \geq 0 \quad \forall i \in [n] \end{cases}$$

correctly separate all data

$$\Leftrightarrow \begin{cases} \max_{w,b} & \frac{1}{\|w\|} \min_{i \in [n]} |w^T x_i + b| = f(w,b) \\ \text{s.t.} & \gamma_i (w^T x_i + b) \geq 0 \quad i \in [n] \end{cases}$$

observation: $f(w,b) = f(kw, kb)$ for arbitrary k !

maintain same
object, but shrink the maximizer set
 \longleftrightarrow

$$\begin{cases} \max_{w,b} & \frac{1}{\|w\|} \\ \text{s.t.} & \gamma_i (w^T x_i + b) \geq 0 \quad i \in [n] \\ & \min_{i \in [n]} \gamma_i (w^T x_i + b) = 1 \end{cases}$$

$$\Leftrightarrow \begin{cases} \min_{w,b} & \frac{1}{2} w^T w \\ \text{s.t.} & \gamma_i (w^T x_i + b) \geq 1 \quad i \in [n] \end{cases} \rightarrow \boxed{\text{Primal Form}}$$

(it must exist some k , s.t. $\underline{\gamma_k (w^T x_k + b) = 1}$)

4. Lagrange Duality / KKT Condition

① Non-Linear Programming

$$\rightarrow \begin{cases} \min_x & f(x) \\ \text{s.t.} & g_i(x) = 0 \quad i \in [m] \\ & h_j(x) \leq 0 \quad j \in [l] \\ & x \in X \quad (X \subseteq \mathbb{R}^p) \end{cases}$$

Generally speaking, we can have KKT optimality condition

a) Necessary part, if \bar{x} is a minimizer, then there exists $\{u_i : i \in [m]\} \{v_j : j \in [l]\}$, such that:

$$\begin{cases} \nabla f(\bar{x}) + \sum_i u_i \nabla g_i(\bar{x}) + \sum_j v_j \nabla h_j(\bar{x}) = 0 \\ g_i(\bar{x}) = 0, \quad h_j(\bar{x}) \leq 0 \\ v_j \geq 0 \\ v_j h_j(\bar{x}) = 0 \quad \forall j \in [l] \end{cases}$$

b) Sufficient Part: under convexity condition of f, g_i, h_j , & there exists $\{u_i : i \in [m]\} \{v_j : j \in [l]\}$ s.t. KKT Condition holds at $\bar{x} \in \mathbb{R}^n$.
 then \bar{x} is a global minimizer

② From the duality perspective

Consider Lagrangian $L(x; u, v)$

$$:= f(x) + \sum_i u_i g_i(x) + \sum_j v_j h_j(x)$$

Lagrangian Dual function $\Theta(u, v)$

$$\Theta(u, v) := \inf_{x \in \mathcal{X}} L(x; u, v)$$

Concave func

Note: $\Theta(u, v) = \inf_x L(x; u, v)$

$$\leq L(\bar{x}; u, v)$$

$$\leq f(\bar{x})$$

$$= \text{primal objective}$$

→ weak duality

$$(D) \leq (P)$$

generally holds

Thus, our interest is $\begin{cases} \max_{u, v} \Theta(u, v) \\ \text{s.t. } v_j \geq 0 \quad j \in [l] \end{cases} \leq f(\bar{x})$

↓
the largest lower bound of primal problem

(Regularity)

Note: under Sufficient Condition (e.g., Slater's condition + convexity)

Strong Duality holds!

That means, there exists $\begin{cases} \hat{x} \in \text{Primal Solution} \\ (\hat{u}, \hat{v}) \in \text{Dual Solution} \end{cases}$

such that $\Theta(\hat{u}, \hat{v}) = f(\hat{x}) \rightarrow (D) = (P)$

Rmk: If Strong Duality holds, then

① directly, it implies the $\begin{cases} \text{primal solution } \hat{x} \\ \text{dual solution } (\hat{u}, \hat{v}) \end{cases}$ will satisfy

KKT condition system (necessary condition)

② if the program is convex program,

then it will imply:

$\{\hat{x}; \hat{u}, \hat{v}\}$ is the primal, dual solution

$\Leftrightarrow \{\hat{x}; \hat{u}, \hat{v}\}$ satisfy KKT condition system

5. Dual of Hard / Soft Margin SVM

a) conclusion:

hard-margin

$$\begin{cases} \max_{\alpha} & -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_{i=1}^n \alpha_i \\ \text{s.t} & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \alpha_i \geq 0 \quad i \in [n] \end{cases}$$

soft-margin

$$\begin{cases} \max_{\alpha} & -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle + \sum_{i=1}^n \alpha_i \\ \text{s.t} & \sum_{i=1}^n \alpha_i y_i = 0 \\ & \underline{0 \leq \alpha_i \leq C} \quad i \in [n] \end{cases}$$

b) derivation for soft-margin case

$$\begin{array}{l} \text{Primal Formulation} \\ \left[\begin{array}{ll} \min_{w, b, \xi} & \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i \\ \text{s.t} & y_i (w^T x_i + b) \geq 1 - \xi_i \quad i \in [n] \\ & \xi_i \geq 0 \quad i \in [n] \end{array} \right. \end{array}$$

$$L(w, b, \zeta; \alpha, \beta) = \frac{1}{2} w^T w + C \sum_{i=1}^n \zeta_i + \sum_{i=1}^n \alpha_i (1 - \zeta_i - \gamma_i (w^T x_i + b)) - \sum_{i=1}^n \beta_i \zeta_i$$

① Stationary

$$\begin{cases} \nabla_w L = w - \sum_{i=1}^n \alpha_i \gamma_i x_i \\ \nabla_b L = \sum_{i=1}^n \alpha_i \gamma_i \\ \nabla_{\zeta_i} L = C - \alpha_i - \beta_i \end{cases}$$

for solution analysis

KKT optimality

② Feasibility

③ Complementary Condition $\begin{cases} \alpha_i (1 - \zeta_i - \gamma_i (w^T x_i + b)) = 0 \\ \beta_i \zeta_i = 0 \end{cases}$

$$\Theta(\alpha, \beta) = \min_{w, b, \zeta} L(w, b, \zeta; \alpha, \beta)$$

$$= \begin{cases} \frac{1}{2} \hat{w}^T \hat{w} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i \gamma_i x_i^T \hat{w}, & \begin{cases} \sum \alpha_i \gamma_i = 0 \\ \alpha_i + \beta_i = C \end{cases} \\ -\infty & , \text{ otherwise} \end{cases}$$

$$\Rightarrow \begin{cases} \max_{\alpha, \beta} \Theta(\alpha, \beta) \\ \text{s.t. } \alpha \geq 0, \beta \geq 0 \end{cases} \Leftrightarrow \begin{cases} \max_{\alpha, \beta} -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \gamma_i \gamma_j \langle x_i, x_j \rangle + \sum_{i=1}^n \alpha_i \\ \text{s.t. } \sum_{i=1}^n \alpha_i \gamma_i = 0 \\ \alpha_i + \beta_i = C \quad i \in [n] \end{cases}$$

$$\Leftrightarrow \begin{cases} \max_{\alpha} -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \gamma_i \gamma_j \langle x_i, x_j \rangle + \sum_{i=1}^n \alpha_i \\ \text{s.t. } \sum_{i=1}^n \alpha_i \gamma_i = 0 \\ 0 \leq \alpha_i \leq C \quad i \in [n] \end{cases}$$

Rmk: (Dual Solution Analysis)

From KKT Condition, we are able to reconstruct (\hat{w}, \hat{b})

from the dual solution $\underline{\hat{\alpha}}$

① $\hat{w} = \sum_i \alpha_i \gamma_i x_i$

② to derive \hat{b} , we have to determine:

→ index sv such that $0 < \alpha_{sv} < C$

→ it forces $\begin{cases} \zeta_{sv} = 0 \\ 1 - \zeta_{sv} - \gamma_{sv} (\hat{w}^T x_{sv} + \hat{b}) = 0 \end{cases}$

$$\Rightarrow \gamma_{sv} (\hat{w}^T x_{sv} + \hat{b}) = 1$$

$$\Rightarrow \hat{w}^T x_{sv} + \hat{b} = \gamma_{sv}$$

$$\Rightarrow \hat{b} = \gamma_{sv} - \sum_{i=1}^n \hat{\alpha}_i \gamma_i \langle x_i, x_{sv} \rangle$$

$$= \gamma_{sv} - \sum_{i \in SV} \hat{\alpha}_i \gamma_i \langle x_i, x_{sv} \rangle$$

c) dual \Rightarrow kernel

just by replace all $\langle x_i, x_j \rangle$ by $k(x_i, x_j)$
 $= \langle \phi(x_i), \phi(x_j) \rangle$

6. Re-think Soft-Margin SVM

\rightarrow Primal Formulation

$$\begin{cases} \min_{w, b, \xi} & \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i \\ \text{s.t} & \gamma_i (w^T x_i + b) \geq 1 - \xi_i \quad i \in [n] \\ & \xi_i \geq 0 \quad i \in [n] \end{cases}$$

$$\Rightarrow \begin{cases} \min_{w, b, \xi} & \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i \\ \text{s.t} & \xi_i \geq \max \{0, 1 - \gamma_i (w^T x_i + b)\} \quad i \in [n] \end{cases}$$

$$\begin{aligned} & \downarrow \\ & \begin{cases} \min_{w, b, \xi} & f(w, b) + C \sum_{i=1}^n \xi_i \\ \text{s.t} & \xi_i \geq g_i(w, b) \quad i \in [n] \end{cases} \\ & = \min_{w, b} \underbrace{f(w, b) + C \sum_{i=1}^n g_i(w, b)} \end{aligned}$$

$$\Rightarrow \min_{w, b} \underbrace{\frac{1}{2} w^T w + C \sum_{i=1}^n \max \{0, 1 - \gamma_i (w^T x_i + b)\}}$$

regularization

loss term

$$\Rightarrow \min_{w, b} \sum_{i=1}^n \text{Loss}_{\text{hinge}}(z_i) + \frac{1}{2C} \|w\|_2^2 \rightarrow \text{Soft-margin SVM}$$

$$\begin{cases} z_i = y_i (w^T x_i + b) \rightarrow \text{agreement} \\ \text{Loss}_{\text{hinge}}(z) = \max\{0, 1 - y_i (w^T x_i + b)\} \end{cases}$$

Re-cap: what is Logistic Regression?

$$\begin{aligned} \hat{y}_i &= \text{sigmoid}(w^T x + b) \\ &= \frac{1}{1 + \exp(-(w^T x + b))} \end{aligned}$$

→ optimization:

Cross-entropy Loss

$$\text{Loss}(w, b) = - \sum_{i=1}^n \left[y_i \log\left(\frac{1}{1 + \exp(-(w^T x_i + b))}\right) + (1 - y_i) \log\left(\frac{\exp(-(w^T x_i + b))}{1 + \exp(-(w^T x_i + b))}\right) \right]$$

$$= - \sum_{i=1}^n \left[\log\left(\frac{1}{1 + \exp(w^T x_i + b)}\right) + y_i (w^T x_i + b) \right]$$

$$= - \sum_{i=1}^n \log\left(\frac{1}{1 + \exp(-\tilde{y}_i (w^T x_i + b))}\right) \quad \tilde{y}_i \in \{-1, 1\}$$

$$= \sum_{i=1}^n \log(1 + \exp(-\tilde{y}_i (w^T x_i + b)))$$

$$:= \sum_{i=1}^n \text{Loss}_{\text{logistic}}(z_i) \quad \begin{cases} z_i = \tilde{y}_i (w^T x_i + b) \\ \text{Loss}_{\text{logistic}}(z) = \log(1 + \exp(-z)) \end{cases}$$

7. Optimization for Dual Problem

→ Block Coordinate Descent

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x_1, \dots, x_\ell) \\ \text{s.t.} \quad & x \in \mathcal{X} \end{aligned}$$

$$\begin{aligned} \text{here } \sum_{j=1}^{\ell} n_j &= n \\ x_j &\in \mathbb{R}^{n_j} \end{aligned}$$

- Algo**:
- 1) set $x^{(0)}$
 - 2) for $j = 1, 2, \dots, \ell$
$$x_j^{(k+1)} = \underset{x_j}{\operatorname{argmin}} f(\underbrace{x_1^{(k+1)}}_{\leftarrow}, \underbrace{x_j}_{\leftarrow}, \underbrace{x_{>j}^{(k)}}_{\leftarrow})$$
 - 3) $k \rightarrow k+1$.

Apply to **Dual SVM** \Rightarrow SMD Algo

Recap: **Dual SVM** \Rightarrow
$$\begin{cases} \min_{\alpha} \quad \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K_{ij} - \sum_{i=1}^n \alpha_i \\ \text{s.t.} \quad 0 \leq \alpha_i \leq C \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

SMD Algo: 1) initialize α

2) choose (i, j) pairs in $[n]$.

3)

$$(\hat{\alpha}_i, \hat{\alpha}_j) = \begin{cases} \min_{\alpha} \quad \frac{1}{2} K_{ii} \alpha_i^2 + \frac{1}{2} K_{jj} \alpha_j^2 - y_i y_j K_{ij} \alpha_i \alpha_j \\ \quad - \alpha_i - \alpha_j \\ \text{s.t.} \quad 0 \leq \alpha_i \leq C, \quad 0 \leq \alpha_j \leq C \\ \quad \alpha_i y_i + \alpha_j y_j = 0 \end{cases}$$

$$(\alpha_i, \alpha_j) \leftarrow (\hat{\alpha}_i, \hat{\alpha}_j)$$

4) Repeat