

Announcement: Office hours every Thurs 4pm. on Zoom (943998564)

Primal SVM

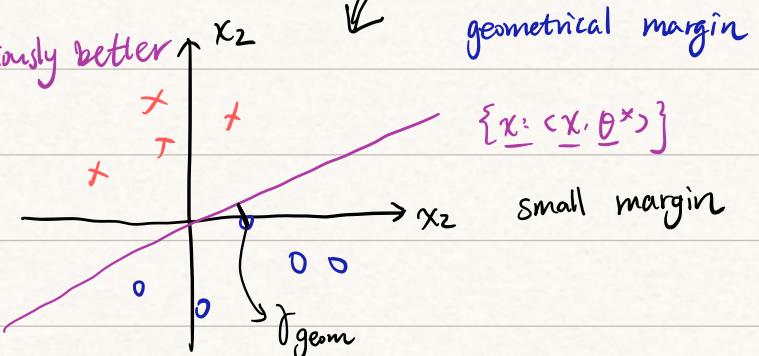
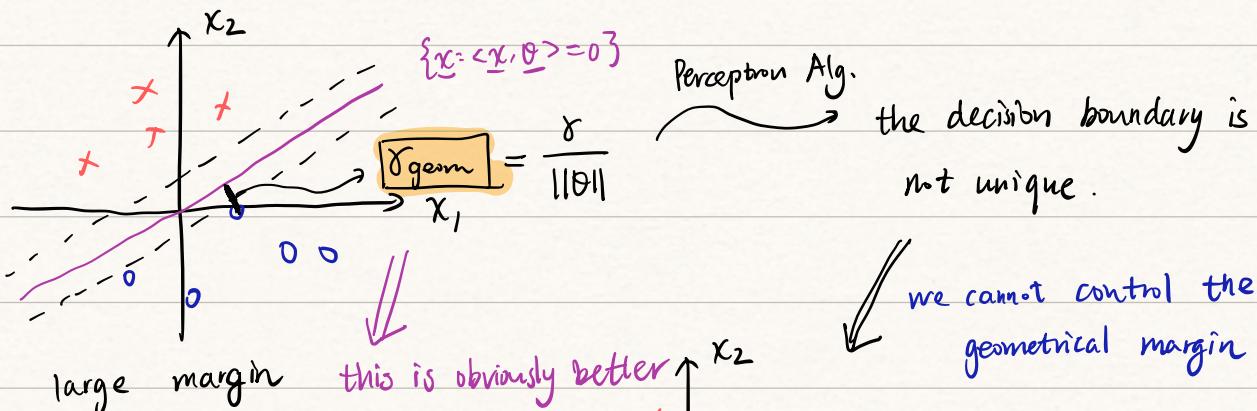
Recall: $\mathcal{D} = \{(\underline{x}_t, y_t) : 1 \leq t \leq n\}$ $\underline{x}_t \in \mathbb{R}^d$, $y_t \in \{-1, +1\}$.

want to find a linear classifier:

$$\mathcal{F}_{lin} = \{f: \mathbb{R}^d \rightarrow \{-1, +1\} : f(\underline{x}) = \text{sign}(\langle \underline{x}, \underline{\theta} \rangle) \text{ for some } \underline{\theta} \in \mathbb{R}^d\}.$$



Perceptron Alg- to learn the linear classifier.



$$J = \min_{1 \leq t \leq n} y_t \langle \underline{x}_t, \underline{\theta} \rangle$$

$\xrightarrow{\text{agreement}}$

Rmk: Perceptron alg. don't give us a unique $\underline{\theta} \in \mathbb{R}^d$.

The result $L = \{x : \langle x, \underline{\theta} \rangle = 0\}$ may have a small margin.

$\underline{\theta}$: parameter for the linear classifier

$$\mathcal{F}_{lin} = \{f: \mathbb{R}^d \rightarrow \{-1, +1\} : f_{\underline{\theta}}(\underline{x}) = \text{sgn}(\langle \underline{x}, \underline{\theta} \rangle) \text{ for some } \underline{\theta} \in \mathbb{R}^d\}$$

Constraint ① Classifies all samples correctly
 $\Leftrightarrow y_t \langle \underline{x}_t, \underline{\theta} \rangle \geq \gamma \text{ for all } t = 1, 2, \dots, n \Rightarrow \text{we want to find } \underline{\theta}$

Objective ② At the same time, $\max \gamma_{geom}$ satisfy this

$$\Leftrightarrow \max_{\|\theta\|} \frac{r}{\|\theta\|}$$

optimization problems!

$\max_{\theta \in \mathbb{R}^d} \frac{1}{\|\theta\|}$ subject to $y_t < \underline{\theta}^T x_t, \underline{\theta} > \geq r \Rightarrow$ its solution is guaranteed by r - linearly separable.

maximize min. distance to decision boundary.

minimize max. correctly classified (r is positive).

$$\min_{\theta \in \mathbb{R}^d} \frac{\|\theta\|}{r} \text{ s.t. } y_t < \frac{\theta}{r}, x_t > \geq 1 \text{ for every } t.$$

Actually r determined by $\underline{\theta}$

$$\text{let } \underline{\theta}' \text{ be } \frac{\underline{\theta}}{r}$$

$$\min_{\theta' \in \mathbb{R}^d} \|\underline{\theta}'\| \text{ s.t. } y_t < \underline{\theta}', x_t > \geq 1 \text{ for every } t.$$



$$\min_{\theta \in \mathbb{R}^d} \|\theta\| \text{ s.t. } y_t < \underline{\theta}, x_t > \geq 1 \text{ for every } t.$$

Primal-SVM (without offset & slack)

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2} \|\theta\|^2 \text{ s.t. } y_t < \underline{\theta}, x_t > \geq 1 \text{ for every } t$$

Uniqueness

Claim: The solution to Primal-SVM is unique. → Contradiction

Pf: $\exists 2$ different solution vectors $\underline{\theta}_1, \underline{\theta}_2$. ($\underline{\theta}_1 \neq \underline{\theta}_2$)

It must hold $\|\underline{\theta}_1\| = \|\underline{\theta}_2\|$ (obviously)

Consider $\underline{\theta} = \underline{\theta}_1 + \underline{\theta}_2$

$$y_t < \underline{\theta}, x_t > = \frac{1}{2} y_t < \underline{\theta}_1, x_t > + \frac{1}{2} y_t < \underline{\theta}_2, x_t >$$

$$\geq \frac{1}{2} + \frac{1}{2} = 1 \Rightarrow \underline{\theta} \text{ is feasible}$$

$$\|\underline{\theta}\| = \left\| \frac{1}{2} (\underline{\theta}_1 + \underline{\theta}_2) \right\| \leq \frac{1}{2} \|\underline{\theta}_1\| + \frac{1}{2} \|\underline{\theta}_2\| = \|\underline{\theta}_1\| = \|\underline{\theta}_2\|$$

i) $\|\underline{\theta}\| < \|\underline{\theta}_1\| = \|\underline{\theta}_2\| \rightarrow$ contradicts to the choice of $\underline{\theta}_1$ & $\underline{\theta}_2$

ii) $\|\theta_1\| = \|\theta_2\| = \|\theta\| \Rightarrow$ triangle inequality holds with equality

$$\hookrightarrow \|a+b\| = \|a\| + \|b\| \Leftrightarrow a = \lambda b$$

$$\|\frac{1}{2}\theta_1 + \frac{1}{2}\theta_2\| = \|\frac{1}{2}\theta_1\| + \|\frac{1}{2}\theta_2\| \Leftrightarrow \frac{1}{2}\theta_1 = \frac{1}{2}\theta_2 \cdot \lambda.$$



$$\|\theta_1\| = \|\theta_2\| \Rightarrow \lambda = 1 \Rightarrow \theta_1 = \theta_2.$$

contradiction!
($\theta_1 \neq \theta_2$)

\Rightarrow The solution is unique #

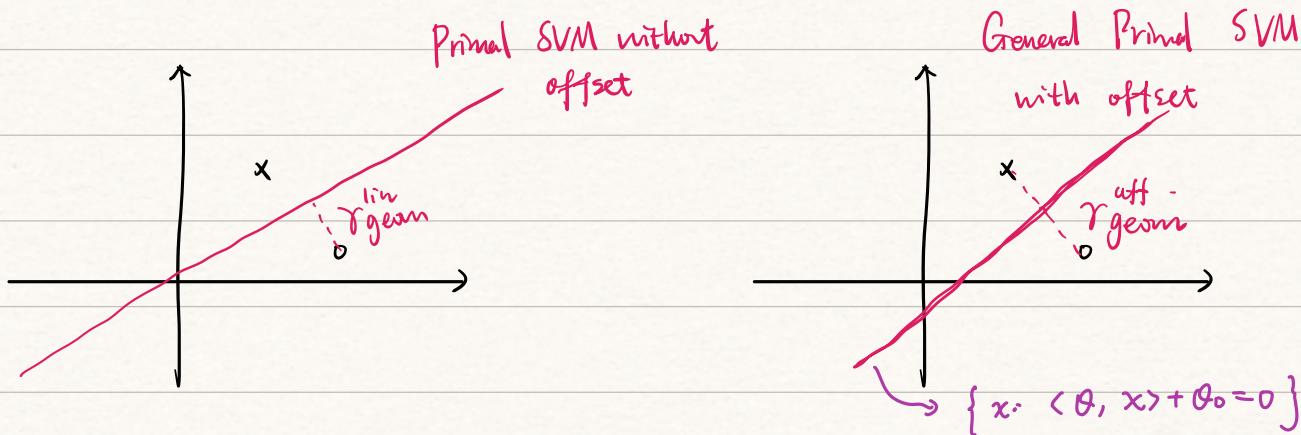
General Primal SVM with offset

① $\mathcal{F}_{aff} := \{ f: \mathbb{R}^d \rightarrow \{+1, -1\} : f(\underline{x}) = \text{sgn}(\langle \underline{x}, \underline{\theta} \rangle + \theta_0) \text{ for some } \underline{\theta} \in \mathbb{R}^d \text{ & } \theta_0 \in \mathbb{R} \}$

② Each $f \in \mathcal{F}_{aff}$ is parametrized & $(\underline{\theta}, \theta_0) \in \mathbb{R}^{d+1}$

③ Decision boundary is set of all \underline{x} satisfying $\{ \underline{x} : \langle \underline{x}, \underline{\theta} \rangle + \theta_0 = 0 \}$.

Rank: offset parameter can lead to larger γ_{geom}



$$\rightarrow r_{\text{geom}}^{\text{lin}} \leq r_{\text{geom}}^{\text{aff}}$$

Primal SVM with offset problem:

$$\min_{(\underline{\theta}, \theta_0) \in \mathbb{R}^{d+1}} \frac{1}{2} \|\underline{\theta}\|^2 \quad \text{s.t. } y_t (\langle \underline{\theta}, \underline{x}_t \rangle + \theta_0) \geq 1 \quad \text{for all } t$$

By writing the constraints, we are assuming that :

dataset is affinely-separable

$$\Rightarrow \exists (\underline{\theta}, \theta_0) \in \mathbb{R}^{d+1} \quad \text{s.t. } y_t (\langle \underline{\theta}, \underline{x}_t \rangle + \theta_0) \geq 0 \quad \text{for all } t$$

Question: Why not convert objective function to $\frac{1}{2} \left\| \begin{bmatrix} \underline{\theta} \\ \theta_0 \end{bmatrix} \right\|^2$?

Rmk : i) θ_0 only appears in constraints of obj. function

ii) Different from modeling each feature vector

\underline{x}_t as $\underline{x}_t = \begin{bmatrix} \underline{x}_t \\ 1 \end{bmatrix}$. (not the right thing to do!)

Ex: \mathcal{D} : γ -affinely separable

$$L := \{ \underline{x} : \langle \underline{x}, \underline{\theta} \rangle + \theta_0 = 0 \}$$

With the offset θ_0 , the min distance over all training examples

to L is still $\frac{\gamma}{\|\underline{\theta}\|}$ (not θ_0) .

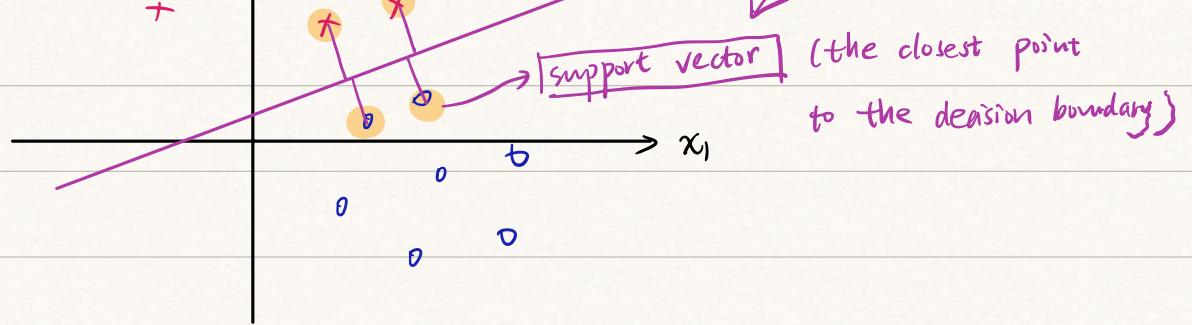
Robustness of Primal SVM

The primal-SVM with offset prob. is quadratic program

{ objective quadratic
constraints linear

Solution to QP depends only on a small # of training samples





Quantify the quality of Primal-SVM with offset

↪ use Leave-one-out Cross-Validation error

$$\text{LOOCV} := \frac{1}{n} \sum_{t=1}^n \underbrace{\text{Loss}(y_t, f(x_t; \underline{\theta}^{(t)}, \theta_0^{(t)}))}_{\text{o-1 loss}}$$

parameters learned when
(y_t, x_t) is removed from \mathcal{D}

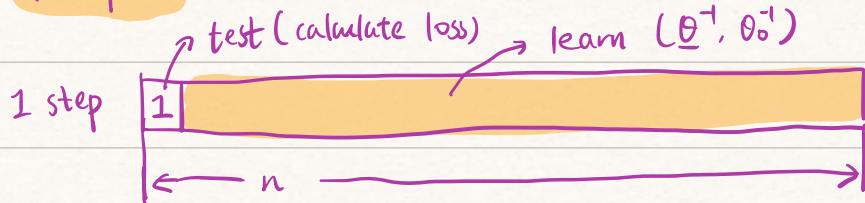
Actually there are n SVMs!

$$(\underline{\theta}^{(t)}, \theta_0^{(t)}) = \underset{(\underline{\theta}, \theta_0) \in \mathbb{R}^{d+1}}{\operatorname{argmin}} \frac{1}{2} \|\underline{\theta}\|^2 \text{ s.t. } y_s (\langle \underline{\theta}, x_s \rangle + \theta_0) \geq 1$$

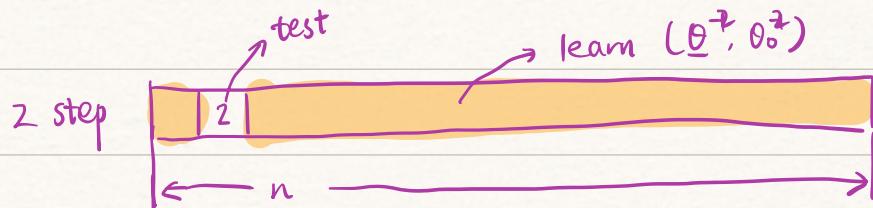
for $s \in \{1, 2, \dots, n\} / \{t\}$

$$\text{Loss}(y, \hat{y}) = \begin{cases} 1, & y \neq \hat{y} \\ 0, & \text{else} \end{cases}$$

The process of LOOCV:



→ try to mimic the scenario



that you are training samples
and test on an independent
set!

Rmk: If LOOCV is small, model generalization well!

Prop: Let the # of support vectors be $N \in \{0, 1, \dots, n\}$.

Then $\text{LOOCV} \leq \frac{N}{n}$ (obviously)

Rmk: If # of support vectors is small,

then PRIMAL-SVM has a good quality of generalization!

Allowing Misclassified Samples

→ dataset is not affinely separable

Real datasets are often not affinely separable



PRIMAL-SVM with offset & slack variables

$$\min_{(\underline{\theta}, \theta_0, \underline{\xi})} \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{t=1}^n \xi_t \quad C > 0.$$

$$\text{s.t. } y_t (\langle \underline{\theta}, \mathbf{x}_t \rangle + \theta_0) \geq 1 - \xi_t \quad t=1,2,\dots,n$$

① if dataset is affinely separable, then $\xi = 0$ can be achieved

② if dataset is not affinely separable, ξ can not be 0

$$\xi \geq 0$$

Rank: i) Margin constraint violated for t^{th} example if $\xi_t > 0$

ii) Penalty to objective $\rightarrow C \xi_t$

iii) $C \rightarrow +\infty$, then $\xi_t \rightarrow 0$ for all t

⇒ we obtain original SVM (without ξ_t slack)

iv) for general case $C > 0$, some margin violation is allowed.

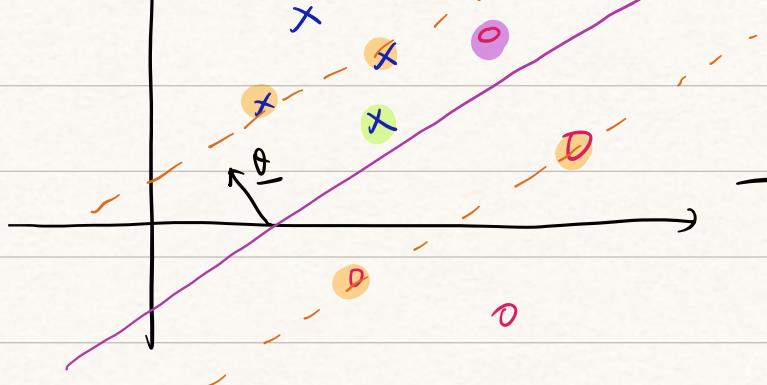
v) trade-off between max-margin objective $\frac{1}{2} \|\underline{\theta}\|^2$

and margin violation $C \sum_t \xi_t$

Q: Which are the support vectors for SVM with slack & offset?



+ --- Decision Boundary.

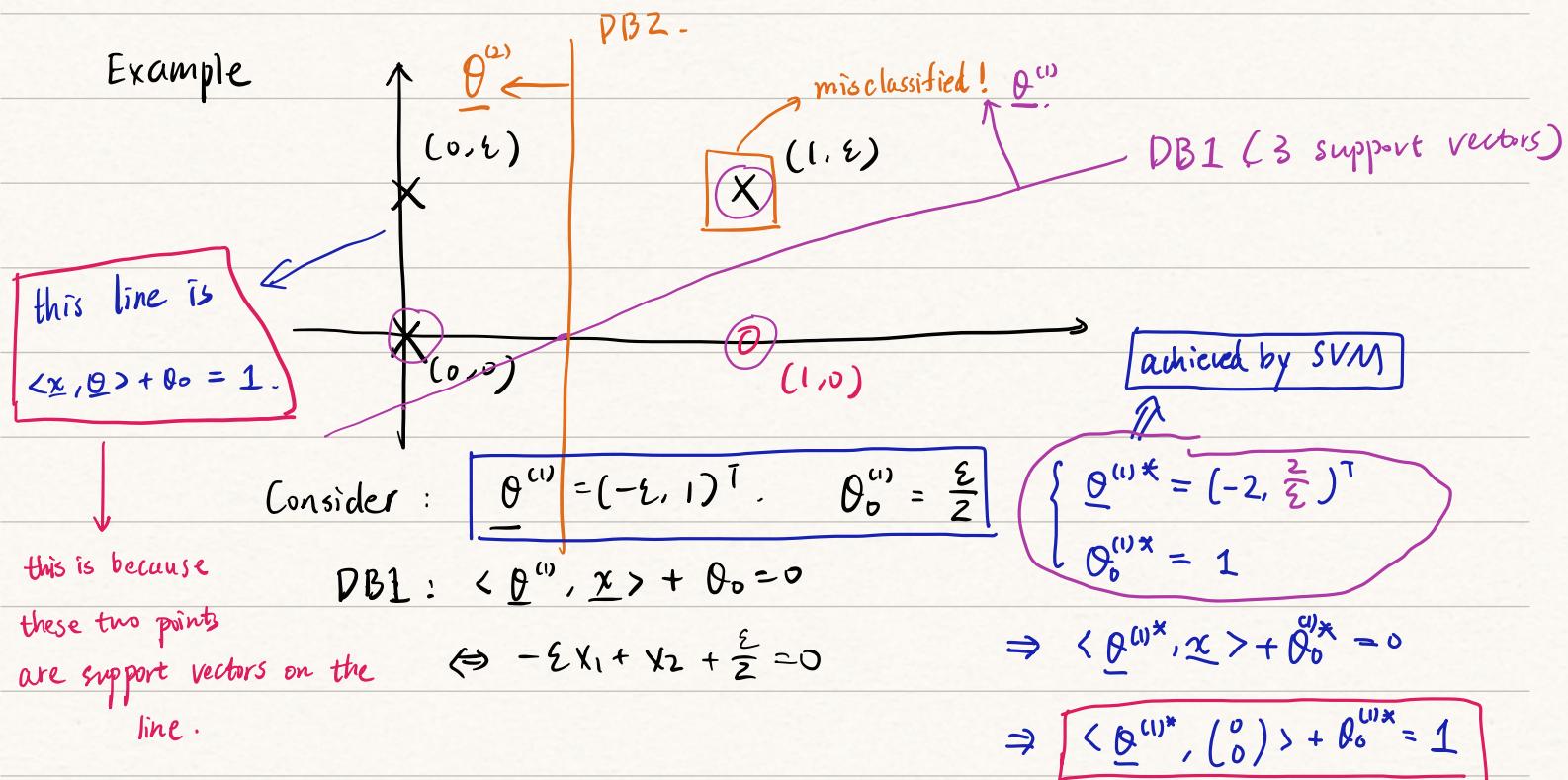


Answer: i) Points that lie on margin \Rightarrow usual case

ii) Those that violate margin constraints but not enough
to be misclassified

iii) Misclassified points

Example



Classify all samples correctly, but margin small.

$$\underline{\theta}^{(2)*} = (-2, 1) \quad \theta_0^{(2)*} = 1$$

for SVM algorithm,
we must guarantee that

$$\underline{\theta}^{(2)} = (-1, 0) \quad \theta_0^{(2)} = \frac{1}{2}$$

$$DB2 : -x_1 + \frac{1}{2} = 0 \Rightarrow x_1 = \frac{1}{2}$$

for those (x_t, y_t) lie in the boundary

Top right point is misclassified $\Rightarrow 1$ tiny error

margin is larger

$r_2 >> r_1$

\Rightarrow training error costs $\sum \mathbb{1}_{\hat{y}_t \neq y_t}$ since $\langle \theta^{\text{opt}}, (\xi_t) \rangle + \theta_0^{\text{opt}}$

explain why $\hat{y}_t = 2$

$$= -2 + 1$$

$$= -1$$

$$= 1 - \textcircled{2}$$

\downarrow
 \hat{y}_t

$$\gamma_1 = \frac{1}{\|\theta^{\text{opt}}\|}$$

$$= \frac{1}{\sqrt{4 + \frac{4}{\xi_t^2}}}$$

$$\frac{1}{2\gamma_1^2} \leq \frac{1}{2\gamma_2^2} + 2C$$

\hookrightarrow this is actually from:

$$\gamma_2 = \frac{1}{\|\theta^{\text{opt}}\|}$$

$$= \frac{1}{2}$$

$$\text{objective: } \frac{1}{2}\|\theta\|^2 + C \sum \hat{y}_t$$

$$\frac{1}{2\gamma_2^2}$$

$\left\{ \begin{array}{l} \text{If } C \rightarrow \infty, \text{ this inequality is satisfied and all points are classified correctly} \\ \text{If } C \approx 0, \text{ this means we are allowed to make some mistakes!} \end{array} \right.$

Outline:

1. Start from Perceptron \rightarrow disadvantage

$\left\{ \begin{array}{l} \text{not unique} \\ \text{not robust} \end{array} \right.$

\downarrow
still linear classifier
first.

SVM · (HARD) \longrightarrow USE FOR LINEARLY
SEPARABLE DATASET!

Model: $\max_{\gamma} \frac{\gamma}{\|\theta\|}$

$$\text{s.t. } y_t \langle \underline{x}_t, \underline{\theta} \rangle \geq \gamma \quad \text{for all } t$$

$$\min_{\underline{\theta}} \frac{\|\underline{\theta}\|}{\gamma} := \|\underline{\theta}\|$$

$$\text{s.t. } y_t \langle \underline{x}_t, \frac{\underline{\theta}}{\|\underline{\theta}\|} \rangle \geq 1 \quad \text{for all } t$$

$$\min_{\underline{\theta}} \|\underline{\theta}\| \Leftrightarrow \min_{\underline{\theta}} \frac{1}{2} \|\underline{\theta}\|^2$$

$$\text{s.t. } y_t \langle \underline{x}_t, \underline{\theta} \rangle \geq 1 \quad \text{for all } t$$

Support vector (\underline{x}_t, y_t)



Satisfy that $y_t \langle \underline{x}_t, \hat{\underline{\theta}} \rangle = 1$.



this implies MARGIN

$$\gamma_{\text{geom}} = \frac{1}{\|\hat{\underline{\theta}}\|}$$

2. SVM with offsets \rightarrow Affine Classifier.

$$f_{\underline{\theta}, \theta_0}(x) = \langle \underline{x}, \underline{\theta} \rangle + \theta_0$$

$$\text{Model: } \min_{\underline{\theta}, \theta_0} \frac{1}{2} \|\underline{\theta}\|^2$$

Note, $\gamma_{\text{geom}} = \frac{1}{\|\hat{\underline{\theta}}\|}$ also holds!

$$\text{s.t. } y_t (\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0) \leq 1$$

3. Advantage of SVM (HARD)

{ Uniqueness of Solution
Robustness

\downarrow
i.e., support vectors will restrict
 $\text{LOOCV} \leq \frac{K}{n}$

4. Soft SVM (SVM with offset & slackness)

① Model: $\min_{\underline{\theta}, \theta_0, \gamma} \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_t \gamma_t$

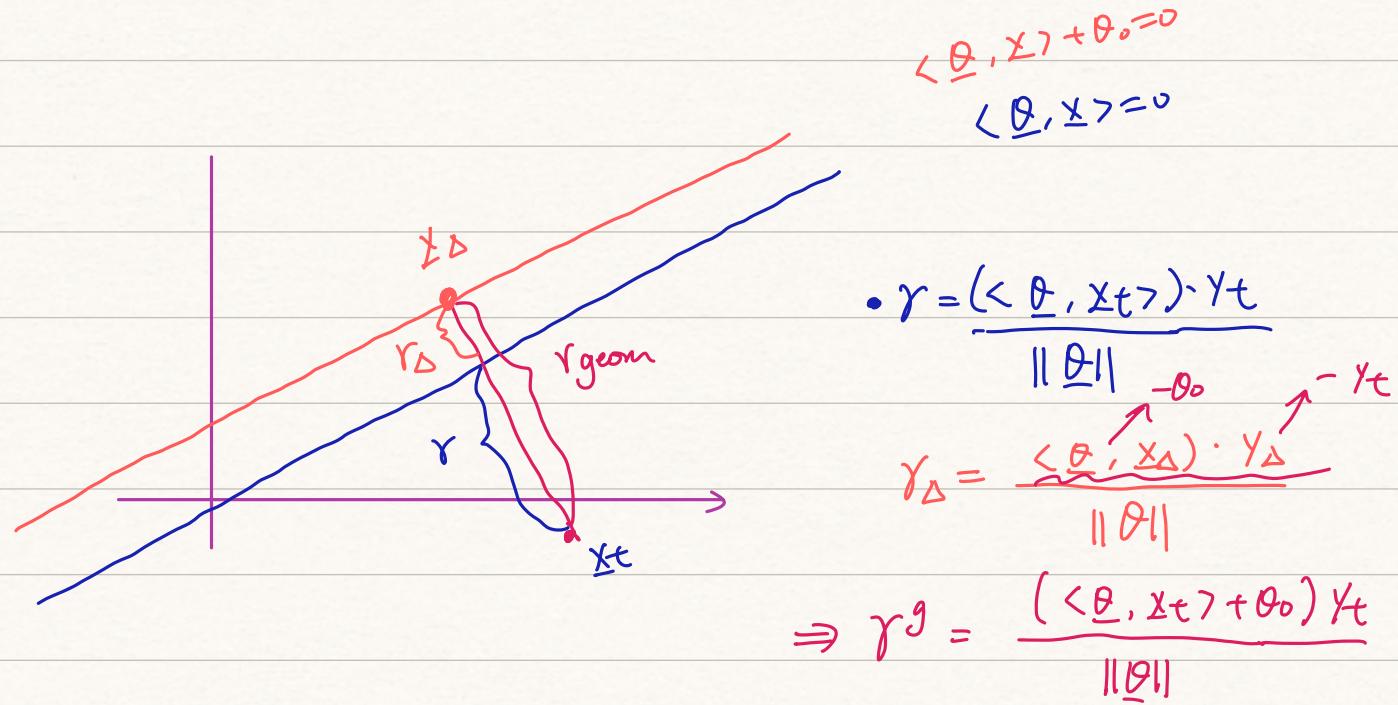
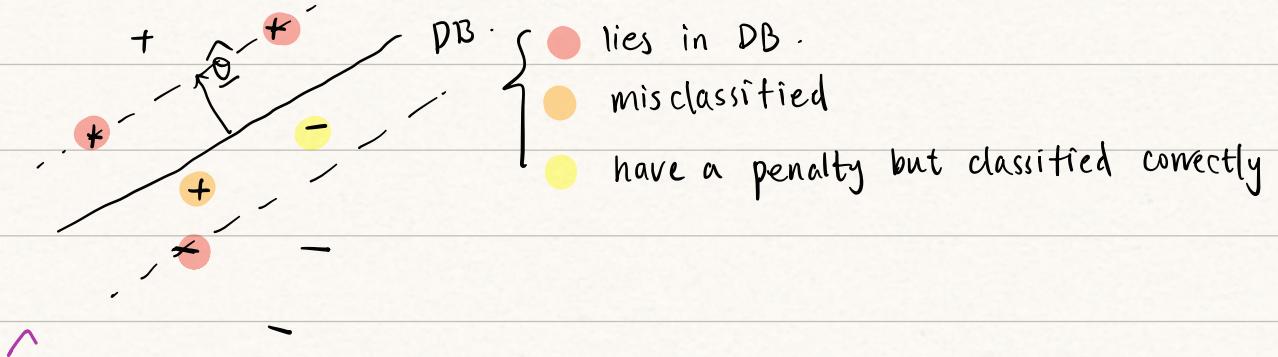
s.t. $y_t (\langle \underline{x}_t, \underline{\theta} \rangle + \theta_0) \leq 1 - \gamma_t \quad t=1, 2, \dots, n$

the meaning of $\underline{\theta}, \hat{\theta}_0$:

we have the { decision boundary
margin boundary

the distance from DB to MB is $\frac{1}{\|\underline{\theta}\|}$

② Three types of Support vectors



$$\text{For another case } \begin{cases} y_\Delta = y_t \\ r^g = r - r_\Delta \end{cases} \quad \checkmark$$

$$\Rightarrow r^g = \frac{(\langle \underline{\theta}, \underline{x}_t \rangle + \theta_0) y_t}{\|\underline{\theta}\|}$$

\Rightarrow for affine-SVM

$$\min \frac{1}{2} \|\underline{\theta}\|^2$$

$$\text{s.t. } y_t (\langle \underline{\theta}, \underline{x}_t \rangle + \theta_0) \geq 1$$

$$\Rightarrow \frac{1}{\|\underline{\theta}^*\|} \rightarrow \underline{r^g}$$