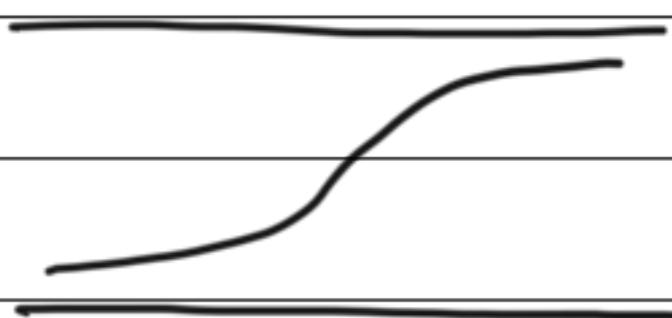


Consider logistic regression

$P(y=1|x; \theta)$ is modeled by $h_\theta(x) = g(\theta^T x)$



: if $\theta^T x \gg 0$, $y = 1$] very
if $\theta^T x \ll 0$. $y = 0$] confident
prediction

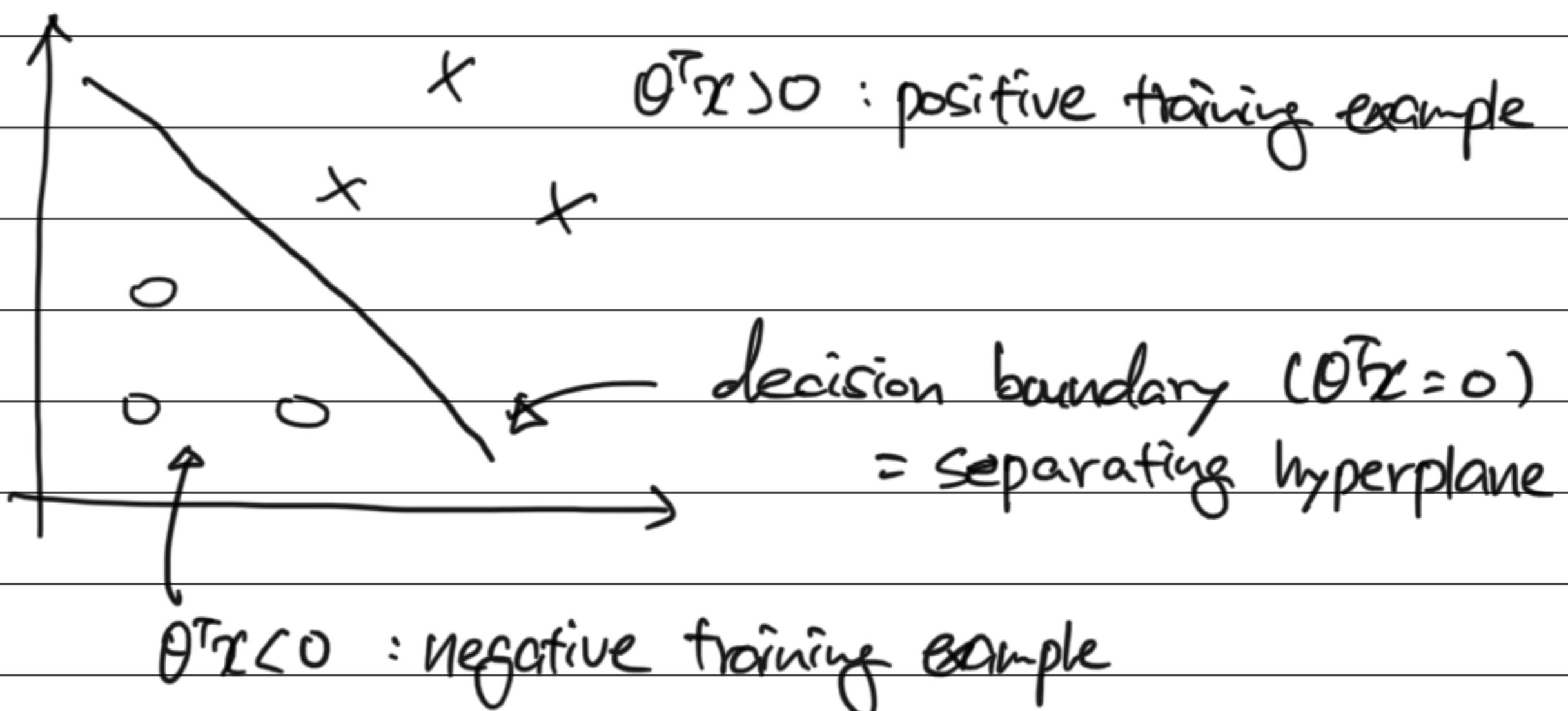
Given a training set.

if we can find θ so that $\theta^T x^{(i)} \gg 0$

whenever $y^{(i)} = 1$

$\theta^T x^{(i)} \ll 0$ whenever $y^{(i)} = 0$

⇒ This seems to be a nice goal to aim for. And we'll soon formalize this idea using the notion of functional margin



* If a point is far from the separating hyperplane, we may be significantly more confident in our predictions

• Notation

We will be considering a linear classifier for a boundary classification problem with labels y and features x

$$y \in \{+1\}$$

$$h_\theta(x) = g(\theta^T x) \rightarrow h_{w,b}(x) = g(w^T x + b)$$

b takes the role of what was previously θ_0 , and w takes the role of $[\theta_1, \dots, \theta_n]^T$

• Functional margin

Given a training example $(x^{(i)}, y^{(i)})$,

We define the functional margin of (w, b)

$$\Rightarrow \hat{\gamma}^{(i)} = y^{(i)}(w^T x + b)$$

if $y^{(i)}=1$, then for the functional margin to be large (very confident), we need $w^T x + b$ to be a large positive number.

if $y^{(i)}=-1$, we need $w^T x + b$ to be a large negative number.

Moreover, if $y^{(i)}(w^T x + b) > 0$, then our prediction on this example is correct.

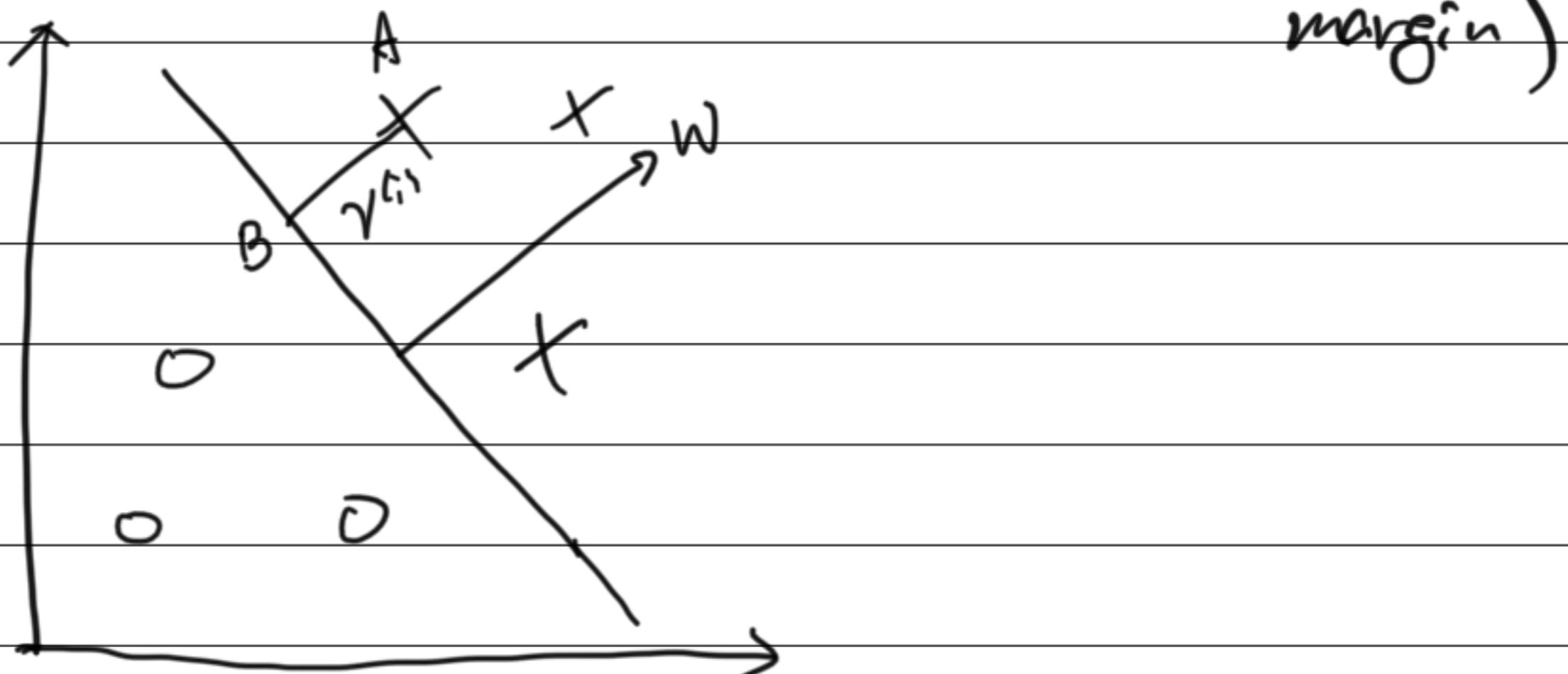
Given a training set $S = \{(x^{(i)}, y^{(i)}) : i=1 \dots m\}$, we define the functional margin of (w, b) with respect to S to be the smallest of the functional margins of the individual training examples.

$$\hat{\gamma} = \min_{i=1 \dots m} \hat{\gamma}^{(i)}$$

* Replacing (w, b) with $(2w, 2b)$ results in multiplying our functional margin by a factor of 2.

\Rightarrow It seems that by exploiting freedom to scale w and b . We can make the functional margin arbitrarily large without really changing anything meaningful.

- Geometric margin (Generalization of the functional margin)



The decision boundary corresponding (w, b) is shown, along with the vector w .

(Note that w is orthogonal to the separating hyperplane)

Consider the point A , which represents the input $x^{(i)}$ of some training example with label $y^{(i)}=1$.

Its distance to the decision boundary, $\gamma^{(i)}$, is given by the line segment AB .

How can we find $\gamma^{(i)}$?

- $w/\|w\|$ is a unit-length vector pointing in the same direction as w
- A represents $x^{(i)}$
- B is given by $x^{(i)} - \gamma^{(i)} \cdot w/\|w\|$

This point lies on the decision boundary

All points x on the decision boundary satisfy $w^T x + b = 0$.

$$w^T \left(x^{(i)} - \gamma^{(i)} \frac{w}{\|w\|} \right) + b = 0$$

$$w^T x^{(i)} + b = \underbrace{\gamma^{(i)} \frac{w^T w}{\|w\|}}_{= \gamma^{(i)} \|w\|^2}$$

$$w^T x^{(i)} + b = \gamma^{(i)} \|w\|$$

$$\Rightarrow \gamma^{(i)} = \frac{w^T x^{(i)} + b}{\|w\|} = \left(\frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|}$$

We define the geometric margin of (w, b) with respect to a training example $(x^{(i)}, y^{(i)})$ to be

$$\gamma^{(i)} = y^{(i)} \left(\left(\frac{w}{\|w\|} \right)^T x^{(i)} + \frac{b}{\|w\|} \right)$$

Note that, if $\|w\| = 1$, the functional margin equals the geometric margin

* Given a training set $S = \{(x^{(i)}, y^{(i)}) ; i=1, \dots, m\}$,

the geometric margin of (w, b) with respect to S to be the smallest of the geometric margins on the individual training examples.

$$\gamma = \min_{x \in S} \gamma$$

• The optimal margin classifier

: Find a decision boundary that maximize the geometric margin

$$\max_{\gamma, w, b} \gamma$$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq \gamma, \quad i=1, \dots, m$$

$$\|w\| = 1 \leftarrow \text{non-convex}$$

⇒ Transform the problem into a nicer one

$$\max_{\gamma, w, b} \frac{\gamma}{\|w\|} \leftarrow \text{non-convex}$$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq \hat{\gamma}, \quad i=1, \dots, m$$

We can add an arbitrary scaling constraint on w and b without changing anything.

⇒ We will introduce the scaling constraint that the functional margin of w, b with respect to the training set must be 1:

$$\hat{\gamma} = 1$$

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

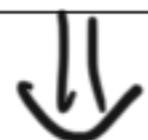
$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1$$

Maximizing $\frac{1}{\|w\|} = \frac{1}{\|w\|}$ is the same thing as

minimizing $\|w\|^2$

$$\left. \begin{array}{l} \min_{w,b} \|w\|^2 \end{array} \right\}$$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1$$



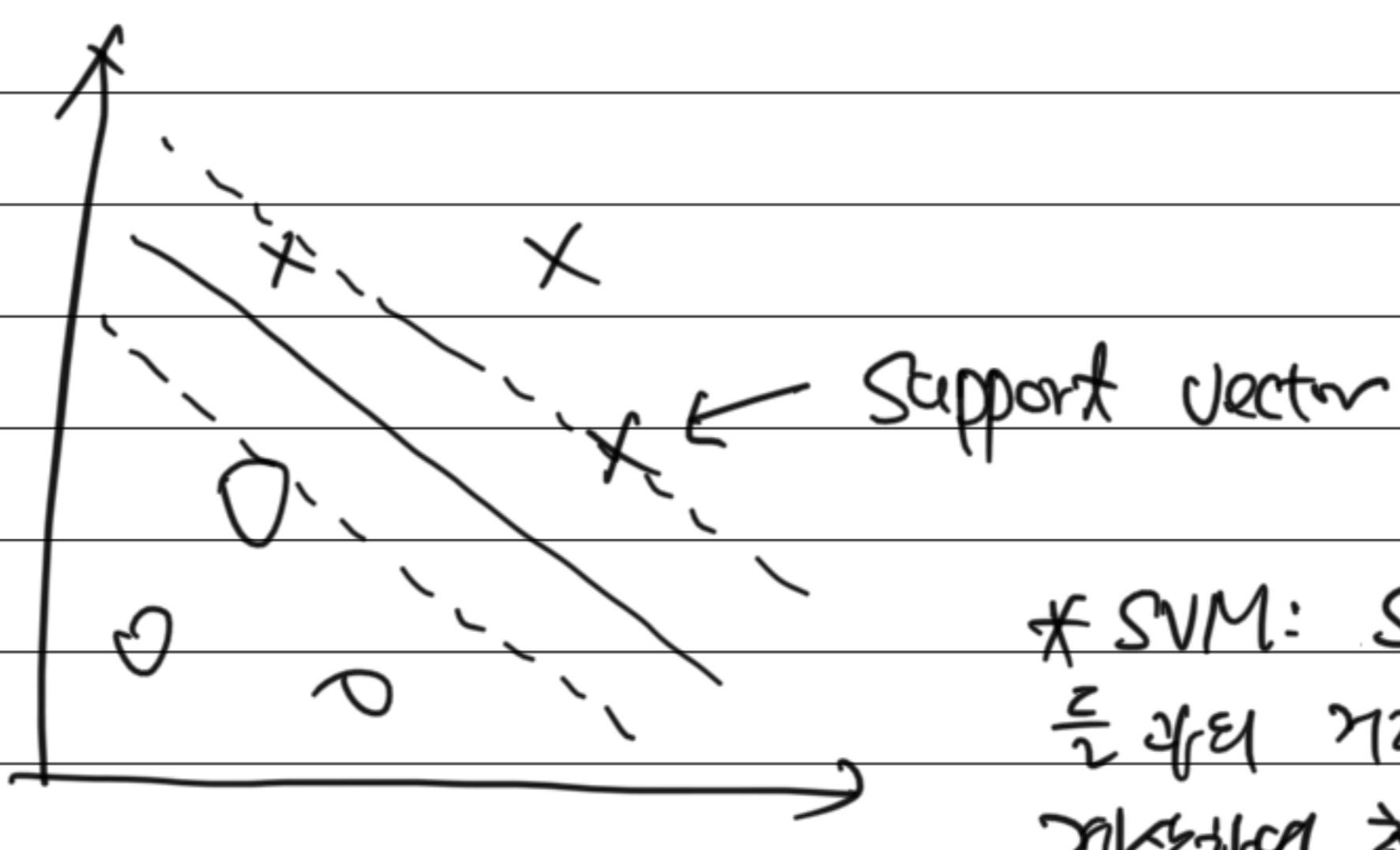
Nicer form:

$$\left. \begin{array}{l} \min_{w,b} \frac{1}{2} \|w\|^2 \end{array} \right\}$$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1, i=1,\dots,m$$

This optimization problem can be solved using quadratic programming (QP).

↳ optimal margin classifier (SVM)



* SVM: Support vector
는 즉각적인 거리(margin)을
가지도록 하는 것을 말함.