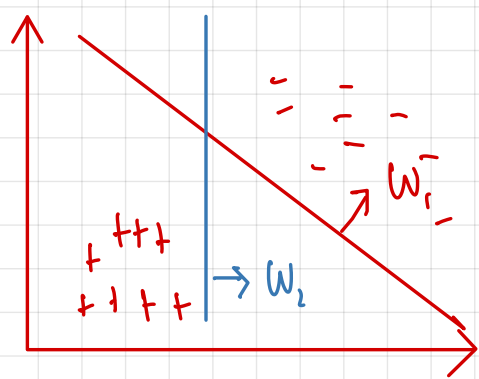# Support vector machine (SVM)
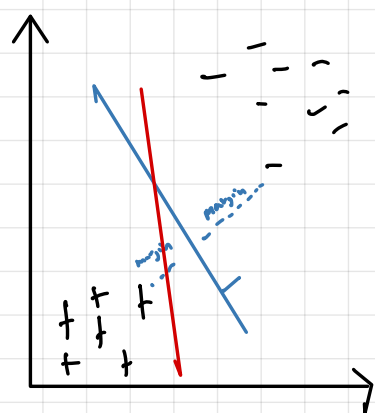
- Extendtion ของ perceptor
    - perceptron : find a hyperphane if it exists
    (how many hyperplane ?)
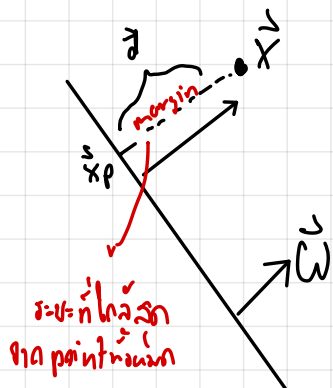


: many hyperplane but which one is the best?
    ↳ SVM

- SVM : find the maximum margin separating hyperplane : margin สูงสุด
    ↳ ระยะทางที่สั้นที่สุดจาก hyperplane ไปยัง point



- The margin $(\gamma)$: is the closest distance from the hyperplane to the closest points in either classes



$\vec{x}_p = \vec{x} - \vec{d}$ ; $\vec{x} = \vec{x}_p + d$

$\vec{w}^T \vec{x}_p + b = 0$  [$\vec{x}_p$ lies on the hyperplane] [$x_p$ อยู่กับ hyperplane]

$\vec{w}^T (\vec{x} - \vec{d}) + b = 0$ ; $d = \alpha \vec{w}$ ; for some $\alpha$ ( $\alpha$ เป็นตัวคูณ/ขยาย $w$ เพื่อให้มันมีความยาว )

$\vec{w}^T (\vec{x} - \alpha \vec{w}) + b = 0$

$$\alpha = \frac{w^T \vec{x} + b}{w^T w}$$

**find distance** $\Rightarrow \vec{d} = \left( \frac{\vec{w}^T \vec{x} + b}{\vec{w}^T \vec{w}} \right) \vec{w}$

$\Rightarrow \| \vec{d} \|_2 = \sqrt{\vec{d}^T \vec{d}} = \sqrt{(\alpha \vec{w})^T (\alpha \vec{w})}$

$= \sqrt{\alpha^2 w^T w}$

$= \alpha \sqrt{\vec{w}^T \vec{w}}$

$= \frac{\vec{w}^T \vec{x} + b}{\vec{w}^T \vec{w}} \sqrt{\vec{w}^T \vec{w}}$

$= \frac{\vec{w}^T \vec{x} + b}{\sqrt{\vec{w}^T \vec{w}}}$

# find margin

$$\gamma = \min_{\vec{x}} \frac{\vec{w}^T\vec{x}+b}{\sqrt{\vec{w}^T\vec{w}}} = \gamma(w,b) = \min_{\vec{x_i}} \frac{w^Tx_i+b}{\sqrt{w^Tw}}$$

$\hookrightarrow ||w||_2$

## Maximum margin hyperplane:

$$\vec{w},b = \max_{\vec{w},b}\left(\underbrace{\min_{\vec{x}} \frac{\vec{w}^T\vec{x}+b}{\sqrt{w^Tw}}}_{\text{margin term}}\right)$$

$\}$ ยังไม่ separate $\{+,-\}$

## - Maximum margin separating hyperplane:

$$\vec{w},b = \max_{\vec{w},b}\left(\min_{\forall\vec{x_i}} \frac{\vec{w}^T\vec{x_i}+b}{\sqrt{\vec{w}^T\vec{w}}}\right) \text{ s.t.}$$

$\underline{\forall i, \ y_i(\vec{w}^Tx_i+b) \geqslant 0}$ $\Big\}$ objective function

$\{$ constraints

- Simplification of finding such $\vec{w},b$: [Classify all point] แปลง 1,-1 ได้

  - $\vec{w},b = \max_{w,b}\left(\frac{1}{\sqrt{\vec{w}^T\vec{w}}} \underbrace{\min_{\forall\vec{x_i}} w^T\vec{x_i}+b}_{\text{margin}}\right)$ s.t. $\forall i, \ y_i(\vec{w}^Tx_i+b)\geqslant 0$
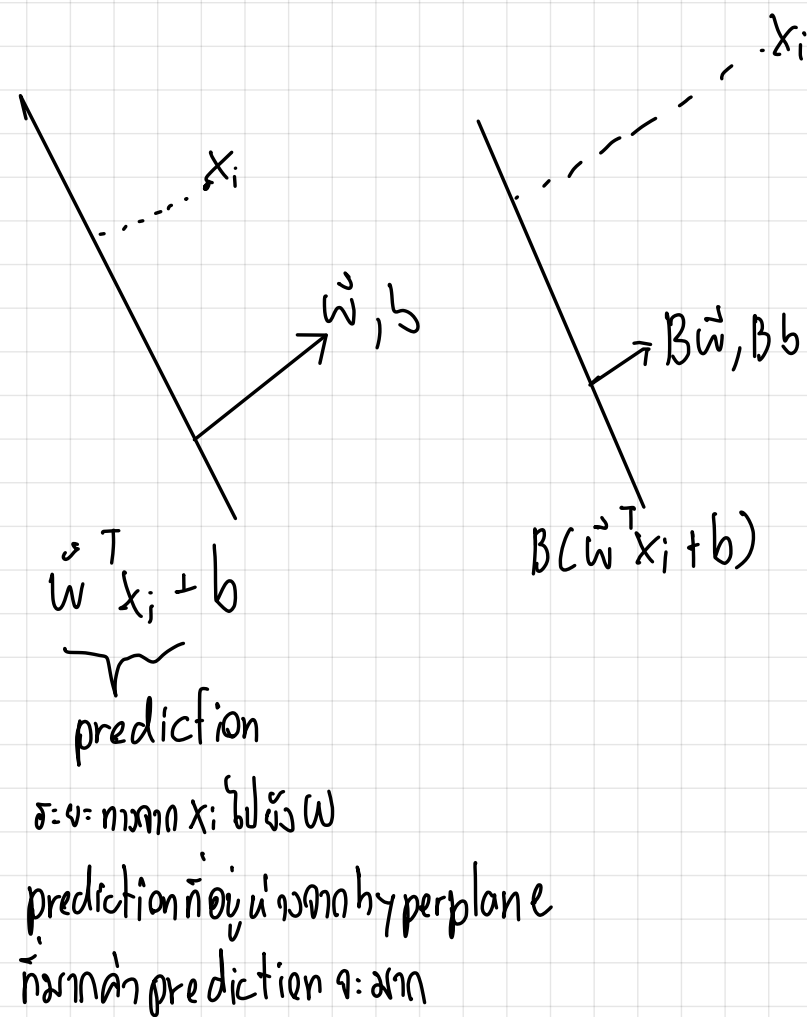
  $\Rightarrow$ Because the hyperplane is scale invariant,

  we can enforce

  $$\min_{\forall\vec{x_i}} |\vec{w}^T\vec{x_i}+b| = 1 \quad \text{(another constraint)}$$

$\Rightarrow$ $\boxed{\vec{w},b = \max_{w,b} \frac{1}{\sqrt{\vec{w}^T\vec{w}}}}$ s.t. $\min_{\forall\vec{x_i}} |\vec{w}^T\vec{x_i}+b|=1$

$\forall i, \ y_i(\vec{w}^Tx_i+b)\geqslant 0$ $\Big\}$ $\forall i, \ y_i(\vec{w}^Tx_i+b)\geq 1$



$\underset{\text{prediction}}{\underbrace{\vec{w}^Tx_i+b}}$

0:ฉ = ทางตาม $x_i$ ไปตาม $w$

prediction ก็อยู่ห่างจาก hyperplane

ก็ยากต่อ prediction 9:มาก

$\boxed{\vec{w}^*,b^* = \min_{\vec{w},b} \underset{||\vec{w}||_2}{\sqrt{\vec{w}^T\vec{w}}} = \min_{\vec{w},b} \vec{w}^T\vec{w}}$

$||w||_2 = \sqrt{w^Tw}$

$||w||_2^2 = w^Tw$

$w^Tw \propto \sqrt{w^Tw}$

$\hookrightarrow$ we can find them by using QCQP

## Final Formulation:

quadratic function

$$\boxed{\vec{w},b = \min_{\vec{w},b} \underset{\text{}}{(w^Tw)} \text{ s.t. } \forall i, \underset{\text{linear inequalities}}{(y_i(\vec{w}^Tx_i+b)} \geq 1}$$

Goal of SVM is to find $\vec{w},b$ according to the formulation

To find $\vec{w},b$ we can use Quadratic Programming Solver / QCQP)

Interpretation: find $\vec{w}, b$ where $\vec{w}$ is of minimum magnitude
such that all points lie at least 1 unit away from
the hyperplane on the correct side. [w is the simplest solution]



$\gamma > 1$   $\rightarrow x^* = \gamma = 1$

$\gamma = 1$

$\gamma = 1$

$\rightarrow x^* = \gamma = 1$

$\gamma > 1$

$(w', b')$

$\vec{w}, b$

$(w', b') = (Bw, Bb)$

Note: There always exist $x_i$ st.
$|\vec{w}^T x_i + b| = 1$ (margin)
($x_i$ is the closest point)

The vector $\vec{w}$ (and $b$) supports the closest point.

Note $y_i (\vec{w}^T \vec{x_i} + b) \geq 1$

ต้องมีค่ามากกว่า 1

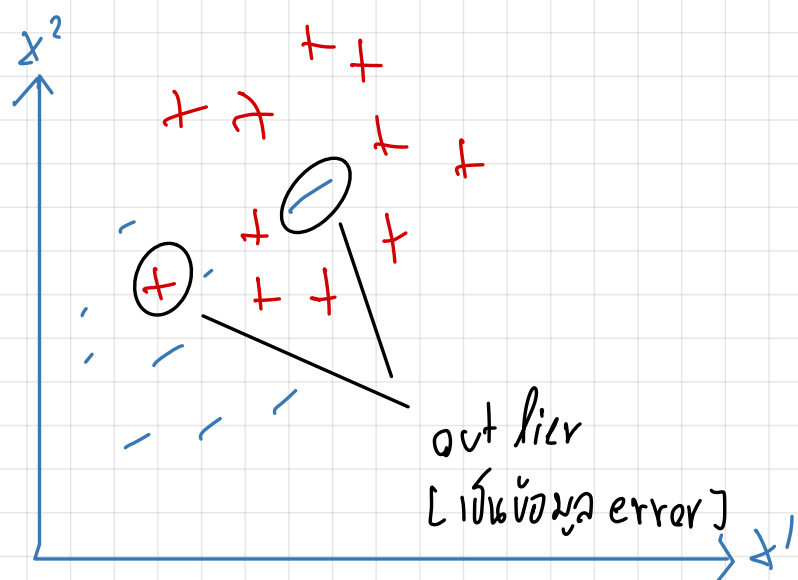$\gamma = 1$ เสมอ [ค่าใกล้เส้น]

$\gamma(w, b) = w^T x^* + b$
$\gamma'(w', b') = w'^T x^* + b'$
$\quad = Bw^T x^* + Bb$
$\quad = B(w^T x^* + b)$

<u>Support Vector</u> : The vector $w$ and $b$ supports the closest point $x^*$

# <u>Problem Percepton</u> [low Dimensional Data set]



$x^2$

out lier
[เป็นข้อมูล error]

$x^1$

$\therefore$ can't use percepton [infinte loop]

# <u>Dealing with non-linearly separable data:</u>

- IDEA : We may sacrifice some outlier(s)
  - in order to place the hyperplane

# Softing the constraints [relax constraint]

$\forall i, \ y_i(w^T x_i + b) \geq 0.7$ [ห่าง atleast 0.7]
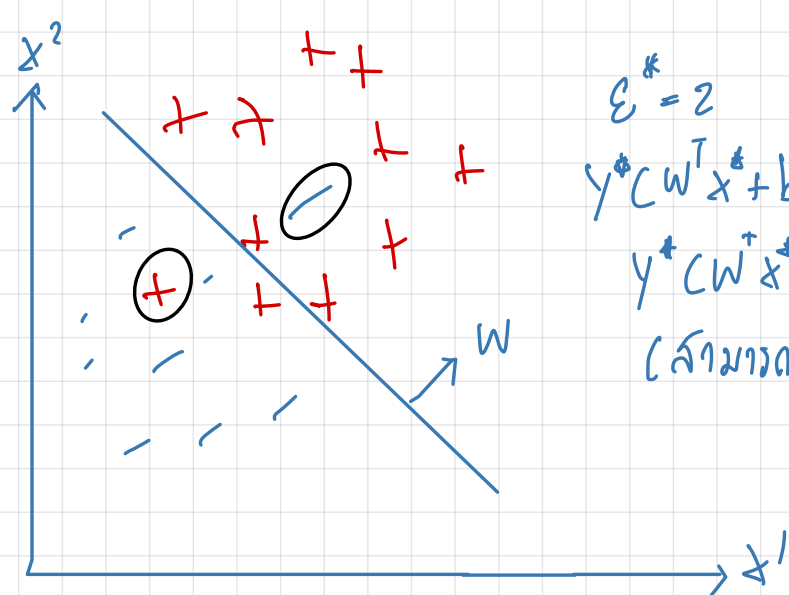
$\qquad\qquad\qquad < 0$ [can classifly wrong side]

- create new variable for constraints.
- **Fix** outlier
    - we allow the constraints to be soften slightly with the introdection of slack variable (หย่อน)

$$\varepsilon_i \geq 0 \ , \forall i$$
(\psi)

$\forall i, \ y_i(w^T x_i + b) \geq 1 - \varepsilon_i$ : corespanse ของแต่ละ $= y_i$ ; $\varepsilon_i \geq 0$

$\qquad\qquad\qquad$ ($\varepsilon_i$ ไม่จำเป็นต้อง เหมือน กัน)



$\varepsilon^* = 2$

$y^*(w^T x^* + b) < 0$

$y^*(w^T x^* + b) \geq -1$

(สามารถดวงไว้ ผิดฝั่ง ได้)

hyper parameter [หน้าที่เรา set]

จากเดิม → New

$$w^*, b^*, \varepsilon^* = \min_{w, b, \varepsilon} \left[ w^T w + C \sum_{i=1}^{n} \varepsilon_i \right] \ ; \ C \geq 0$$

$\qquad$ ($\varepsilon = (\varepsilon_1, \varepsilon_2, \ldots, \varepsilon_n)$)

if set $C = \infty$ [no relax] $\forall i, \ y_i(w^T x_i + b) \geq 1 - \varepsilon_i$

if set $C = 0$ [free]

- The slack variable $\varepsilon_i$ allow $x_i$ to be closer the hyperplane or even be on the wrong side but there is a penalty in the objective function for such slack. (exchange cost)

Penalty $C \to +\infty$, svm will try to make all the points to be on correct side.

$\qquad\quad$ $C \to 0$, svm will sacrifice some points

# Unconstrainted formulation:

- we set $\varepsilon_i$ as follows:

$$\varepsilon_i = \begin{cases} 1 - y_i(w^Tx_i + b) & \text{if } y_i(w^Tx_i + b) < 1 \\ 0 & \text{if } y_i(w^Tx_i + b) \geq 1 \end{cases}$$

- In other words

$$\varepsilon_i = \max\left(1 - y_i(w^Tx_i + b), 0\right)$$

- Hence, we can rewrite

$1 - h_{w,b}(x_i)$

$$w^*, b^* = \min_{w, b} \left( w^Tw + C \sum_{i=1}^{n} \max\left(1 - y_i(w^Tx + b), 0\right) \right)$$

$\underbrace{\phantom{w^Tw}}$ regularizer $\ell_2$

$\underbrace{\phantom{max}}$ loss function (hinge loss)

SVM with soft constraints

Many ML algorithms can be expressed via the optimization problem of the form

$$w^*, b^* = \min_{w} \left( \frac{1}{n} \sum_{i=1}^{n} \ell(h_w(x_i), y_i) + \lambda r(w) \right)$$

average loss

loss function of h() with w as parameter

regularizer

**Similar**

put $\lambda$ in regularizer and cancel $C$ cause $\lambda$, $C$ is balance in each other.

Note: regularizer like pior so it used for MAP

if $k = \frac{1}{C}$

$$w^*, b^* = \min_{w, b} \left( w^Tw + C \sum_{i=1}^{n} \max\left(\overline{1 - y_i(w^Tx + b), 0}\right) \right)$$

$$= \min\left( \lambda w^Tw + \sum_{i=1}^{n} \max(1 - y_i(w^Tx_i + b), 0) \right)$$

$\ell_2$-regularizer          hinge loss

if $C = \frac{1}{\lambda}$   $\hookrightarrow \|w\|_2^2 = w^Tw = (w^1)^2 + (w^2)^2 + \ldots + (w^d)^2$

$$\min_{w, b} \left( w^Tw + C \sum_{i=1}^{n} \max\left(\overline{1 - y_i(w^Tx + b), 0}\right) \right)$$

$n$

# Empirical risk minimization

- Many learning algorithms can be written in a form of an Optimization problem with objective to minimize some loss function $l$ and a regularizer $r()$
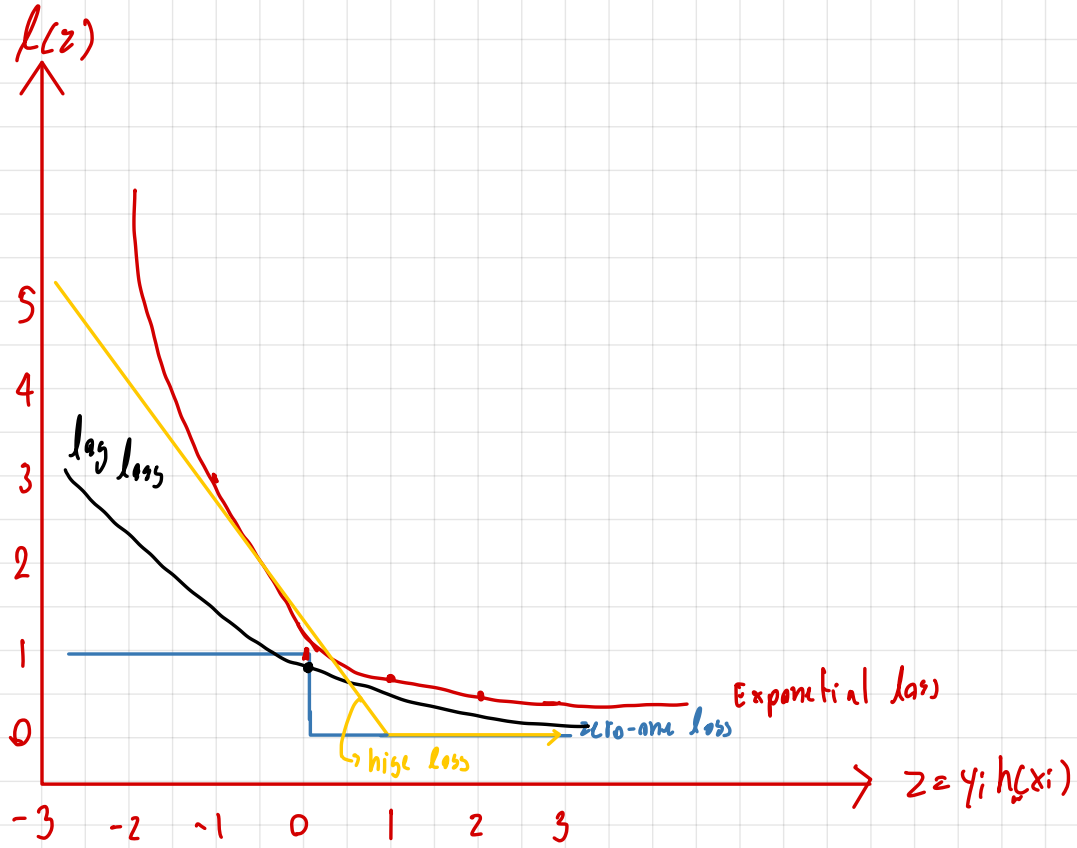
$$w^* = \min_{w} \frac{1}{n} \sum_{i=1}^{n} \underbrace{l(h_w(x_i), y_i)}_{\text{loss}} + \underbrace{\lambda r(w)}_{\text{regularizer}} \quad \text{like pior}$$

$\underbrace{\qquad\qquad\qquad\qquad}_{\text{average loss}}$

**Example**

$$W_{Map} = \boxed{\underset{w}{\arg\min} \; \frac{1}{n} \sum_{i=1}^{n} \underbrace{(w^T x_i - y_i)^2}_{Sq-loss} + \underbrace{\lambda \; w^T w}_{\text{regularizer}}}$$

# Commonly Used Binary Classification loss functions. (not for regression)

| loss $l(h_w(x_i, y_i))$ | Usage | Comments |
|---|---|---|
| Hige-loss $[max(1-h_w(x_i)y_i, 0)]^P$ | - Standard [P=1] <br> - Hige less SVM (P=2) C differentiable | When used for standard SVM, the loss function denotes the size of the margin between the linear separator and its closest points in either class |
| log-loss $(log(1 + e^{-h_w(x_i)y_i})$ | - logistic regression | - One of the most popular loss functions in ML, since its output are well calibrated probs. |
| Exponential loss $e^{-h_w(x_i)y_i}$ | - Ada Boost | - This loss function is very aggressive it increases exponentialy with Vale of $-h_w(x_i)y_i$ it's sensitive to noisy data |
| zero - one loss $\delta(sign \; h_w(x_i) \neq y_i)$ | - Actual classification loss | - non continous cannot optimize in practice |

$\ell(z)$

lag loss

Exponetial loss

zero-one loss

hige loss

$z = y_i h(x_i)$

- Assume $y_i \in \{-1, 1\}$

hige loss
- miss when data in margin or high margin / loss will be linear

zero-one loss

if $z > 0$  loss = 0 : classify

if $z < 0$  loss = 1 : miss classify [count 1]

exponetina loss

if $z > 0$ loss approch to 0 : classify

if $z < 0$ loss will jump high loss and try to decrease it : miss classify

lag-loss
if $z < 0$ loss will likely linear
: miss classify

# Regression

| loss $\ell(h_w(x_i), y_i)$ | Comment |
|---|---|
| Squared loss $(h_w(x_i) - y_i)^2$ | - most popular regression loss function<br><br>- Estimate mean label<br><br>- Pros: Differentiable everywhere<br>- Cons: Sensitive to outliers / nosie |
| Absolute loss $\lfloor h_w(x_i) - y_i \rfloor$ | - Also, very popular loss function<br>- estimetes median label<br>- Pros : less sensitive to neise<br>- Cons: Not differntiable at 0 |
| Huber loss $\frac{1}{2}(h_w(x_i) - y_1)^2$ if $\lvert h_w(x_i) - y_i \rvert < \delta$ - otherwise $\delta(\lvert h(x_i) - y_i \rvert - \frac{\delta}{2})$ | - A.K.A. smooth absolute loss<br> - Pros Best of "Both worlds" (squared + Absolute loss)<br>- Once - differentiable<br>- Takes on behavior of sq-loss when loss is small, and absolute loss when loss is large. |

log - Cosh

$log(Cosh(h(x_i)-y_i)),$

$Cosh(x) = \dfrac{e^x + e^{-x}}{2}$

- Pros Similar to Huber loss, but twice differtiable every where

# Regularizer

$w^* = \min\limits_{w} \dfrac{1}{n} \sum\limits_{i=1}^{n} \ell(h_w(x_i), y_i) + \lambda r(w)$

like pior

regularizer

loss

average loss

∴ data ที่เรามีมาจาก การเก็บของ population
ทำให้ไม่รู้ว่า data นั้น มีการกระจายตัวปกติหรือไม่เลยต้อง มี
ตัว regularizer (over fit ใน collect data แต่ไม่ใช่นอก data)
        ↳ เป็นต้องการไม่ให้ overfit

- without regularizer, we always end up with only minimizing loss function on traning data. This often *leads* to "overfitting".

- In general, regularizer corresponds to the notion of simplicity / complexity of the solution to the optimization problem
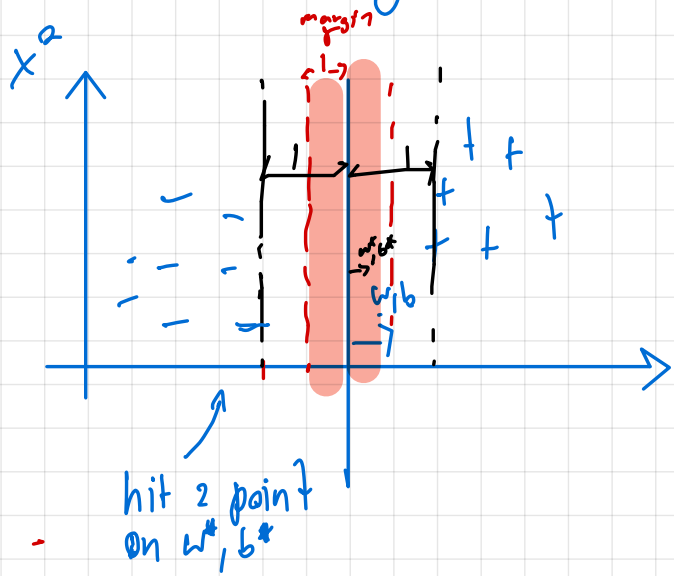
SVM: $\min\limits_{w, b} w^T w$  → regularizer

subject to $\forall_i, y_i (w^T x_i + b) \geqslant 1$

loss

# Maximum margin Solution $w^*, b^*$

$x^a$

norm $w^*, b^*$ ลดลง เมื่อ 8 เพิ่มขึ้น

hit 2 point
on $w^*, b^*$

---

$$\min_{w} \frac{1}{n} \sum_{i=1}^{n} l(h_w(x_i), y_i) + \lambda r(w)$$

for any value $\lambda \geq 0$

there exists $B \geqslant 0$

and vice varse

$\Longleftrightarrow$

$$\min_{w} \frac{1}{n} \sum_{i=1}^{n} l(h_w(x_i), y_i)$$

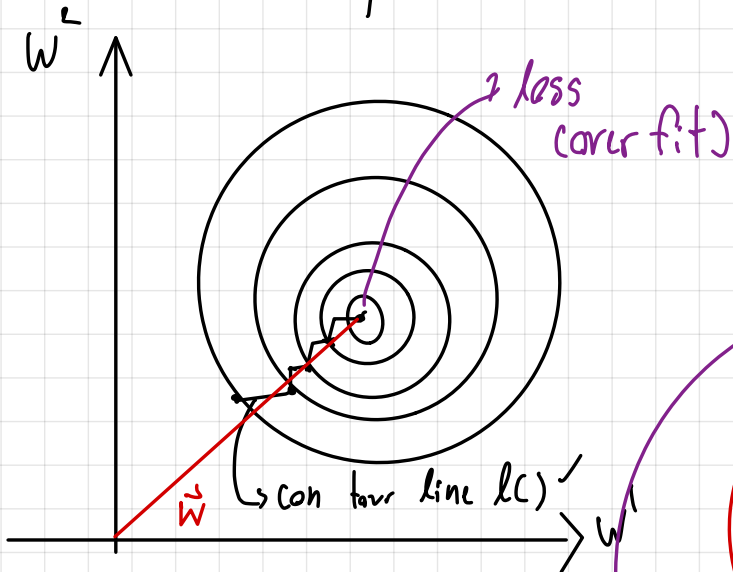$$\text{subject to } r(w) \leq B$$

$$\min_{w} \frac{1}{n} \sum_{i=1}^{n} l(h_w(x_i), y_i)$$

$$\text{subject to } r(w) \leq B$$

$$r(w) = w^T w = \|w\|_2^2 = (w^1)^2 + (w^2)^2 \leq B$$

## <u>l2 regularizer</u>

พิจารณาตัว loss only

$\to x^2 + y^2 \leq B \to$ สมการวงกลม

กรณีพิจารณา B ด้วย

Note: ตัว B จะเป็นตัว ลดความ overfit
จาก point loss ยิ่ง radius B มาก
ต่ายิ่ง overfit ( R : B มาก overfit)
(R : B น้อย น้อยลง)



$w^2$

loss (over fit)

con tour line $l(C)$

$\tilde{w}$

$w^1$

$w^2$

con tour line $l(C)$

$\tilde{w}$

$w^1$

$\sqrt{B}$

พิจารณา ใน วงกลม

put constant (decrease over fit)

# $l_1$ - regularizer

สมการ เส้นตรง

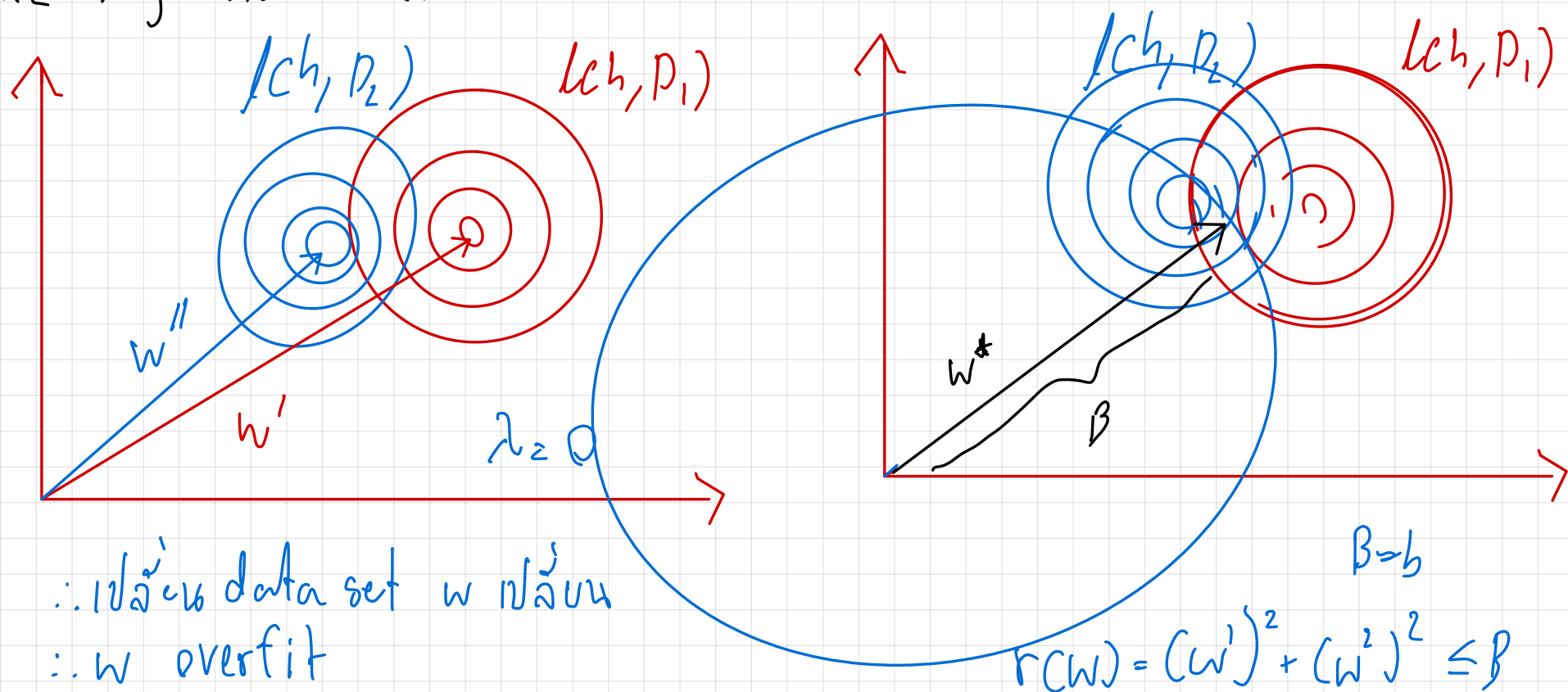$r(w) = \|w\|_1 = |w^1| + |w^2| \le B$



$\begin{bmatrix} w^1 = 0 \\ w^2 = 100 \end{bmatrix}$

solution (w)

เป็น เว้ feature เดียว (Coordinate ที่สำคัญ น้อยตัดออก)

w will be sparse solution

-set coordinate ที่ไม่สำคัญ ออก

# $l_2$ - Regularizer Review



$l(h, D_2)$    $l(h, D_1)$

$w''$   $w'$   $\lambda = 0$

$l(h, D_2)$    $l(h, D_1)$

$w^*$   $B$

$B \to b$

$r(w) = (w^1)^2 + (w^2)^2 \le B$

∴ เปลี่ยน data set w เปลี่ยน
∴ w overfit

$w^*$ is optimal solution within the constraint

$w^*$ จะไม่อยู่ ที่ minimum ของ loss
ถ้า เพิ่ม ค่า B จะสามารถ ไปที่ minimum ของ loss
ทั้ง $D_1$, $D_2$ ได้