# 6. Support Vector Machines (SVM)

Given a training set $(x_1, y_1), \ldots, (x_n, y_n)$

$x_i \in \mathbb{R}^p$ : data points

$y_i \in \{-1, +1\}$ : class membership (2 classes/groups)

Assume: Exists a separating hyperplane $H$

s.t. $\{x_i / y_i = 1\}$ and $\{x_i / y_i = -1\}$

are separated by $H$.     (will be released later)

$\boxed{\text{Fig 1}}$

## 6.1. Hyperplanes and Margins.

Representing hyperplanes in $\mathbb{R}^p$:

a) Given $a \in \mathbb{R}^p$

$\{x \in \mathbb{R}^p / a^T x = 0\}$ is the $(p-1)$-dim

(linear space orth. to $a$.

$\boxed{\text{Fig 2}}$

b) Given $a \in \mathbb{R}^p, b \in \mathbb{R}$. Consider

$\{x \in \mathbb{R}^p / a^T x - b = 0\}$

It holds

$a^T x - b = 0 \iff a^T x - \dfrac{a^T a}{\|a\|^2} b = 0$

$\iff a^T \left(x - \dfrac{b}{\|a\|^2} a\right) = 0$   $\boxed{\text{Fig 3}}$

Hence: $\{x / a^T x - b = 0\}$ is a lin. subspace
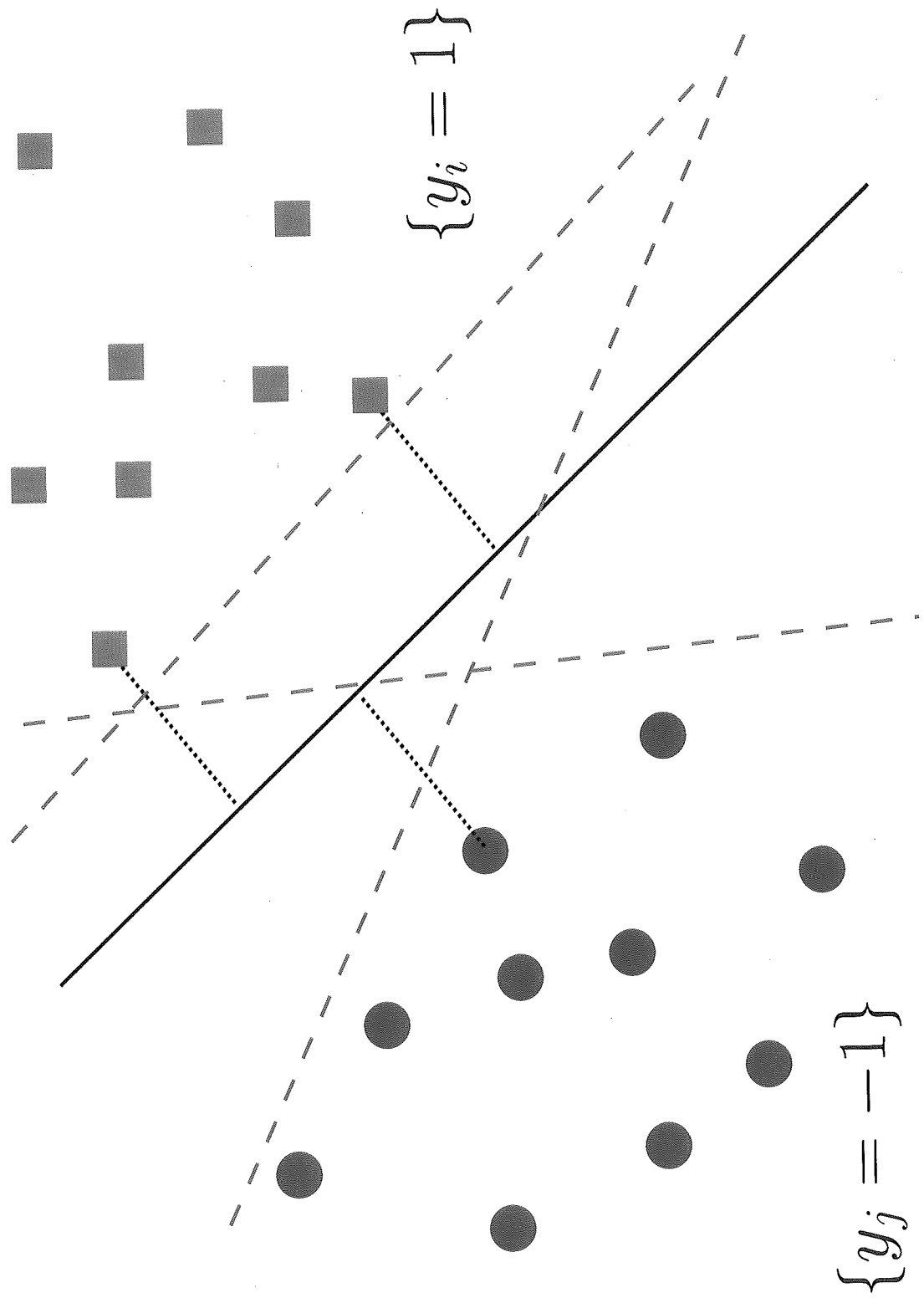
shifted by $\dfrac{b}{\|a\|^2} a$, a hyperplane.

— 1 —

Fig. 1

$\{y_i = 1\}$

$\{y_j = -1\}$

Fig 2



$$\{x \mid a^T x = 0\}$$

$a$

Fig. 3

$$\frac{b}{\|a\|^2}a$$

$$a$$

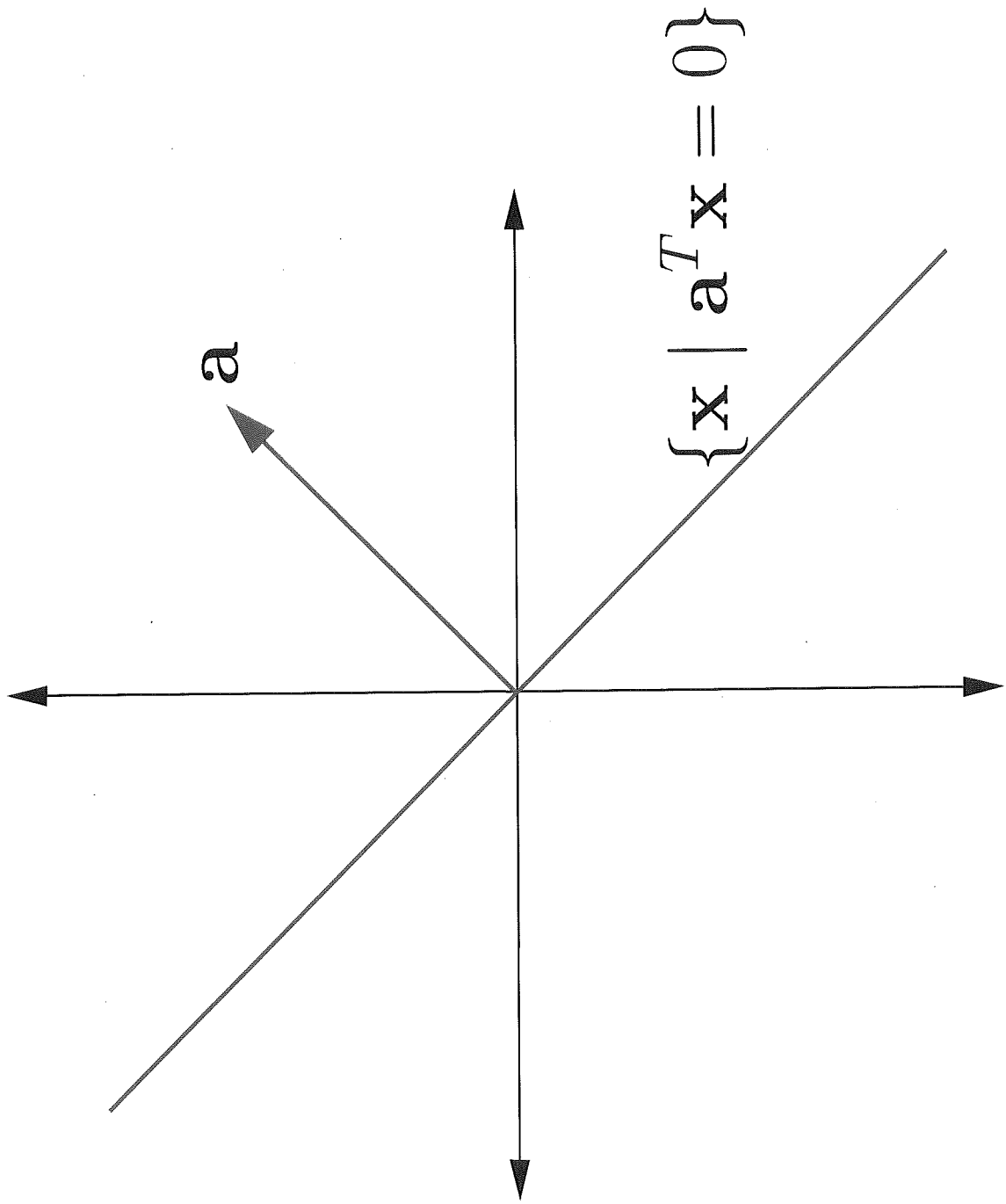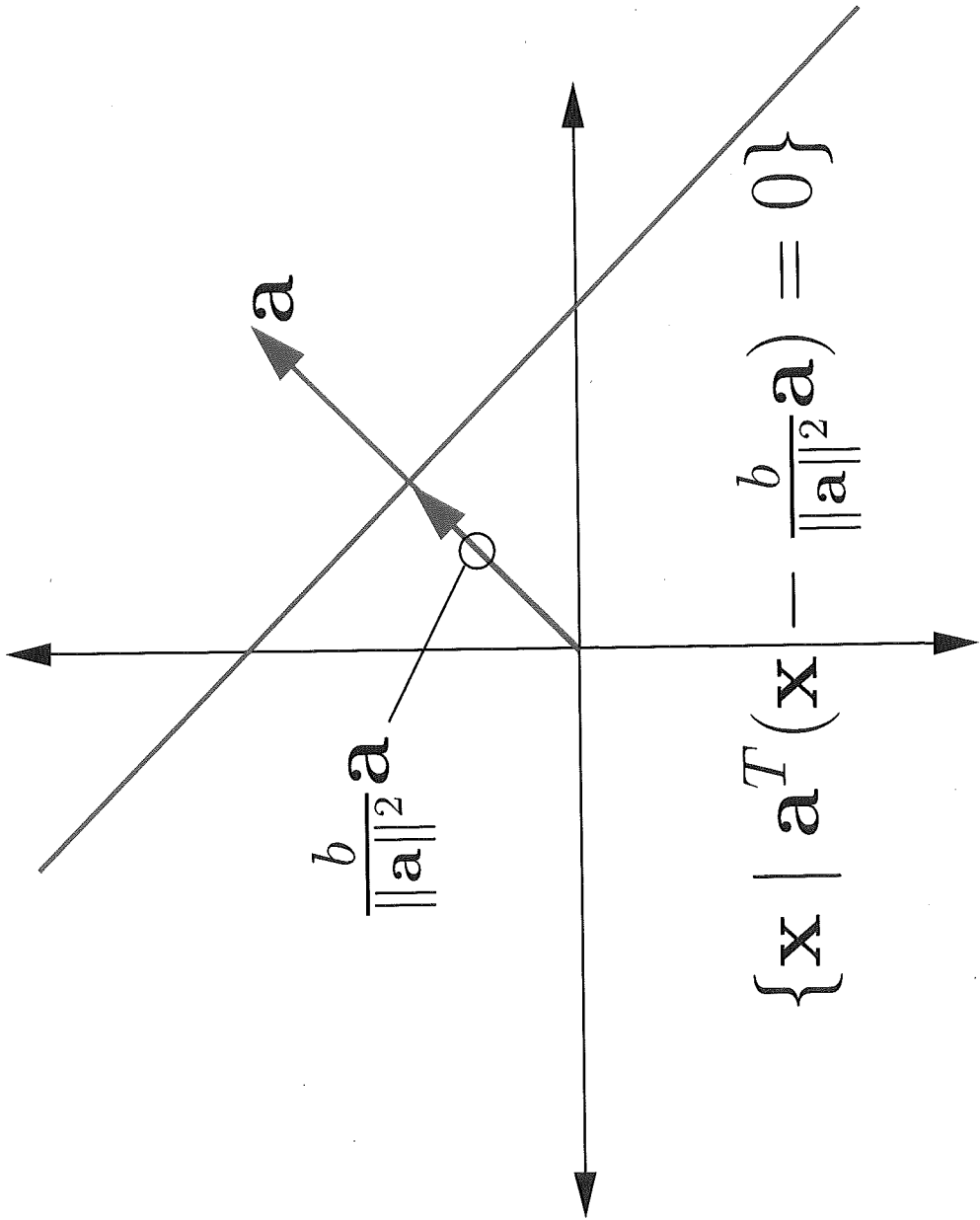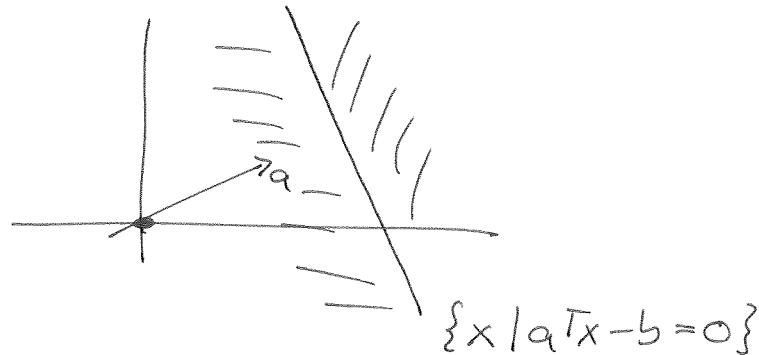$$\{x \mid a^T(x - \frac{b}{\|a\|^2}a) = 0\}$$

$\{x \in \mathbb{R}^P \mid a^T x \gtreqless b\}$ is called half-space :



$\{x \mid a^T x - b = 0\}$

c) Given $a \in \mathbb{R}^P$, $b_1, b_2 \in \mathbb{R}$

Distance between $H_1 = \{a^T x - b_1 = 0\}$, $H_2 = \{a^T x - b_2 = 0\}$



$\{a^T x - b_1 = 0\}$    $\{a^T x - b_2 = 0\}$

Both hyperplanes are parallel and orthogonal to $a$.

Pick $x_1, x_2$ such that

$x_1 = \lambda_1 a$       $x_2 = \lambda_2 a$

$a^T x_1 - b_1 = 0$       $a^T x_2 - b_2 = 0$

Then

$\lambda_1 a^T a - b_1 = 0$       $\lambda_2 a^T a - b_2 = 0$

$\lambda_1 \|a\|^2 - b_1 = 0$       $\lambda_2 \|a\|^2 - b_2 = 0$

$\lambda_1 = \dfrac{b_1}{\|a\|^2}$       $\lambda_2 = \dfrac{b_2}{\|a\|^2}$

and

$$\|x_2 - x_1\| = \|\lambda_2 a - \lambda_1 a\| = |\lambda_2 - \lambda_1| \|a\|$$

$$= \left| \frac{b_2}{\|a\|^2} - \frac{b_1}{\|a\|^2} \right| \|a\| = \frac{1}{\|a\|} |b_2 - b_1|$$

Hence the distance between $H_1$ and $H_2$

is $\qquad \frac{1}{\|a\|} |b_2 - b_1|$

d) Given $a \in \mathbb{R}^P$, $b \in \mathbb{R}$, $x_0 \in \mathbb{R}^P$

Distance between $H = \{x \mid a^T x - b = 0\}$ and point $x_0$.

Consider auxiliary hyperplane containing $x_0$.

$$H_0 = \{x \mid a^T x - b_0 = 0\} = \{x \mid a^T x - a^T x_0\}$$

(since $b_0 = a^T x_0$ since $a^T x_0 - b_0 = 0$)

By c), the distance between $H$ and $H_0$ is

$$\frac{1}{\|a\|} |b - a^T x_0|.$$

This distance is called <u>margin</u> of $x_0$.
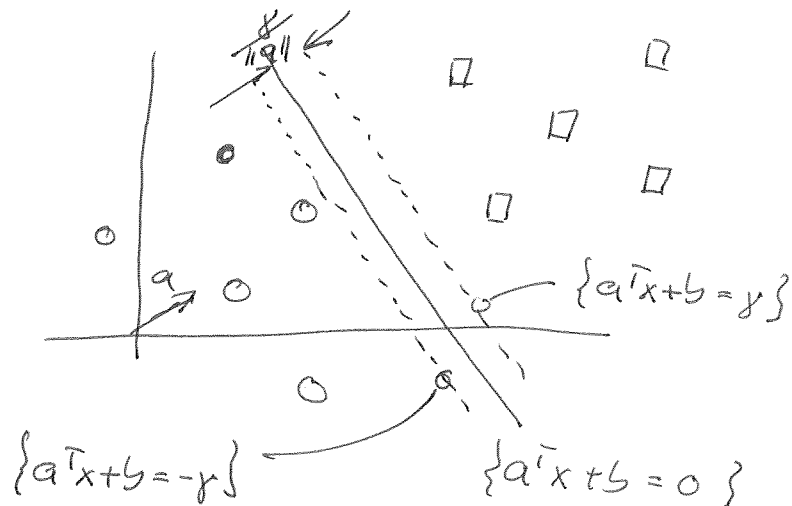
## 6.2 The optimal margin classifier

Given a training set $(x_1, y_1), \dots, (x_n, y_n)$, $x_i \in \mathbb{R}^p, y_i \in \{-1, 1\}$

Assume there exists a separating hyperplane.

$$\{x \mid a^T x + b = 0\}$$

Then
$$\left. \begin{array}{l} y_i = +1 \implies a^T x_i + b \geq \gamma \\ y_i = -1 \implies a^T x_i + b \leq -\gamma \end{array} \right\} \text{ for some } \gamma \geq 0$$

Hence
$$y_i (a^T x_i + b) \geq \gamma \quad \text{for some } \gamma \geq 0 \text{ for all } i = 1, \dots, n$$



Objective : Find a hyperplane $\{x \mid a^T x + b = 0\}$

such that the minimum margin is maximum.

$$\max_{\substack{\gamma \geq 0 \\ a \in \mathbb{R}^p, b \in \mathbb{R}}} \frac{\gamma}{\|a\|} \qquad \text{s.t.} \quad y_i (a^T x_i + b) \geq \gamma$$

(not scale invariant)

$\Leftrightarrow \min\limits_{\gamma, a, b} \dfrac{\gamma}{\|a\|} \dfrac{\|a\|}{\gamma} \quad$ s.t. $\quad y_i\left(\dfrac{a^T}{\gamma}x_i + \dfrac{b}{\gamma}\right) \geq 1$

$\Leftrightarrow \min\limits_{\substack{a \in \mathbb{R}^p \\ b \in \mathbb{R}}} \|a\| \quad$ s.t. $\quad y_i(a^T x_i + b) \geq 1$

$\Leftrightarrow \min\limits_{a \in \mathbb{R}^p, b \in \mathbb{R}} \dfrac{1}{2}\|a\|^2 \quad$ s.t. $\quad y_i(a^T x_i + b) \geq 1 \quad i = 1, \dots, n$

In summary

(OMC)

(opt. margin classifier)

$\boxed{\begin{array}{l} \text{Given } (x_1, y_1), \dots, (x_n, y_n), x_i \in \mathbb{R}^p, y_i \in \{-1, 1\} \\[2mm] \min\limits_{a \in \mathbb{R}^p, b \in \mathbb{R}} \dfrac{1}{2}\|a\|^2 \quad \text{s.t.} \quad y_i(a^T x_i + b) \geq 1, i = 1, \dots, n \end{array}}$

Quadratic optimization problem with linear constraints, special case of a convex optimization problem.

○ Assume $a^*$ is an optimum solution of (OMC) and $x_k$ some point with minimum margin. Then

$$y_k(a^{*T}x_k + b^*) = 1$$

$\Leftrightarrow \quad (a^{*T}x_k + b^*) = y_k \quad (\text{since } y_k^2 = 1)$

$\Leftrightarrow \quad b^* = y_k - a^{*T}x_k$

Hence, $b^* = y_k - a^{*T}x_k$ is the optimum $b$-value.

○ The solution ($\underline{w}^*$, $a^*$, $b^*$) is called the
  optimal margin classifier.  $\boxed{\text{Fig 4}}$

○ Use commercial or public domain software
  to solve (OMC).

Problem solved?    Yes and no!

Consider: ● Smarter way to solve (OMC)

● Non-separability

## 6.3. SVM and Lagrange Duality

Brief excursion on convex optimization

○ Convex optimization problem:

$$(P) \quad \text{minimize} \quad f_0(x)$$
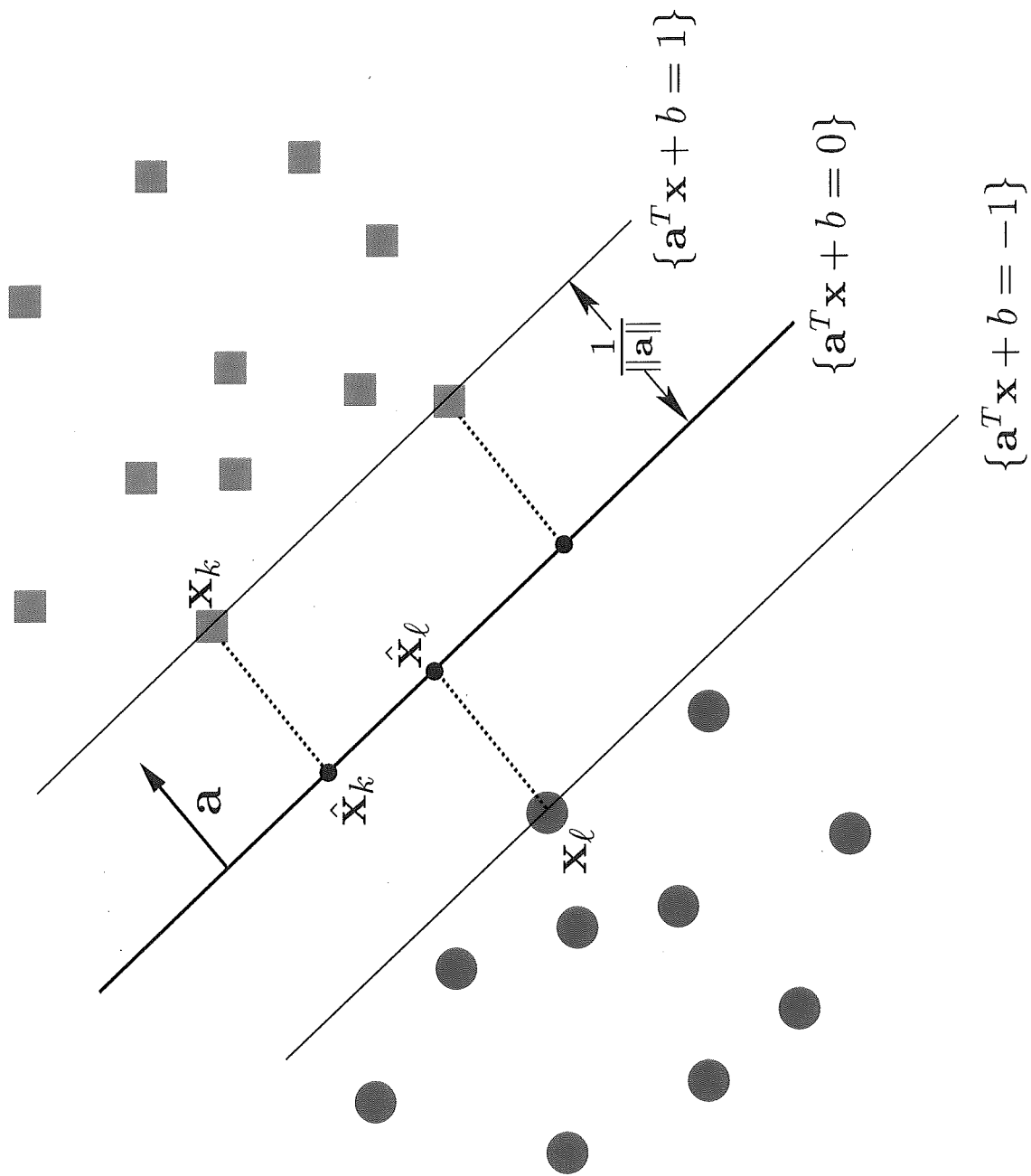$$\text{s.t.} \quad f_i(x) \le 0, \quad i = 1, \ldots, m$$
$$h_i(x) = 0, \quad i = 1, \ldots, p$$

$f_0, f_i$ are convex, $h_i$ are linear.

○ Lagrangian: (prime function)

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) + \sum_{i=1}^{p} \nu_i h_i(x)$$

Fig 4

$\{a^T x + b = 1\}$

$\{a^T x + b = 0\}$

$\{a^T x + b = -1\}$

$\dfrac{1}{\|a\|}$

$x_k$

$\hat{x}_\ell$

$a$

$\hat{x}_k$

$x_\ell$

o Lagrangian dual function

$$g(\lambda, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \nu)$$

$$\mathcal{D} = \bigcap_{i=0}^{m} \text{dom}(f_i) \cap \bigcap_{i=1}^{p} \text{dom}(h_i)$$

o Lagrangian dual problem:

$$(\mathcal{D}) \qquad \max g(\lambda, \nu)$$

$$\text{s.t. } \lambda \geq 0$$

o <u>Weak duality theorem</u>:

$$g(\lambda^*, \nu^*) \leq f_0(x^*)$$

$\lambda^*, \nu^*$ opt. solutions of $(\mathcal{D})$, $x^*$ opt. solution of $(P)$.

o Strong duality:

$$g(\lambda^*, \nu^*) = f_0(x^*)$$

o If the constraints are linear the "Slater's condition" holds, which implies that $g(\lambda^*, \nu^*) = f_0(x^*)$, "strong duality" holds, the "duality gap is $\emptyset$"