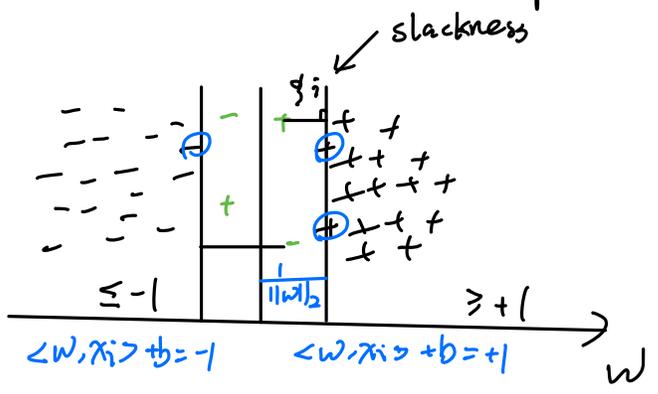


Wrap up SVM

General case: Non-separable case



Previously, $\min \frac{1}{2} \|w\|_2^2$

s.t. $y_i (\langle w, x_i \rangle + b) \geq 1, i=1, \dots, n$

Now, relax the constraints:

$$\min \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \eta_i$$

← tuning constant
← penalty

s.t. $y_i (\langle w, x_i \rangle + b) \geq 1 - \eta_i \dots \alpha_i \geq 0$ (Primal)

$\eta_i \geq 0, i=1, \dots, n \dots \mu_i \geq 0$

Solved in
close form

$$\text{primal} \Rightarrow \min_{(w,b,\eta)} \max_{(\alpha,\mu)} L \Rightarrow \max_{(\alpha,\mu)} \min_{(w,b,\eta)} L \Rightarrow \max Q \text{ (Dual)}$$

Primal parameters: w, b, η

Dual parameters: α, μ (Lagrange multipliers)

$$L(w, b, \eta, \alpha, \mu) = \underbrace{\frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \eta_i}_{\text{Primal objective}} + \underbrace{\sum_{i=1}^n \alpha_i (1 - \eta_i - y_i \langle w, x_i \rangle + b)}_{\leq 0 \text{ in original form}} + \sum_{i=1}^n \mu_i (-\eta_i)$$

Still concave - convex

free the constraints

min w, b, ξ $\mathcal{L} \longrightarrow Q(\alpha)$

← Searching for the saddle point

(same for all 3 rounds)

Representer

$\frac{\partial \mathcal{L}}{\partial w} = 0 \implies w = \sum_{i=1}^n \alpha_i y_i x_i$
 Alternative: $\|w - \hat{w}\|_2^2$

$\frac{\partial \mathcal{L}}{\partial b} = 0 \implies \sum_{i=1}^n \alpha_i y_i = 0$ (caused by b)
 Alter: $+\infty / -\infty$ trick

$\frac{\partial \mathcal{L}}{\partial \xi_i} = 0 \implies \alpha_i = C - \mu_i \leq C, i=1, \dots, n$ (caused by non-separability)
 (box constraints)

$Q(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i y_i x_i \right\|_2^2$ (Same as that in the 1st round)

Kernelize:

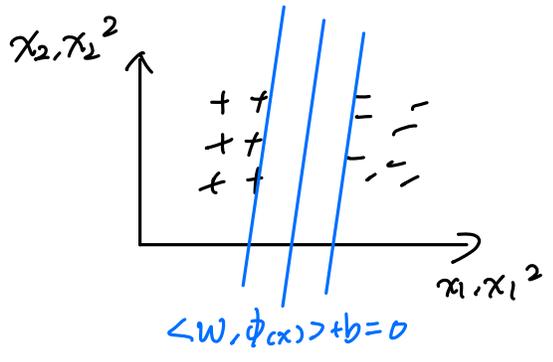
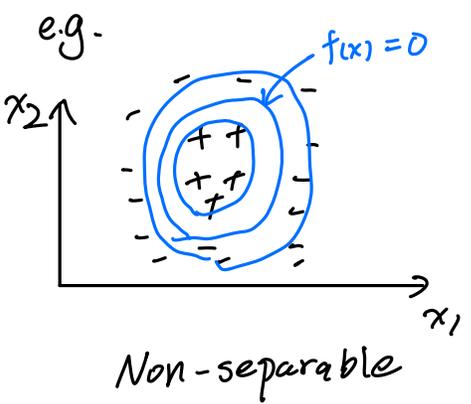
$\hat{w} = \sum_{i=1}^n \alpha_i y_i x_i$

← interpolation

$f(x) = \langle \hat{w}, x \rangle + b = \sum_{i=1}^n \alpha_i y_i \langle x_i, x \rangle + b$ (learned classifier)

$Q(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$ (training)

$x \longrightarrow \phi(x) \quad \langle x, x' \rangle \longrightarrow \langle \phi(x), \phi(x') \rangle = k(x, x')$



Dual optimization :

$$\max_{\alpha_i \in [0, c], \forall i} Q(\alpha)$$

$$\sum_i \alpha_i y_i = 0$$

If remove $\sum_{i=1}^n \alpha_i y_i = 0$ (assuming $b=0$)

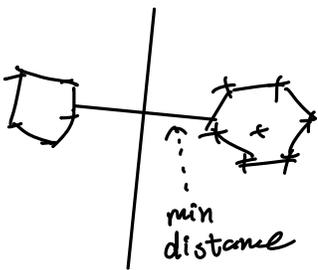
$\max_{\alpha_i \in [0, c]} Q(\alpha)$ (Coordinate ascent to solve it)
Each time update one α_i

However, if we allow b , one more constraint : $\sum_{i=1}^n \alpha_i y_i = 0$

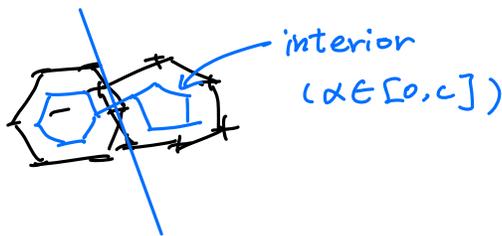
Sol: Each time, change $Q(\alpha_i, \alpha_j)$ (Sequential minimal optimization)
SMO (Linear case, b important)

↕
Kernel case, b negligible
(Dual Coordinate ascent)

Geometry of Dual optimization :



Separable



Non-separable

Primal: Max margin

Dual: Min distance

Stats aspect / Connected to logistic regression

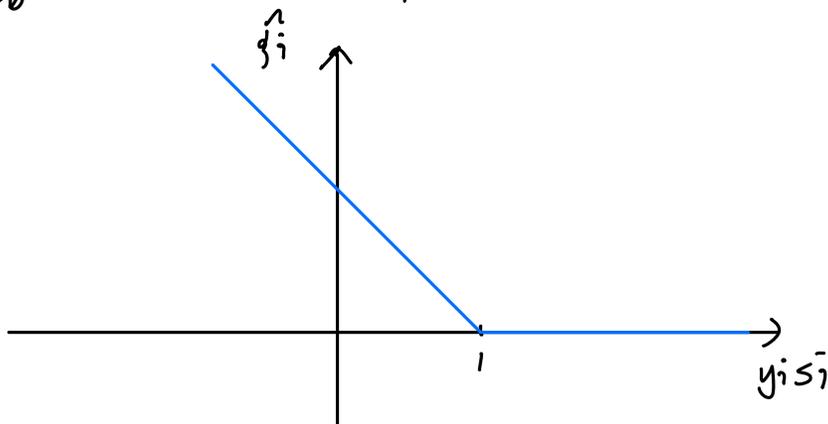
Another way to translate constraints:

$$y_i (\langle w, x_i \rangle + b) \geq 1 \implies \hat{\xi}_i = 0 \quad \leftarrow \text{penalty}$$

$$\dots < 1 \implies \hat{\xi}_i = 1 - y_i (\langle w, x_i \rangle + b)$$

$$\hat{\xi}_i = \max(0, 1 - y_i (\underbrace{\langle w, x_i \rangle + b}_{s_i})) \quad \text{hinge loss}$$

$$\min_{w, b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \max(0, 1 - y_i (\langle w, x_i \rangle + b))$$



Logistic regression:

$$s_i = \langle w, x_i \rangle + b$$

$$p_i = \frac{e^{s_i}}{1 + e^{s_i}} = \Pr(y_i = +1 | s_i), \quad 1 - p_i = \frac{1}{1 + e^{s_i}} = \Pr(y_i = -1 | s_i)$$

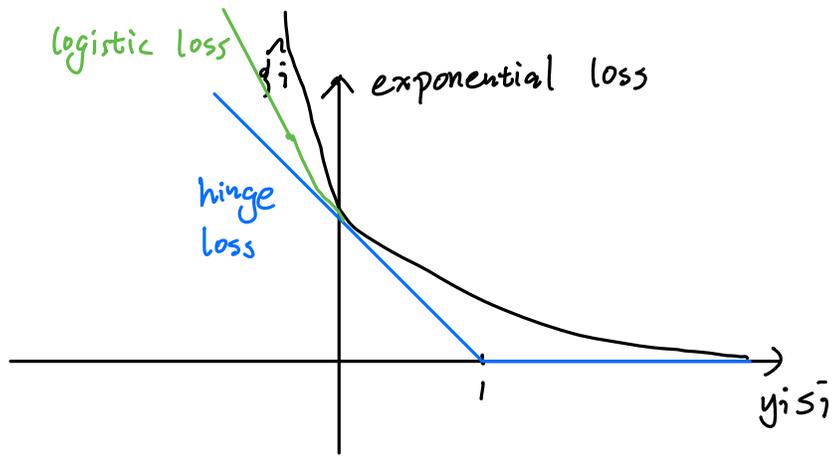
Note: y_i is 0/1

$$P(y_i | s_i) = \frac{1}{1 + e^{-y_i s_i}}$$

$$\begin{aligned} \text{Loss} &= -\log P(y_i | s_i) = \log(1 + e^{-y_i s_i}) \quad \text{logistic loss} \\ &= \begin{cases} -y_i s_i, & y_i s_i \rightarrow -\infty \\ e^{-y_i s_i}, & y_i s_i \rightarrow +\infty \end{cases} \end{aligned}$$

(Similar to $1 - y_i s_i$, SVM)
(exponential loss, Adaboost)

similar to lasso forward selection



Rethink

$$SVM : \sum_{i=1}^n \max(0, 1 - y_i \langle w, x_i \rangle + b) + \frac{1}{2C} \|w\|_2^2 \quad \text{margin consideration}$$

$$\text{logistic} : \sum_{i=1}^n \log(1 + \exp(\langle w, x_i \rangle + b)) + \lambda \|w\|_2^2 \quad \text{shrinkage consideration}$$

ridge logistic regression