

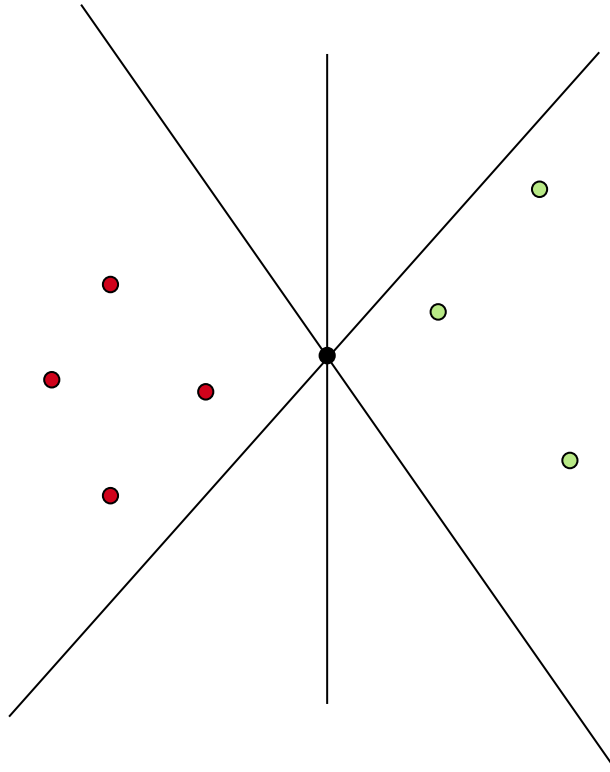
# Support Vector Machines

Machine Learning Techniques

# Outline

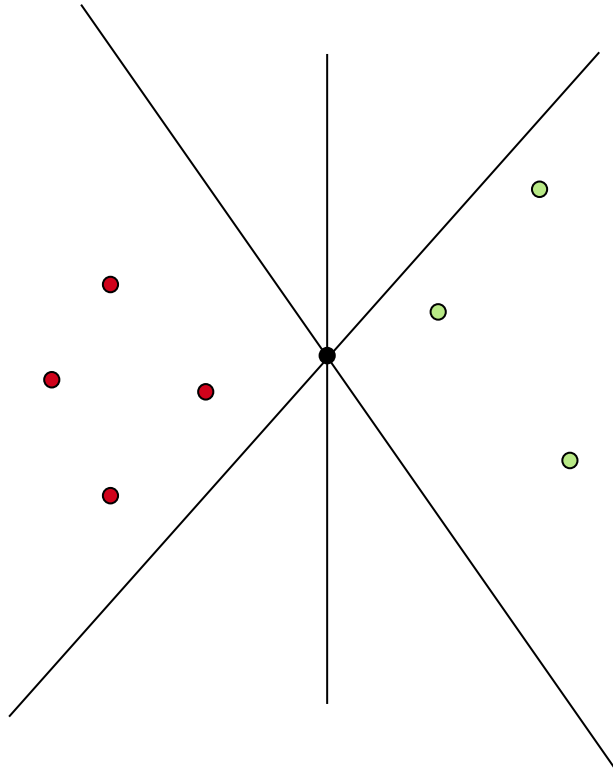
- Margin
- Max-margin classifier
- Duality
- Weight vector
- Support vectors
- Hard-margin Linear-SVM
- Soft-margin SVM

# From Perceptrons to SVM



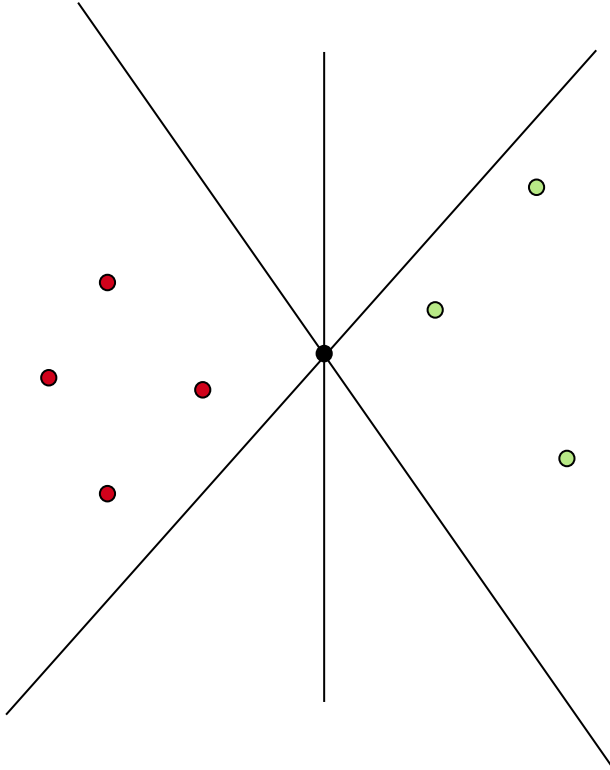
For linearly separable data with  $\gamma$  margin:

# From Perceptrons to SVM



For linearly separable data with  $\gamma$  margin:  
(1) Infinite number of valid linear classifiers.

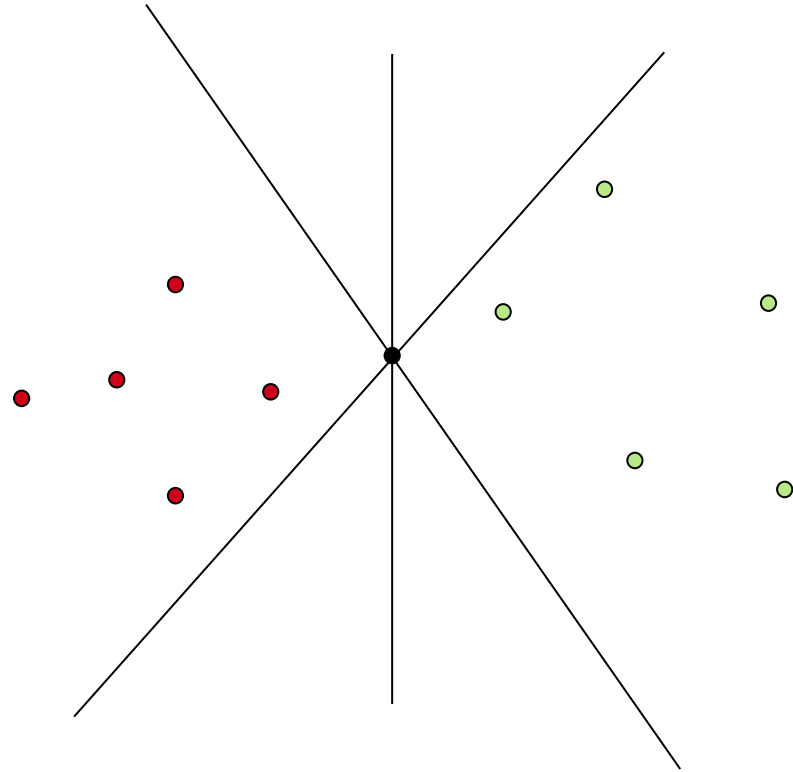
# From Perceptrons to SVM



For linearly separable data with  $\gamma$  margin:

- (1) Infinite number of valid linear classifiers.
- (2) Perceptron returns a valid linear classifier.

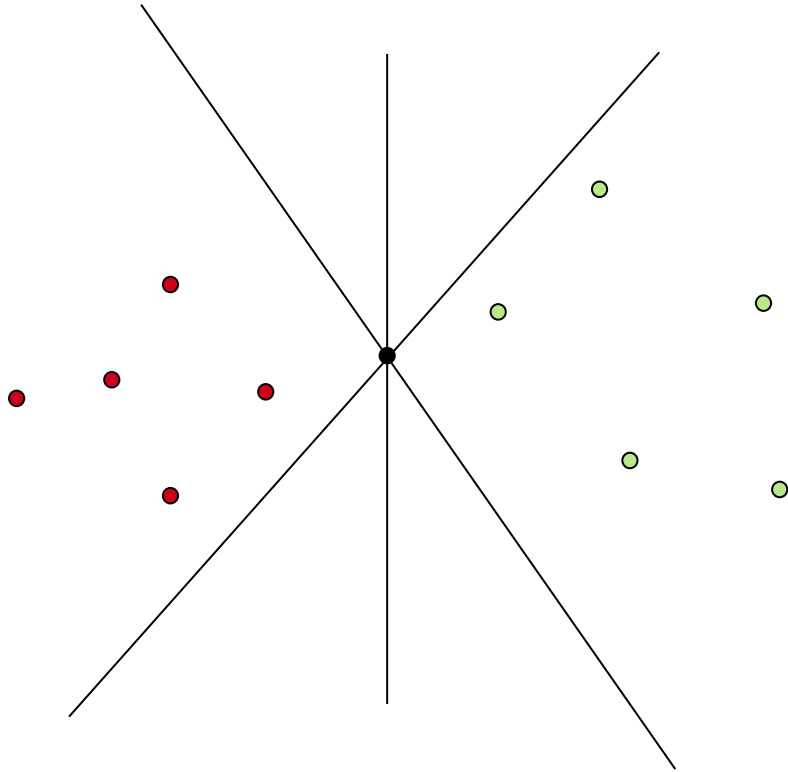
# From Perceptrons to SVM



For linearly separable data with  $\gamma$  margin:

- (1) Infinite number of valid linear classifiers.
- (2) Perceptron returns a valid linear classifier.
- (3) Is it the "best"?

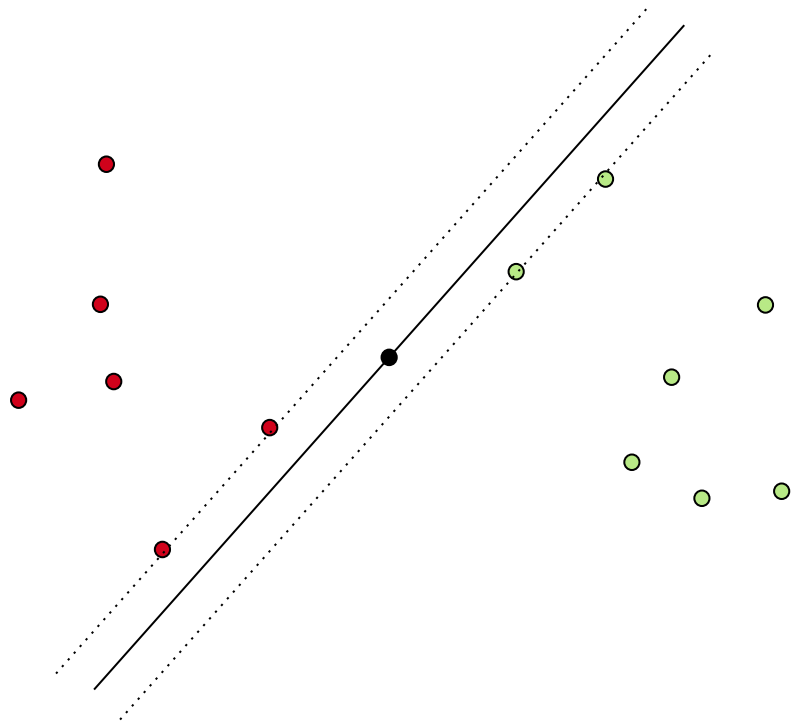
# From Perceptrons to SVM



For linearly separable data with  $\gamma$  margin:

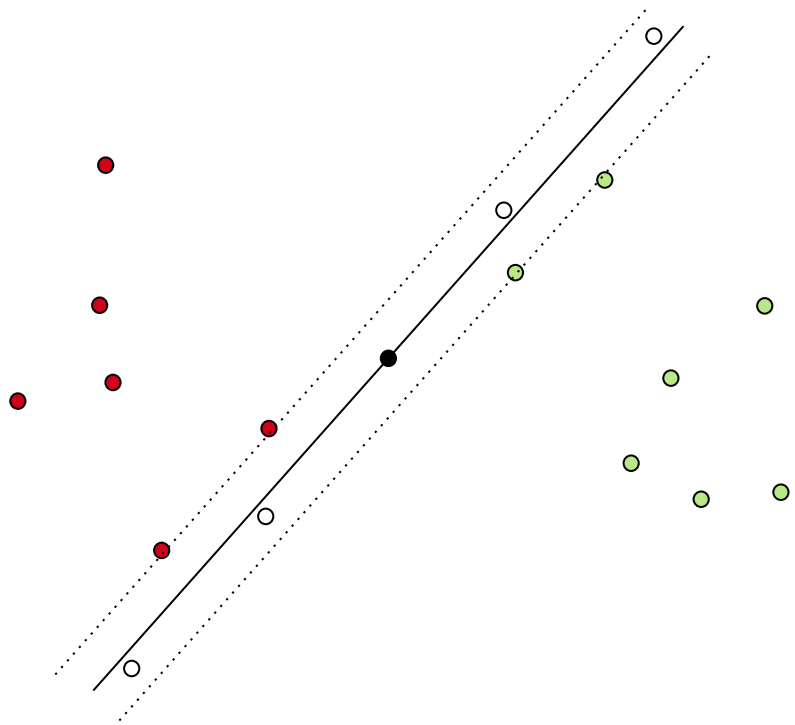
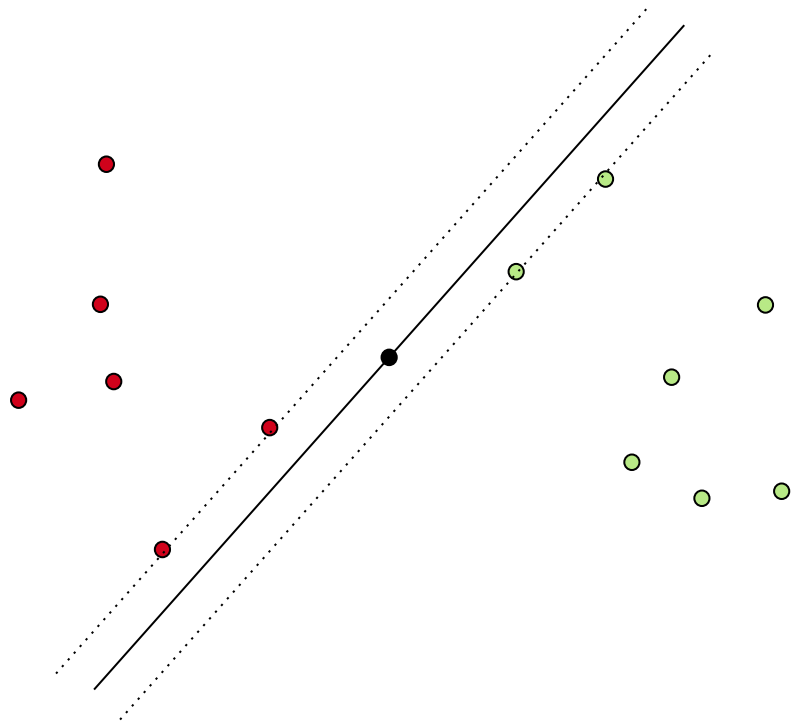
- (1) Infinite number of valid linear classifiers.
- (2) Perceptron returns a valid linear classifier.
- (3) Is it the "best"?
- (4) What is a good notion of "best"?

Margin

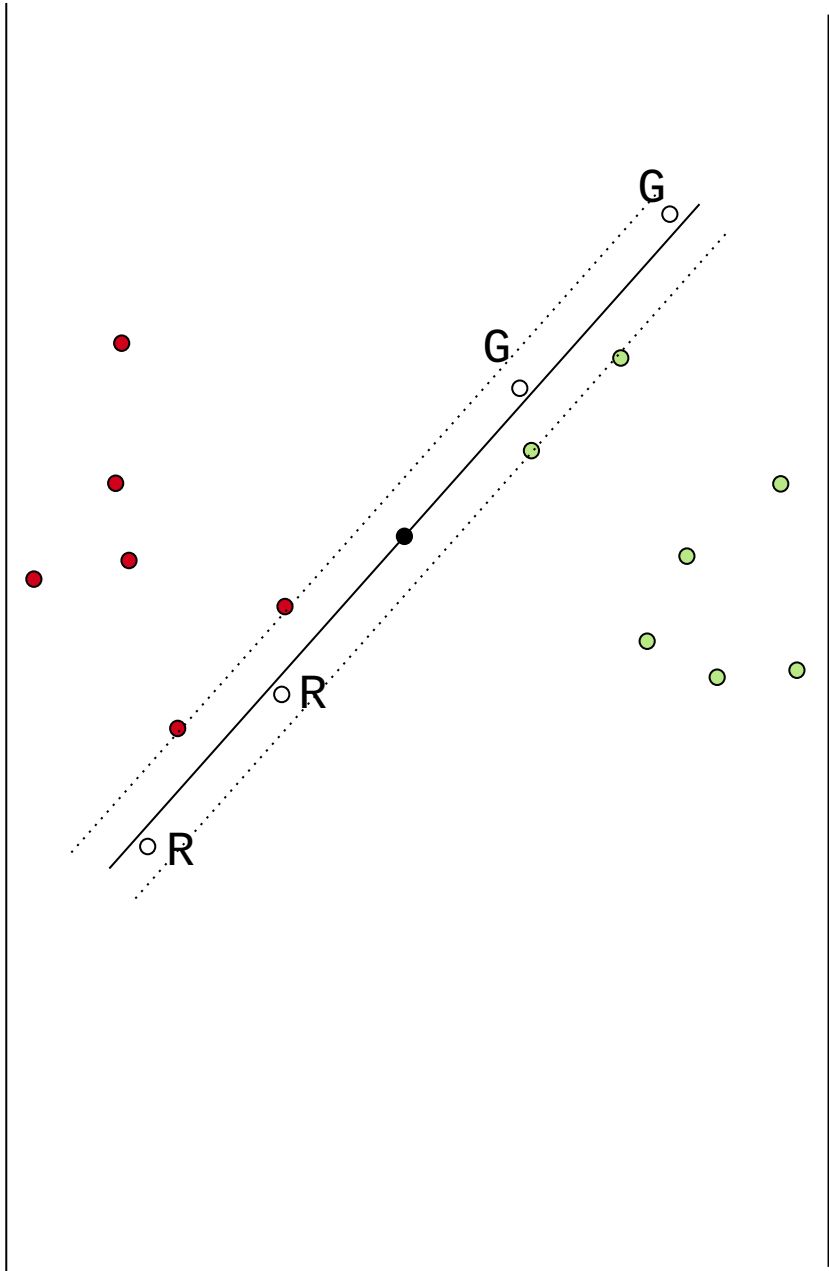
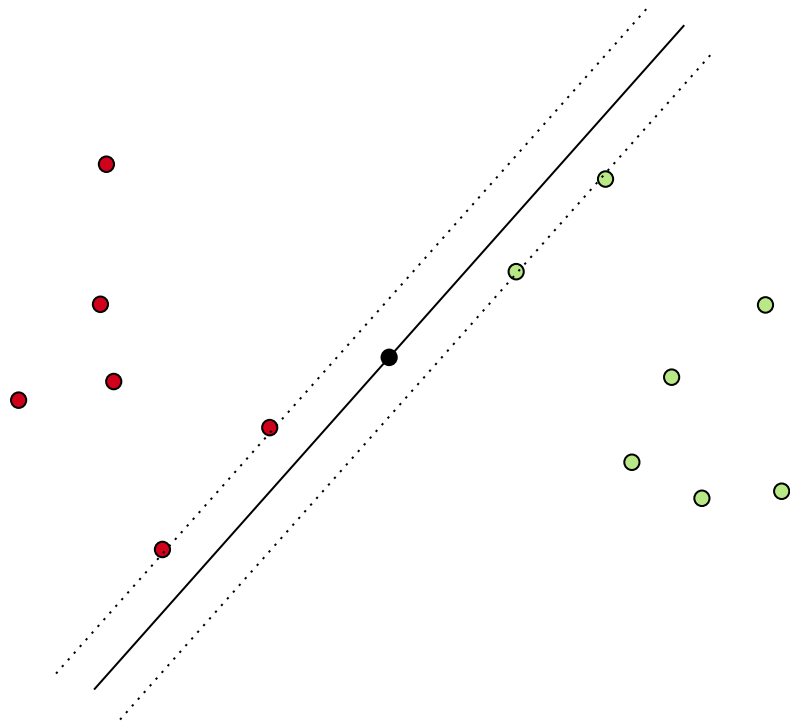




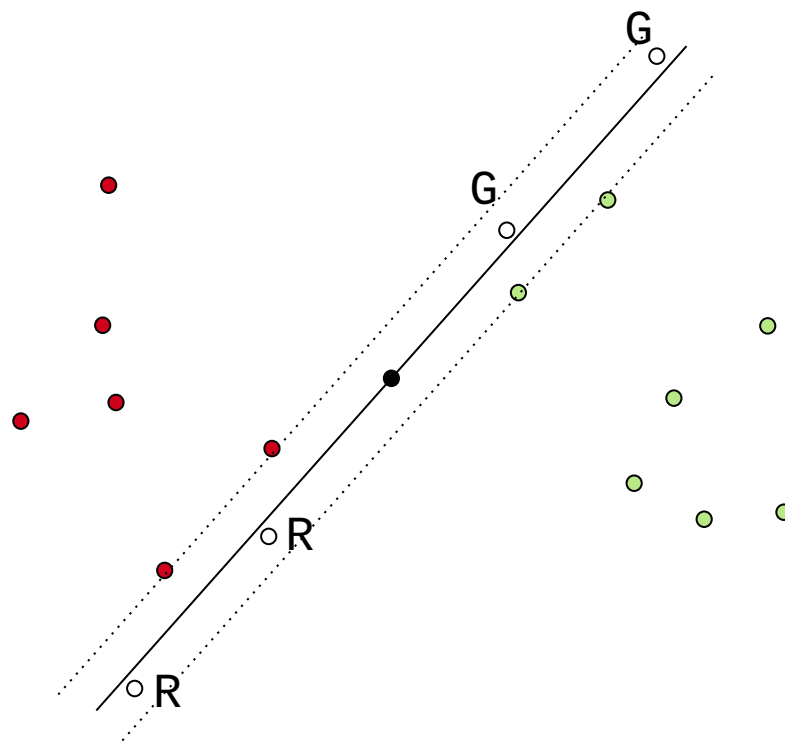
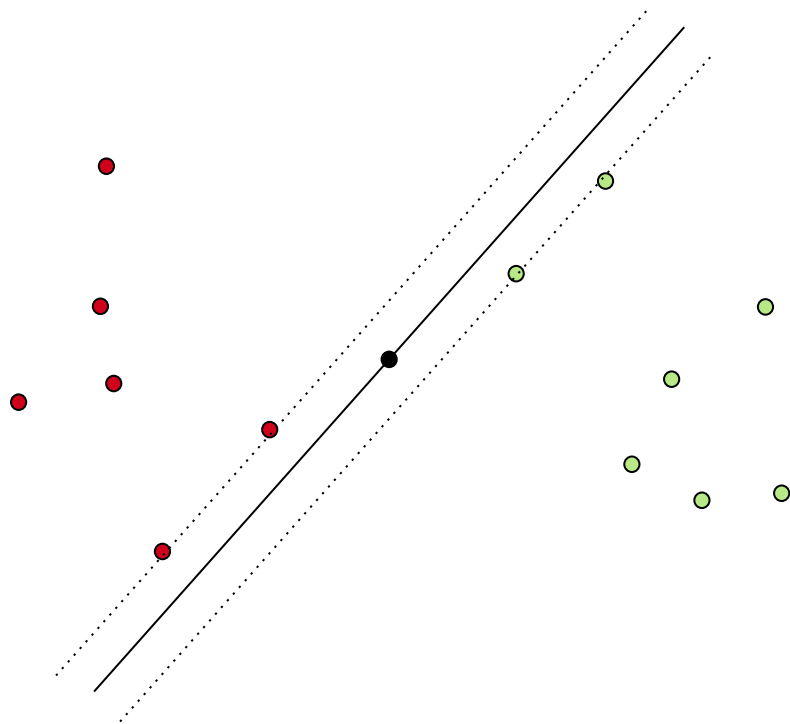
# Margin



Margin

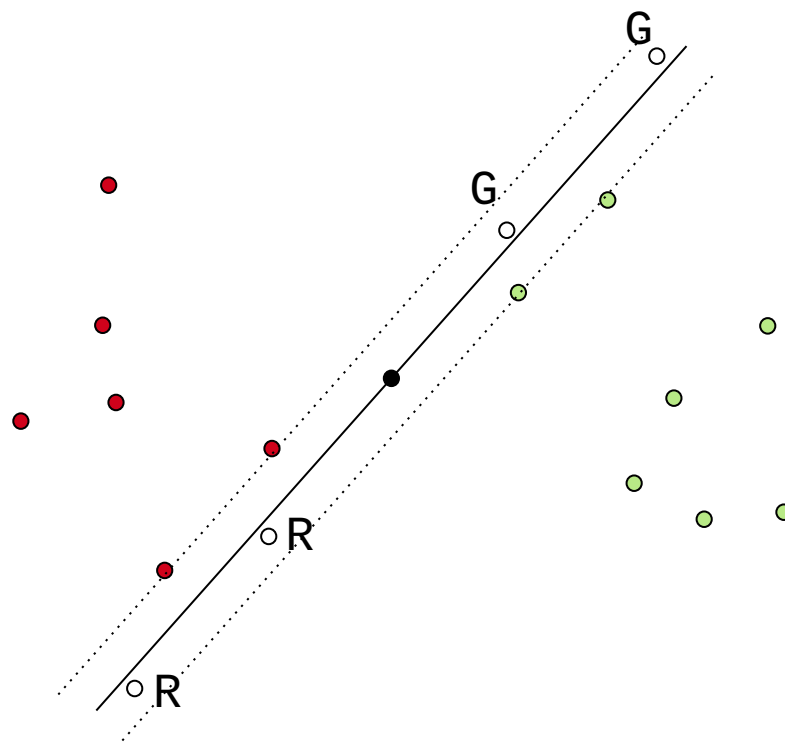
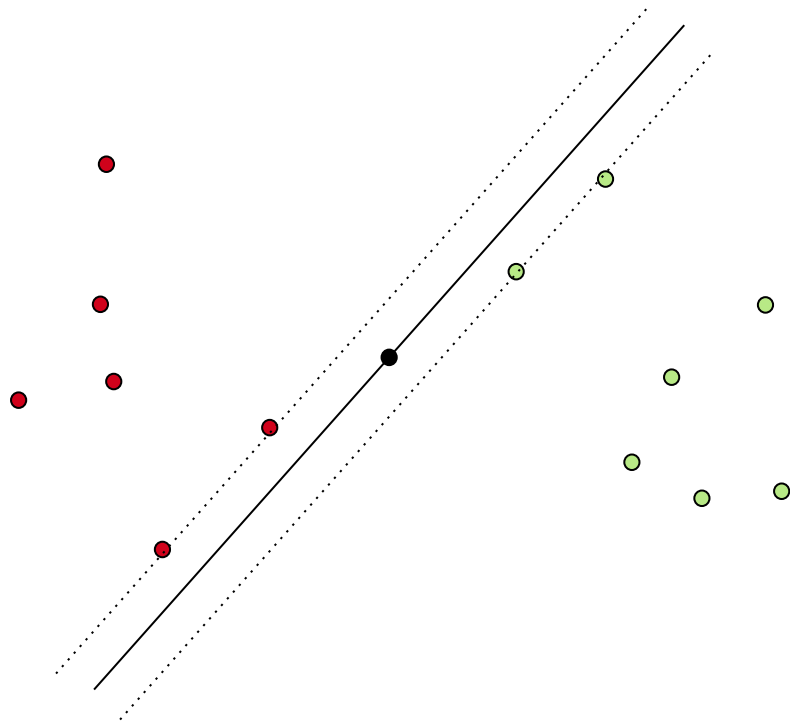


# Margin

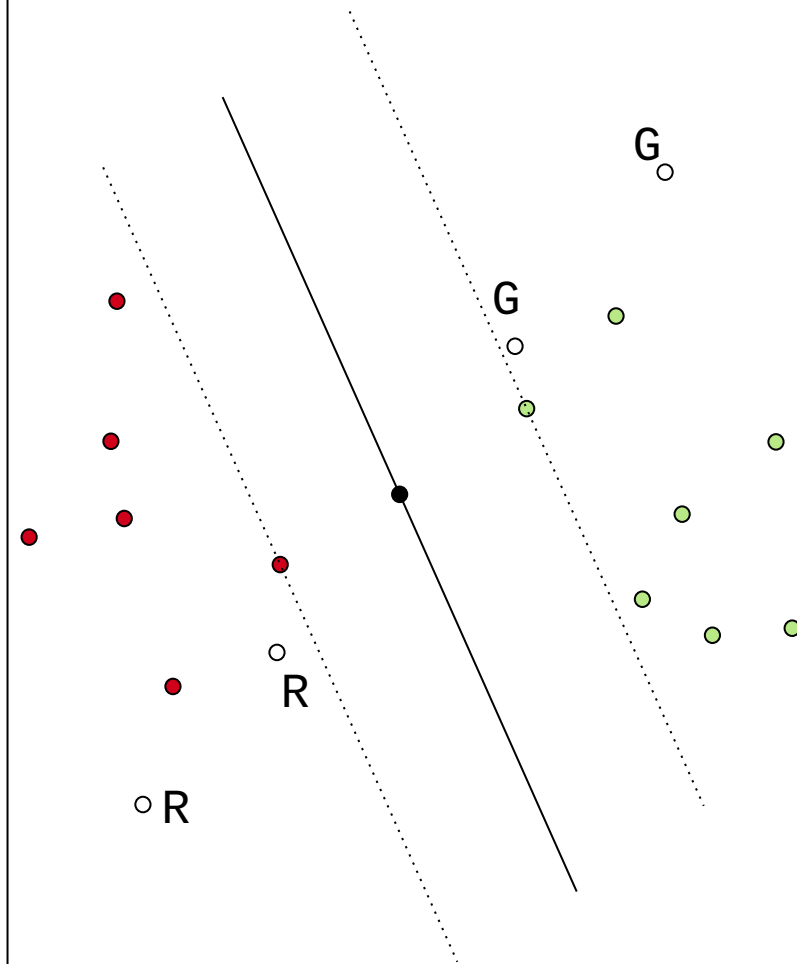


Small margin  
Doesn't generalize well

# Margin



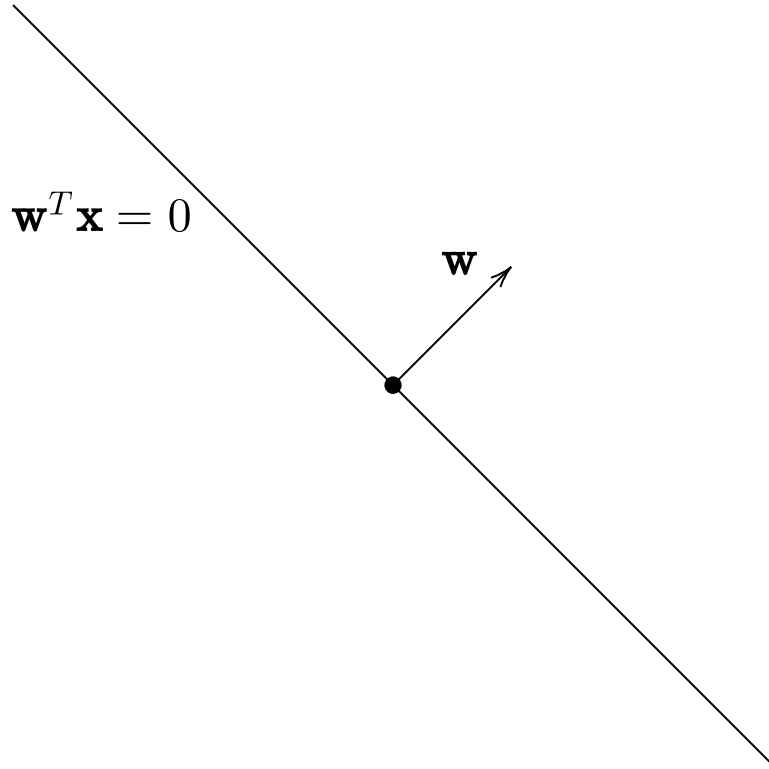
Small margin  
Doesn't generalize well



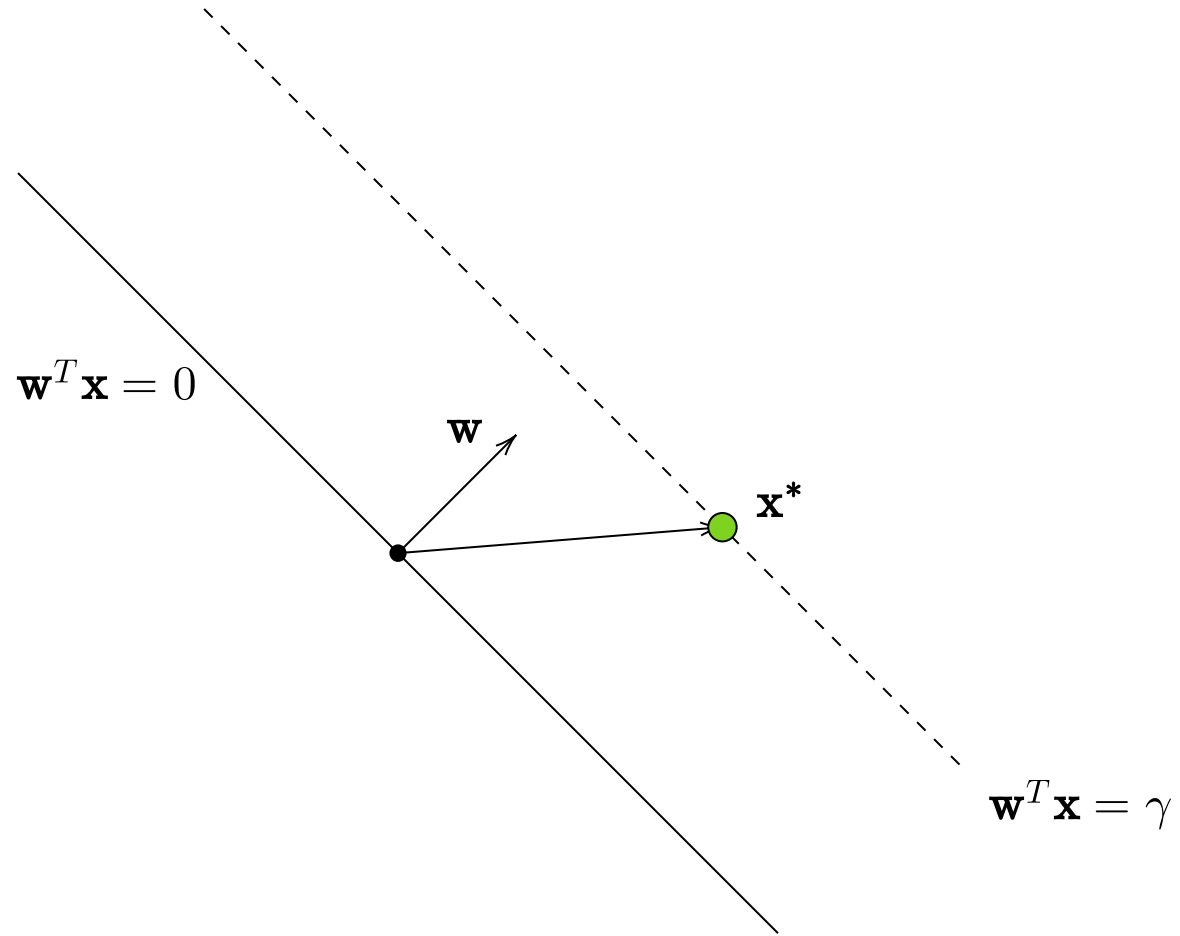
Large margin  
Better generalization

# Computing the Margin

For any linear classifier represented by  $\mathbf{w}$ :



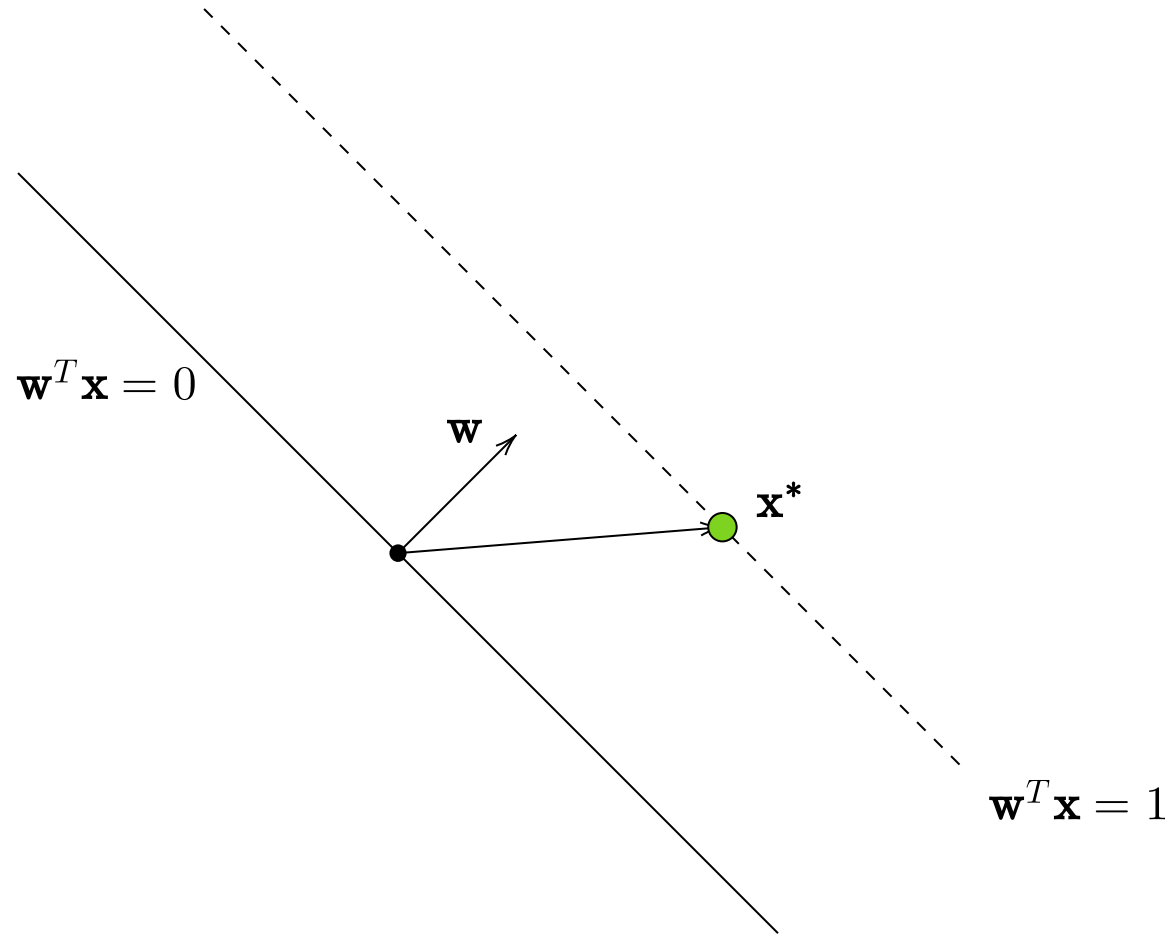
# Computing the Margin



For any linear classifier represented by  $\mathbf{w}$ :

- (1) Find the point closest to it  $\rightarrow \mathbf{x}^*$

# Computing the Margin

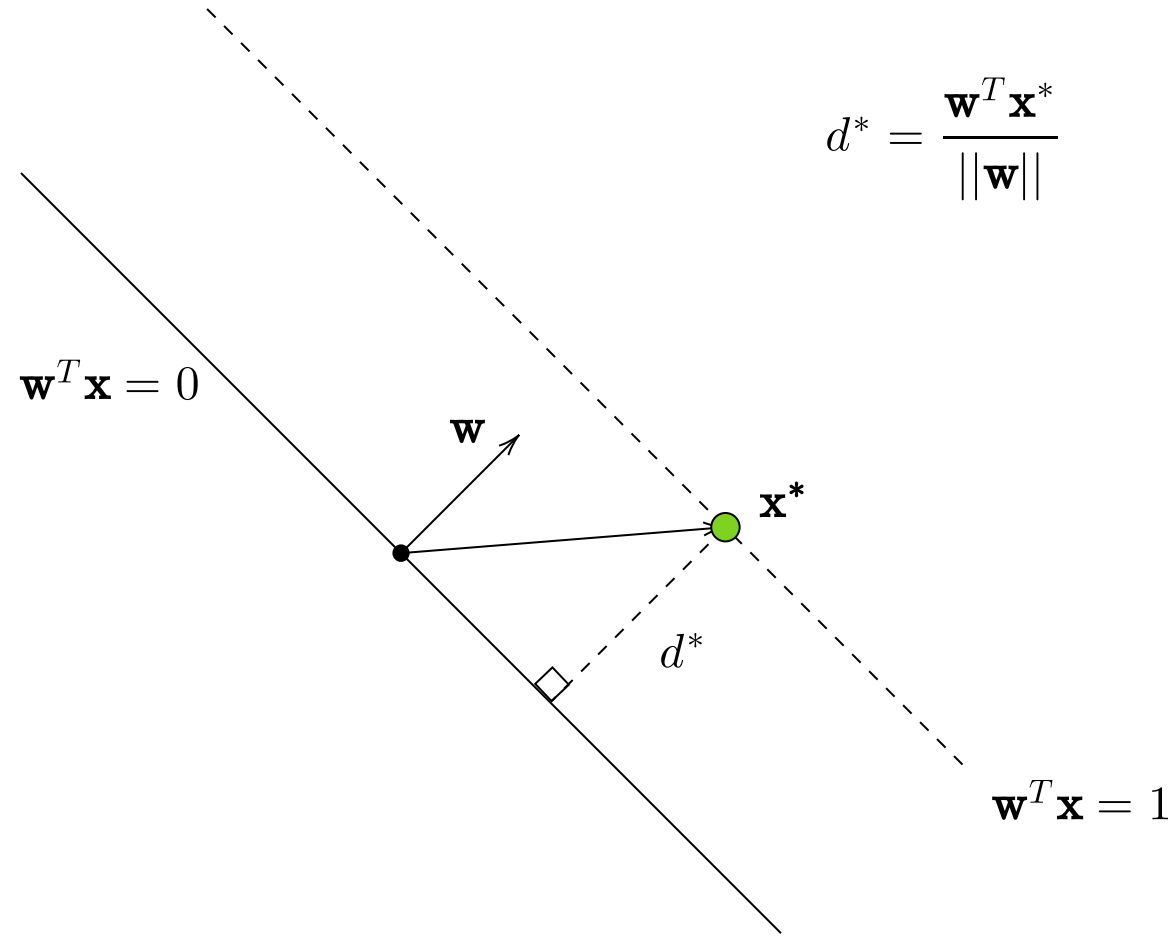


For any linear classifier represented by  $\mathbf{w}$ :

(1) Find the point closest to it  $\rightarrow \mathbf{x}^*$

(2) Scale  $\mathbf{w}$  such that  $\mathbf{x}^*$  lies on  $\mathbf{w}^T \mathbf{x} = 1$

# Computing the Margin



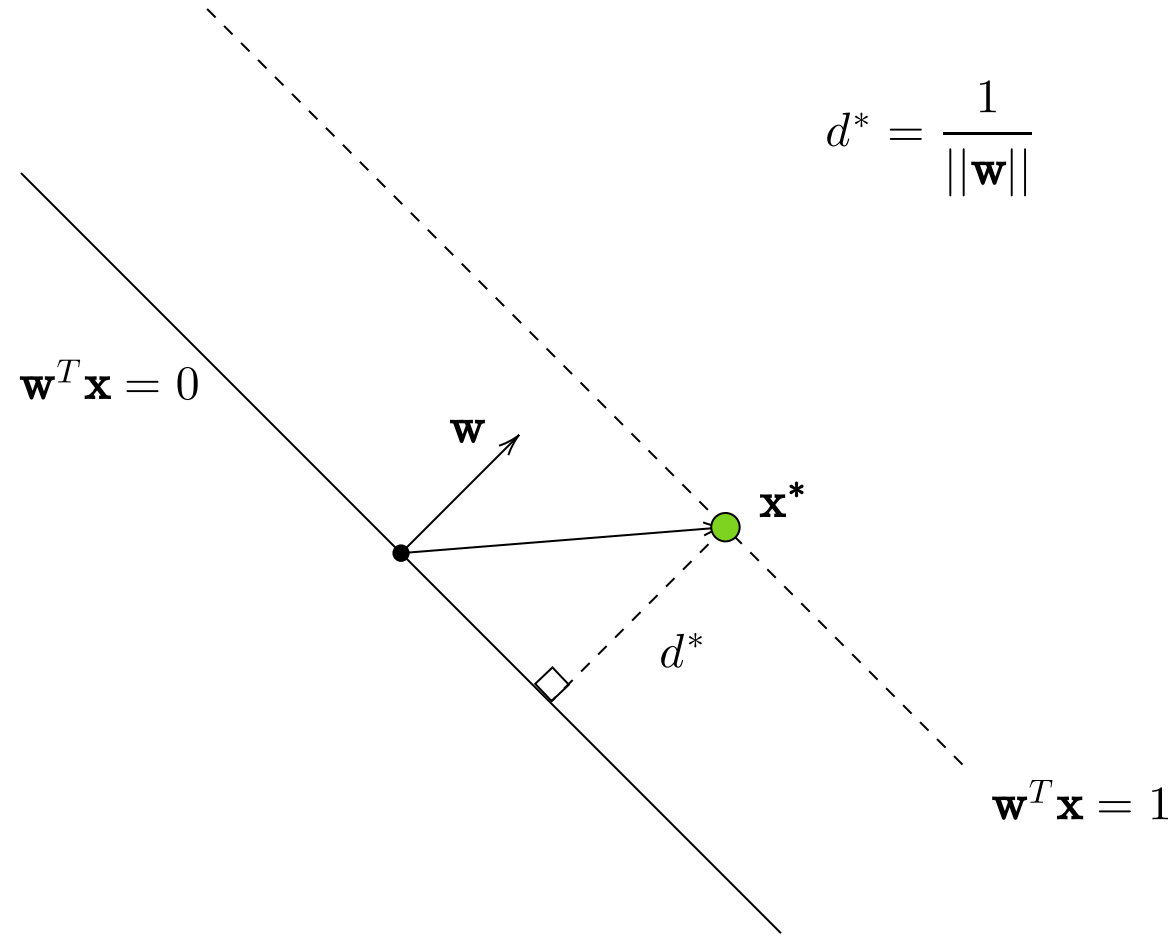
$$d^* = \frac{\mathbf{w}^T \mathbf{x}^*}{\|\mathbf{w}\|}$$

For any linear classifier represented by  $\mathbf{w}$ :

- (1) Find the point closest to it  $\rightarrow \mathbf{x}^*$
- (2) Scale  $\mathbf{w}$  such that  $\mathbf{x}^*$  lies on  $\mathbf{w}^T \mathbf{x} = 1$
- (3) Distance of  $\mathbf{x}^*$  from the line is  $\frac{1}{\|\mathbf{w}\|}$



# Computing the Margin

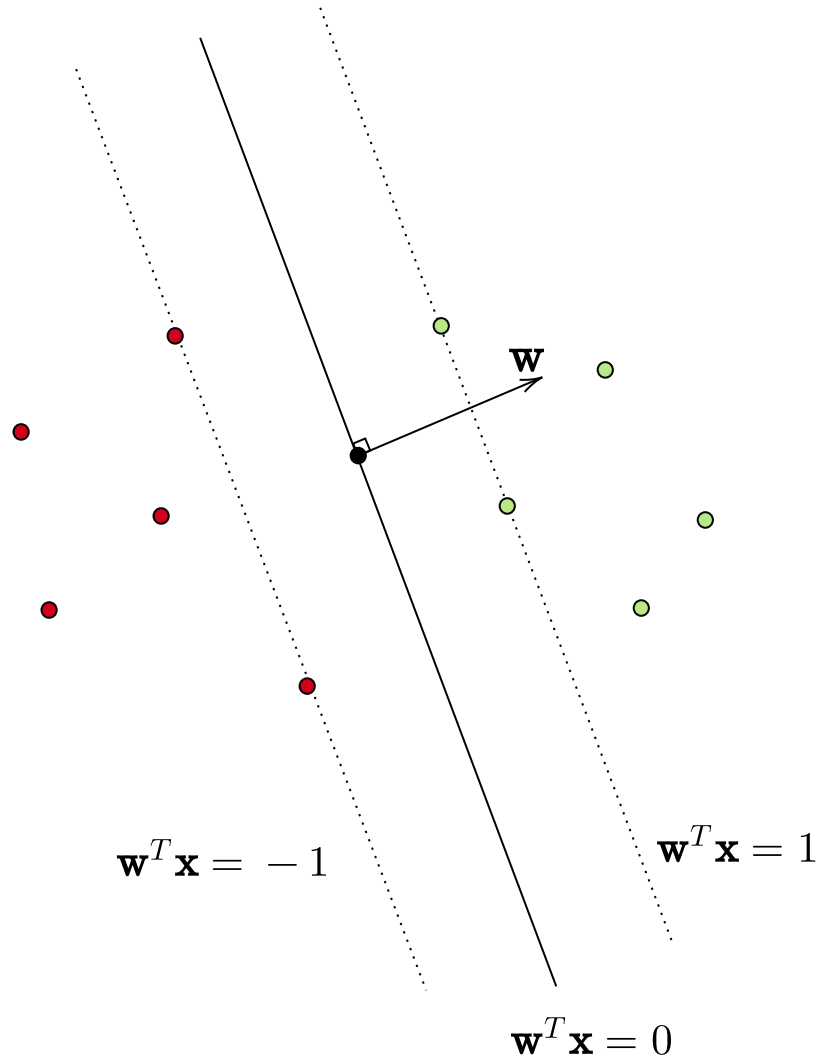


$$d^* = \frac{1}{\|\mathbf{w}\|}$$

For any linear classifier represented by  $\mathbf{w}$ :

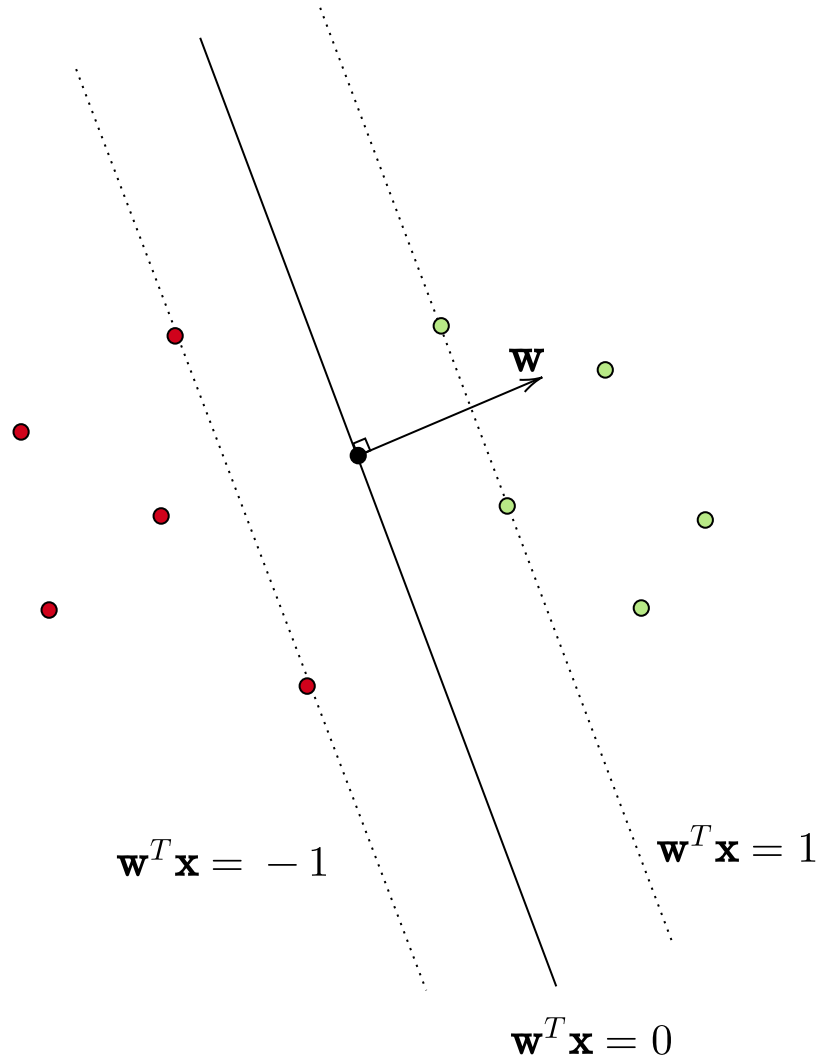
- (1) Find the point closest to it  $\rightarrow \mathbf{x}^*$
- (2) Scale  $\mathbf{w}$  such that  $\mathbf{x}^*$  lies on  $\mathbf{w}^T \mathbf{x} = 1$
- (3) Distance of  $\mathbf{x}^*$  from the line is  $\frac{1}{\|\mathbf{w}\|}$
- (4) This is the (geometric) margin for this linear classifier.

# Beyond the "margin"



$$(\mathbf{w}^T \mathbf{x}_i) y_i \geq 1, \quad 1 \leq i \leq n$$

# Max-Margin Classifier

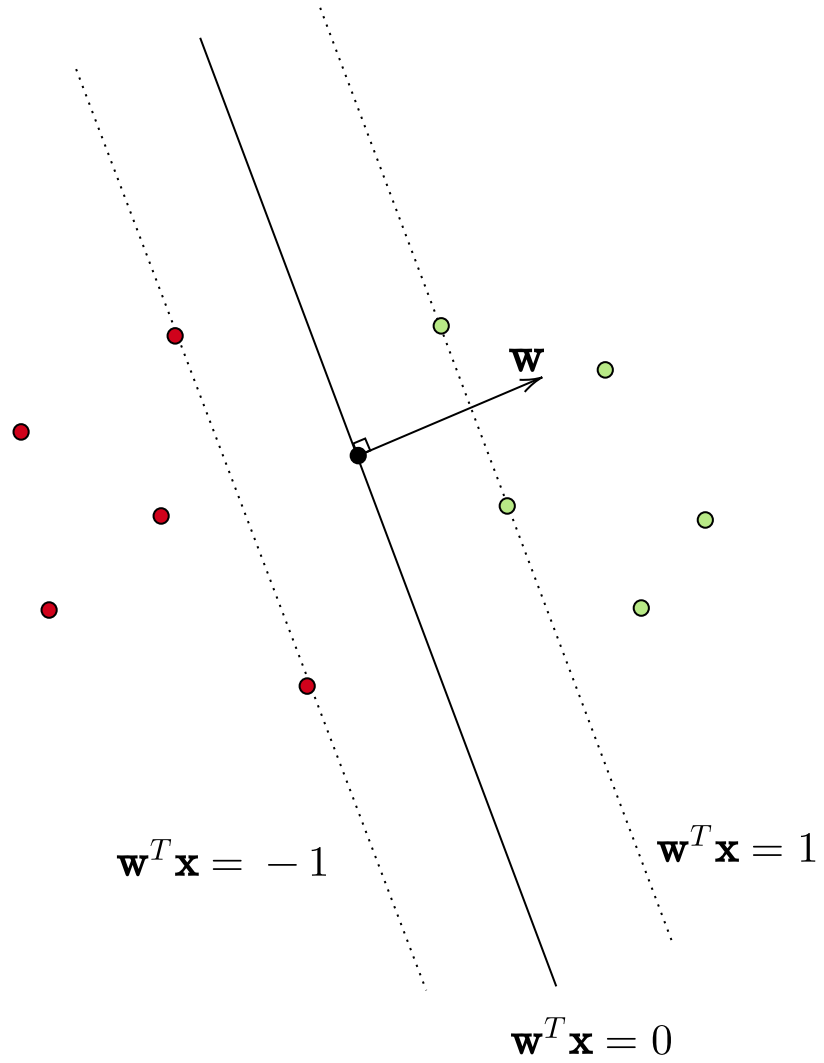


$$\max_{\mathbf{w}} \quad \frac{1}{\|\mathbf{w}\|}$$

sub. to

$$(\mathbf{w}^T \mathbf{x}_i) y_i \geq 1, \quad 1 \leq i \leq n$$

# Max-Margin Classifier



$$\min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2}$$

sub. to

$$(\mathbf{w}^T \mathbf{x}_i) y_i \geq 1, \quad 1 \leq i \leq n$$

# Primal and Dual

$$\min_{\mathbf{w}} f(\mathbf{w})$$

sub. to

$$g(\mathbf{w}) \leq 0$$

# Primal and Dual

$$\min_{\mathbf{w}} \quad f(\mathbf{w})$$

sub. to

$$g(\mathbf{w}) \leq 0$$

$$\max_{\alpha \geq 0} \quad f(\mathbf{w}) + \alpha g(\mathbf{w}) = \begin{cases} f(\mathbf{w}), & g(\mathbf{w}) \leq 0 \\ \infty, & g(\mathbf{w}) > 0 \end{cases}$$

# Primal and Dual

$$\max_{\alpha \geq 0} \quad f(\mathbf{w}) + \alpha g(\mathbf{w}) = \begin{cases} f(\mathbf{w}), & g(\mathbf{w}) \leq 0 \\ \infty, & g(\mathbf{w}) > 0 \end{cases}$$

$$\min_{\mathbf{w}} \quad f(\mathbf{w})$$

sub. to  $\equiv$

$$\min_{\mathbf{w}} \left[ \max_{\alpha \geq 0} \quad f(\mathbf{w}) + \alpha g(\mathbf{w}) \right]$$

$$g(\mathbf{w}) \leq 0$$

# Primal and Dual

$$\max_{\alpha \geq 0} f(\mathbf{w}) + \alpha g(\mathbf{w}) = \begin{cases} f(\mathbf{w}), & g(\mathbf{w}) \leq 0 \\ \infty, & g(\mathbf{w}) > 0 \end{cases}$$

$$\min_{\mathbf{w}} f(\mathbf{w})$$

sub. to

$\equiv$

$$\min_{\mathbf{w}} \left[ \max_{\alpha \geq 0} f(\mathbf{w}) + \alpha g(\mathbf{w}) \right]$$

$\equiv$

$$\max_{\alpha \geq 0} \min_{\mathbf{w}} f(\mathbf{w}) + \alpha g(\mathbf{w})$$



Strong Duality

$$g(\mathbf{w}) \leq 0$$



# Primal and Dual

$$\max_{\alpha \geq 0} f(\mathbf{w}) + \alpha g(\mathbf{w}) = \begin{cases} f(\mathbf{w}), & g(\mathbf{w}) \leq 0 \\ \infty, & g(\mathbf{w}) > 0 \end{cases}$$

$$\min_{\mathbf{w}} f(\mathbf{w})$$

sub. to

$$g(\mathbf{w}) \leq 0$$

$\equiv$

$$\min_{\mathbf{w}} \left[ \max_{\alpha \geq 0} f(\mathbf{w}) + \alpha g(\mathbf{w}) \right]$$

$\equiv$

$$\max_{\alpha \geq 0} \min_{\mathbf{w}} f(\mathbf{w}) + \alpha g(\mathbf{w})$$



Strong Duality

$$\min_{\mathbf{w}} \frac{||\mathbf{w}||^2}{2}$$

sub. to

$$(\mathbf{w}^T \mathbf{x}_i) y_i \geq 1, \quad 1 \leq i \leq n$$

# Primal and Dual

$$\max_{\alpha \geq 0} f(\mathbf{w}) + \alpha g(\mathbf{w}) = \begin{cases} f(\mathbf{w}), & g(\mathbf{w}) \leq 0 \\ \infty, & g(\mathbf{w}) > 0 \end{cases}$$

$$\min_{\mathbf{w}} f(\mathbf{w})$$

sub. to

$\equiv$

$$\min_{\mathbf{w}} \left[ \max_{\alpha \geq 0} f(\mathbf{w}) + \alpha g(\mathbf{w}) \right]$$

$\equiv$

$$\max_{\alpha \geq 0} \min_{\mathbf{w}} f(\mathbf{w}) + \alpha g(\mathbf{w})$$



Strong Duality

$$g(\mathbf{w}) \leq 0$$

$$\min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2}$$

$\equiv$

$$\min_{\mathbf{w}} \max_{\alpha \geq 0} \frac{\|\mathbf{w}\|^2}{2} + \sum_{i=1}^n \alpha_i \left[ 1 - (\mathbf{w}^T \mathbf{x}_i) y_i \right]$$

sub. to

$$(\mathbf{w}^T \mathbf{x}_i) y_i \geq 1, \quad 1 \leq i \leq n$$

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$$

# Primal and Dual

$$\max_{\alpha \geq 0} f(\mathbf{w}) + \alpha g(\mathbf{w}) = \begin{cases} f(\mathbf{w}), & g(\mathbf{w}) \leq 0 \\ \infty, & g(\mathbf{w}) > 0 \end{cases}$$

$$\min_{\mathbf{w}} f(\mathbf{w})$$

sub. to

$$g(\mathbf{w}) \leq 0$$

$\equiv$

$$\min_{\mathbf{w}} \left[ \max_{\alpha \geq 0} f(\mathbf{w}) + \alpha g(\mathbf{w}) \right]$$

$\equiv$

$$\max_{\alpha \geq 0} \min_{\mathbf{w}} f(\mathbf{w}) + \alpha g(\mathbf{w})$$



Strong Duality



$$\min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2}$$

sub. to

$\equiv$

$$\min_{\mathbf{w}} \max_{\alpha \geq 0} \frac{\|\mathbf{w}\|^2}{2} + \sum_{i=1}^n \alpha_i \left[ 1 - (\mathbf{w}^T \mathbf{x}_i) y_i \right]$$

$\equiv$

$$\max_{\alpha \geq 0} \min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} + \sum_{i=1}^n \alpha_i \left[ 1 - (\mathbf{w}^T \mathbf{x}_i) y_i \right]$$

$$(\mathbf{w}^T \mathbf{x}_i) y_i \geq 1, \quad 1 \leq i \leq n$$

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$$

# Formulating the Dual

$$\max_{\boldsymbol{\alpha} \geq 0} \min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} + \sum_{i=1}^n \alpha_i [1 - (\mathbf{w}^T \mathbf{x}_i) y_i]$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}_i y_i$$

# Formulating the Dual

$$\max_{\boldsymbol{\alpha} \geq 0} \min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} + \sum_{i=1}^n \alpha_i [1 - (\mathbf{w}^T \mathbf{x}_i) y_i]$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}_i y_i$$

$d \times n$

$$\mathbf{X} = \begin{bmatrix} | & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ | & & | \end{bmatrix}$$

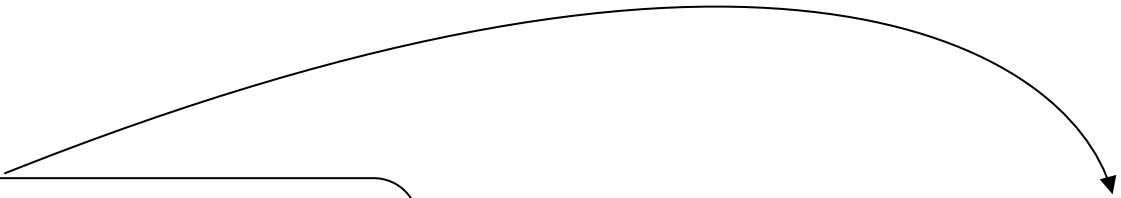
$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$$

$n \times n$

$$\mathbf{Y} = \begin{bmatrix} y_1 & & 0 \\ & \ddots & \\ 0 & & y_n \end{bmatrix}$$

$$\mathbf{XY}\boldsymbol{\alpha} = \begin{bmatrix} | & & | \\ y_1 \mathbf{x}_1 & \cdots & y_n \mathbf{x}_n \\ | & & | \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \sum_{i=1}^n \alpha_i (y_i \mathbf{x}_i) = \mathbf{w}$$

# Formulating the Dual

$$\max_{\boldsymbol{\alpha} \geq 0} \min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} + \sum_{i=1}^n \alpha_i [1 - (\mathbf{w}^T \mathbf{x}_i) y_i]$$


$$\mathbf{w} = \mathbf{X}\mathbf{Y}\boldsymbol{\alpha}$$

$d \times n$

$$\mathbf{X} = \begin{bmatrix} | & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n \\ | & & | \end{bmatrix}$$

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$$

$n \times n$

$$\mathbf{Y} = \begin{bmatrix} y_1 & & 0 \\ & \ddots & \\ 0 & & y_n \end{bmatrix}$$

$$\max_{\boldsymbol{\alpha} \geq 0} \boldsymbol{\alpha}^T \mathbf{1} - \frac{\boldsymbol{\alpha}^T (\mathbf{Y}^T \mathbf{X}^T \mathbf{X} \mathbf{Y}) \boldsymbol{\alpha}}{2}$$

# Advantages of the Dual

$$\max_{\boldsymbol{\alpha} \geq 0} \boldsymbol{\alpha}^T \mathbf{1} - \frac{\boldsymbol{\alpha}^T (\mathbf{Y}^T \mathbf{X}^T \mathbf{X} \mathbf{Y}) \boldsymbol{\alpha}}{2}$$

- (1)  $\boldsymbol{\alpha} \in \mathbb{R}^n$  and  $\mathbf{w} \in \mathbb{R}^d$ , if  $n \ll d$ , we are solving for fewer variables in the dual
- (2) Simpler constraints, just bounds.
- (3) The appearance of  $\mathbf{X}^T \mathbf{X} \Rightarrow$  kernels.

# Support Vectors

$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix}$$

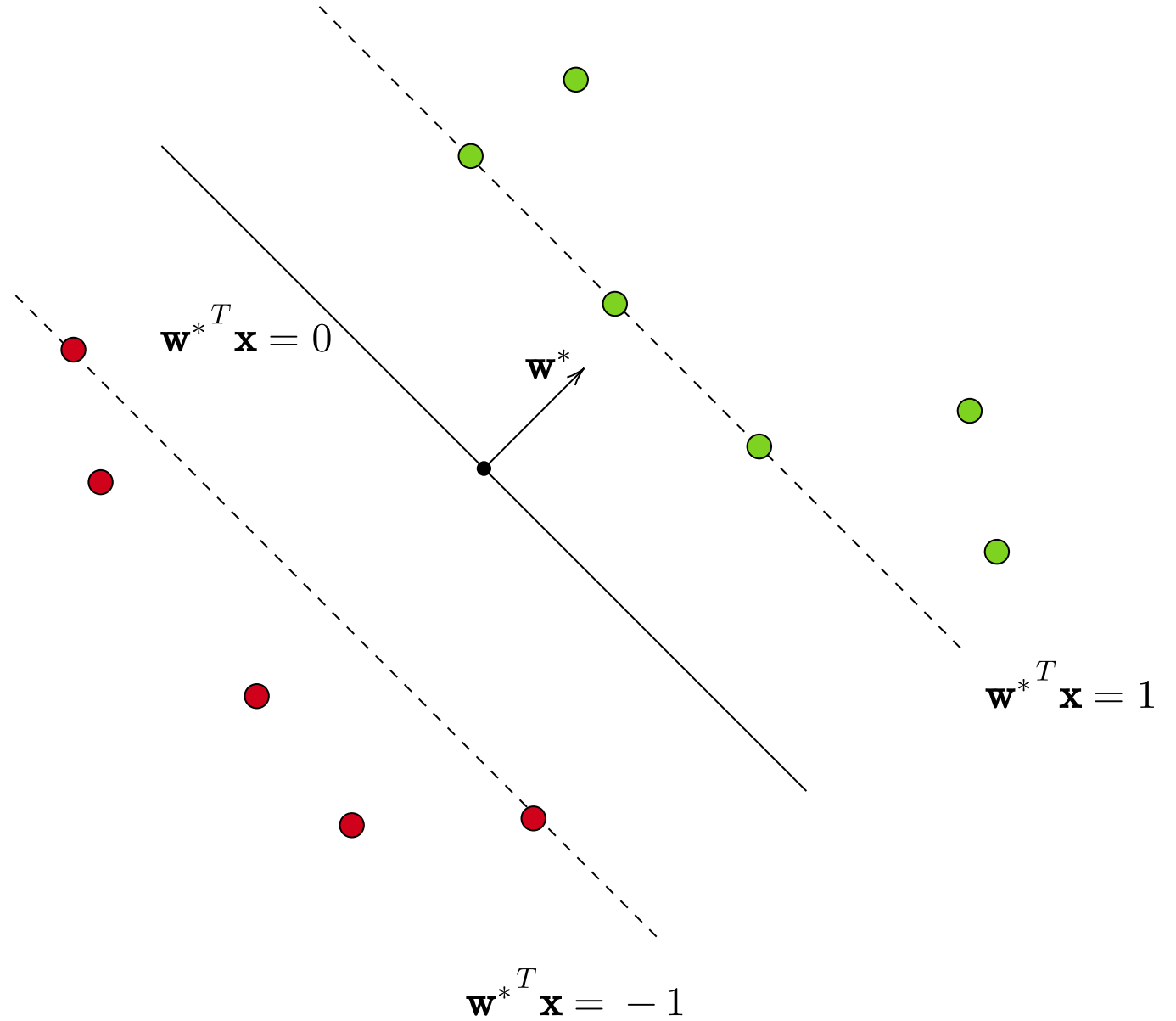
$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$



# Support Vectors

$$\mathbf{a}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix}$$

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

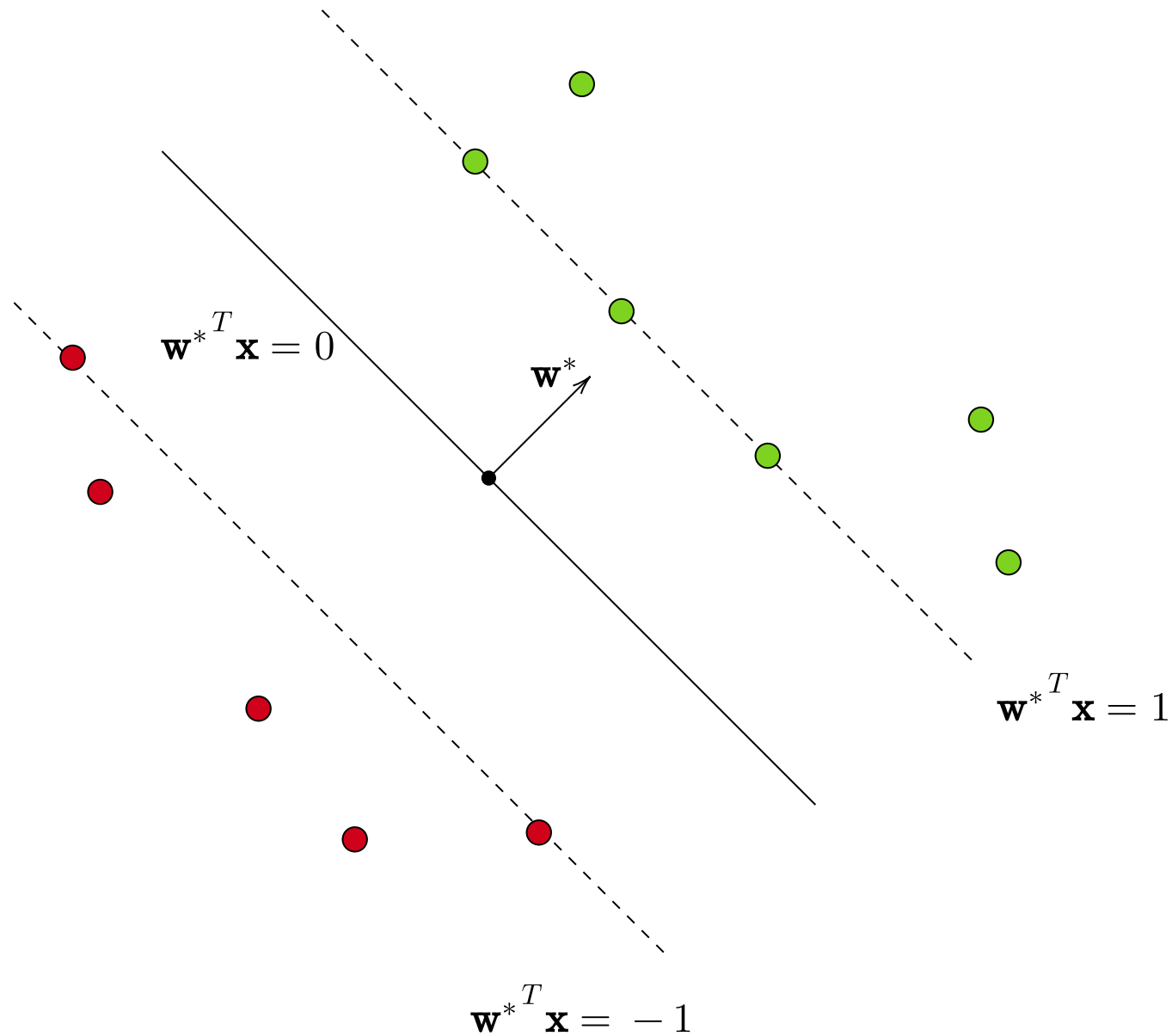


# Support Vectors

$$\mathbf{a}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix}$$

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

$$\alpha_i^* \left[ 1 - (\mathbf{w}^{*T} \mathbf{x}_i) y_i \right] = 0$$



# Support Vectors

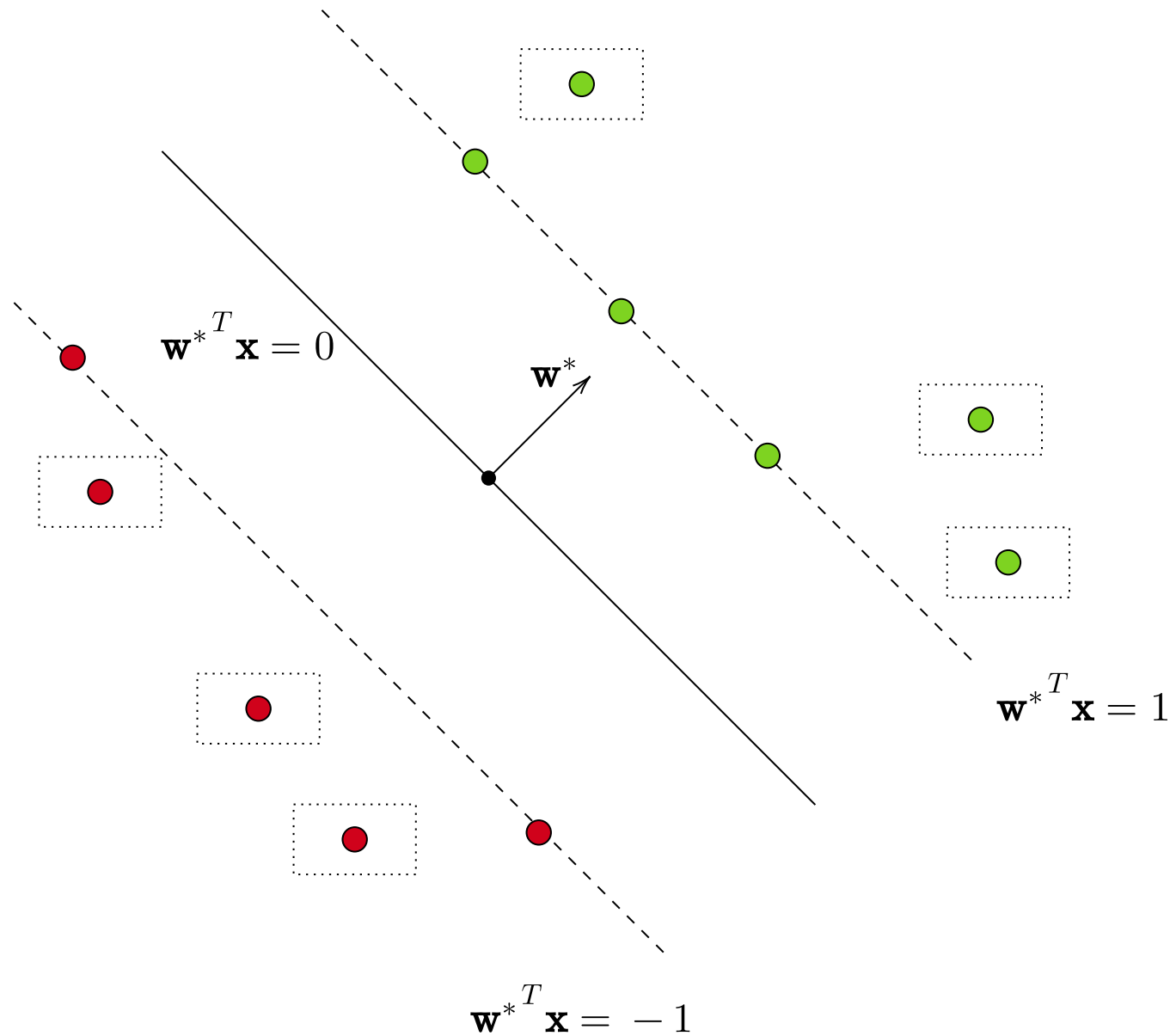
$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix}$$

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

$$\alpha_i^* \left[ 1 - (\mathbf{w}^{*T} \mathbf{x}_i) y_i \right] = 0$$

$$y_i (\mathbf{w}^{*T} \mathbf{x}_i) \neq 1$$

$$\alpha_i^* = 0$$



# Support Vectors

$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix}$$

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

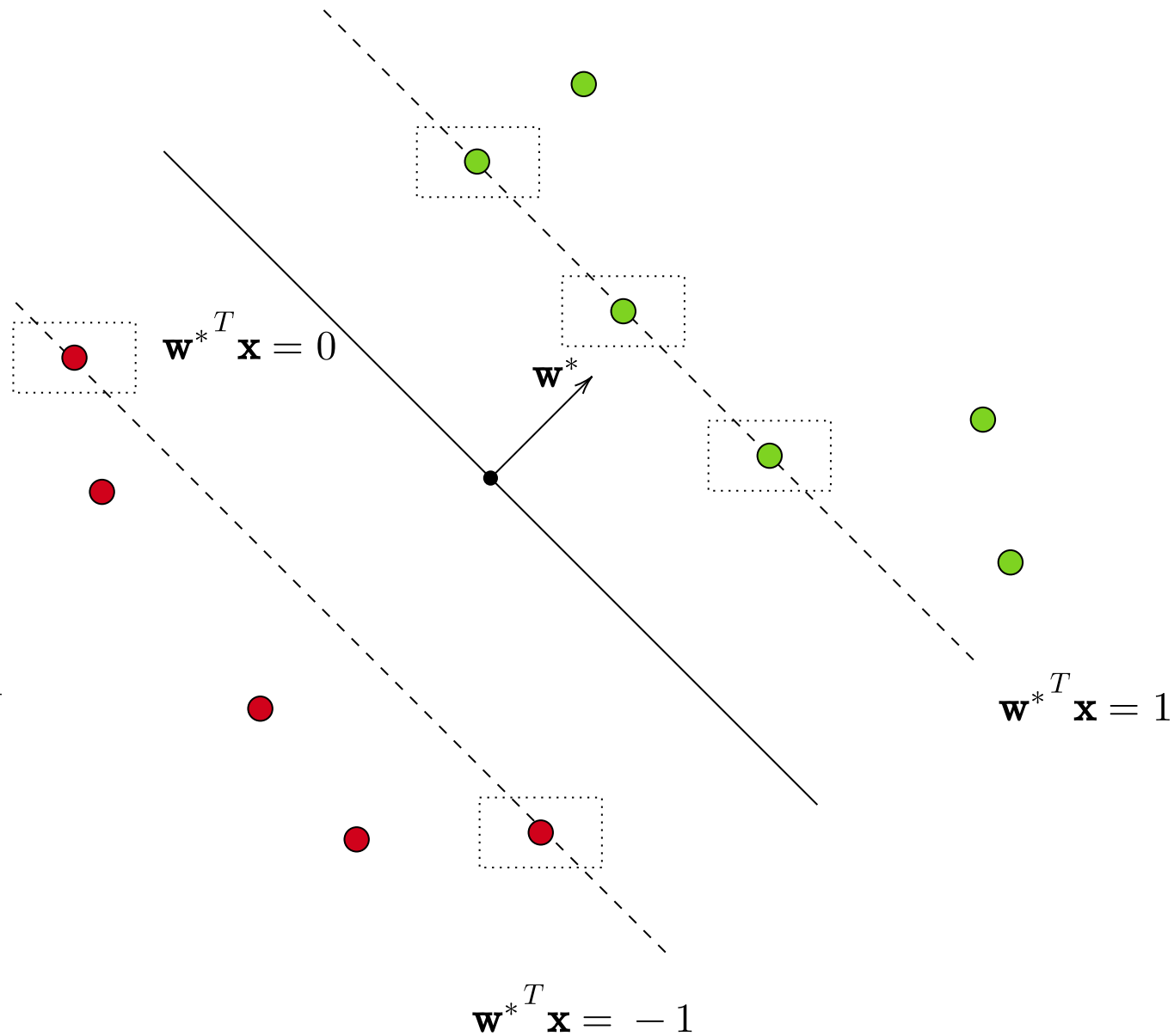
$$\alpha_i^* \left[ 1 - (\mathbf{w}^{*T} \mathbf{x}_i) y_i \right] = 0$$

$$y_i (\mathbf{w}^{*T} \mathbf{x}_i) \neq 1$$

$$\alpha_i^* = 0$$

$$y_i (\mathbf{w}^{*T} \mathbf{x}_i) = 1$$

$$\alpha_i^* \geq 0$$



# Support Vectors

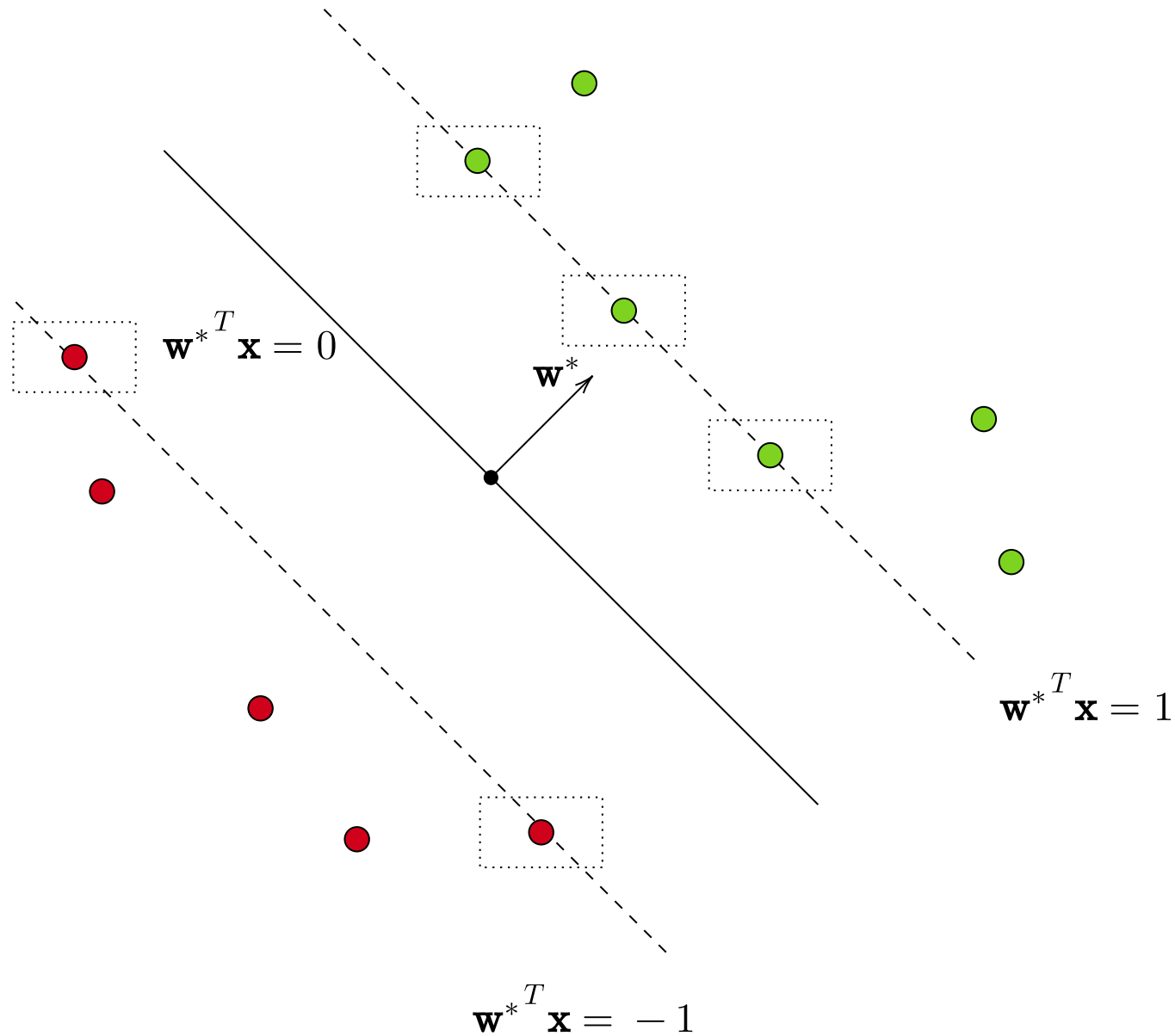
$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix}$$

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

$$\alpha_i^* \left[ 1 - (\mathbf{w}^{*T} \mathbf{x}_i) y_i \right] = 0$$

$$\alpha_i^* > 0$$

$$y_i (\mathbf{w}^{*T} \mathbf{x}_i) = 1$$



# Support Vectors

$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix}$$

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

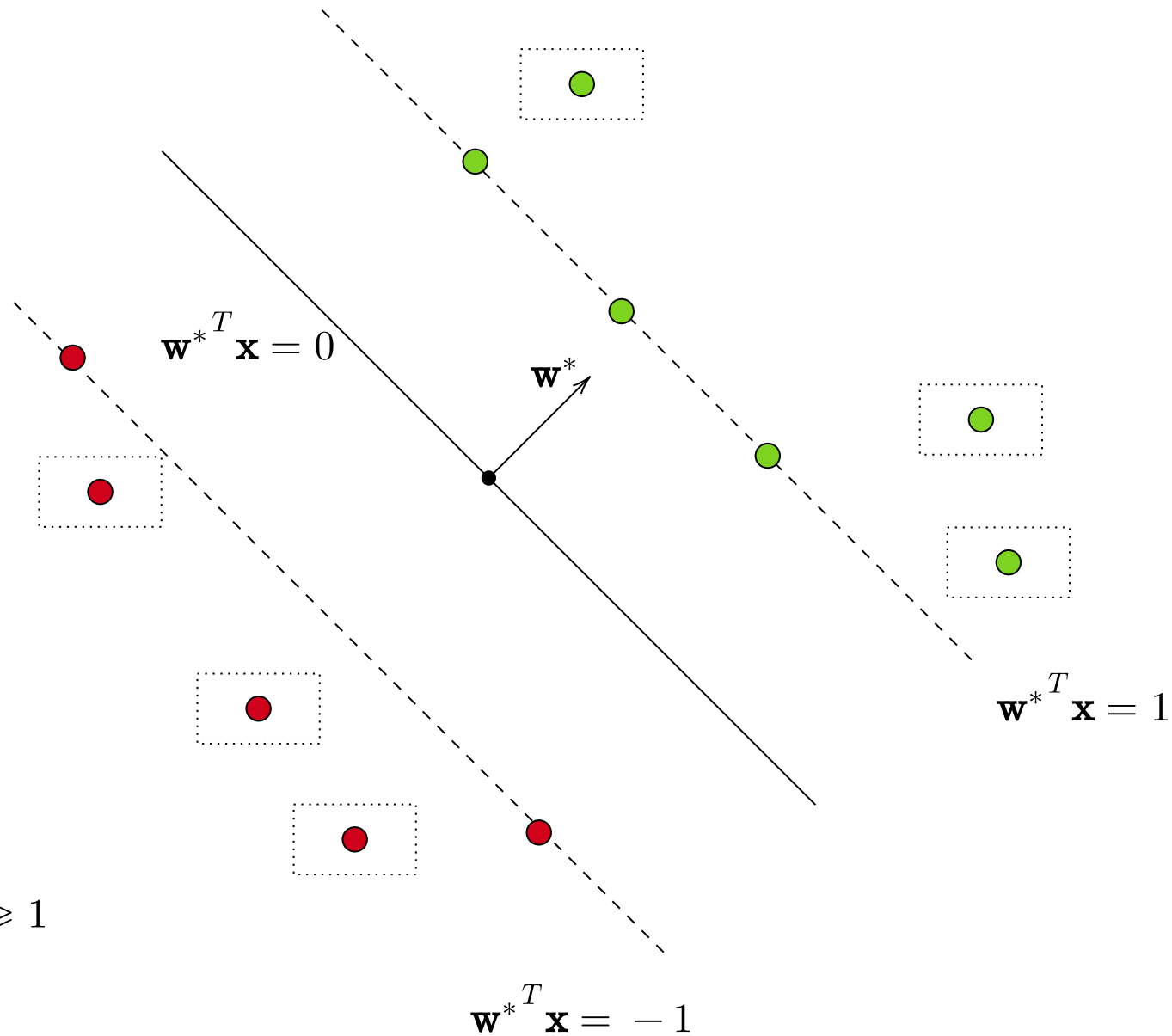
$$\alpha_i^* \left[ 1 - (\mathbf{w}^{*T} \mathbf{x}_i) y_i \right] = 0$$

$$\alpha_i^* > 0$$

$$\alpha_i^* = 0$$

$$y_i (\mathbf{w}^{*T} \mathbf{x}_i) = 1$$

$$y_i (\mathbf{w}^{*T} \mathbf{x}_i) \geq 1$$



# Support Vectors

$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix}$$

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

**Definition:** A support vector is a point for which  $\alpha_i^* > 0$

# Support Vectors

$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix} \quad \mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

**Definition:** A support vector is a point for which  $\alpha_i^* > 0$

$$\alpha_i^* g_i(\mathbf{w}^*) = 0 \implies \alpha_i^* \left[ 1 - \left( (\mathbf{w}^*)^T \mathbf{x}_i \right) y_i \right] = 0$$

Complementary Slackness



# Support Vectors

$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix} \quad \mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

$$\alpha_i^* g_i(\mathbf{w}^*) = 0 \implies \alpha_i^* \left[ 1 - \left( (\mathbf{w}^*)^T \mathbf{x}_i \right) y_i \right] = 0$$

Complementary Slackness

**Definition:** A support vector is a point for which  $\alpha_i^* > 0$

Every support vector lies on one of the two supporting hyperplanes  $(\mathbf{w}^*)^T \mathbf{x} = \pm 1$

# Support Vectors

$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix} \quad \mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

$$\alpha_i^* g_i(\mathbf{w}^*) = 0 \implies \alpha_i^* \left[ 1 - \left( (\mathbf{w}^*)^T \mathbf{x}_i \right) y_i \right] = 0$$

Complementary Slackness

**Definition:** A support vector is a point for which  $\alpha_i^* > 0$

Every support vector lies on one of the two supporting hyperplanes  $(\mathbf{w}^*)^T \mathbf{x} = \pm 1$

Every point that is **not** on one of the two supporting hyperplanes has  $\alpha_i^* = 0$ .

# Support Vectors

$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix} \quad \mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

$$\alpha_i^* g_i(\mathbf{w}^*) = 0 \implies \alpha_i^* \left[ 1 - \left( (\mathbf{w}^*)^T \mathbf{x}_i \right) y_i \right] = 0$$

## Complementary Slackness

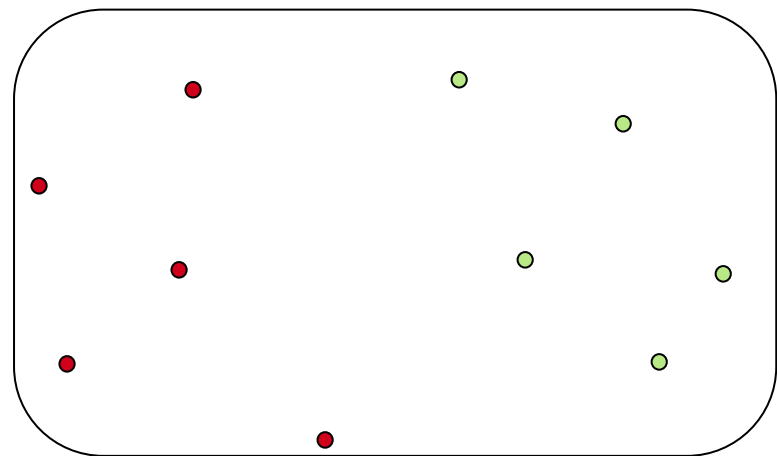
Is every point on one of the two supporting hyperplanes a support vector?

**Definition:** A support vector is a point for which  $\alpha_i^* > 0$

Every support vector lies on one of the two supporting hyperplanes  $(\mathbf{w}^*)^T \mathbf{x} = \pm 1$

Every point that is **not** on one of the two supporting hyperplanes has  $\alpha_i^* = 0$ .

# Hard-Margin, Linear-SVM

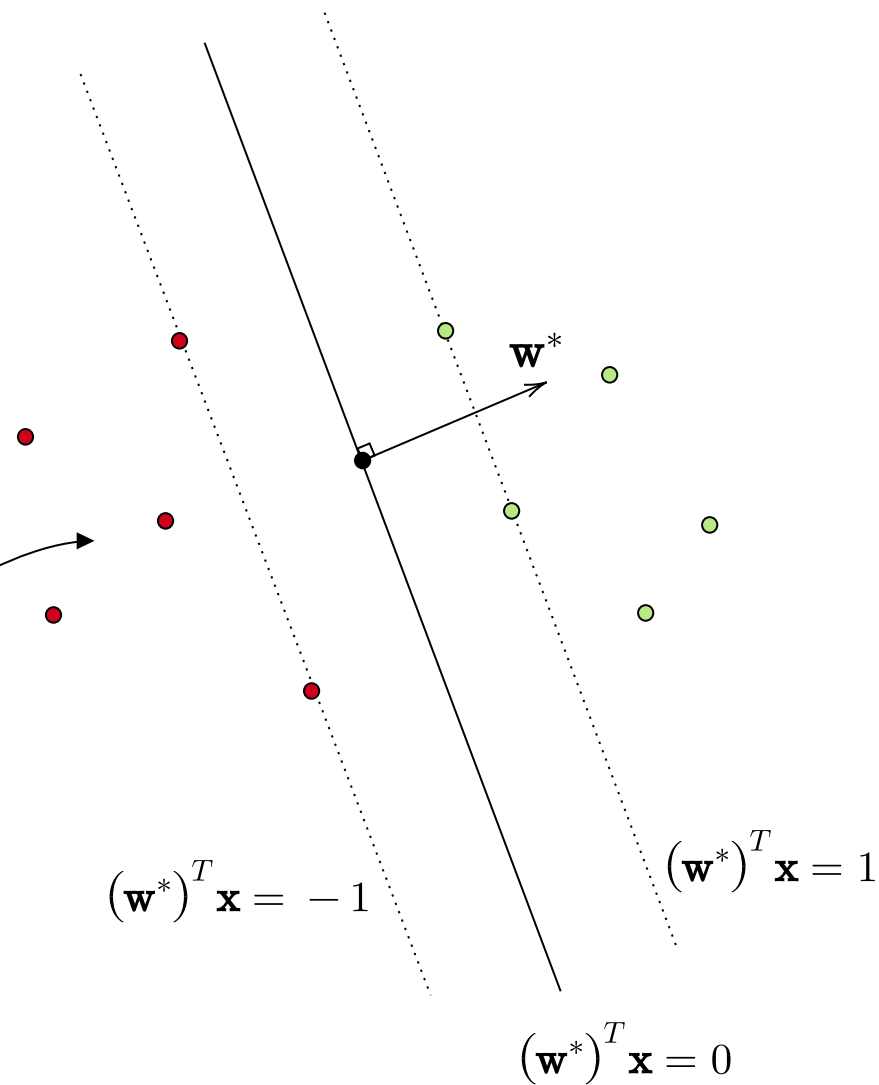


$$\hat{y} = \begin{cases} 1, & (\mathbf{w}^*)^T \mathbf{x} \geq 0 \\ -1, & (\mathbf{w}^*)^T \mathbf{x} < 0 \end{cases}$$

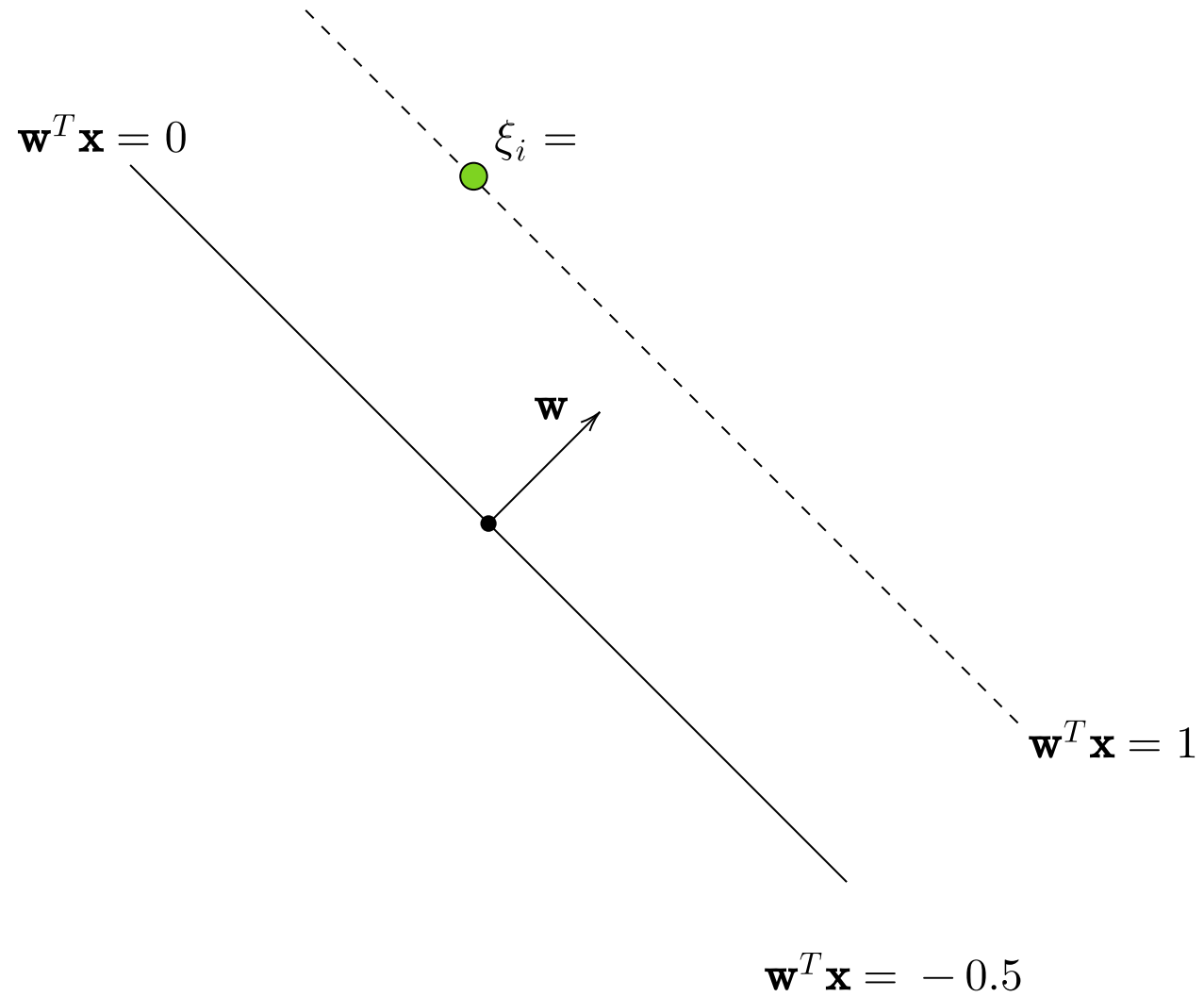
SVM  
Solver

$$\boldsymbol{\alpha}^* = \begin{bmatrix} \alpha_1^* \\ \vdots \\ \alpha_n^* \end{bmatrix}$$

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* \mathbf{x}_i y_i$$

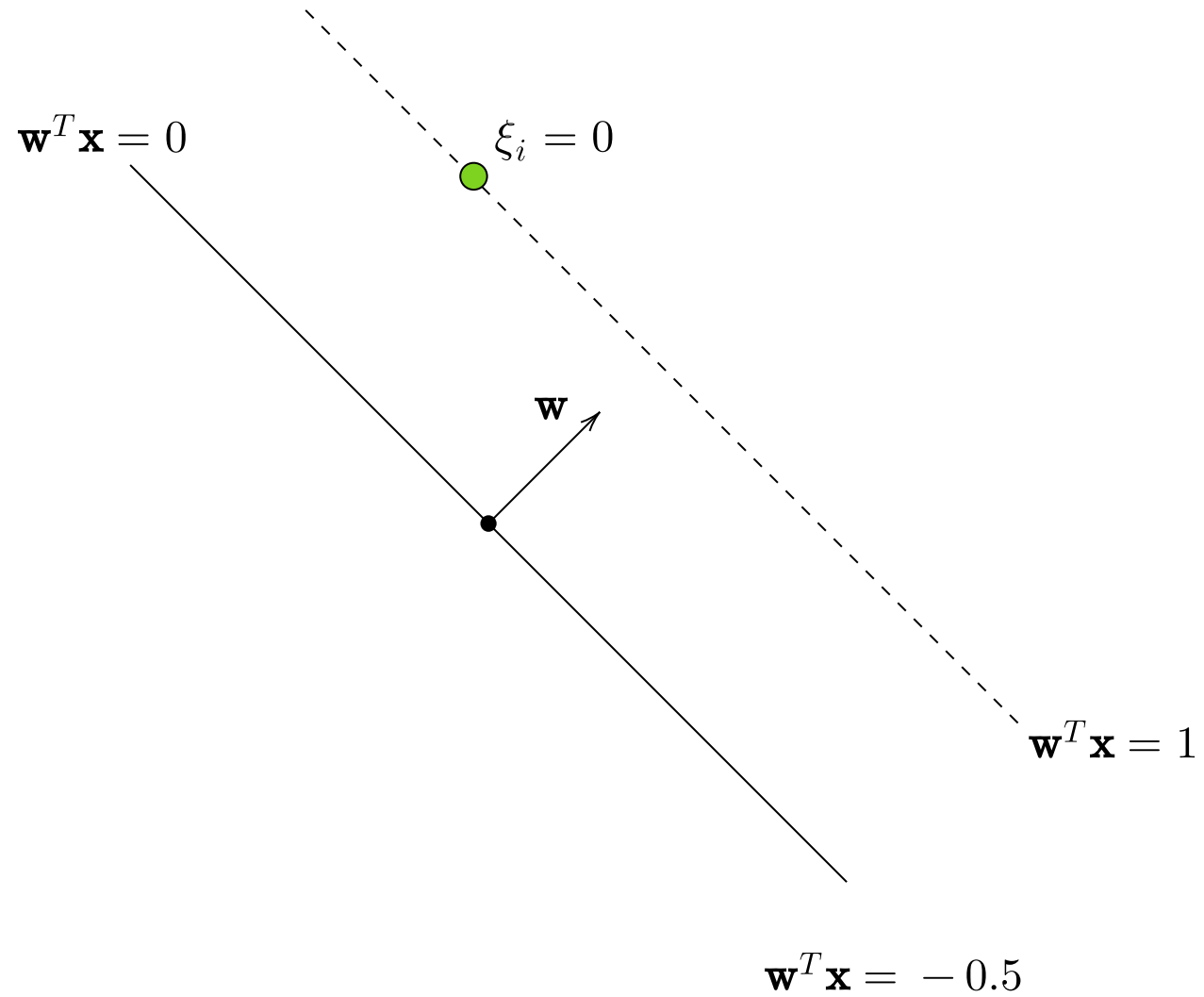


# Soft-Margin



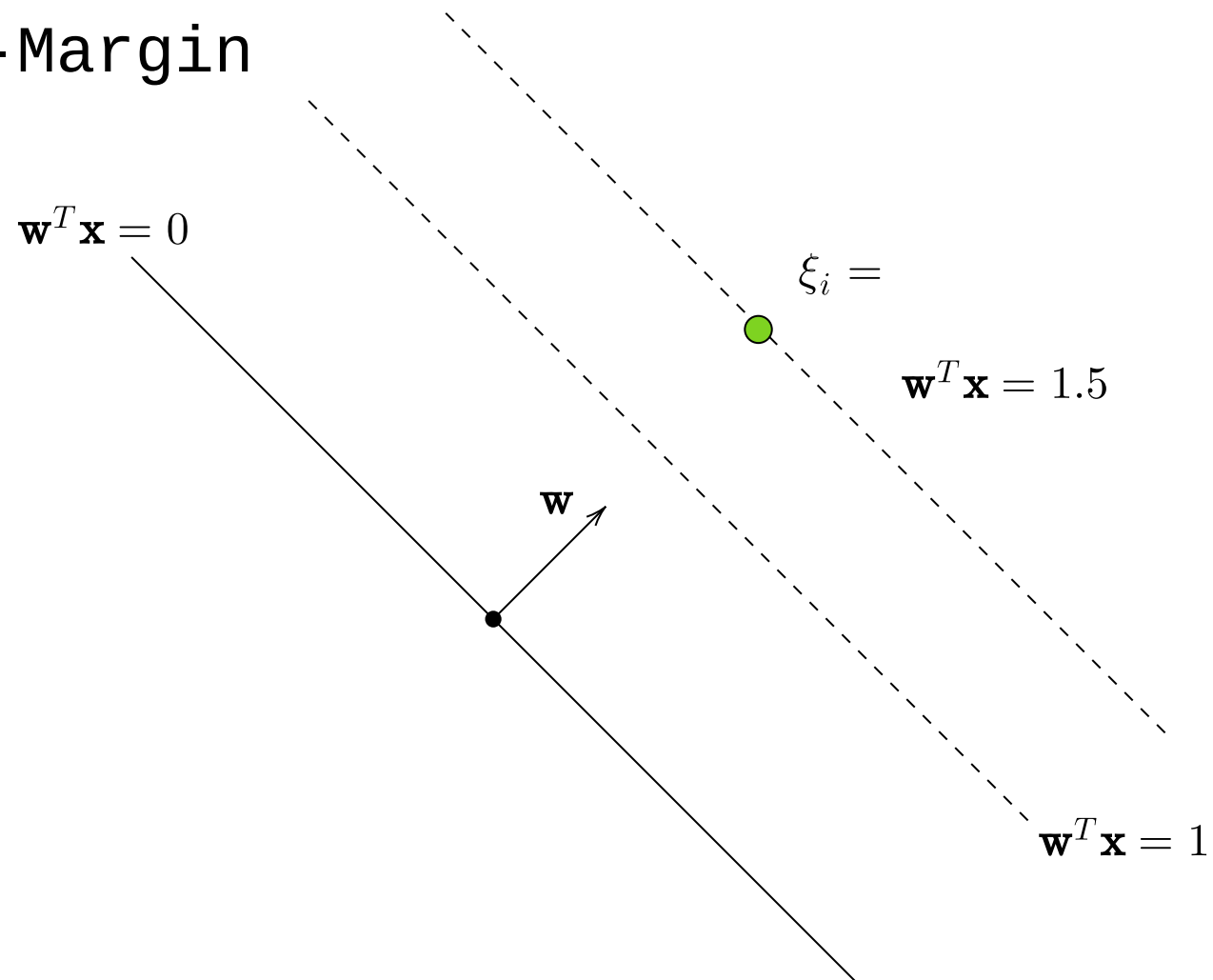
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

# Soft-Margin



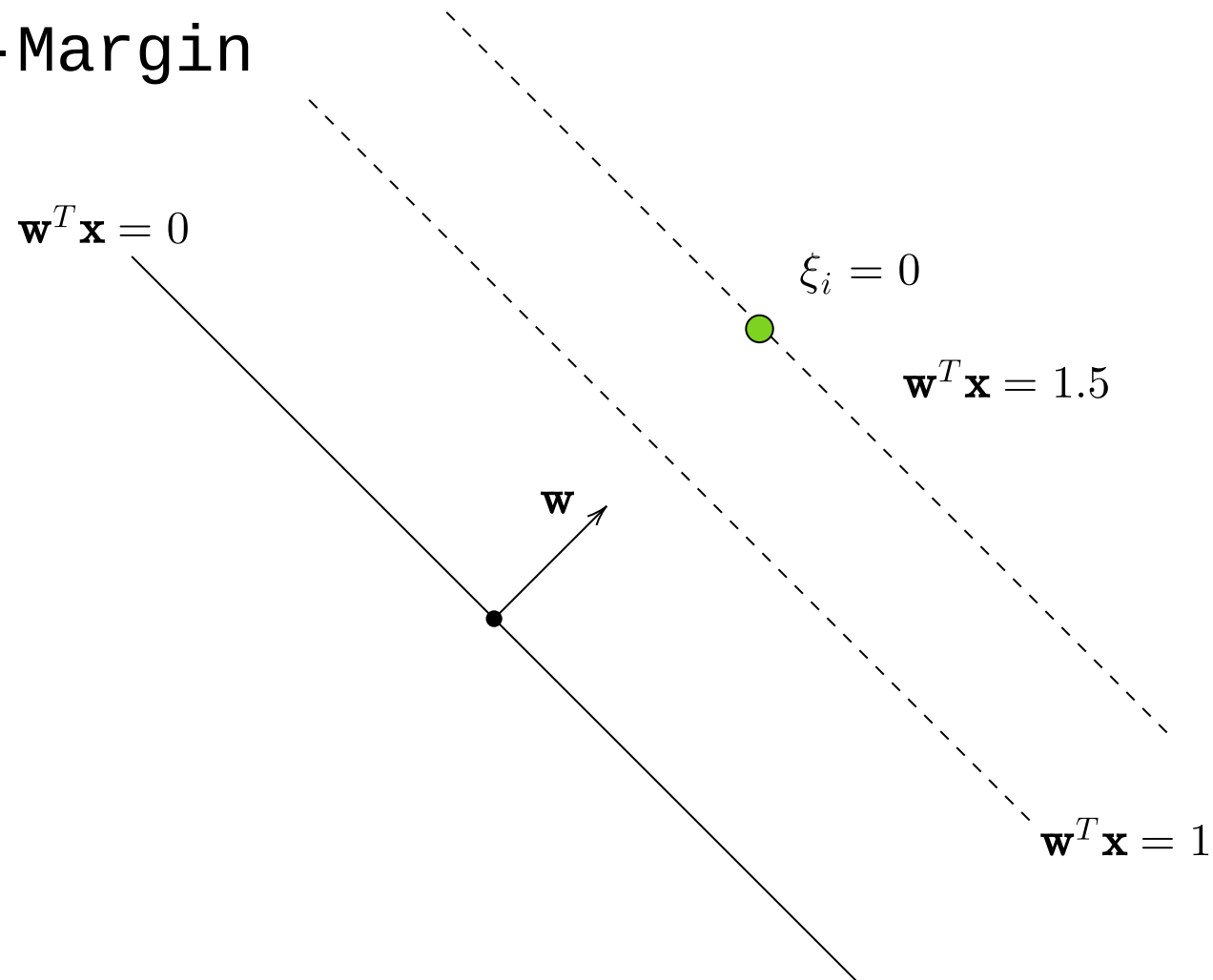
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

# Soft-Margin



$$(\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i \geq 1$$
$$\xi_i \geq 0$$

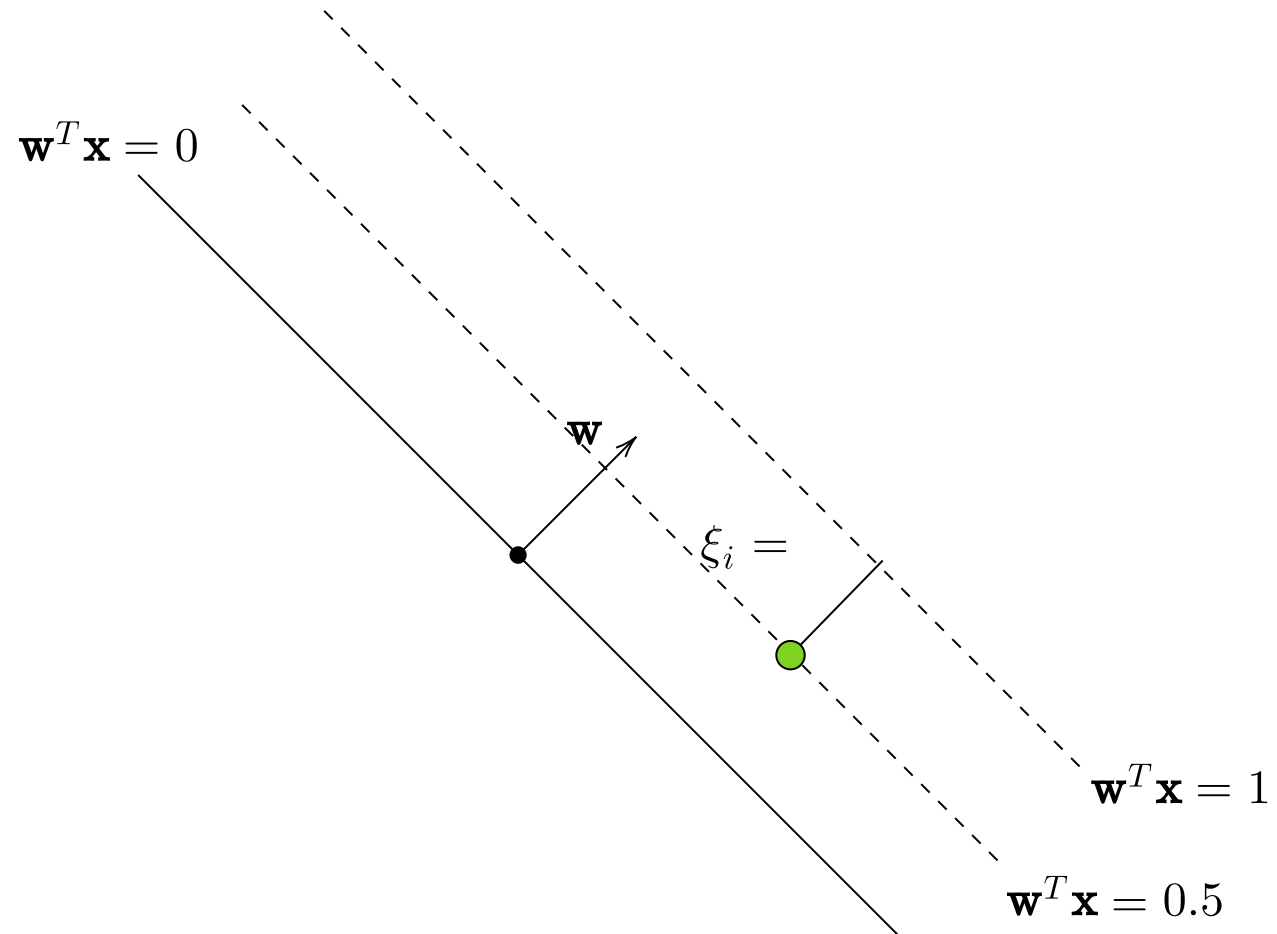
# Soft-Margin



$$(\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i \geq 1$$
$$\xi_i \geq 0$$

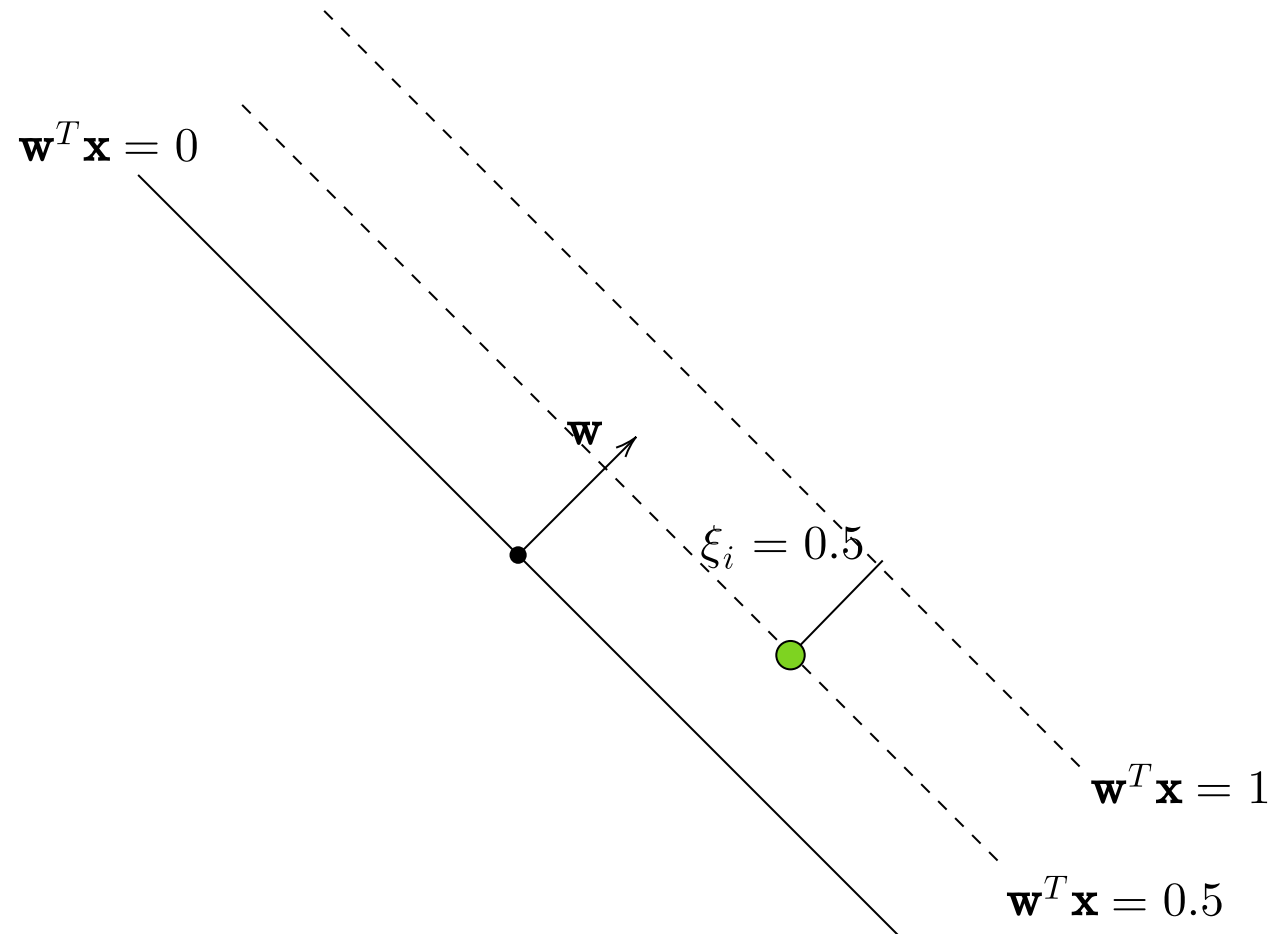


# Soft-Margin



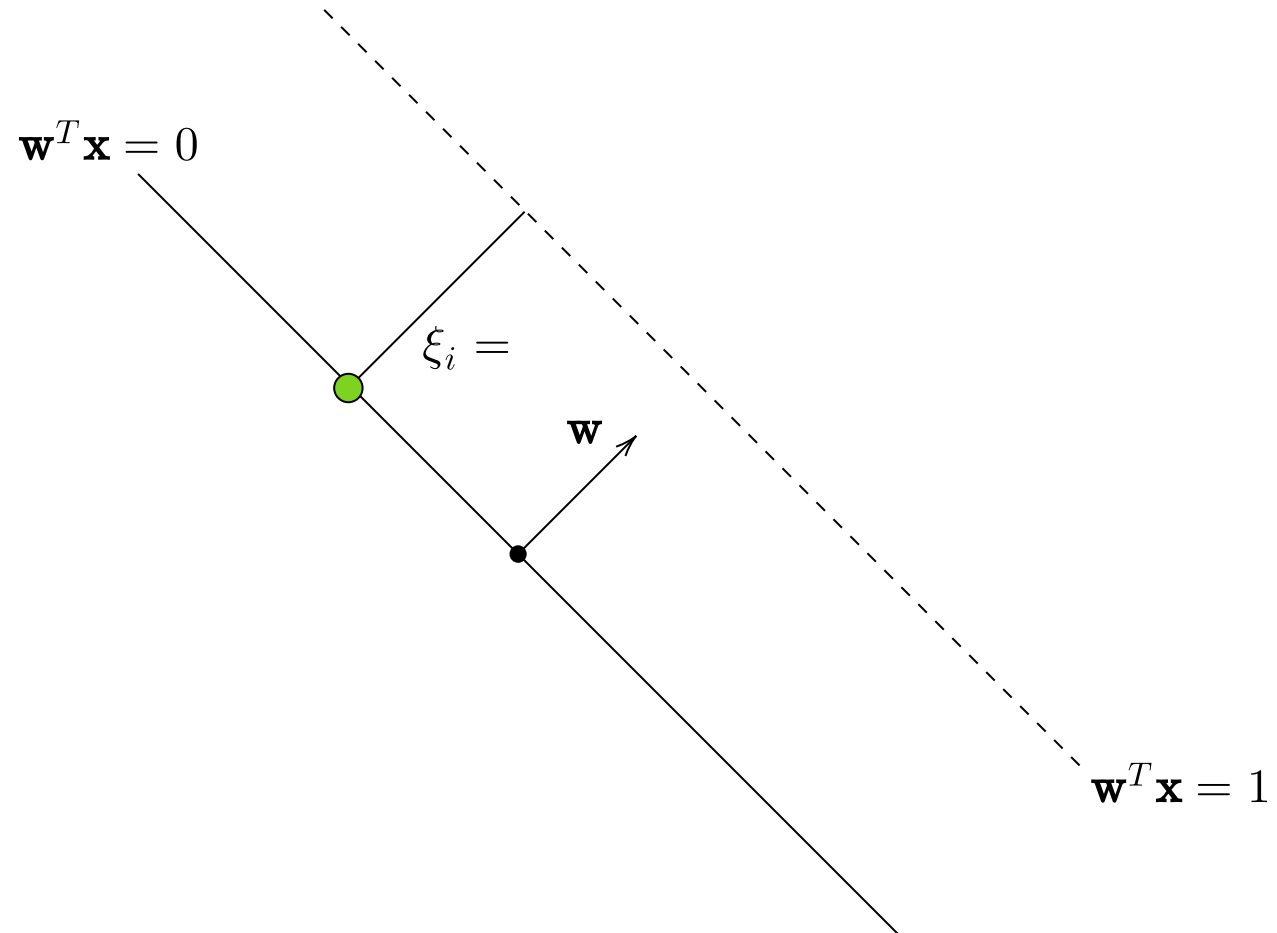
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

# Soft-Margin



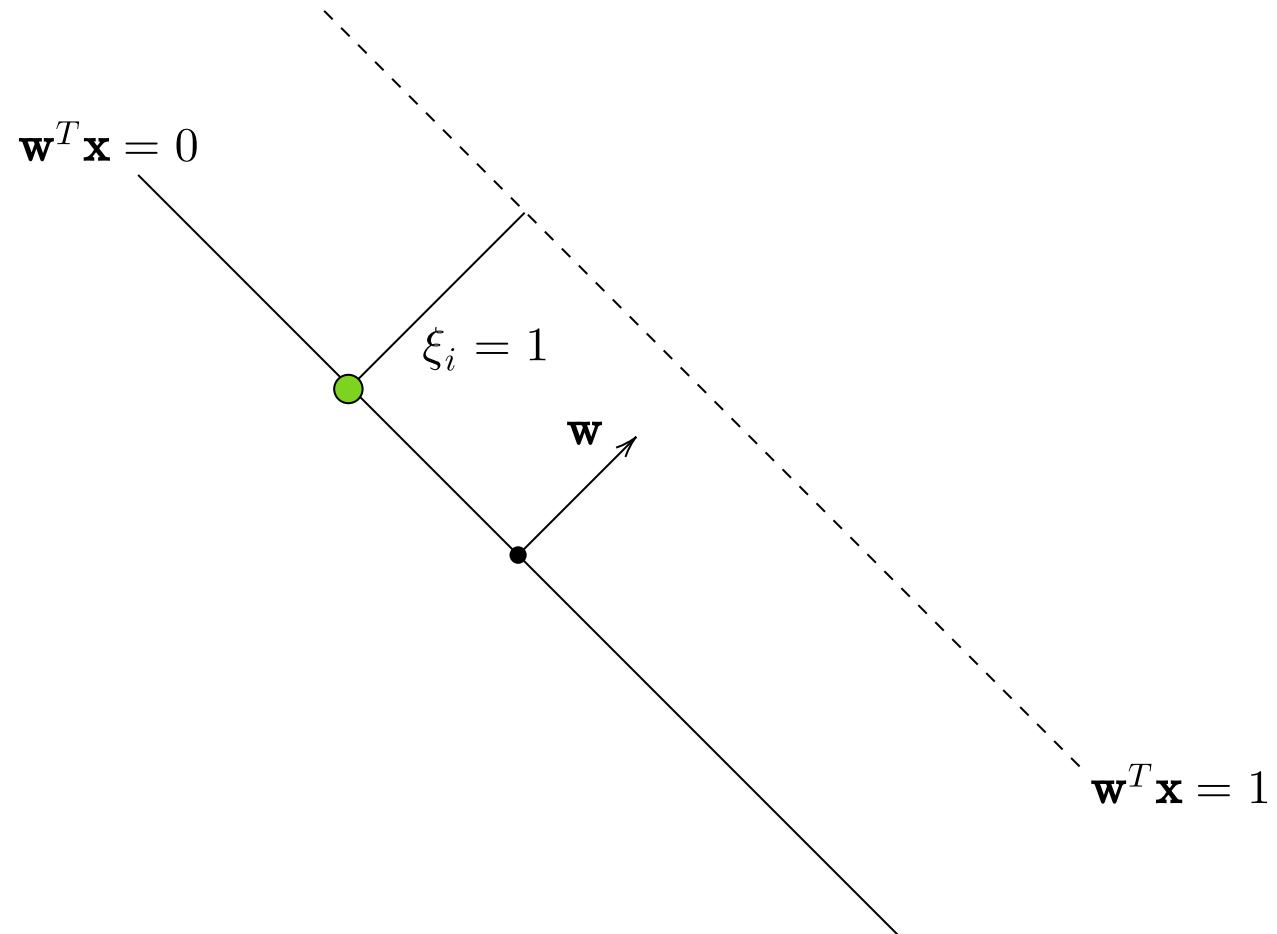
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

# Soft-Margin



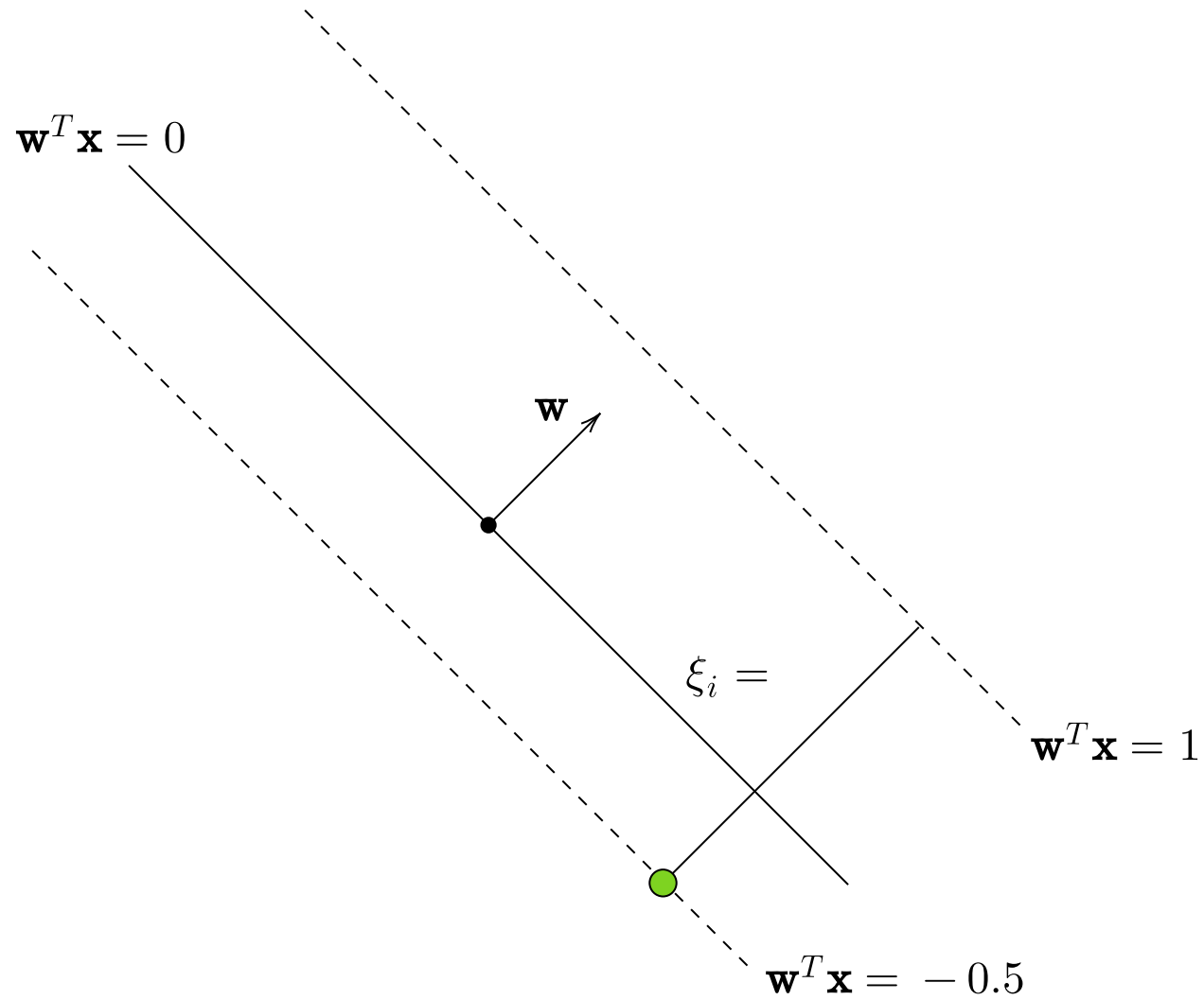
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

# Soft-Margin



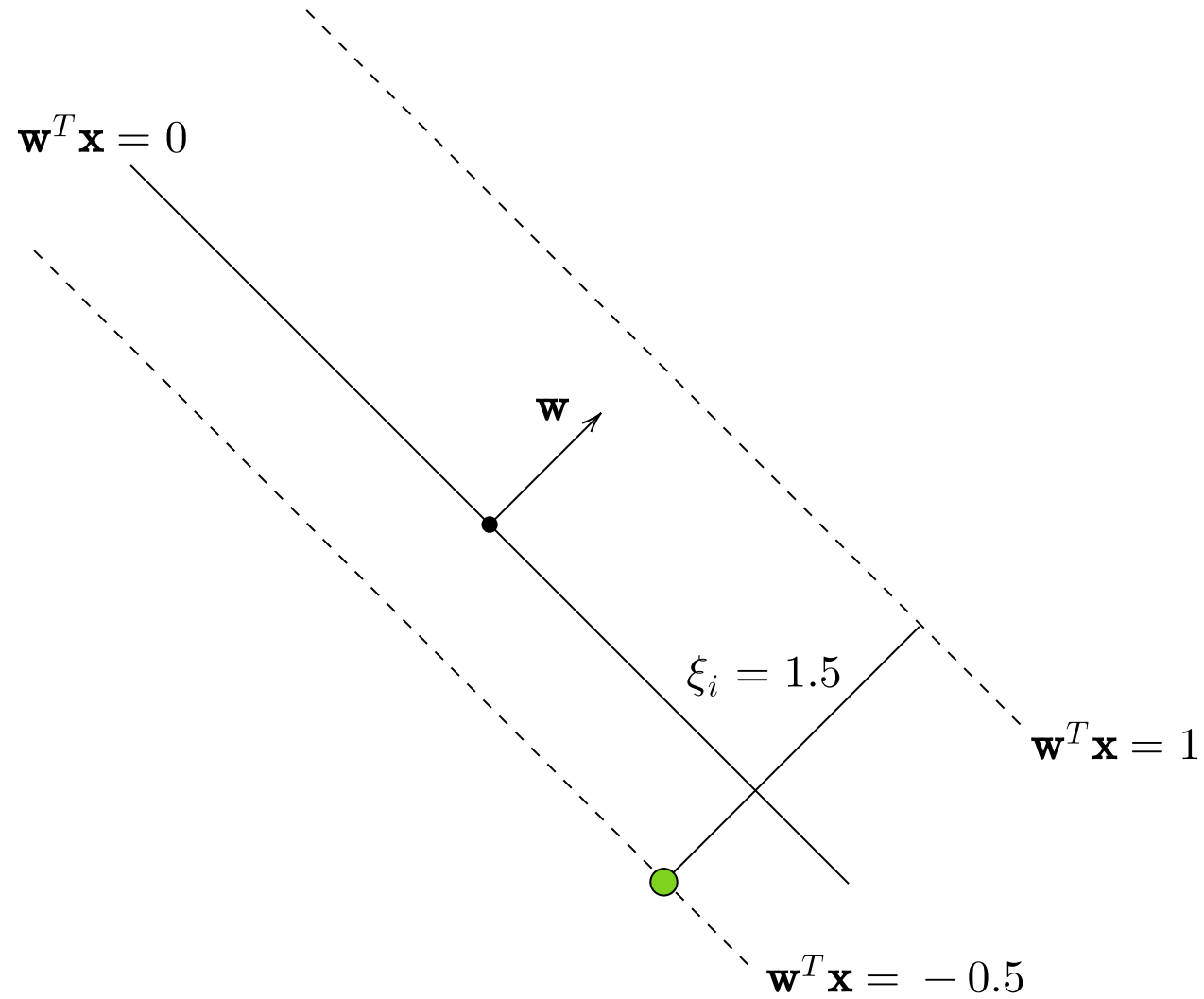
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

# Soft-Margin



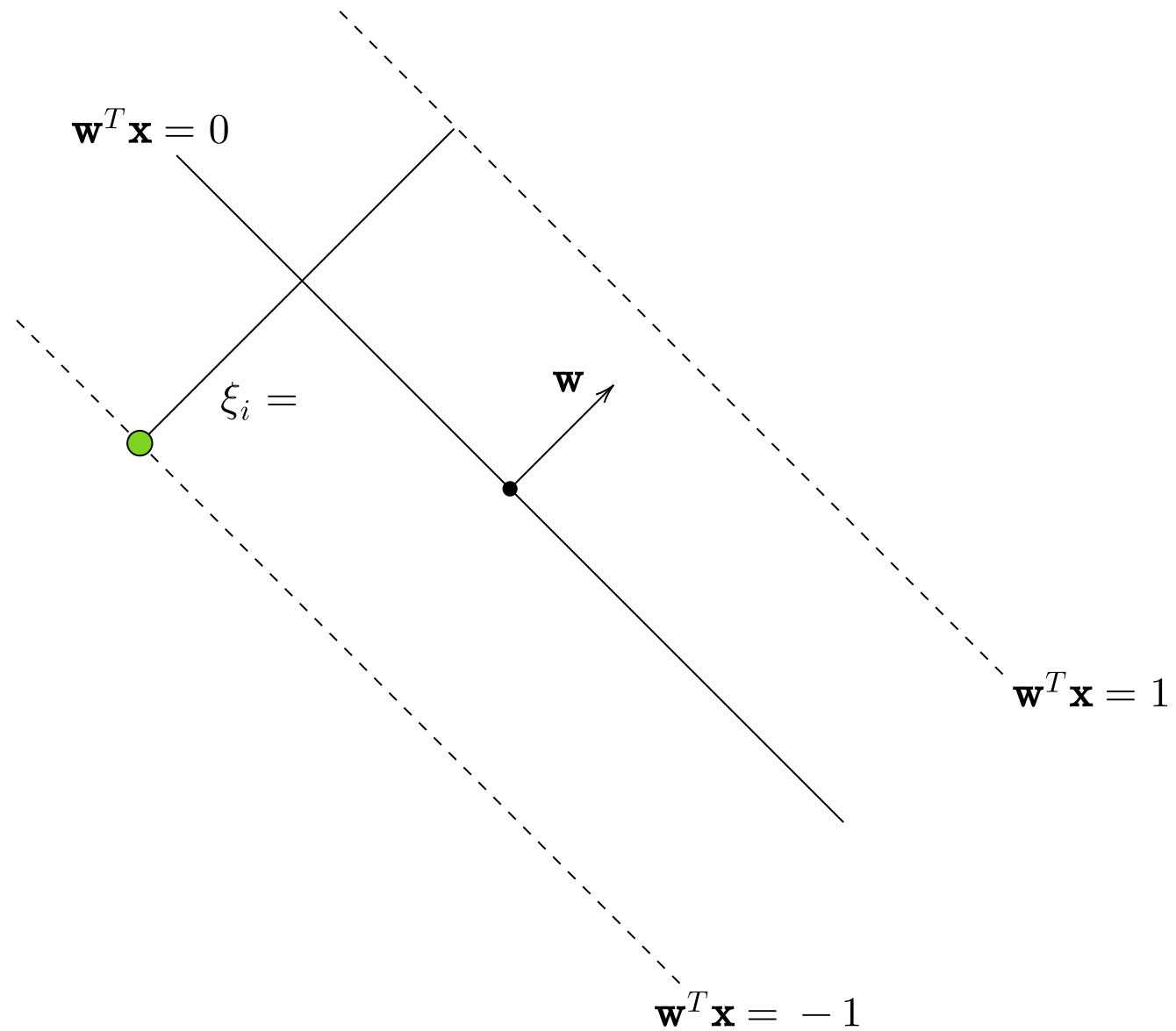
$$(\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i \geq 1$$
$$\xi_i \geq 0$$

# Soft-Margin



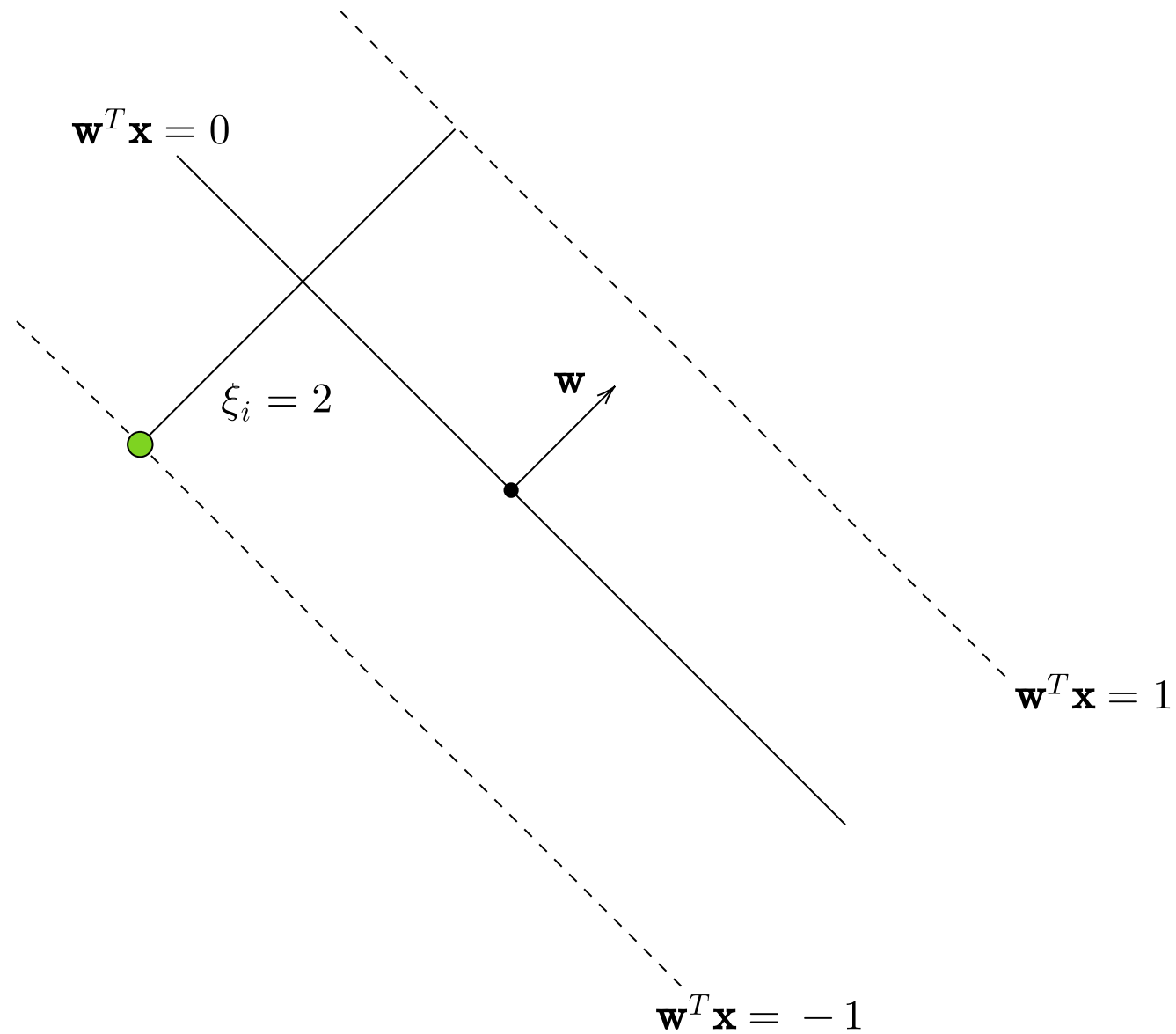
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

# Soft-Margin



$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

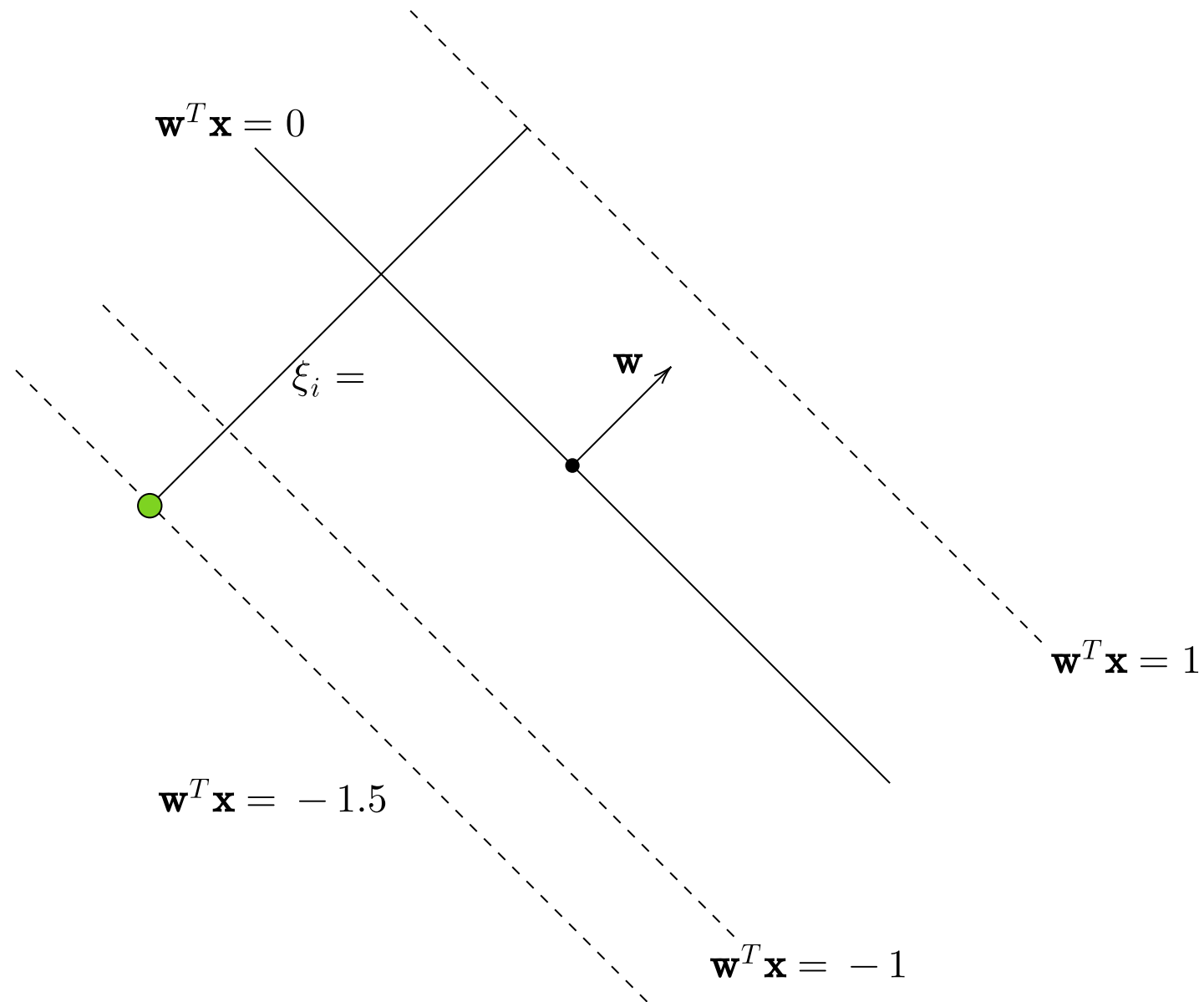
# Soft-Margin



$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

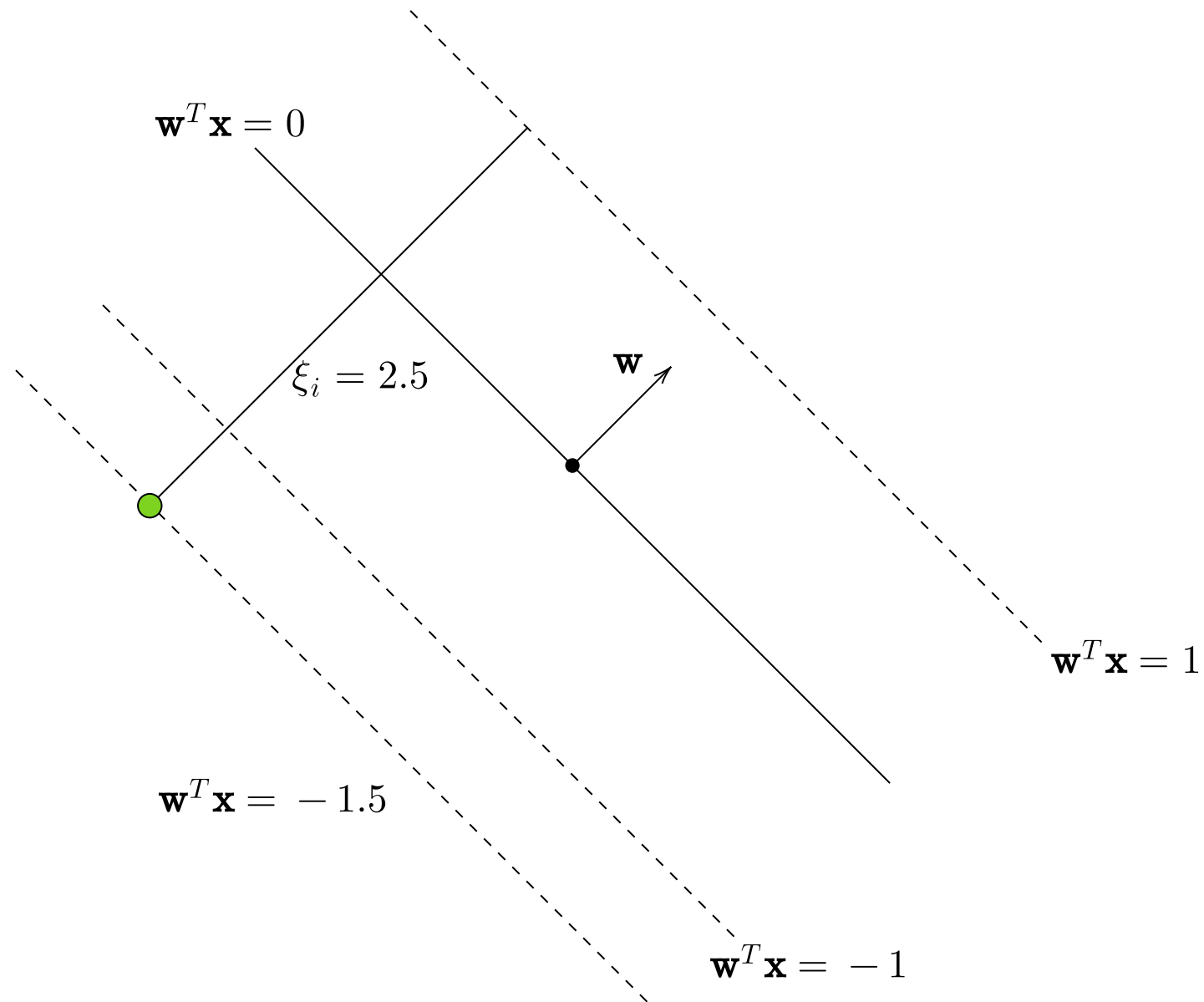


# Soft-Margin



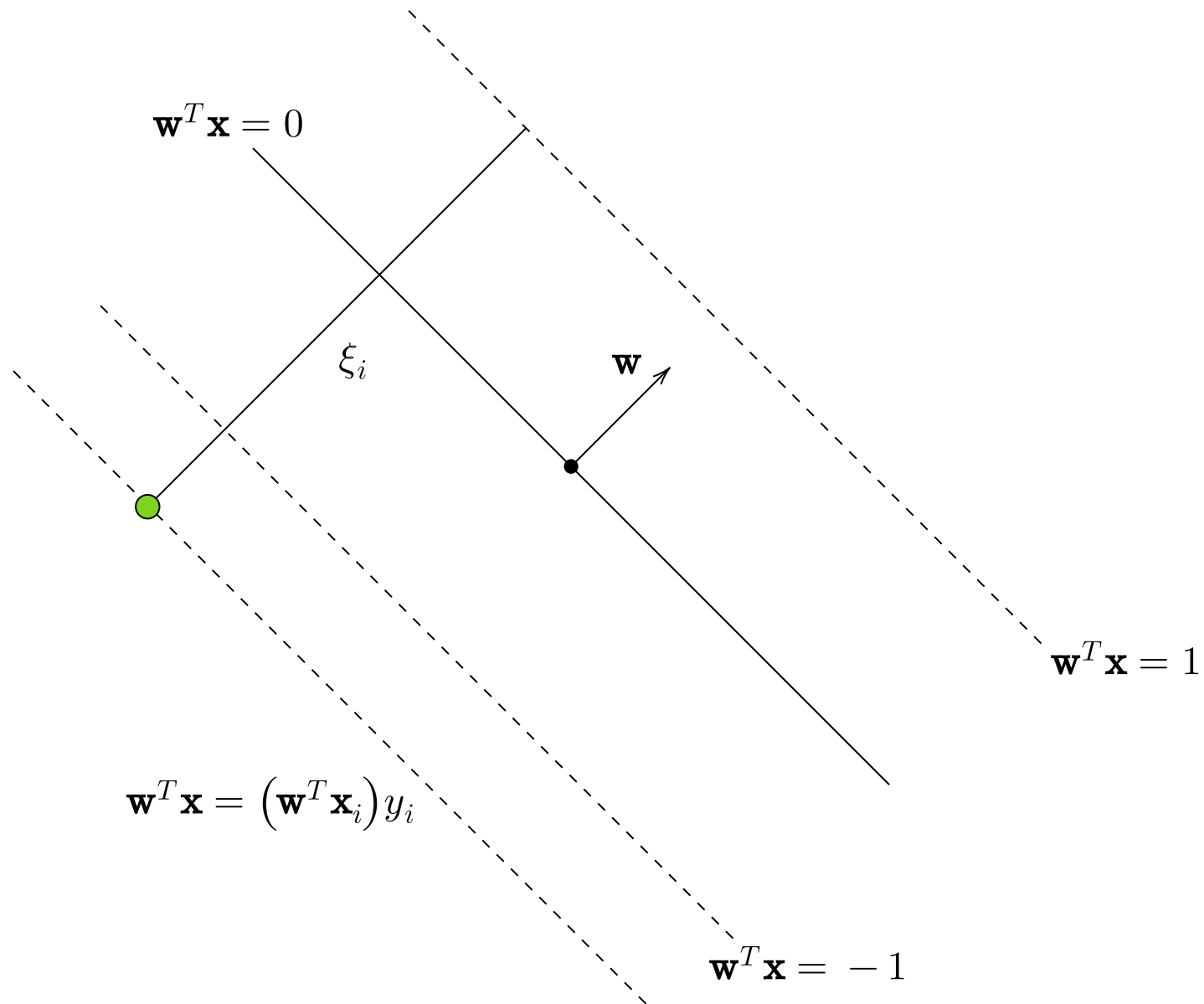
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

# Soft-Margin



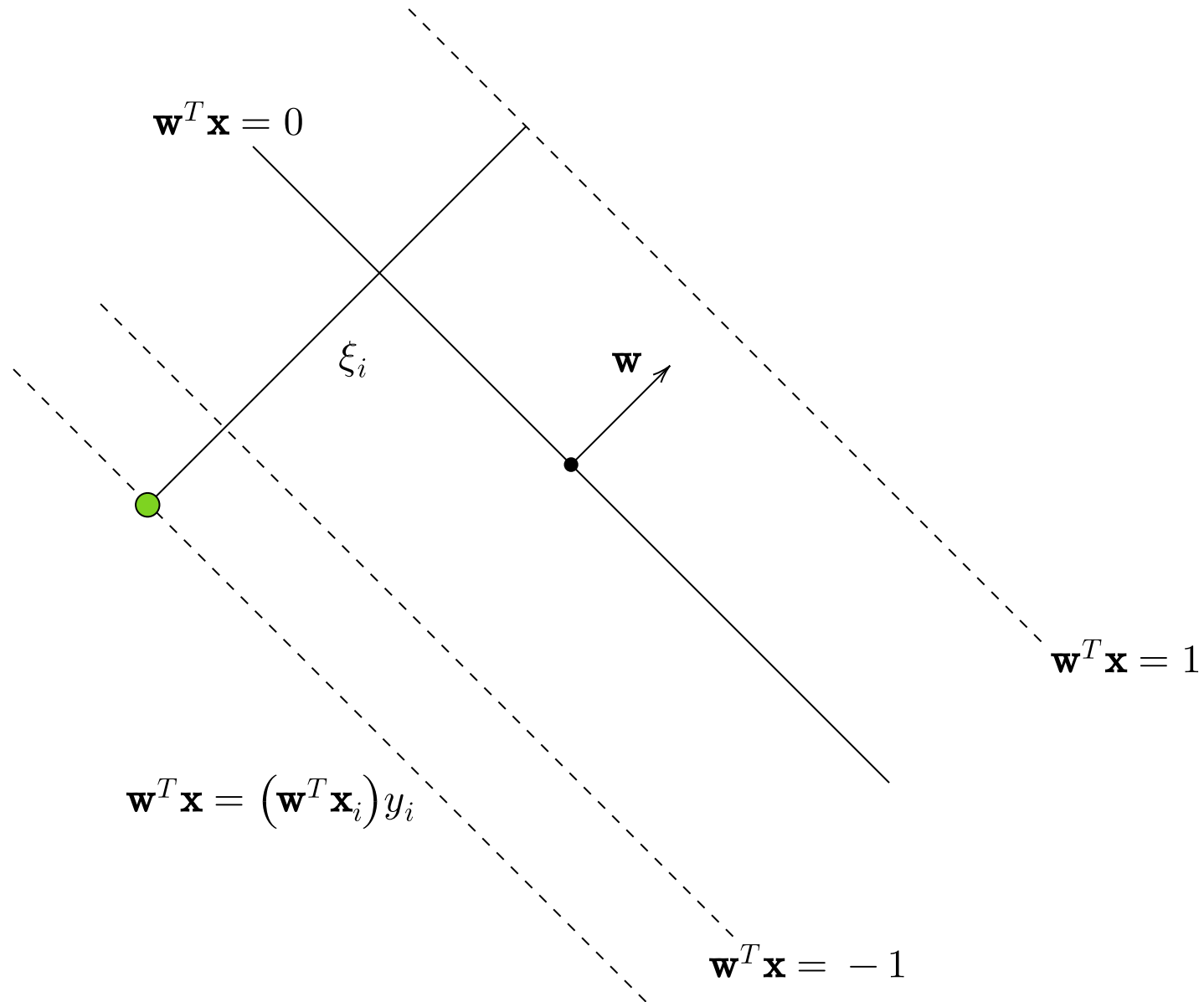
$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

# Soft-Margin



$$(\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i \geq 1$$
$$\xi_i \geq 0$$

# Soft-Margin



$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1 \\ \xi_i &\geq 0 \end{aligned}$$

$$\xi_i = \max(1 - (\mathbf{w}^T \mathbf{x}_i) y_i, 0)$$

# Soft-Margin, Linear-SVM

$$\min_{\mathbf{w}} \quad \frac{||\mathbf{w}||^2}{2} + C \cdot \sum_{i=1}^n \xi_i$$

sub. to

$$(\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i \geq 1, \quad 1 \leq i \leq n$$

$$\xi_i \geq 0, \quad 1 \leq i \leq n$$

# Soft-Margin, Linear-SVM: Hinge-loss formulation

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{\|\mathbf{w}\|^2}{2} + C \cdot \sum_{i=1}^n \xi_i \\ \text{sub. to} \quad & (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i \geq 1, \quad 1 \leq i \leq n \\ & \xi_i \geq 0, \quad 1 \leq i \leq n \end{aligned} \quad \equiv \quad \min_{\mathbf{w}} \quad \frac{\|\mathbf{w}\|^2}{2} + C \cdot \sum_{i=1}^n \max(0, 1 - (\mathbf{w}^T \mathbf{x}_i) y_i)$$

# Soft-Margin, Linear-SVM: Hinge-loss formulation

$$\min_{\mathbf{w}} \quad \frac{\|\mathbf{w}\|^2}{2} + C \cdot \sum_{i=1}^n \xi_i$$

sub. to

$$\begin{aligned} (\mathbf{w}^T \mathbf{x}_i) y_i + \xi_i &\geq 1, & 1 \leq i \leq n \\ \xi_i &\geq 0, & 1 \leq i \leq n \end{aligned}$$

$\equiv$

$$\min_{\mathbf{w}} \quad \frac{\|\mathbf{w}\|^2}{2} + C \cdot \sum_{i=1}^n \max(0, 1 - (\mathbf{w}^T \mathbf{x}_i) y_i)$$

Regularization

Hinge Loss

Model

Data

## Soft-Margin, SVM: Hinge-loss formulation

$$\min_{\mathbf{w}} \quad \frac{\|\mathbf{w}\|^2}{2} + C \cdot \sum_{i=1}^n \max(0, 1 - (\mathbf{w}^T \mathbf{x}_i) y_i) \quad (1)$$

Terms in the objective:

(1)  $\frac{\|\mathbf{w}\|^2}{2}$  controls the width of the margin  
Smaller the value of  $\|\mathbf{w}\|$ , wider the margin

(2)  $\sum_{i=1}^n \max(0, 1 - (\mathbf{w}^T \mathbf{x}_i) y_i)$  is the hinge-loss. Wider the margin, larger the loss.

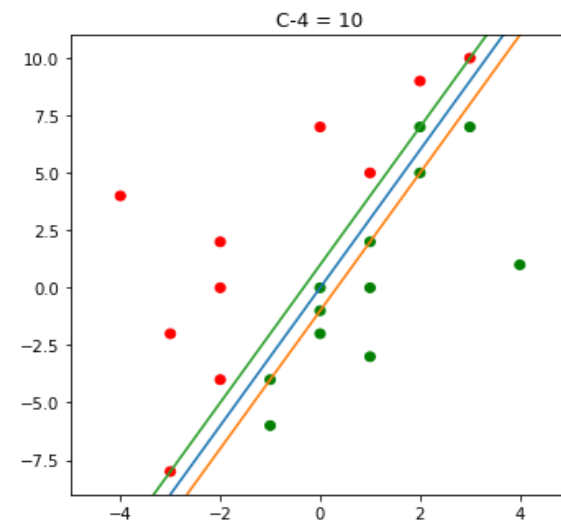
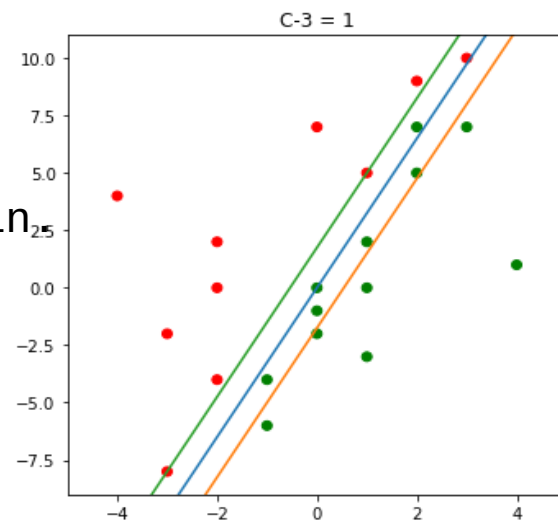
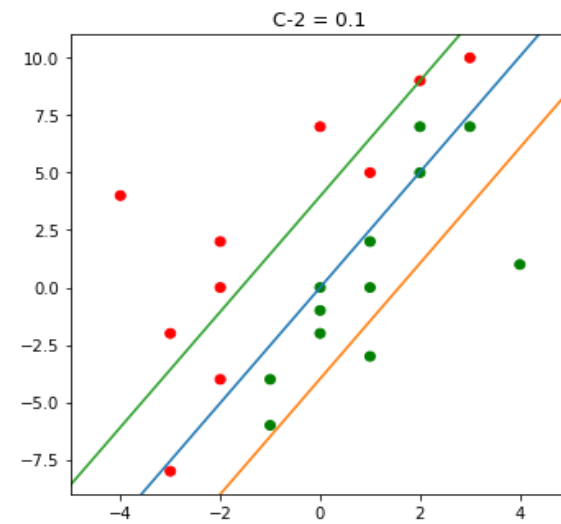
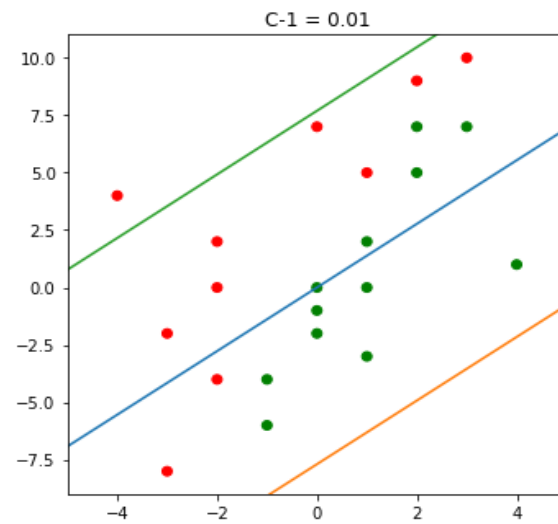


# Soft-Margin, SVM: Hinge-loss formulation

$$\min_{\mathbf{w}} \frac{\|\mathbf{w}\|^2}{2} + C \cdot \sum_{i=1}^n \max(0, 1 - (\mathbf{w}^T \mathbf{x}_i) y_i)$$

(1) (2)

- (1) and (2) work in opposite directions
- If  $\|\mathbf{w}\|$  decreases, the margin becomes wider, which increases the hinge-loss.
- $C$  controls the tradeoff between (1) and (2):
  - If  $C$  is small, we are fine with a wide margin.
  - If  $C$  is large, we prefer a narrow margin.
  - If  $C \rightarrow \infty$ , we do not tolerate bribery at all.



# Miscellaneous

## Terminology used

- (1) Hard-margin, Linear-SVM
- (2) Hard-margin, Kernel-SVM
- (3) Soft-margin, Linear-SVM
- (4) Soft-margin, Kernel-SVM

## Additional points

- Discriminative model
- Weight vector is a sparse linear combination of data-points
- The dual is a quadratic programming problem