

Summary

→ reformulate Ridge Reg as Kernel Ridge Reg.

without specifying feature map

only works for Regression
(not classification)

SVM

(here we use interchangeably)

① Linear f^{\ddagger} : $f(x) = w^T x + b$ (actually is Affine function)

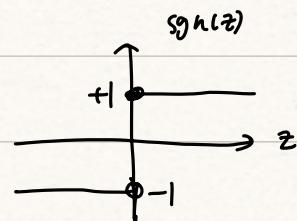
hyperplane $\Rightarrow f(x) = 0$

$$\mathcal{H}_{\text{SVM}} := \{f : f(x) = \text{sgn}(w^T x + b)\}$$

② Binary Classification \leftarrow **SVM**

$$\rightarrow \mathcal{D} = \{(x_i, y_i) : x_i \in \mathbb{R}^d, y_i \in \{-1, +1\}\}_{i=1}^N$$

$$\rightarrow \text{linear decision function: } f(x) = \text{sgn}(w^T x + b)$$



Linearly Separable Assumption:

$$\exists f^* \text{ (linear decision } f^{\ddagger}) \text{ s.t. } f^*(x_i) = y_i$$

$$\Leftrightarrow \exists (\hat{w}, \hat{b}) \text{ s.t. } y_i(\hat{w}^T x_i + \hat{b}) \geq 0 \text{ for } i=1, 2, \dots, N$$

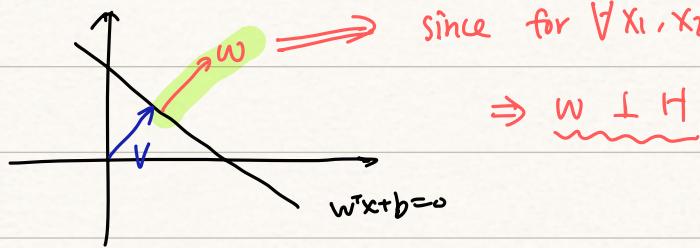
③ There are many choices of 'Separable Hyperplane' (Non-unique)

Maximize Margin
to find the optimal hyperplane

{robust generalization (IDEA)}

Hyperplane

[E.g.]



since for $\forall x_1, x_2 \in H$, then $w^T(x_1 - x_2) = 0$

$$\Rightarrow w \perp H$$

How to calculate the distance from H to origin?

$$\textcircled{1} \quad V \parallel w$$

$$\textcircled{2} \quad V \in H$$

$$\Rightarrow w^T V + b = 0 \quad \& \quad V = d w$$

$$\Rightarrow d = -\frac{b}{\|w\|^2}$$

$$\Rightarrow \text{distance} = \|V\| = \left\| -\frac{b}{\|w\|^2} w \right\| = \frac{|b|}{\|w\|}$$

$$\text{generalization: distance } (x_0, H) = \underbrace{\frac{|w^T x_0 + b|}{\|w\|}}$$

Proof

$$\textcircled{1} \quad x - x_0 \parallel w$$

$$\textcircled{2} \quad x \in H \Rightarrow w^T x + b = 0$$

$$\textcircled{3} \quad \text{distance} = \left| \langle x - x_0, \frac{w}{\|w\|} \rangle \right| = \frac{|w^T x_0 + b|}{\|w\|}$$

(4)

Formulation

$$\textcircled{1} \quad \max_{w,b} \gamma^g(w,b)$$

s.t. (w,b) can separate correctly

the ideally formulation

$$\Leftrightarrow \max_{w,b} \min_{i=1,2,\dots,N} \gamma_i^g(w,b) \rightarrow \gamma_i^g = \frac{|w^T x_i + b|}{\|w\|}$$

s.t. $y_i(w^T x_i + b) > 0$ for all $i=1,2,\dots,N$

$$\Leftrightarrow \max_{w,b} \frac{1}{\|w\|} \min_{i=1,2,\dots,N} (w^T x_i + b)$$

s.t. $y_i(w^T x_i + b) > 0 \quad \forall i$

} equivalent

$$\Leftrightarrow \max_{w,b,\lambda} \frac{1}{\|w\|} \cdot \lambda$$

s.t. $y_i(w^T x_i + b) \geq \lambda \quad \forall i$

} (The trick step)

$$\max_{w,b} \frac{1}{2} \|w\|^2$$

s.t. $y_i(w^T x_i + b) \geq 1 \quad \forall i$

equivalent in the sense that it will share the same objective value but may have different optimizer

⑥

$$\min_{w,b} \frac{1}{2} w^T w$$

s.t. $y_i(w^T x_i + b) \geq 1 \quad \forall i$

Primal Form of SVM

2 Ques [equivalence convexity]

Dual Form (Lagrangian Multiplier)

optimization problem

$$\begin{aligned} & \min_z F(z) \\ & \text{s.t. } z \in P \rightarrow \begin{cases} g_i(x) = 0 \\ h_j(x) \leq 0 \end{cases} \end{aligned}$$

Lagrangian Multiplier Framework (KKT necessary condition)

Lagrangian Function

$$L(x, \lambda, \mu) = f(x) + \sum \lambda_i g_i(x) + \sum \mu_j h_j(x)$$

KKT Condition [Necessary]

① Stationarity

$$\nabla_x L(\hat{x}; \hat{\lambda}, \hat{\mu}) = 0$$

there exists $\hat{\lambda}_i, \hat{\mu}_j$, s.t. $\hat{x} \in Z$

Lagrangian Multipliers

② Dual feasibility $\hat{\mu} \geq 0$

③ Primal feasibility $g_i(\hat{x}) = 0$

$$h_j(\hat{x}) \leq 0$$

④ Complementary Slackness

$$\hat{\lambda}_i h_j(\hat{x}) = 0$$

Should under some constraint qualification condition

Application of KKT

i) distance between x_0 and H

$$\text{dist}(x_0, H) = \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2$$

s.t. $\mathbf{w}^\top \mathbf{x} + b = 0$

→ convex problem

$$L(x; \mu) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2 + \mu(\mathbf{w}^\top \mathbf{x} + b)$$

$$\nabla_{\mathbf{x}} L(\hat{\mathbf{x}}; \hat{\mu}) = 0 \Rightarrow (\hat{\mathbf{x}} - \mathbf{x}_0) + \hat{\mu} \mathbf{w} = 0$$

$$\Rightarrow \boxed{\begin{matrix} \mathbf{w}^\top \mathbf{x} \\ \downarrow \\ -b \end{matrix}} - \mathbf{w}^\top \mathbf{x}_0 + \hat{\mu} \mathbf{w}^\top \mathbf{w} = 0$$

→ primary feasibility

$$\Rightarrow \hat{\mu} = \frac{b + \mathbf{w}^\top \mathbf{x}_0}{\|\mathbf{w}\|_2^2}$$

$$\Rightarrow \boxed{\text{dist}} = |\hat{\mathbf{x}} - \mathbf{x}_0| = \|\hat{\mu} \mathbf{w}\|_2$$

$$= \frac{|w^\top x_0 + b|}{\|\mathbf{w}\|_2}$$

ii) Dual Form of SVM

Recall: Primal form

$$\boxed{\begin{matrix} \min_{\mathbf{w}} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \mathbf{w} \\ \text{s.t.} & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad \text{for } \forall i \end{matrix}}$$

$$L(w, b; \mu) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^N \mu_i [1 - y_i(w^\top \mathbf{x}_i + b)]$$

Stationarity

$$\textcircled{1} \quad \nabla_w L(\hat{w}, \hat{b}; \hat{\mu}) = 0 \Rightarrow \hat{w} = \sum_{i=1}^N \hat{\mu}_i y_i \mathbf{x}_i$$

$$\nabla_b L(\hat{w}, \hat{b}, \hat{\mu}) = 0 \Rightarrow \sum \hat{\mu}_i y_i = 0$$

\textcircled{2} Dual Feasibility $\Rightarrow \hat{\mu}_i \geq 0 \quad \forall i = 1, \dots, N$

\textcircled{3} Primal Feasibility $\Rightarrow y_i (\hat{w}^\top \mathbf{x}_i + \hat{b}) \geq 1 \quad \forall i = 1, \dots, N$

\textcircled{4} Complementary Slackness $\Rightarrow \mu_i [1 - y_i (\hat{w}^\top \mathbf{x}_i + \hat{b})] = 0 \quad \forall i = 1, \dots, N$

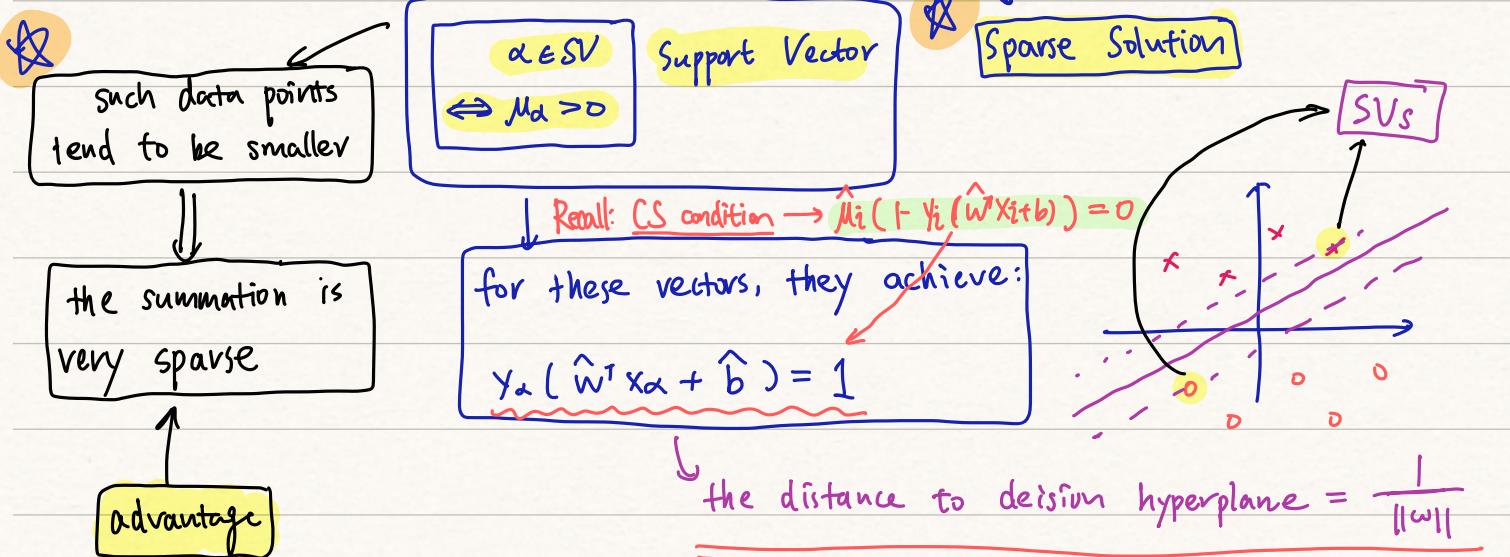
Strong Duality guarantee that, these 2 are equivalent

$$\begin{array}{ll} \max_{\boldsymbol{\mu}} & \min_{\boldsymbol{w}} L(\boldsymbol{w}; \boldsymbol{\mu}) \rightarrow \text{for the convexity of } L(\boldsymbol{w}; \boldsymbol{\mu}) \\ \text{s.t.} & \boldsymbol{\mu} \geq 0 \\ & \min_{\boldsymbol{w}} L(\boldsymbol{w}; \boldsymbol{\mu}) = L(\sum \mu_i \mathbf{x}_i; \boldsymbol{\mu}) \\ & = -\frac{1}{2} \sum \mu_i \mu_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum \mu_i \end{array}$$

$$\leftarrow \begin{array}{l} \max_{\boldsymbol{\mu}} -\frac{1}{2} \sum \mu_i \mu_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum \mu_i \\ \text{s.t. } \boldsymbol{\mu} \geq 0 \end{array}$$

Suppose that we have $\hat{\boldsymbol{\mu}} \geq 0$, then the decision function \hat{f} becomes:

$$\begin{aligned} \hat{\boldsymbol{w}} &= \sum_{i=1}^N \hat{\mu}_i \mathbf{y}_i \mathbf{x}_i \\ \hat{f}(\mathbf{x}) &= \text{sign} (\hat{\boldsymbol{w}}^T \mathbf{x} + \hat{b}) \quad \text{from CS condition, we have, some of } \hat{\mu}_i = 0 \\ &= \text{sign} (\sum_{i=1}^N \hat{\mu}_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + \hat{b}) \\ &= \text{sign} (\sum_{\text{desv}} \hat{\mu}_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + \hat{b}) \end{aligned}$$



Kernel SVM

Recall:

$$\begin{aligned} \hat{f}(\mathbf{x}) &= \text{sgn} (\hat{\boldsymbol{w}}^T \mathbf{x} + \hat{b}) \\ &= \text{sgn} (\sum_i \hat{\mu}_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + \hat{b}) \end{aligned}$$

only depends on $x_i^T x$ (or $\langle x_i, x \rangle$)

IDEA: extend to Kernel trick!

(like Ridge Reg \rightarrow Kernel Ridge Reg)

deal with non-linearly
separable dataset

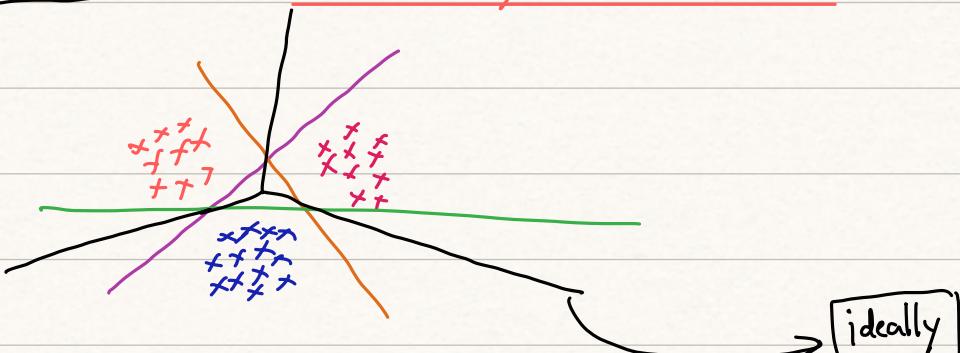
Kernel SVM

$x \xrightarrow{\phi(\cdot)} \phi(x) \rightarrow$ linearly separable in feature space
non-linearly separable

$$\hat{f}(x) = \text{sgn} \left(\sum_i \hat{w}_i y_i k(x_i, x) + \hat{b} \right) \quad k(x_i, x) = \langle \phi(x_i), \phi(x) \rangle$$

generalize to multi-class classification

1 vs all \leadsto several binary-classification tests



Q: How to aggregate these 3 sub-models

Recap of Lagrangian Duality

① KKT condition

given optimization problem:

$$\begin{cases} \min_{x \in \mathbb{R}^n} & f(x) \\ \text{s.t.} & g_i(x) \leq 0 \quad i=1, 2, \dots, I \\ & h_j(x) = 0 \quad j=1, 2, \dots, J \end{cases}$$

(Necessary Part)

$$\hat{x} \in \underset{\substack{x \in \mathbb{R}^n \\ x \in C}}{\operatorname{argmin}} f(x)$$

under some
constraint qualification

① $\hat{x} \in C$
 ② $\exists \lambda_i \geq 0, \mu_i \in \mathbb{R}$ s.t.
 $\nabla f(\hat{x}) + \sum_I \lambda_i \nabla g_i(x) + \sum_J \mu_j \nabla h_j(x) = 0$
 $\lambda_i g_i(\hat{x}) = 0 \quad \forall i \in I$

Define $L(x; \lambda, \mu) = f(x) + \sum \lambda_i g_i(x) + \sum \mu_j h_j(x)$

\longleftrightarrow

① $g_i(\hat{x}) \leq 0, h_j(\hat{x}) = 0$
 ② $\hat{\lambda}_i \geq 0, \hat{\mu}_i \in \mathbb{R}$
 ③ $\nabla_x L(\hat{x}; \hat{\lambda}, \hat{\mu}) = 0$
 ④ $\hat{\lambda}_i g_i(\hat{x}) = 0$

(Sufficient Part)

a) $(\hat{x}; \hat{\lambda}, \hat{\mu})$ satisfy ① ② ③ ④

b) Convex constraints for $\begin{cases} f(x) \\ g_i(x) \\ h_j(x) \end{cases} \Rightarrow \hat{x} \in \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} f(x) \text{ s.t } x \in C$

Note: this part is not related to duality

KEY POINT

focus on the connection Between

② Duality Introduction

{ primal minimizer \hat{x}
dual maximizer $\hat{\lambda}, \hat{\mu}$

given

$$\begin{cases} \min_{x \in \mathbb{R}^n} & f(x) \\ \text{s.t.} & g_i(x) \leq 0 \quad i=1, 2, \dots, I \\ & h_j(x) = 0 \quad j=1, 2, \dots, J \end{cases}$$

Define $L(x; \lambda, \mu) = f(x) + \sum_{i \in I} \lambda_i g_i(x) + \sum_{j \in J} \mu_j h_j(x)$

a)

then optimization problem $\Leftrightarrow \min_x \max_{\substack{\lambda \geq 0 \\ \mu \in \mathbb{R}^J}} L(x; \lambda, \mu)$

b) Weak Duality ($\hat{x} \rightarrow$ optimizer for constrained prob)

$$\text{general result} \leftarrow \min_x L(x; \lambda, \mu) \leq f(x) \quad \text{for } \begin{cases} \text{feasible } x \\ \lambda \geq 0 \\ \mu \end{cases}$$

$$\Rightarrow \max_{\lambda \geq 0, \mu} \min_x L(x; \lambda, \mu) \leq f(\hat{x}) := \min_{x \in C} f(x)$$

dual \leq **PRIMAL**

weaker slater: ① g_i affine, ② feasible Region $\neq \emptyset$

slater: $\exists x \in \text{relint } \mathcal{D}, \text{ s.t. } g_i(x) < 0 \& h_j(x) = 0$

c) Strong Duality

① Slater's / Weak Slater's condition + $\begin{cases} f \text{ convex} \\ g_i \text{ convex} \\ h_j \text{ affine} \end{cases} \Rightarrow$ **strong duality**

② given strong duality holds, then

$(\hat{x}; \hat{\lambda}, \hat{\mu})$ is the optimal solution for $(P; D)$

$$\Rightarrow \begin{cases} ① g_i(\hat{x}) \leq 0 \& h_j(\hat{x}) = 0 \quad \text{for } \forall i, j \\ ② \hat{\lambda} \geq 0 \\ ③ \nabla_x L(\hat{x}; \hat{\lambda}, \hat{\mu}) = 0 \\ ④ \hat{\lambda}_i g_i(\hat{x}) = 0 \quad \underline{\forall i=1, 2, \dots, I} \end{cases}$$

③ given strong duality holds, if (P) is a convex program

(if slater's condition + convex program holds)

$\begin{cases} f \text{ convex} \\ g_i \text{ convex} \\ h_j \text{ affine} \end{cases}$

then $(\hat{x}; \hat{\lambda}, \hat{\mu})$ is the optimal solution for $(P; D)$

$$\Leftrightarrow \begin{cases} ① g_i(\hat{x}) \leq 0 \& h_j(\hat{x}) = 0 \quad \text{for } \forall i, j \\ ② \hat{\lambda} \geq 0 \\ ③ \nabla_x L(\hat{x}; \hat{\lambda}, \hat{\mu}) = 0 \\ ④ \hat{\lambda}_i g_i(\hat{x}) = 0 \quad \underline{\forall i=1, 2, \dots, I} \end{cases}$$

④ **Saddle Point Defn**

$(\hat{x}; \hat{\lambda}, \hat{\mu})$ is saddle point of $L(x; \lambda, \mu)$

$$\Leftrightarrow L(\hat{x}; \hat{\lambda}, \hat{\mu}) = \sup_{\lambda \geq 0} \inf_x L(x; \lambda, \mu) = \inf_{\mu} \sup_{\lambda \geq 0} L(x; \lambda, \mu)$$

Then, $(\hat{x}; \hat{\lambda}, \hat{\mu})$ is saddle point for $L(x; \lambda, \mu)$

\Leftrightarrow strong duality holds & $(\hat{x}; \hat{\lambda}, \hat{\mu})$ is optimizer to $(P; D)$

SVM formulation: (Soft-Margin)

primal :
$$\begin{cases} \min_{w, b, \xi} & \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i \\ \text{s.t.} & y_i(w^T x_i + b) \geq 1 - \xi_i \quad i=1, 2, \dots, N \\ & \xi_i \geq 0 \quad i=1, 2, \dots, N \end{cases}$$

dual: $L(w, b, \xi; \alpha, \beta) = \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i + \sum_{i=1}^N \alpha_i (1 - \xi_i - y_i(w^T x_i + b))$

$$- \sum_{i=1}^N \beta_i \xi_i$$

From KKT Condition, if $(\hat{w}, \hat{b}, \hat{\xi}, \hat{\alpha}, \hat{\beta})$ are optimizer to $(P; D)$

then it holds :
$$\left\{ \begin{array}{l} \textcircled{1} \quad y_i(\hat{w}^T x_i + \hat{b}) \geq 1 - \hat{\xi}_i \\ \hat{\xi}_i \geq 0 \\ \textcircled{2} \quad \hat{\alpha}_i \geq 0, \hat{\beta}_i \geq 0 \\ \textcircled{3} \quad \nabla_w L(\hat{w}, \hat{b}, \hat{\xi}; \hat{\alpha}, \hat{\beta}) = \hat{w} - \sum_{i=1}^N \hat{\alpha}_i y_i x_i = 0 \\ \nabla_b L(\hat{w}, \hat{b}, \hat{\xi}; \hat{\alpha}, \hat{\beta}) = \sum_{i=1}^N \hat{\beta}_i y_i = 0 \\ \nabla_\xi L(\hat{w}, \hat{b}, \hat{\xi}; \hat{\alpha}, \hat{\beta}) = C - \hat{\alpha} - \hat{\beta} = 0 \\ \textcircled{4} \quad \left\{ \begin{array}{l} \hat{\alpha}_i (1 - \hat{\xi}_i - y_i(\hat{w}^T x_i + \hat{b})) = 0 \\ \hat{\beta}_i \hat{\xi}_i = 0 \end{array} \right. \end{array} \right.$$

then. Dual Form $\Rightarrow \left\{ \begin{array}{l} \max_{\alpha, \beta} \min_{w, b, \xi} L(w, b, \xi; \alpha, \beta) \\ \text{s.t. } \alpha \geq 0, \beta \geq 0 \end{array} \right.$

$$\Rightarrow \begin{cases} \max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{i=1}^N \alpha_i \\ \text{s.t. } \alpha \geq 0, \quad C - \alpha \geq 0, \quad \sum_{i=1}^N \alpha_i y_i = 0 \end{cases}$$

$$\Rightarrow \begin{cases} \max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{i=1}^N \alpha_i \\ \text{s.t. } 0 \leq \alpha \leq C, \quad \sum_{i=1}^N \alpha_i y_i = 0 \end{cases}$$

↓

Solve for $\hat{\alpha}$ \Rightarrow

$$\begin{cases} \hat{\beta} = C - \hat{\alpha} \\ \hat{w} = \sum_{i=1}^N \hat{\alpha}_i y_i x_i \\ \text{for those } 0 < \hat{\alpha}_i < C. \end{cases}$$

from CS condition

we have $y_s (\hat{w}^T x_s + \hat{b}) = 1$

↓

$$\begin{aligned} \hat{b} &= y_s - \hat{w}^T x_s \\ &= y_s - \sum_{i=1}^N \hat{\alpha}_i y_i x_i^T x_s \end{aligned}$$

prediction $\hat{f}(x) = \operatorname{sgn}(\hat{w}^T x + \hat{b})$

$$= \operatorname{sgn}\left(\sum_{i=1}^N \hat{\alpha}_i y_i x_i^T x + \hat{b}\right)$$

where \hat{b} is given by $\hat{b} = y_s - \sum_{i=1}^N \hat{\alpha}_i y_i x_i^T x_s$

Another formulation

$$\begin{cases} \min_{w, b, \gamma} \frac{1}{2} w^T w + C \sum_{i=1}^N \gamma_i \\ \text{s.t. } y_i (w^T x_i + b) \geq 1 - \gamma_i \quad i=1, 2, \dots, N \\ \gamma_i \geq 0 \quad i=1, 2, \dots, N \end{cases}$$

Unconstrained Optimization Problem Formulation

\Leftrightarrow

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

loss function (hinge loss)

(Model): $\hat{y} = \operatorname{sgn}(\hat{w}^T x + \hat{b})$

$\operatorname{Loss}_{\text{hinge}}(z) = \max \{0, 1 - z\}$

