



# Support Vector Machines

IF-3270 Pembelajaran Mesin

Teknik Informatika ITB

## Modul 5: Support Vector Machine



**Fariska Z. Ruskanda, S.T., M.T.**  
(fariska@informatika.org)

KK IF -Teknik Informatika - STEI ITB

### 02 SVM for Linearly Separable Data

IF3270 - Pembelajaran Mesin  
(Machine Learning)

# Outline

Support Vectors

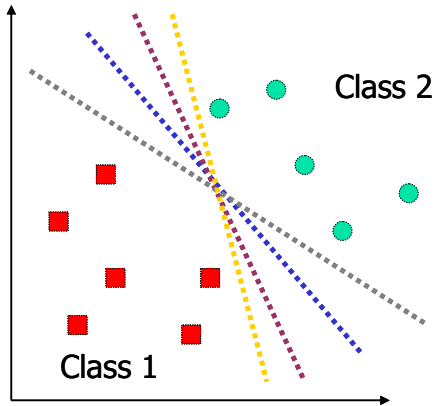
Optimal  
Hyperplane

Quadratic  
Optimization  
Problem

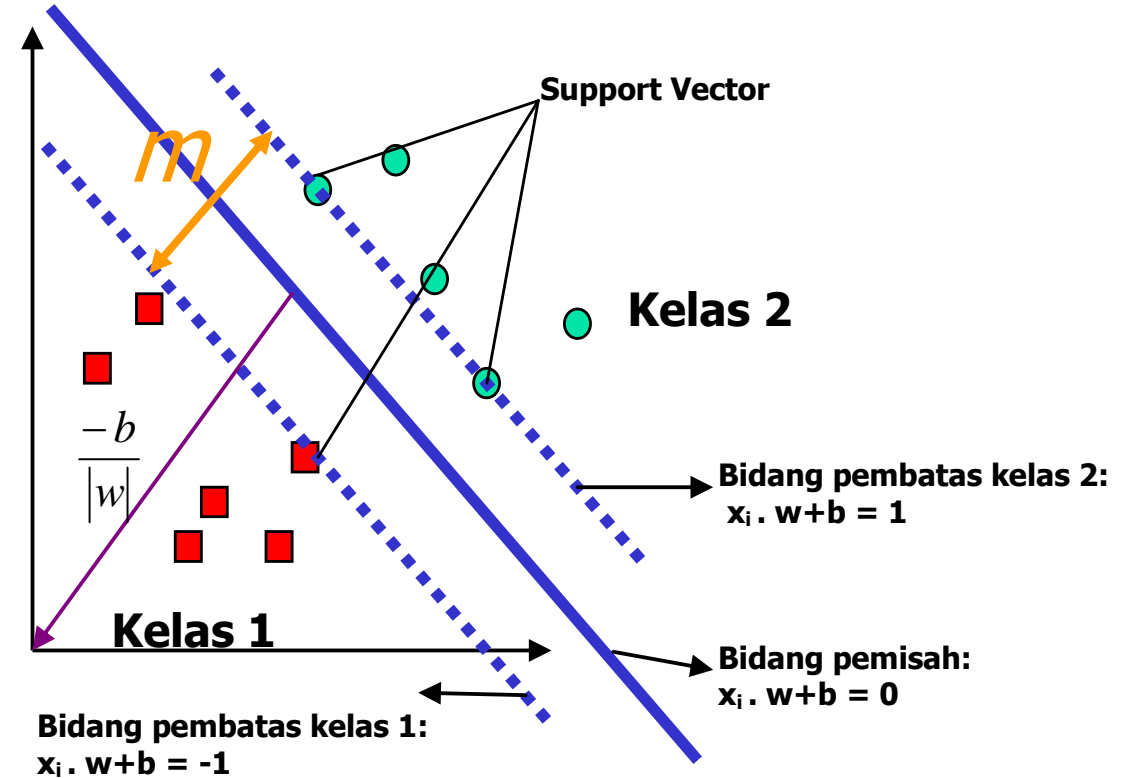
Bidang Pemisah  
Terbaik

Calculation Example

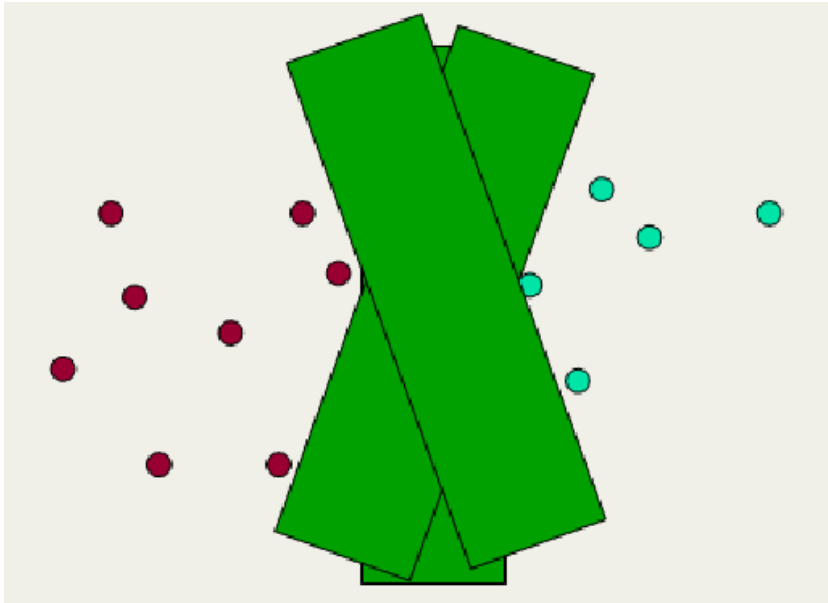
# Mencari Bidang Pemisah Terbaik (1)



Mencari bidang pemisah dengan margin ( $m$ ) terbesar



# Mencari Bidang Pemisah Terbaik (2)



- 2 kelas dapat dipisahkan oleh sepasang bidang pembatas yang sejajar.
  - Bidang pembatas pertama membatasi kelas pertama
  - Bidang pembatas kedua membatasi kelas kedua
- ***Support Vector***: Vector pada training data yang *men-support* bidang pemisah

# Mengapa mencari margin terbesar?

- Bidang pemisah terbaik dengan margin terbesar memiliki generalisasi yang lebih baik
  - Titik yang dekat dengan bidang pemisah merepresentasikan ketidakyakinkan klasifikasi : 50% peluang pengambilan keputusan oleh classifier
- Kapasitas memori untuk menyimpan model menjadi lebih sedikit
  - Hanya support vector yang mempengaruhi pengambilan keputusan

# Support Vectors

- The **data points** that lie closest to the decision surface (or hyperplane)
- They are the data points **most difficult to classify**
- They have direct bearing on the optimum location of the decision surface

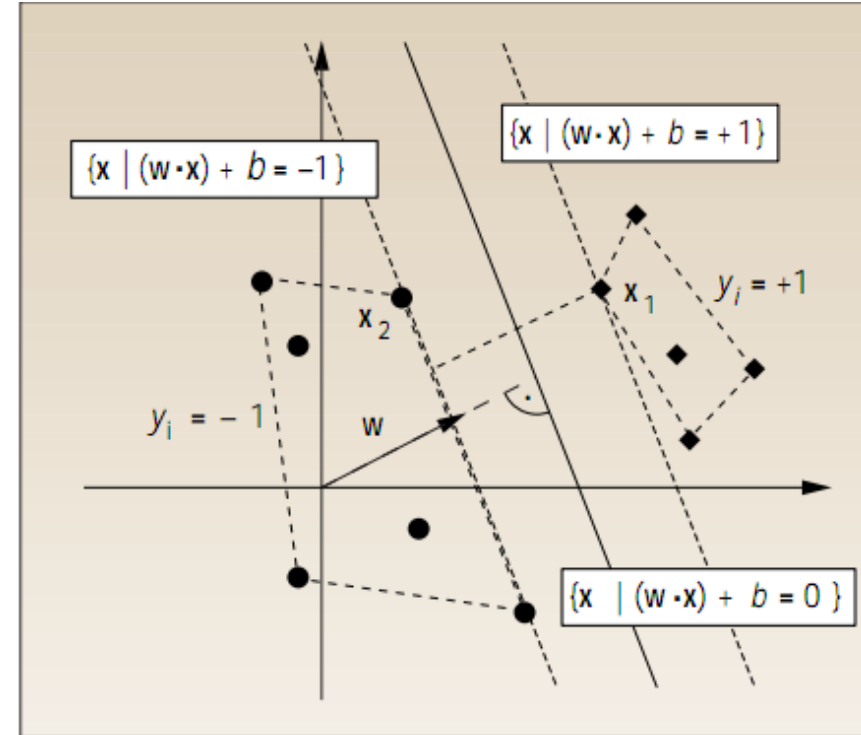
# Support Vectors

- Support vectors memenuhi:

$$|\vec{w} \cdot \vec{x}_i + b| = 1$$

- Untuk itu, semua training data berlabel “+” dan “-” memenuhi juga:

$$y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$$





# Optimal Hyperplane

- Optimal hyperplane:
  - *maximum-margin hyperplane*
  - atau hyperplane with maximal margin separation between two classes.
- Maximal margin ( $2/||w||$ )  $\approx$  minimize  $||w||$  dengan batasan konsistensi training data tetap terjaga (semua klasifikasi benar).

- Minimize: 
$$V(\vec{w}, b) = \frac{1}{2} \vec{w} \cdot \vec{w}$$
- subject to: 
$$y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$$

Recall the distance from a point  $(x_0, y_0)$  to a line:  $Ax + By + c = 0$  is:  $|Ax_0 + By_0 + c| / \sqrt{A^2 + B^2}$ , so,  
 The distance between  $H_0$  and  $H_1$  is then:  
 $|w \cdot x + b| / ||w|| = 1 / ||w||$ , so  
 The total distance between  $H_1$  and  $H_2$  is thus:  $2 / ||w||$

# Optimal Hyperplane (2)

- Pencarian bidang pemisah terbaik dapat dirumuskan menjadi:

$$\begin{array}{ll} \min & \frac{1}{2}|w|^2 \\ \text{s.t} & y_i(x_i \cdot w + b) - 1 \geq 0 \end{array}$$

- Supaya lebih mudah diselesaikan ubah ke formula Lagrangian

$$\begin{array}{ll} \min_{w,b} & L_p(w,b,\alpha) \equiv \frac{1}{2}|w|^2 - \sum_{i=1}^n \alpha_i y_i (x_i \cdot w + b) + \sum_{i=1}^n \alpha_i \\ & \alpha \geq 0 \end{array}$$

- Minimumkan  $L_p$  terhadap  $b$  dan  $w$  (saddle point), diperoleh:

$$\frac{\partial}{\partial b} L_p(w,b,\alpha) = 0 \quad \rightarrow \quad \sum_{i=1}^n \alpha_i y_i = 0$$

$$\frac{\partial}{\partial w} L_p(w,b,\alpha) = 0 \quad \rightarrow$$

$$\frac{\partial L_p}{\partial w} = w - \sum_{i=1}^n \alpha_i y_i x_i = 0$$

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

- Want to look for solution point  $p$  where
 
$$\nabla f(p) = \nabla \lambda g(p)$$

$$g(x) = 0$$
- Or, combining these two as the *Langrangian*  $L$  & requiring derivative of  $L$  be zero:
 
$$L(x,a) = f(x) - ag(x)$$

$$\nabla(x,a) = 0$$



# Langrangian Dual Problem

The Lagrangian Dual Problem: instead of minimizing over  $\mathbf{w}$ ,  $b$ , subject to constraints involving  $a$ 's, we can maximize over  $a$  (the dual variable) subject to the relations obtained previously for  $\mathbf{w}$  and  $b$

Our solution must satisfy these two relations:

$$\mathbf{w} = \sum_{i=1}^l a_i y_i \mathbf{x}_i, \quad \sum_{i=1}^l a_i y_i = 0$$

By substituting for  $\mathbf{w}$  and  $b$  back in the original eqn we can get rid of the dependence on  $\mathbf{w}$  and  $b$ .

Note first that we already now have our answer for what the weights  $\mathbf{w}$  must be: they are a linear combination of the training inputs and the training outputs,  $x_i$  and  $y_i$  and the values of  $a$ . We will now solve for the  $a$ 's by differentiating the dual problem wrt  $a$ , and setting it to zero. Most of the  $a$ 's will turn out to have the value zero. The non-zero  $a$ 's will correspond to the support vectors

# Substitute the Primal Problem

Primal problem:

$$\min L_P = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^l a_i y_i (\mathbf{x}_i \cdot \mathbf{w} + b) + \sum_{i=1}^l a_i$$

$$\text{s.t. } \forall i \ a_i \geq 0$$

$$\mathbf{w} = \sum_{i=1}^l a_i y_i \mathbf{x}_i, \quad \sum_{i=1}^l a_i y_i = 0$$

Dual problem:

$$\max L_D(a_i) = \sum_{i=1}^l a_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l a_i a_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\text{s.t. } \sum_{i=1}^l a_i y_i = 0 \ \& \ a_i \geq 0$$

(note that we have removed the dependence on  $\mathbf{w}$  and  $b$ )

# Quadratic Optimization Problem

- Minimize: 
$$W(\vec{\alpha}) = -\sum_{i=1}^n \alpha_i + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j (\vec{x}_i \cdot \vec{x}_j)$$

subject to: 
$$\sum_{i=1}^n y_i \alpha_i = 0; \forall i \in [1..n]: \alpha_i \geq 0$$

- Support vectors:  $\alpha_i > 0$

- Hyperplane: 
$$\vec{w} \cdot \vec{x} = \left( \sum_{i=1}^n \alpha_i y_i \vec{x}_i \right) \cdot \vec{x} = \sum_{i=1}^n \alpha_i y_i (\vec{x}_i \cdot \vec{x})$$

$$b = y_{sv} - \vec{w} \cdot \vec{x}_{sv}$$

# Optimization Problem Example

Suppose we have two 2D data points:  $[(x_{11}, x_{12}), (x_{21}, x_{22})]$  with labels  $[y_1, y_2]$

- Minimize:  $W(a) = -(a_1 + a_2) + 0.5 \cdot$

$$[y_1 \cdot y_1 \cdot a_1 \cdot a_1 \cdot (x_{11} \cdot x_{11} + x_{12} \cdot x_{12}) + \\ y_1 \cdot y_2 \cdot a_1 \cdot a_2 \cdot (x_{11} \cdot x_{21} + x_{12} \cdot x_{22}) + \\ y_2 \cdot y_1 \cdot a_2 \cdot a_1 \cdot (x_{21} \cdot x_{11} + x_{22} \cdot x_{12}) + \\ y_2 \cdot y_2 \cdot a_2 \cdot a_2 \cdot (x_{21} \cdot x_{21} + x_{22} \cdot x_{22})]$$

subject to:  $(y_1 \cdot a_1 + y_2 \cdot a_2) = 0$

- Support vectors:  $a_1 > 0, a_2 > 0$

- Hyperplane:

$$w \cdot x = a_1 \cdot y_1 \cdot ((x_{11}, x_{12}) \cdot x) + a_2 \cdot y_2 \cdot ((x_{21}, x_{22}) \cdot x)$$

$$b = y_{sv} - w \cdot x_{sv}$$

$$W(\vec{\alpha}) = -\sum_{i=1}^n \alpha_i + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j (\vec{x}_i \cdot \vec{x}_j)$$

$$\sum_{i=1}^n y_i \alpha_i = 0; \forall i \in [1..n]: \alpha_i \geq 0$$

$$\vec{w} \cdot \vec{x} = \left( \sum_{i=1}^n \alpha_i y_i \vec{x}_i \right) \cdot \vec{x} = \sum_{i=1}^n \alpha_i y_i (\vec{x}_i \cdot \vec{x})$$

$$b = y_{sv} - \vec{w} \cdot \vec{x}_{sv}$$

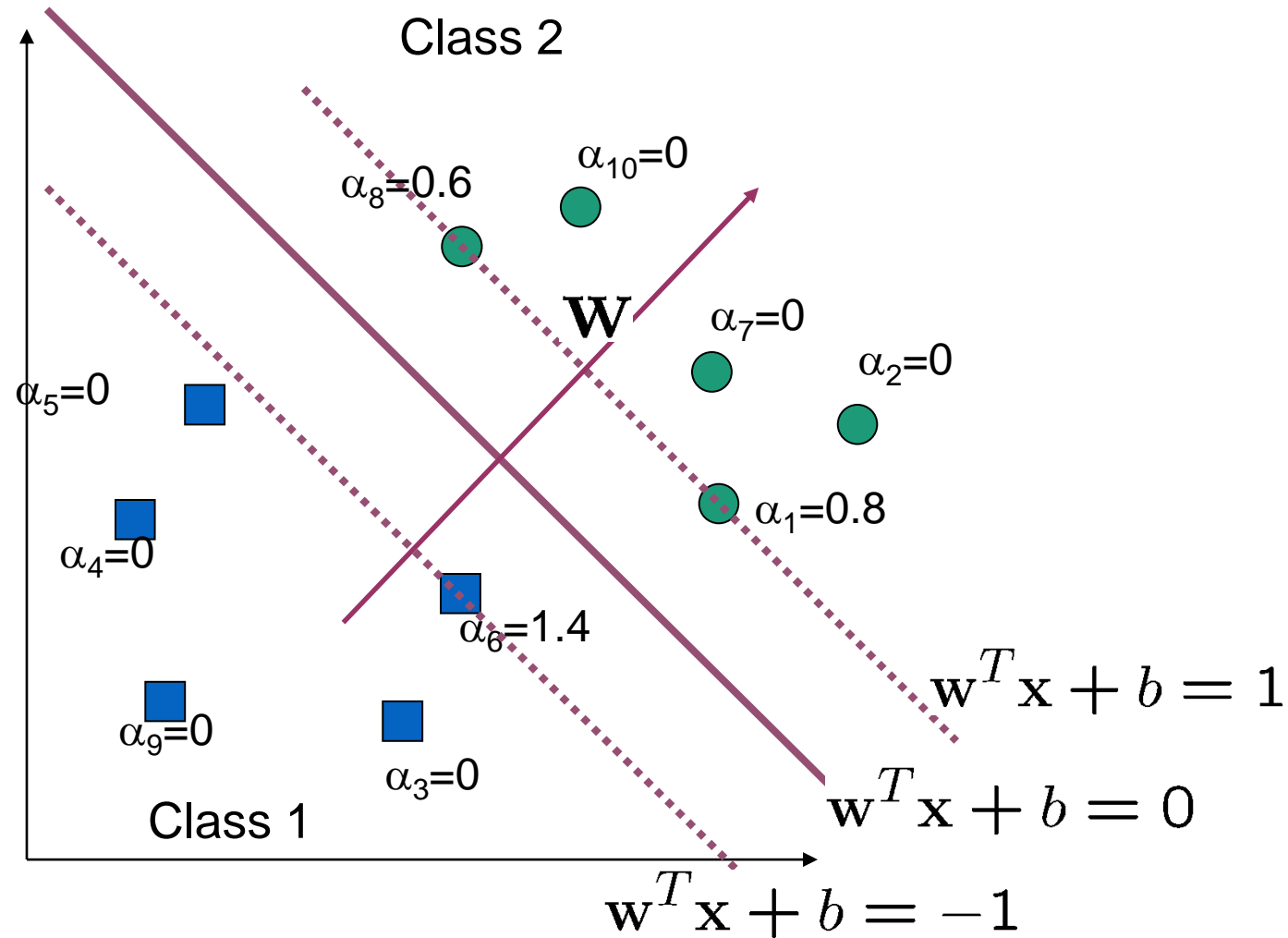


# Quadratic Programming

- Quadratic programming (QP) is the problem of optimizing a quadratic objective function and is one of the simplests form of non-linear programming.
- The objective function can contain bilinear or up to second order polynomial terms, and the constraints are linear and can be both equalities and inequalities.

$$\begin{aligned} \min f(x) &= q^T x + \frac{1}{2} x^T Q x \\ \text{s.t. } Ax &= a \\ Bx &\leq b \\ x &\geq \underline{0} \end{aligned}$$

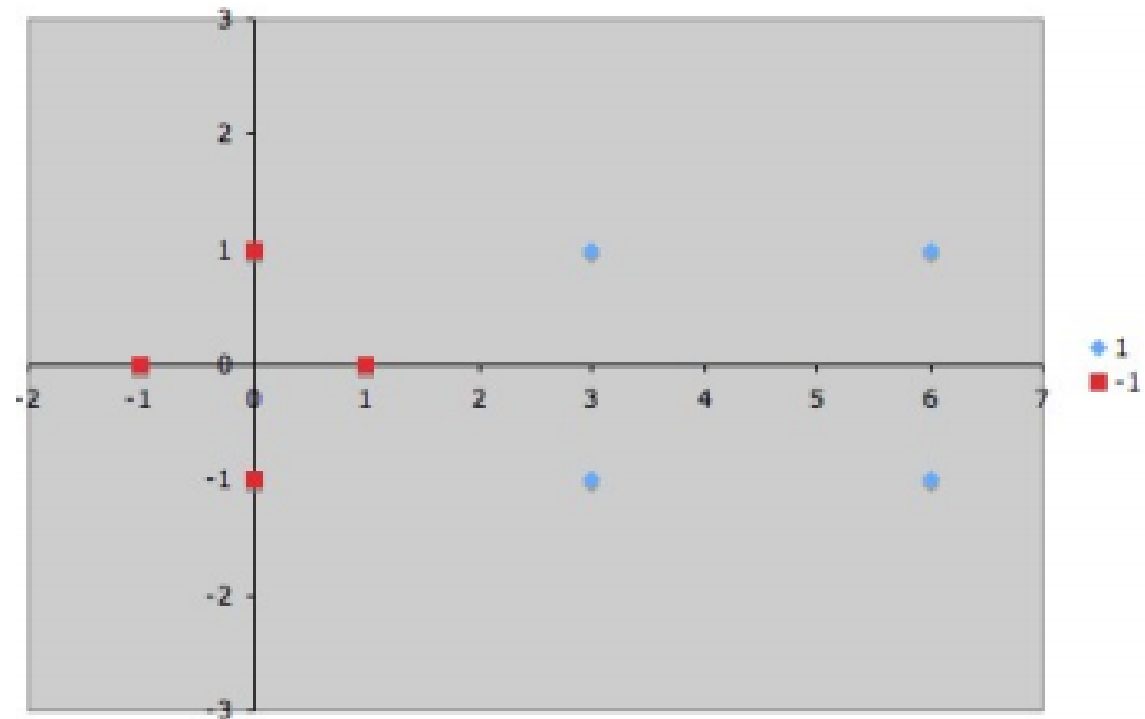
# A Geometrical Interpretation



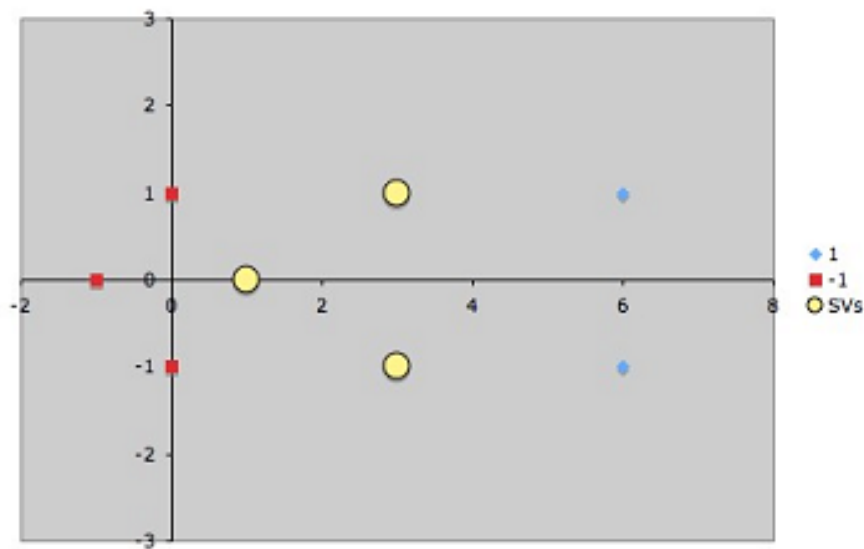


# Contoh

x1	x2	Kelas
3	1	+1
3	-1	+1
6	1	+1
6	-1	+1
1	0	-1
0	1	-1
0	-1	-1
-1	0	-1



- Taken from:  
<http://axon.cs.byu.edu/Dan/678/miscellaneous/SVM.example.pdf>



# Support Vectors

$$f(\vec{x}) = \sum_{i=1}^{nsv} (\alpha_i y_i \vec{x}_i \vec{x}) + b$$

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} \Rightarrow -1 = \alpha_1 \cdot -1 \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha_2 \cdot 1 \cdot \begin{pmatrix} 3 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha_3 \cdot 1 \cdot \begin{pmatrix} 3 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} + b = -\alpha_1 + 3\alpha_2 + 3\alpha_3 + b \dots (1)$$

$$\begin{pmatrix} 3 \\ 1 \end{pmatrix} \Rightarrow 1 = \alpha_1 \cdot -1 \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 1 \end{pmatrix} + \alpha_2 \cdot 1 \cdot \begin{pmatrix} 3 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 1 \end{pmatrix} + \alpha_3 \cdot 1 \cdot \begin{pmatrix} 3 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ 1 \end{pmatrix} + b = -3\alpha_1 + 10\alpha_2 + 8\alpha_3 + b \dots (2)$$

$$\begin{pmatrix} 3 \\ -1 \end{pmatrix} \Rightarrow 1 = \alpha_1 \cdot -1 \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ -1 \end{pmatrix} + \alpha_2 \cdot 1 \cdot \begin{pmatrix} 3 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ -1 \end{pmatrix} + \alpha_3 \cdot 1 \cdot \begin{pmatrix} 3 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} 3 \\ -1 \end{pmatrix} + b = -3\alpha_1 + 8\alpha_2 + 10\alpha_3 + b \dots (3)$$

$$-\alpha_1 + 3\alpha_2 + 3\alpha_3 + b = -1 \dots (1)$$

$$-3\alpha_1 + 10\alpha_2 + 8\alpha_3 + b = 1 \dots (2)$$

$$-3\alpha_1 + 8\alpha_2 + 10\alpha_3 + b = 1 \dots (3)$$

$$-\alpha_1 + \alpha_2 + \alpha_3 = 0 \dots (4)$$

$$(2) - (3): 2\alpha_2 - 2\alpha_3 = 0 \rightarrow \alpha_2 = \alpha_3 \dots (5)$$

$$(5) \text{subs}(1): -\alpha_1 + 6\alpha_2 + b = -1 \dots (6)$$

$$(5) \text{subs}(2): -3\alpha_1 + 18\alpha_2 + b = 1 \dots (7)$$

$$(6) - (7): 2\alpha_1 - 12\alpha_2 = -2$$

$$\alpha_1 - 6\alpha_2 = -1$$

$$\alpha_1 = 6\alpha_2 - 1 \dots (8)$$

$$(8,5) \text{subs}(4): -6\alpha_2 + 1 + \alpha_2 + \alpha_2 = 0$$

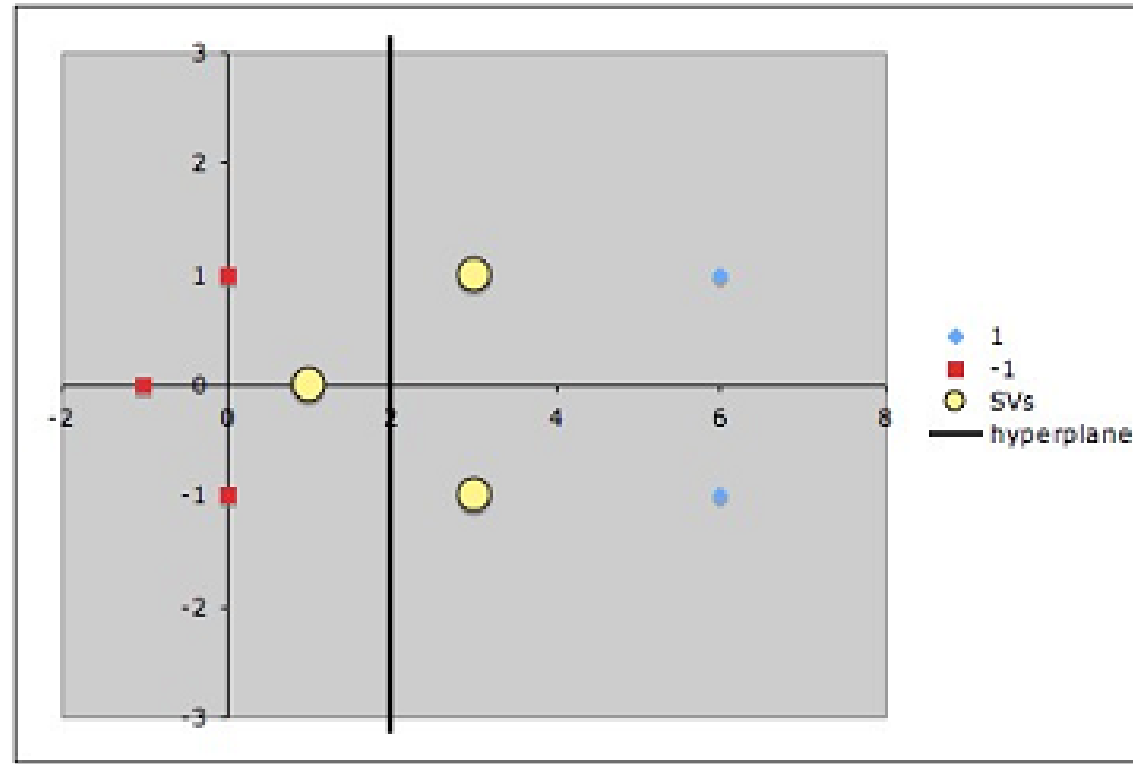
$$-4\alpha_2 = -1 \rightarrow \alpha_2 = \frac{1}{4} = 0.25$$

$$\alpha_1 = 0.5; \alpha_2 = 0.25; \alpha_3 = 0.25$$

$$\text{substitusi}(6): -0.5 + 6 * 0.25 + b = -1$$

$$b = -2$$

# Hipotesis

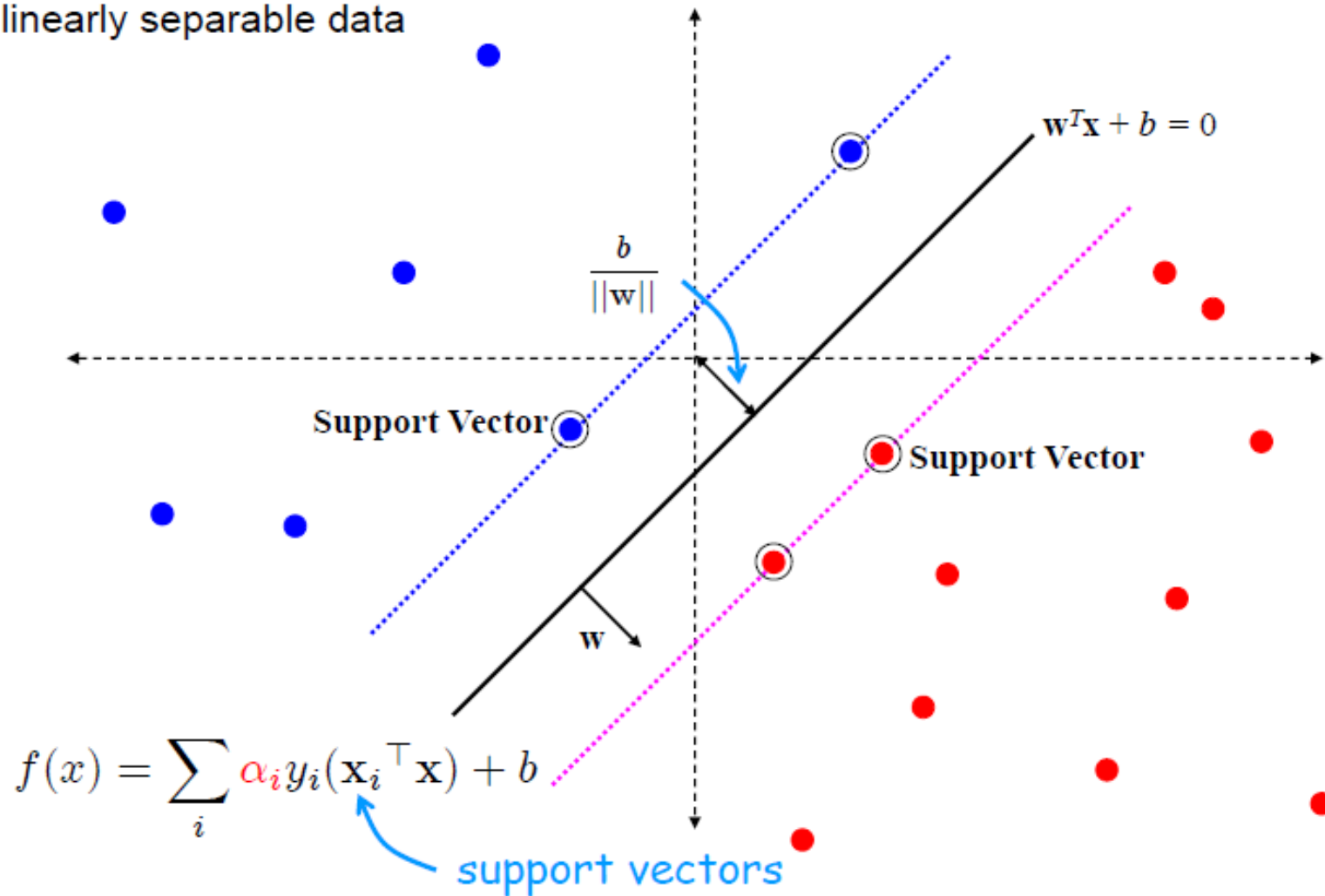


$$f(\vec{x}) = \sum_{i=1}^{nsv} (\alpha_i y_i \vec{x}_i \cdot \vec{x}) - 2; \alpha_1 = 0.5; \alpha_2 = \alpha_3 = 0.25$$

$$f\left(\begin{pmatrix} 6 \\ 1 \end{pmatrix}\right) = \text{sign}\left((0.5)(-1)\begin{pmatrix} 1 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 6 \\ 1 \end{pmatrix} + (0.25)(1)\begin{pmatrix} 3 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 6 \\ 1 \end{pmatrix} + (0.25)(1)\begin{pmatrix} 3 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} 6 \\ 1 \end{pmatrix} - 2\right) = \text{sign}(-3 + 4.75 + 4.25 - 2) = \text{sign}(4) = 1$$

# Summary

linearly separable data



# Inner Product and Similarity

Why should inner product kernels be involved in pattern recognition using SVMs, or at all?

- Intuition is that inner products provide some measure of ‘similarity’
- Inner product in 2D between 2 vectors of unit length returns the cosine of the angle between them = how ‘far apart’ they are

e.g.  $\mathbf{x} = [1, 0]^T$ ,  $\mathbf{y} = [0, 1]^T$

i.e. if they are parallel their inner product is 1 (completely similar)

$$\mathbf{x}^T \mathbf{y} = \mathbf{x} \bullet \mathbf{y} = 1$$

If they are perpendicular (completely unlike) their inner product is 0 (so should not contribute to the correct classifier)

$$\mathbf{x}^T \mathbf{y} = \mathbf{x} \bullet \mathbf{y} = 0$$



## 03 SVM for Non-linearly Separable Data

IF3270 Pembelajaran Mesin