



THE UNIVERSITY OF  
**MELBOURNE**

# Support Vector Machines

**COMP90051 Statistical Machine Learning**

Semester 2, 2020

**QiuHong Ke**

Copyright: University of Melbourne

# Before we start...

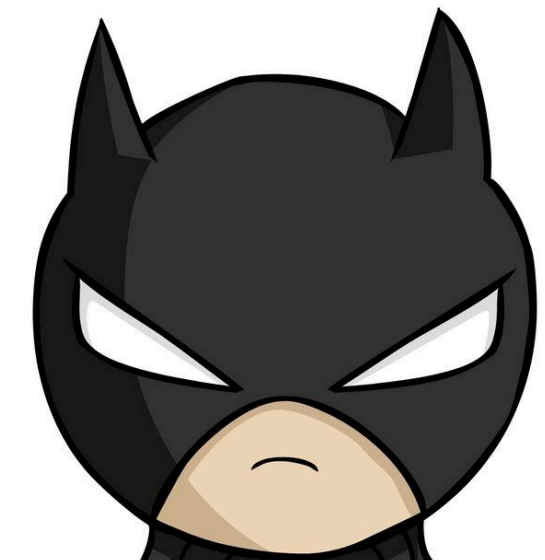
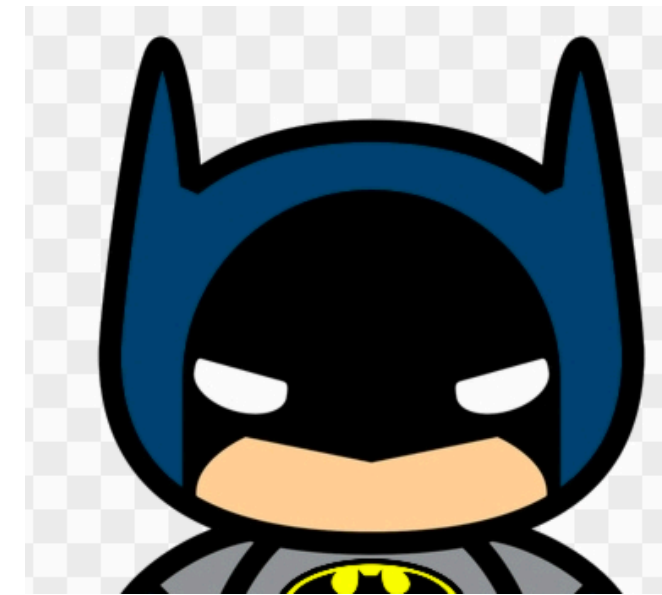
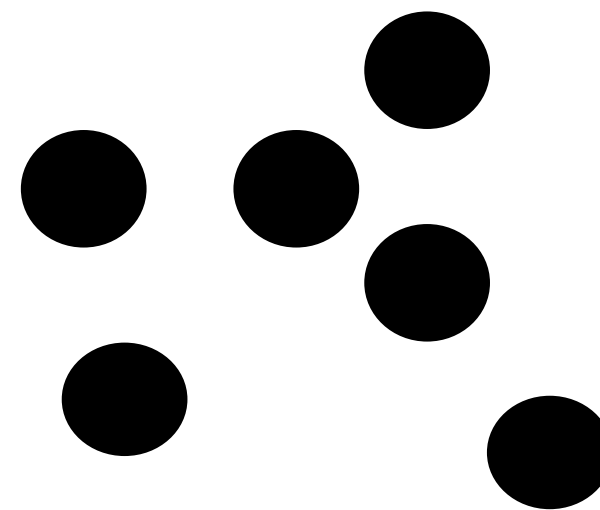
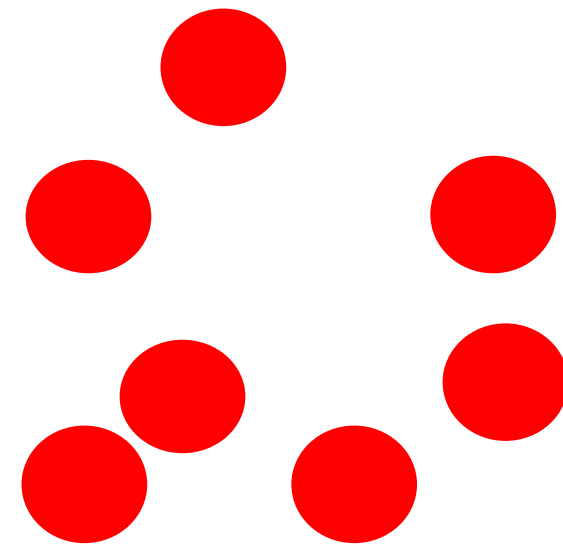
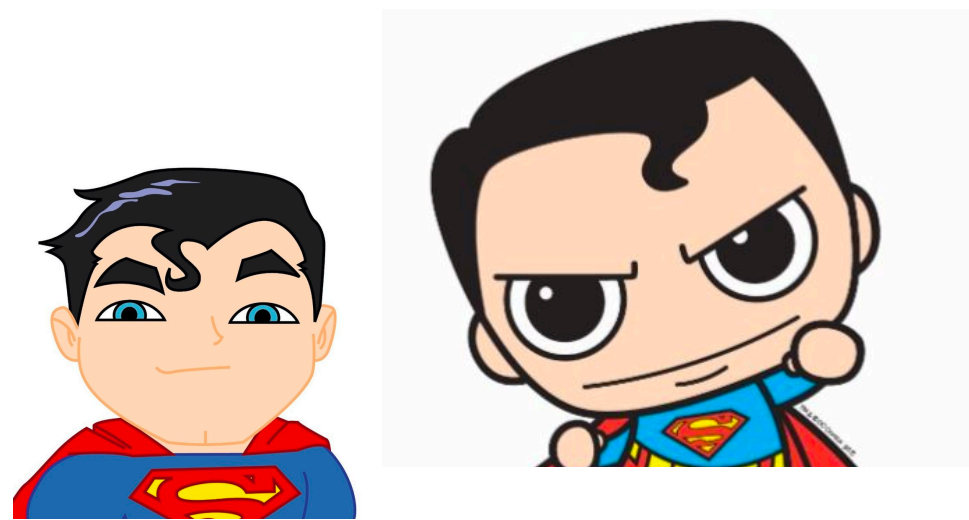
## About me

- 2015.02-2018.04: PhD in UWA
- 2018.05-2019.12: Post-doc in MPII
- From 2020.01: Lecturer in UniMelb
- Research: Action recognition and prediction using machine learning
- Contact:
  - [qiu hong.ke@unimelb.edu.com](mailto:qiu hong.ke@unimelb.edu.com);
  - [comp90051-2020s2-staff@lists.unimelb.edu.au](mailto:comp90051-2020s2-staff@lists.unimelb.edu.au)

# Superman vs Batman



# Superman vs Batman



<https://kidzartworx.com.au/product/2-cartoons-for-kids-5-8yrs-monday-1-30pm-4-30pm-6th-july-2020/>

<https://pngio.com/PNG/a51210-baby-superman-cartoon.html>

<https://www.pinterest.com.au/pin/64809682117268678/>

<https://www.dailymotion.com/video/x6dzn8t>

[https://www.pngitem.com/middle/Jhxmhh\\_batman-cute-png-transparent-png/](https://www.pngitem.com/middle/Jhxmhh_batman-cute-png-transparent-png/)

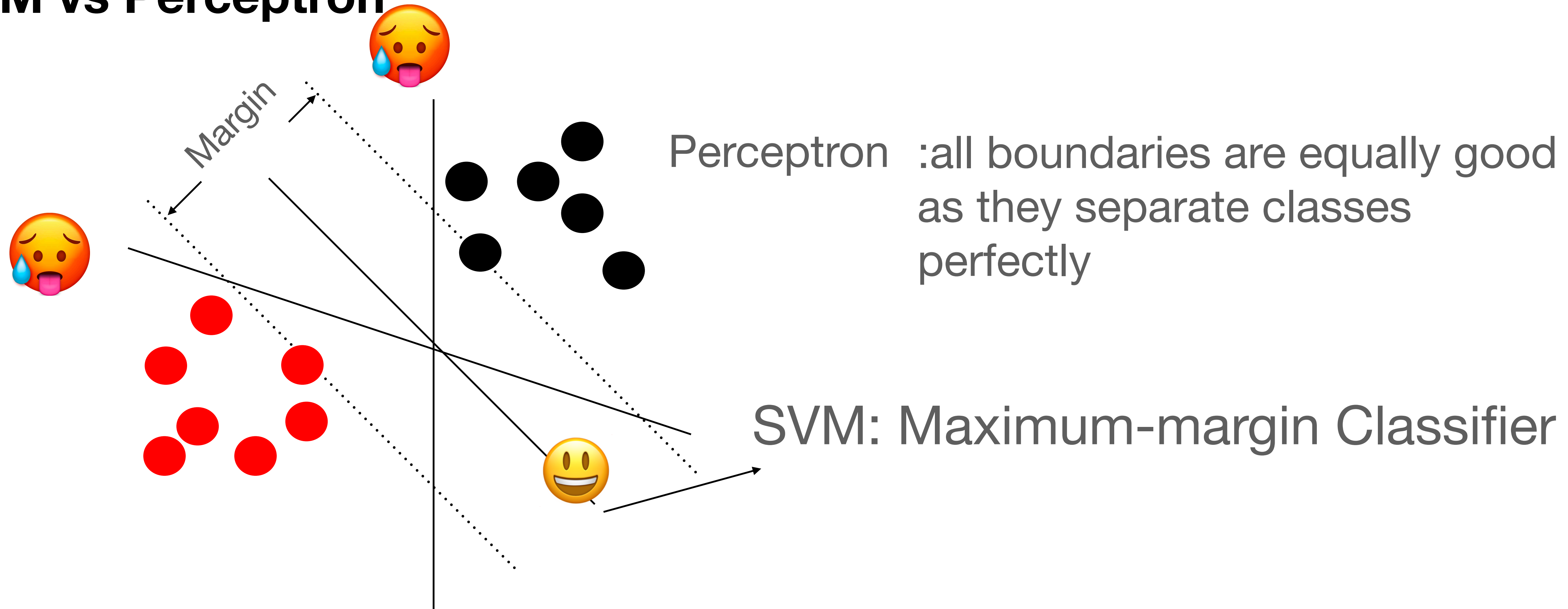
[https://www.clipartmax.com/middle/m2i8H7d3d3b1d3Z5\\_baby-batman-image-baby-batman/](https://www.clipartmax.com/middle/m2i8H7d3d3b1d3Z5_baby-batman-image-baby-batman/)

<https://www.pinterest.com.au/pin/232850243216374160/>

<https://www.pinterest.com.au/pin/849702654667128852/>

# Binary Linear Classifier

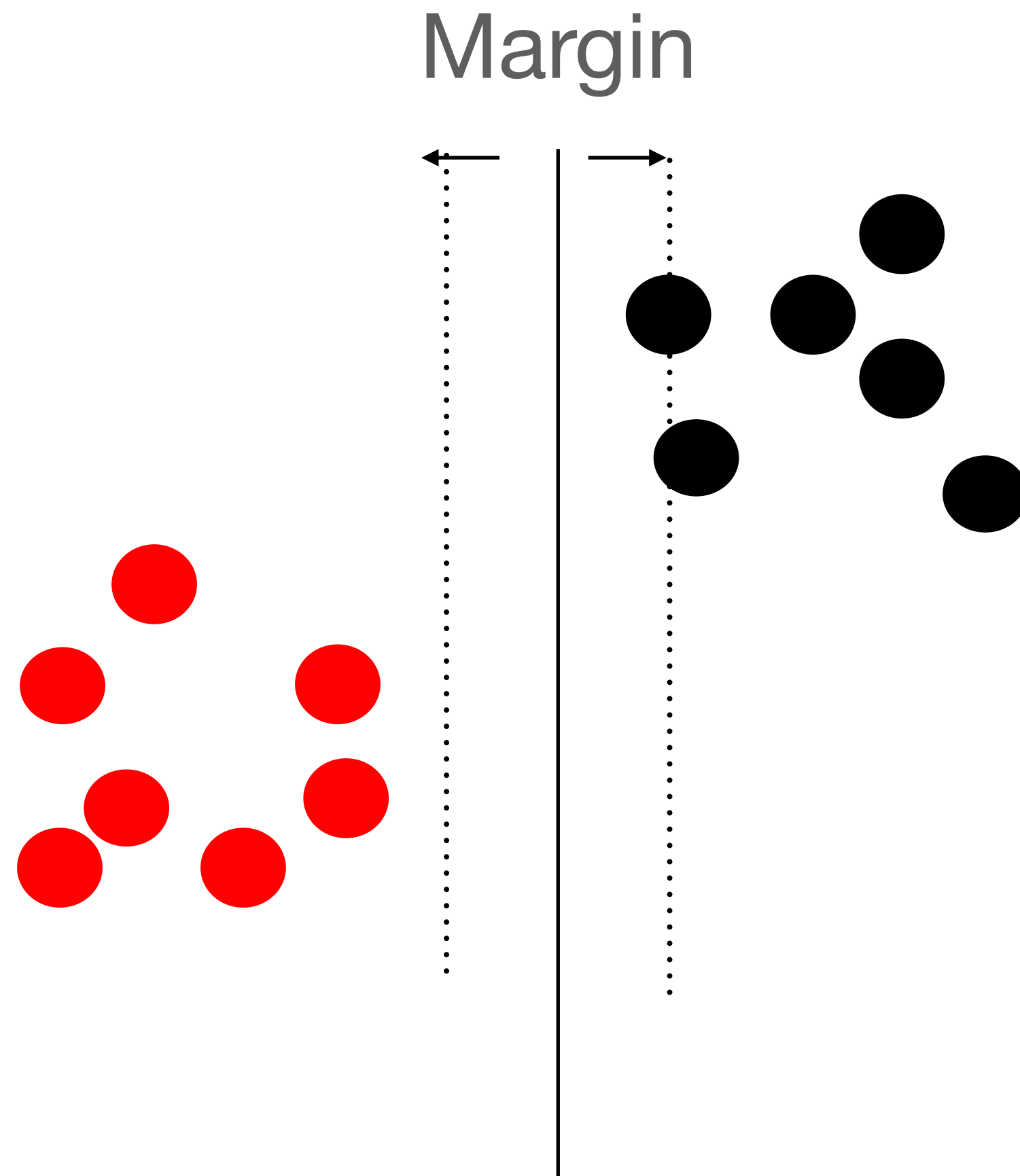
## SVM vs Perceptron



Margin: 2x minimum distance (boundary, data points)

# Binary Linear Classifier

## SVM vs Perceptron



Margin: 2x minimum distance (boundary, data points)

# Outline

- Margin
- Lagrange Duality
- Soft-margin SVM
- Kernels



**Linear classifier**

$$f(x) = w^T x + b$$

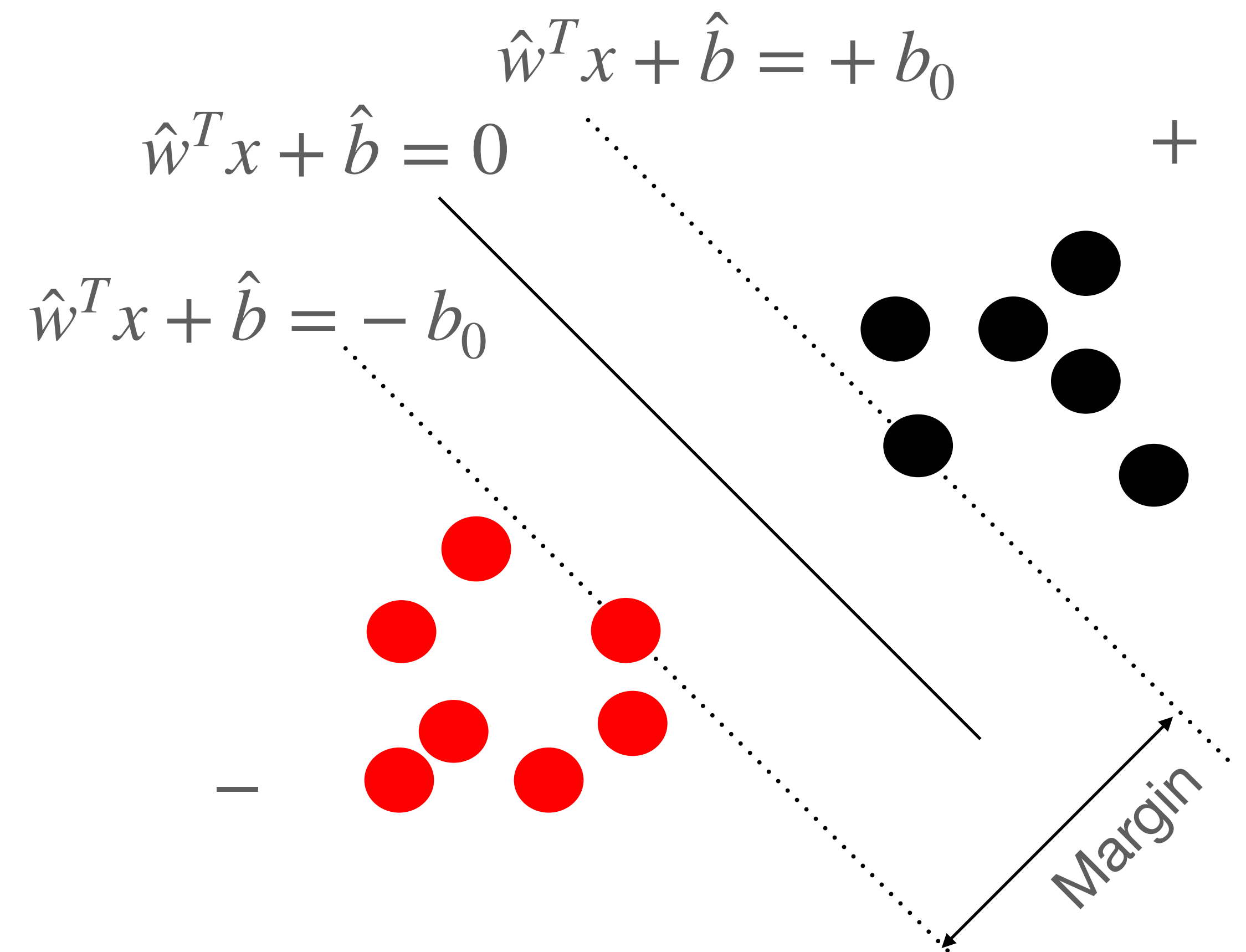
$x$  : Feature vector (column)

$w$  : Weight vector (column)

$T$  : Transpose

$b$  : Bias

$$w^T x = \|w\| \|x\| \cos \theta$$





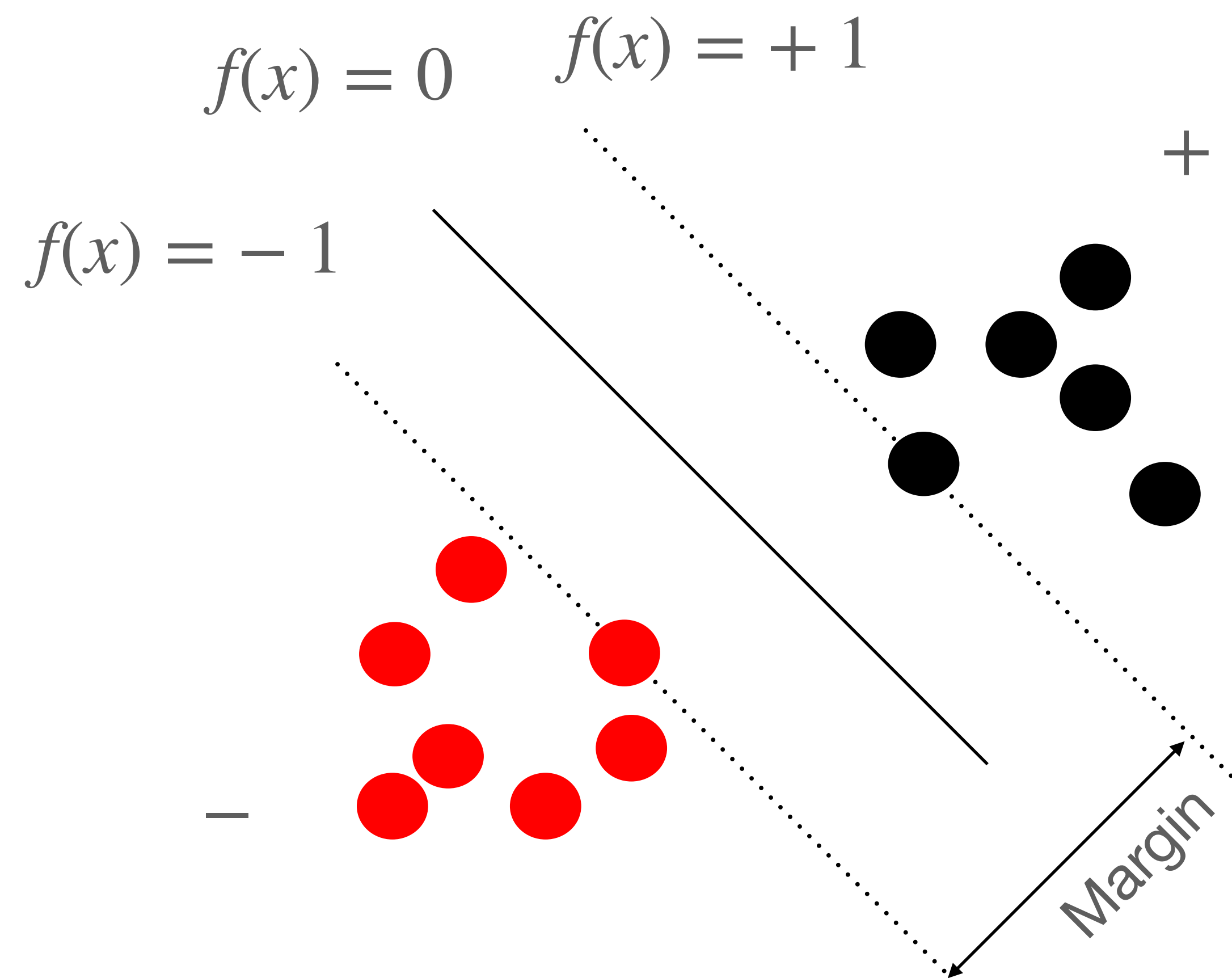
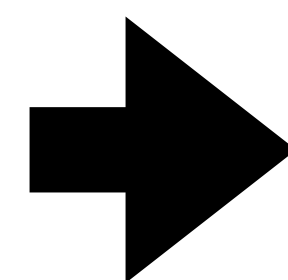
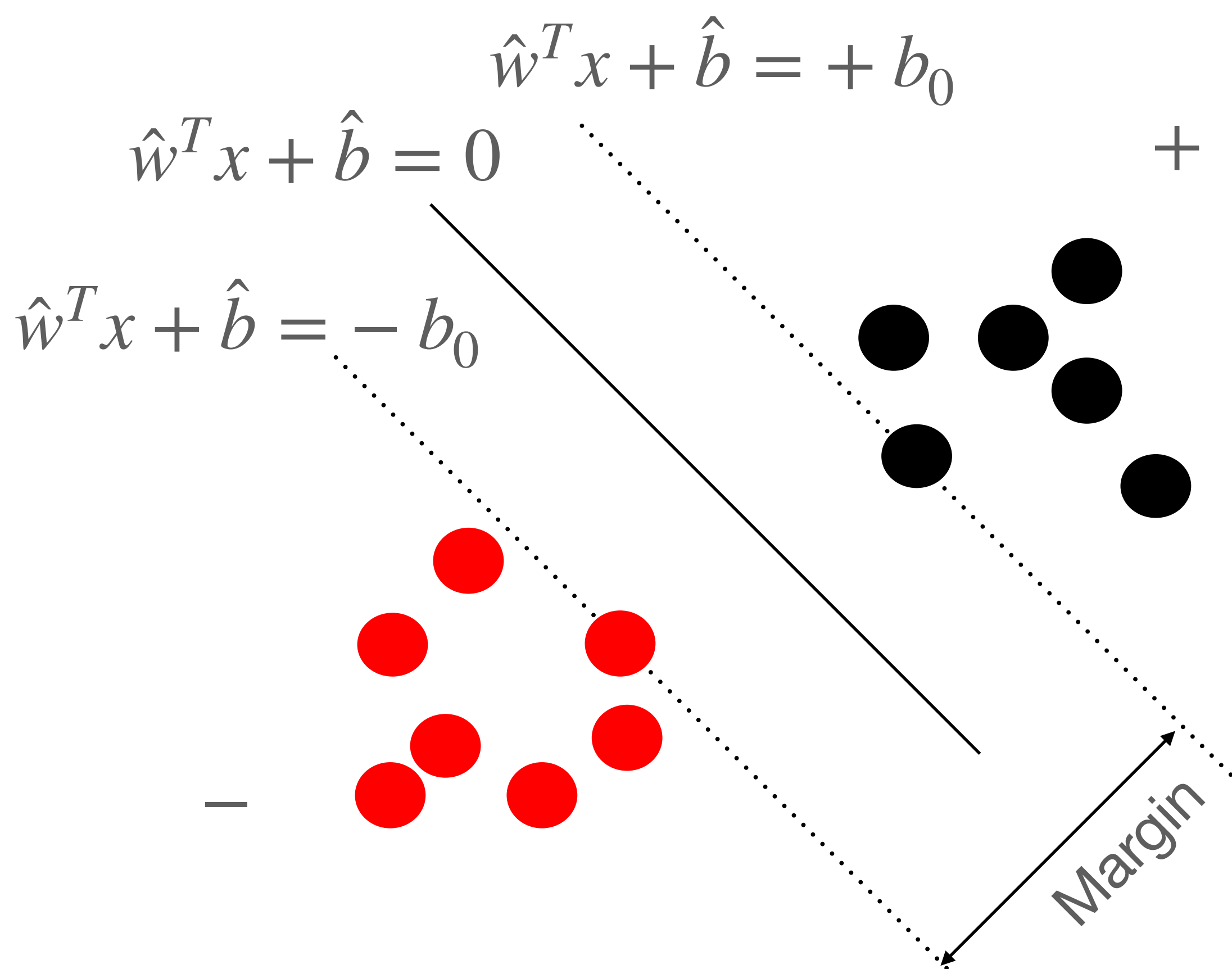
Margin

Lagrange  
Duality

Soft-margin  
SVM

Kernels

$$f(x) = w^T x + b \quad w = \frac{\hat{w}}{b_0} \quad b = \frac{\hat{b}}{b_0}$$



**Margin formula**

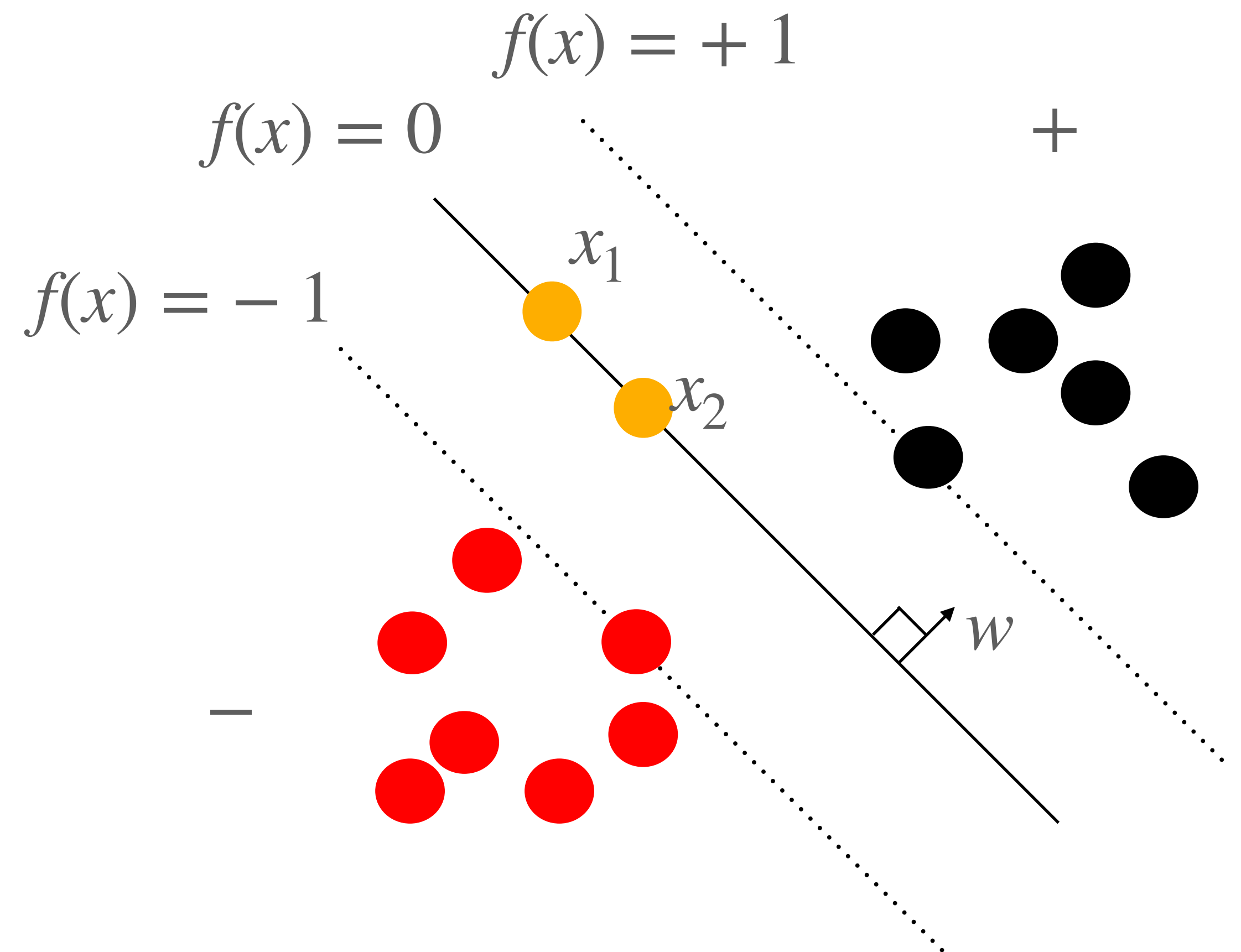
$$f(x) = w^T x + b$$

$$w^T x_1 + b = 0$$

$$w^T x_2 + b = 0$$

$$w^T (x_1 - x_2) = 0$$

$$\|w\| \|x_1 - x_2\| \cos\theta = 0$$



**Margin formula**

$$w^T x^{-1} + b = -1$$

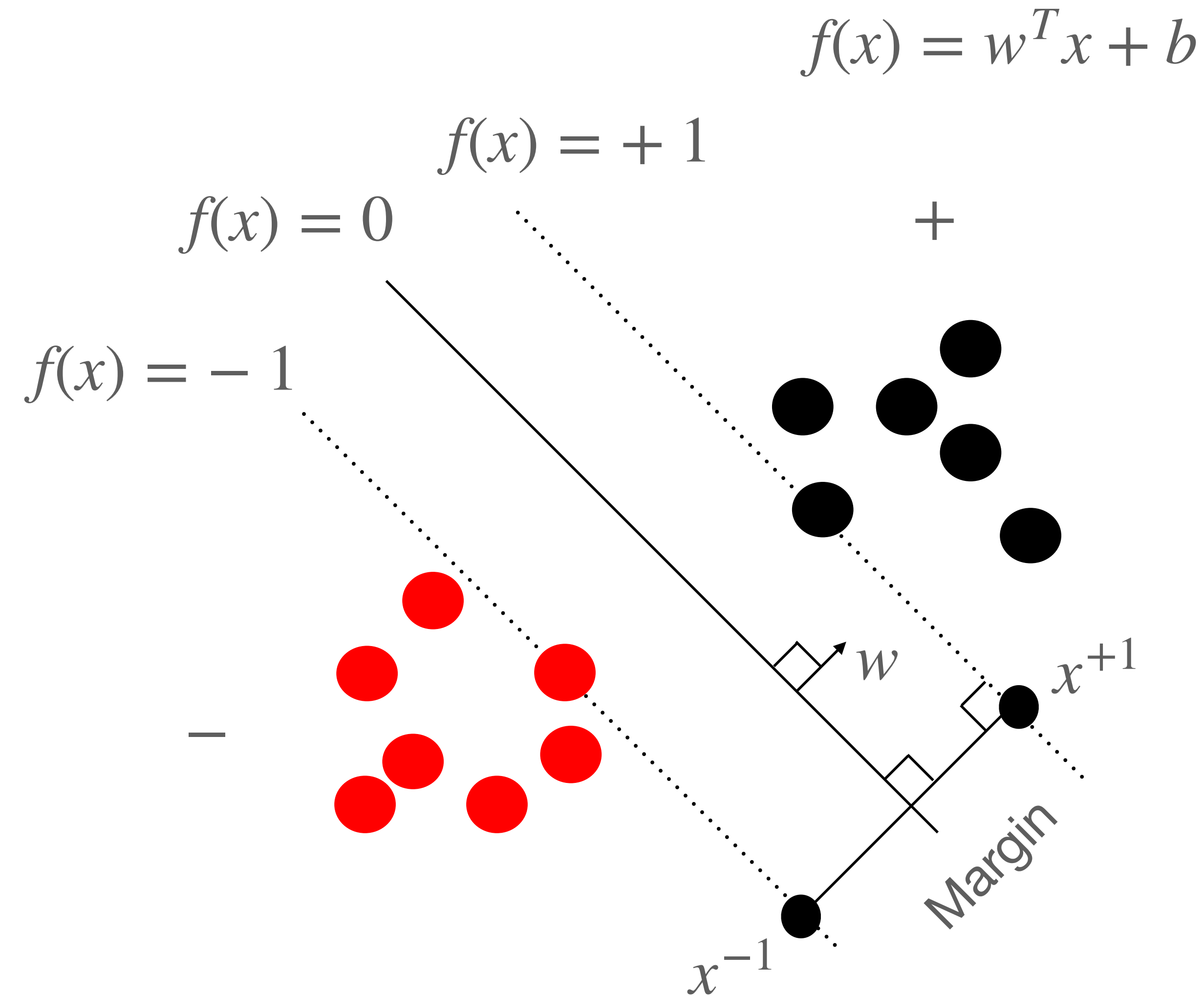
$$w^T x^{+1} + b = 1$$

$$w^T (x^{+1} - x^{-1}) = 2$$

$$\|w\| \cdot \text{Margin} \cdot \cos\theta = 2$$

$$\text{Margin} = \frac{2}{\|w\|}$$

we want to maximise the margin

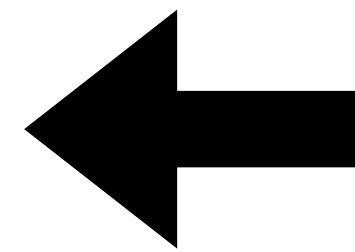


**SVM: Constrained optimisation problem**

$$\min_w \frac{\|w\|^2}{2}$$

s.t

$$1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, \quad i = 1, \dots, n \text{ (data points)}$$



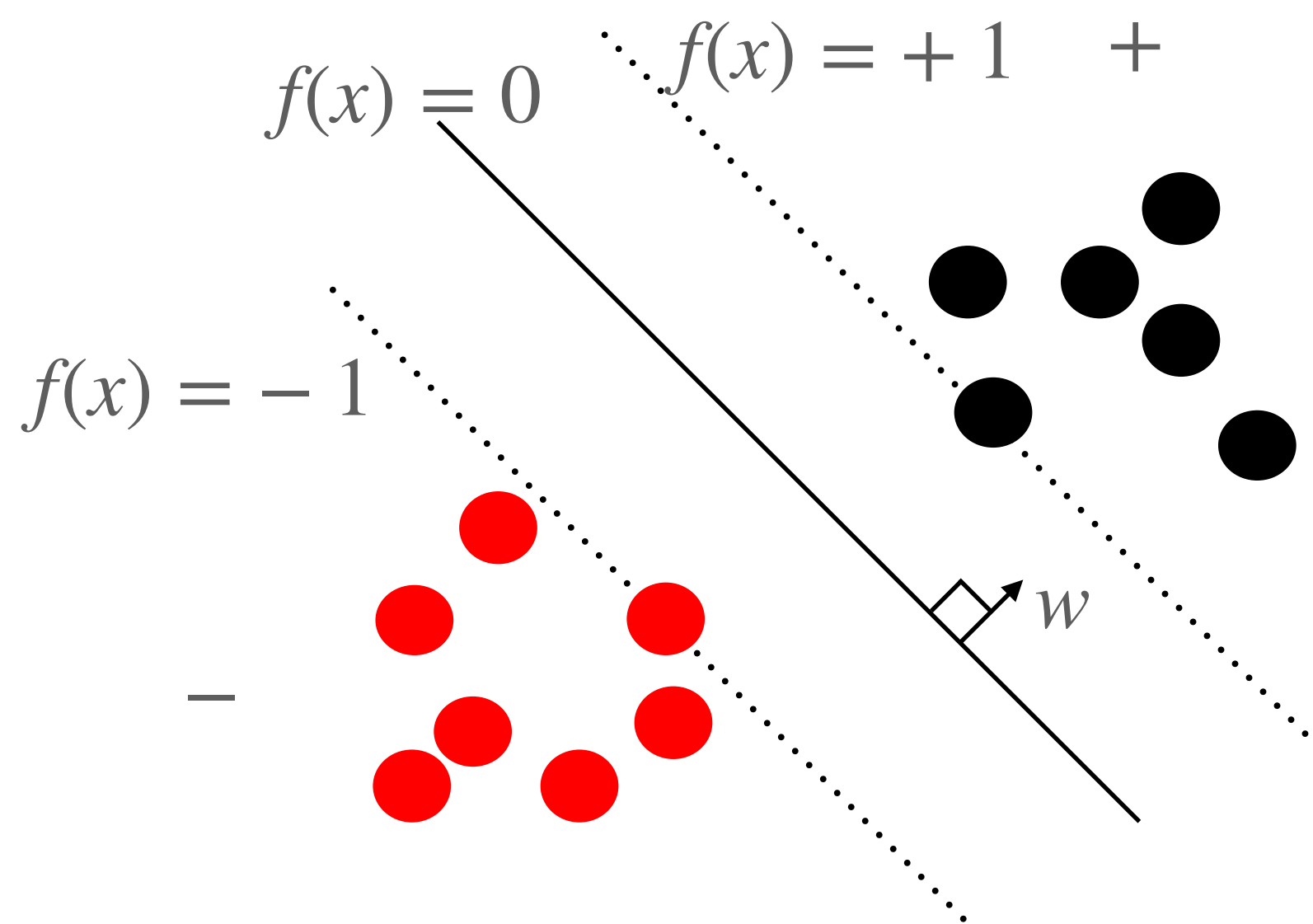
$$\max_w \frac{2}{\|w\|}$$

subject to

$$\text{if } y^{(i)} = +1 : f(x^{(i)}) = w^T x^{(i)} + b \geq +1$$

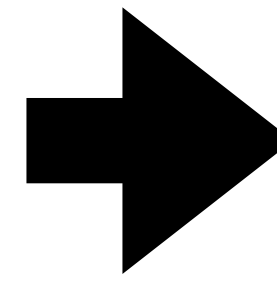
$$\text{if } y^{(i)} = -1 : f(x^{(i)}) = w^T x^{(i)} + b \leq -1$$

( $i = 1, \dots, n$  data points)

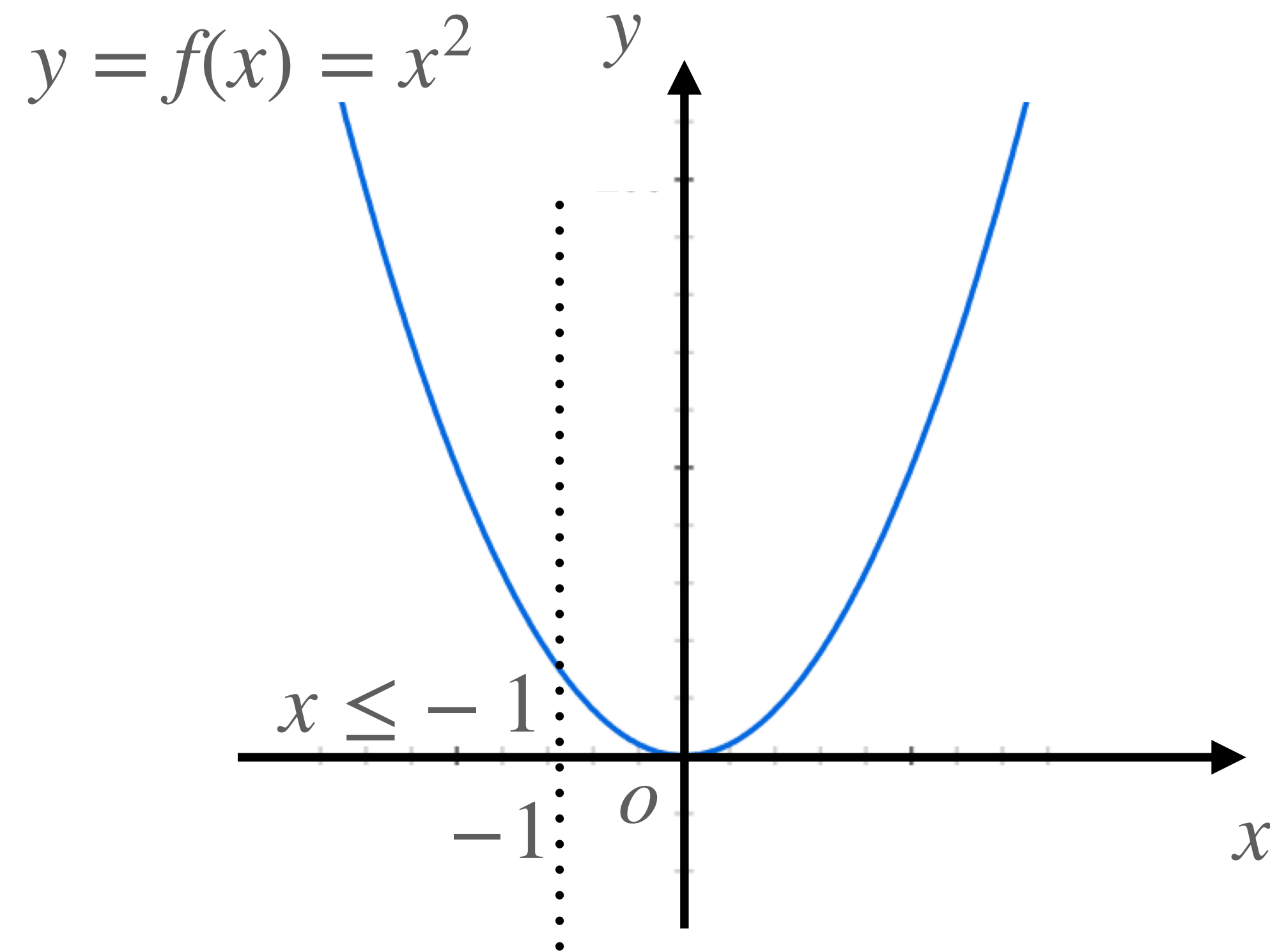


**Primal problem**

$$\begin{aligned} \min_w \quad & \frac{\|w\|^2}{2} \\ \text{s.t} \quad & 1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, \quad i = 1, \dots, n \text{ (data points)} \end{aligned}$$

**Dual problem**

**What's the dual problem?**  
**Why solving primal by  
solving dual problem?**

**Simple example**

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \leq -1 \end{aligned}$$

**Primal problem**

$$\min_x f(x)$$

$$\text{s.t. } g(x) = x + 1 \leq 0$$

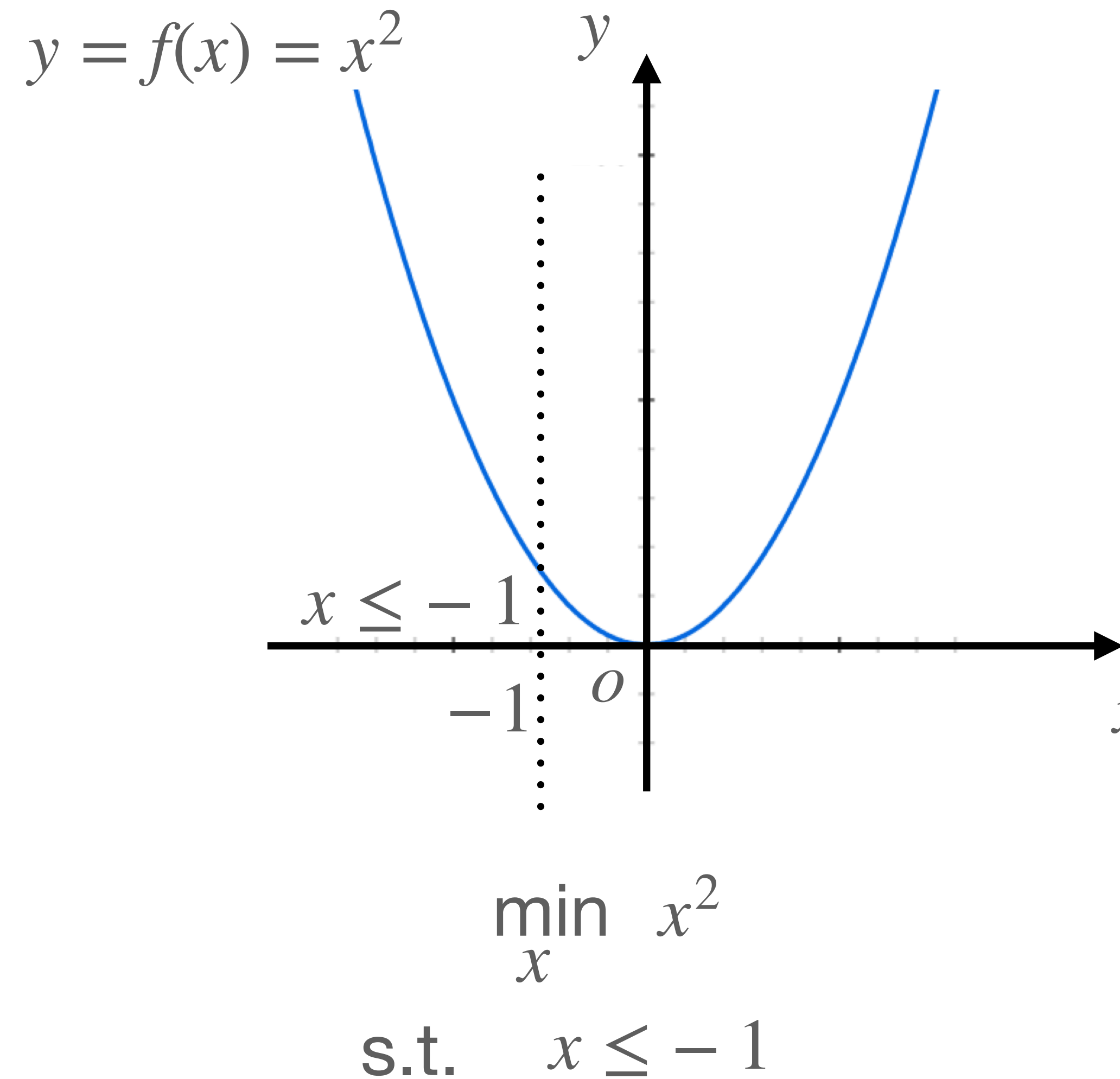
- Construct a function:

$$L(x, \lambda) = f(x) + \lambda g(x)$$

- Set  $\lambda \geq 0$ , calculate  $\max_{\lambda} L(x, \lambda)$

$$g(x) > 0 : \max_{\lambda} L(x, \lambda) = \infty \text{ when } \lambda = \infty$$

$$g(x) \leq 0 : \max_{\lambda} L(x, \lambda) = f(x) \text{ when } \lambda = 0$$





## Primal problem

$$\min_x f(x)$$

$$\text{s.t. } g(x) = x + 1 \leq 0$$

- Construct a function

$$L(x, \lambda) = f(x) + \lambda g(x) : \quad \text{Lagrangian function}$$

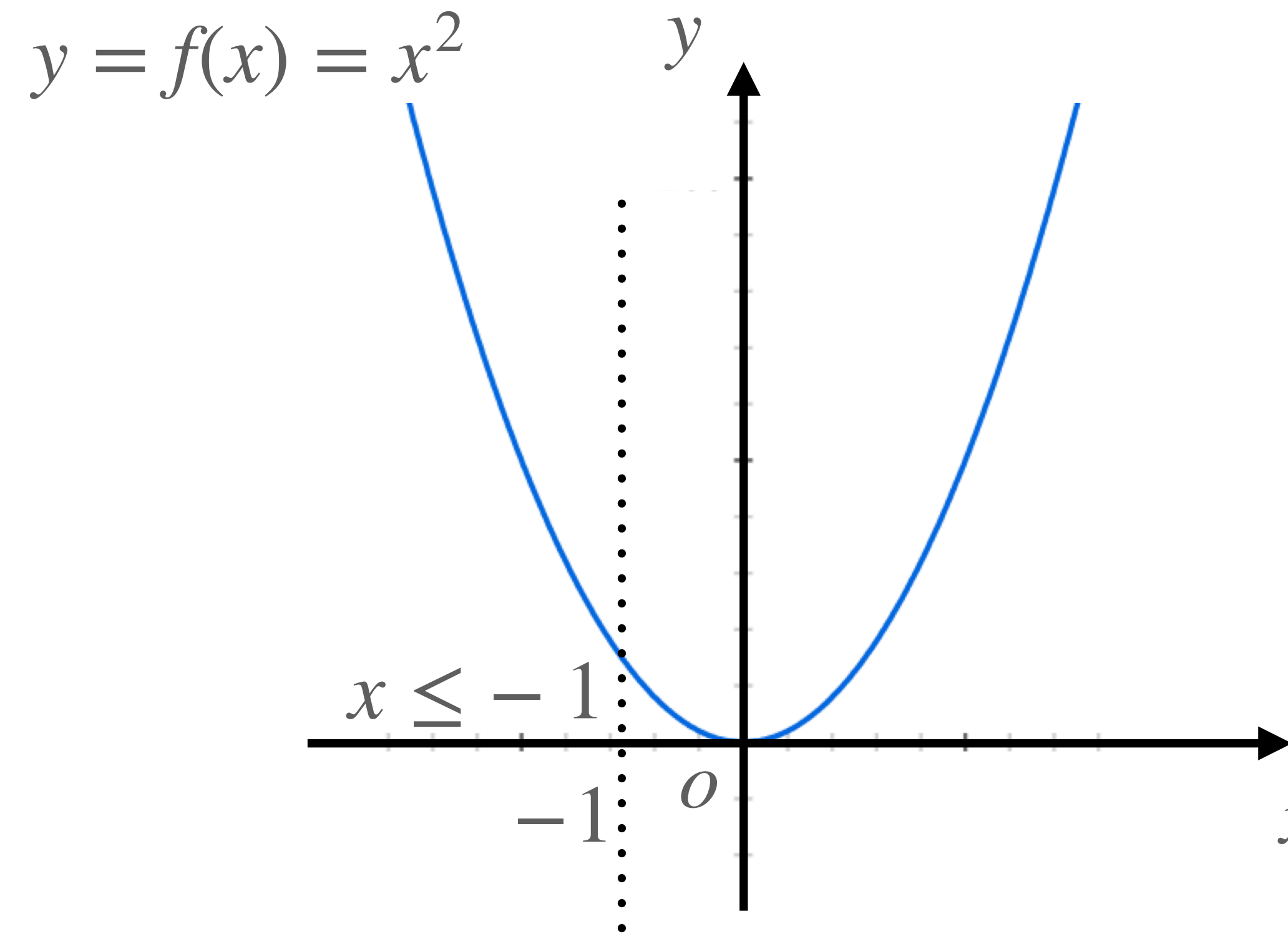
$$\lambda \geq 0 : \quad \text{Lagrange multiplier}$$

- Primal function

$$\theta_p(x) = \max_{\lambda} L(x, \lambda) = f(x) \text{ if } g(x) \leq 0$$

minimise the primal problem is equivalent to minimise the maximum of Lagrangian function

$$\text{So: } \min_x f(x) = \min_x \theta_p(x) = \min_x \max_{\lambda} L(x, \lambda)$$



$$\min_x x^2$$

$$\text{s.t. } x \leq -1$$

**From primal to dual problem**

$$L(x, \lambda) = f(x) + \lambda g(x)$$

$$\lambda \geq 0 \quad g(x) \leq 0$$

- Primal problem:

$$\min_x f(x) = \min_x \max_{\lambda} L(x, \lambda)$$

- Dual problem:

$$\max_{\lambda} \min_x L(x, \lambda) = \max_{\lambda} \theta_d(\lambda)$$

Dual function:  $\theta_d(\lambda) = \min_x L(x, \lambda)$

$$\theta_d(\lambda) = \min_x L(x, \lambda) \leq L(x, \lambda) = \underline{f(x) + \lambda g(x)} \leq f(x)$$

This only hold if  $g(x) < 0$

## From primal to dual problem

- Primal problem:

$$\min_x f(x) = \min_x \max_{\lambda} L(x, \lambda)$$

- Dual problem:

$$\max_{\lambda} \min_x L(x, \lambda) = \max_{\lambda} \theta_d(\lambda)$$

Solutions:

$x^*$  makes  $f(x)$  minimum :  $f(x^*) = p^*$

$\lambda^*$  makes  $\theta_d(\lambda)$  maximum :  $\theta_d(\lambda^*) = d^*$

$$\theta_d(\lambda) = \min_x L(x, \lambda) \leq L(x, \lambda) = f(x) + \lambda g(x) \leq f(x)$$

## From primal to dual problem

- Primal problem:

$$\min_x f(x) = \min_x \max_{\lambda} L(x, \lambda)$$

- Dual problem:

$$\max_{\lambda} \min_x L(x, \lambda) = \max_{\lambda} \theta_d(\lambda)$$

$$\theta_d(\lambda) = \min_x L(x, \lambda) \leq L(x, \lambda) = f(x) + \lambda g(x) \leq f(x)$$

$$d^* = \theta_d(\lambda^*) = \min_x L(x, \lambda^*) \leq L(x^*, \lambda^*) = f(x^*) + \lambda^* g(x^*) \leq f(x^*) = p^*$$

Under some conditions:  $d^* = p^*$

?

Solutions:

$$f(x^*) = p^* = \min_x f(x)$$

$$\theta_d(\lambda^*) = d^* = \max_{\lambda} \theta_d(\lambda)$$

## From primal to dual problem



$$d^* = \theta_d(\lambda^*) = \min_x L(x, \lambda^*) \leq L(x^*, \lambda^*) = f(x^*) + \lambda^* g(x^*) \leq f(x^*) = p^*$$

$$\text{if } \min_x L(x, \lambda^*) = L(x^*, \lambda^*) \text{ and } f(x^*) + \lambda^* g(x^*) = f(x^*)$$

$$d^* = p^*$$

KKT (Karush-Kuhn-Tucker) conditions :

$$g(x) \leq 0 \quad (\text{Primal feasibility})$$

$$\lambda \geq 0 \quad (\text{Dual feasibility})$$

$$\lambda g(x) = 0 \quad (\text{Complementary slackness})$$

$$\frac{\partial L}{\partial x} = 0 \quad (\text{Stationarity})$$

**Dual problem of SVM**

**Primal problem**

$$\min_w \frac{\|w\|^2}{2}$$

s.t.  $g_i(w, b) = 1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, i = 1, \dots, n$  data points

subject to

$$\begin{array}{l} \text{if } y^{(i)} = +1 : f(x^{(i)}) = w^T x^{(i)} + b \geq +1 \\ \text{if } y^{(i)} = -1 : f(x^{(i)}) = w^T x^{(i)} + b \leq -1 \end{array}$$

(1) Lagrangian function:

We need to find the maximum in terms of lambda, and minimum in terms of w and b

$$L(w, b, \lambda) = \frac{\|w\|^2}{2} + \sum_{i=1}^n \lambda_i (1 - y^{(i)}(w^T x^{(i)} + b))$$

(2) dual function  $\theta_d(\lambda) = \min_{w, b} L(w, b, \lambda) :$

$$\frac{\partial L}{\partial w_j} = 0 : w_j = \sum_{i=1}^n \lambda_i y^{(i)} x_j^{(i)}$$

$$\frac{\partial L}{\partial b} = 0 : \sum_{i=1}^n \lambda_i y^{(i)} = 0$$

## Dual problem of SVM

### Dual function

$$\theta_d(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \lambda_i \lambda_k y^{(i)} y^{(k)} (x^{(i)})^T x^{(k)}$$

### Dual problem

$$\begin{aligned} & \max_{\lambda} \theta_d(\lambda) \\ \text{s.t.} \quad & \lambda_i \geq 0 \quad \text{and} \quad \sum_{i=1}^n \lambda_i g_i(x) = 0 \end{aligned}$$

We need this condition since we need to satisfy the complementary slackness in KKT condition



## Support vectors

$g_i(w, b)$  will be negative for the point beyond margin, so lambda will be 0 for those points. The support vectors are the only points that lambda may greater than 0 since  $g$  is 0 in that case to satisfy the complementary slackness

$$\min \frac{\|w\|^2}{2} \quad \text{s.t.} \quad g_i(w, b) = 1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, \quad i = 1, \dots, n \text{ data points}$$

Lagrangian function: 
$$L(w, b, \lambda) = \frac{\|w\|^2}{2} + \sum_{i=1}^n \lambda_i (1 - y^{(i)}(w^T x^{(i)} + b))$$

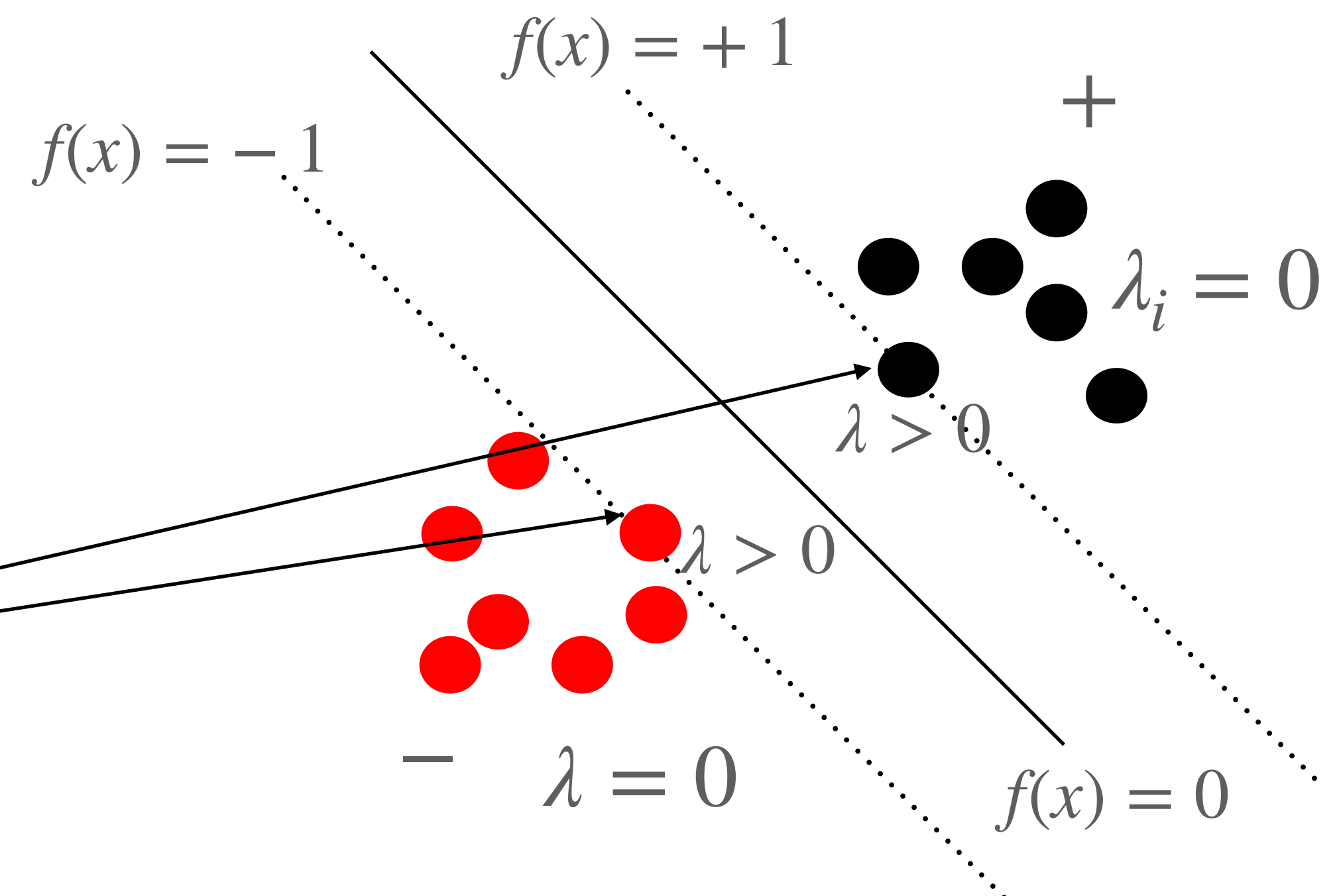
$$\lambda_i \geq 0$$

(Dual feasibility)

$$\lambda_i g_i(w, b) = 0$$

(Complementary slackness)

support vectors:



## Primal vs Dual (Training)

- Primal problem: solve  $d+1$  variables ( $w_j$  and  $b$ ) ( $d$ : dimension of weight vector  $w$ )

$$\text{s.t.} \quad \min_w \frac{\|w\|^2}{2}$$

$$g_i(w, b) = 1 - y^{(i)}(w^T x^{(i)} + b) \leq 0, \quad i = 1, \dots, n \text{ data points}$$

- Dual problem: solve  $n$  variables ( $\lambda_i$ )

$$\begin{aligned} & \max_{\lambda} \theta_d(\lambda) \quad \theta_d(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \lambda_i \lambda_k y^{(i)} y^{(k)} (x^{(i)})^T x^{(k)} \\ & \text{s.t.} \quad \lambda_i \geq 0 \text{ and } \sum_{i=1}^n \lambda_i g_i(x) = 0 \end{aligned}$$

If data size  $n$  is large, ( $n \gg d$ ) solving dual problem is slower than solving primal problem, and vice versa.

## Primal vs Dual (Prediction)

- Primal form:

$$f(x) = w^T x + b \quad \begin{array}{l} f(x) > 0 : \text{positive class} \\ f(x) < 0 : \text{negative class} \end{array}$$

- Dual form:

$$w_j = \sum_{i=1}^n \lambda_i y^{(i)} x_j^{(i)}$$
$$f(x) = \sum_{i=1}^n \lambda_i y^{(i)} (x^{(i)})^T x + b$$

(b can be solved using support vectors:  $f(x) = \pm 1$ )

## Why bother solving dual problem to solve primal problem

Training, solve:

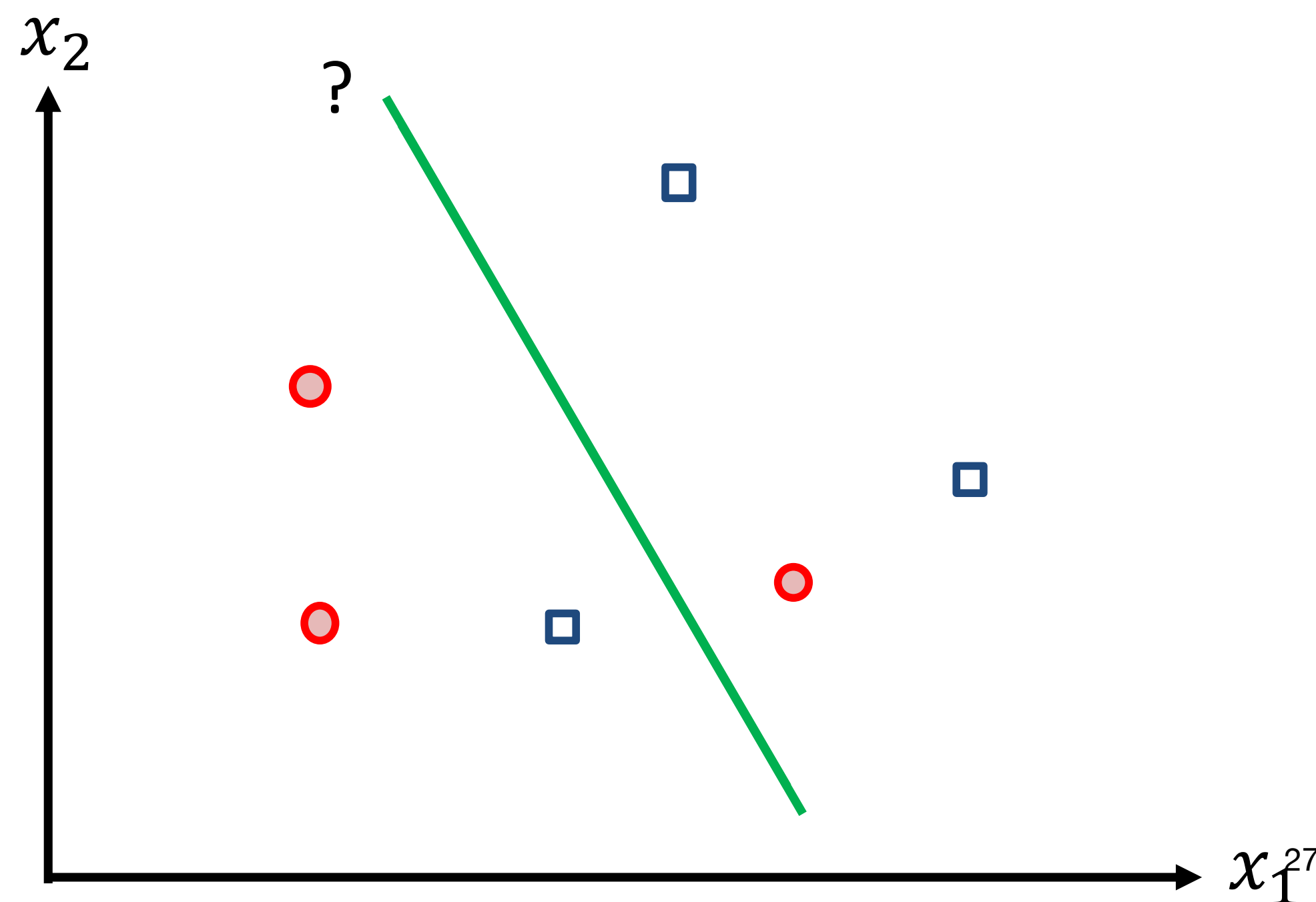
$$\begin{aligned} \max_{\lambda} \quad & \theta_d(\lambda) \\ \text{s.t.} \quad & \lambda_i \geq 0 \text{ and } \sum_{i=1}^n \lambda_i g_i(x) = 0 \end{aligned} \quad \theta_d(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \lambda_i \lambda_k y^{(i)} y^{(k)} (x^{(i)})^T x^{(k)}$$

$$\text{Prediction: } f(x) = \sum_{i=1}^n \lambda_i y^{(i)} (x^{(i)})^T x + b$$

- **Use only support vectors for prediction:** Efficient in prediction
- Inner product: Kernel trick can be used to efficiently handle non-linearly separable data

## Data not linearly separable

- Hard-margin loss is too stringent (*hard!*)
- Real data is unlikely to be linearly separable
- If the data is not separable, hard-margin SVMs are in trouble

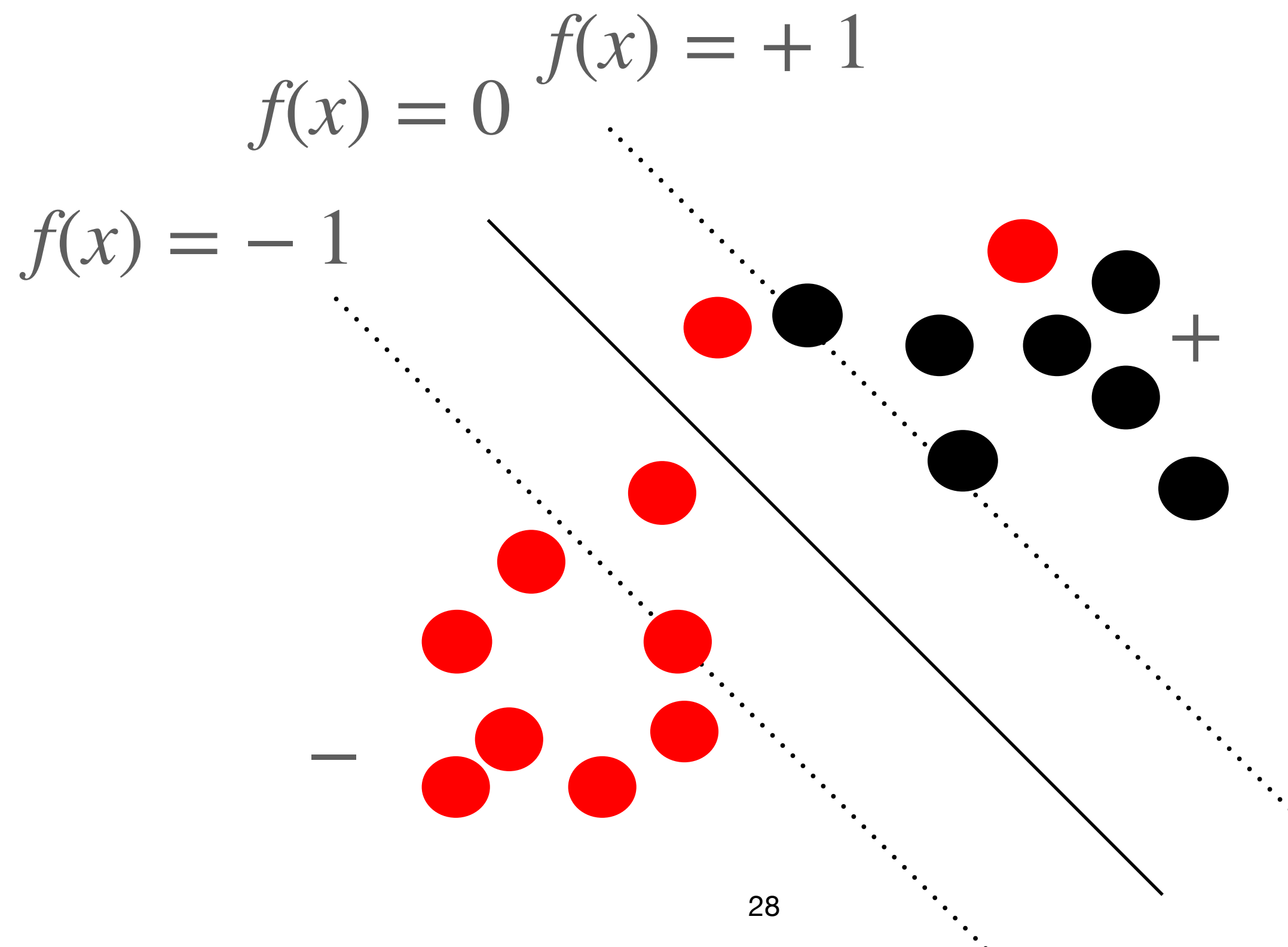


SVMs offer 3 approaches to address this problem:

1. *Relax* the constraints (soft-margin)
2. Still use hard-margin SVM, but *transform* the data (kernel)
3. The combination of 1 and 2 😊

## Soft-margin SVM: 'soft' constraint

- Relax constraints to allow points to be **inside the margin** or even on the **wrong side** of the boundary





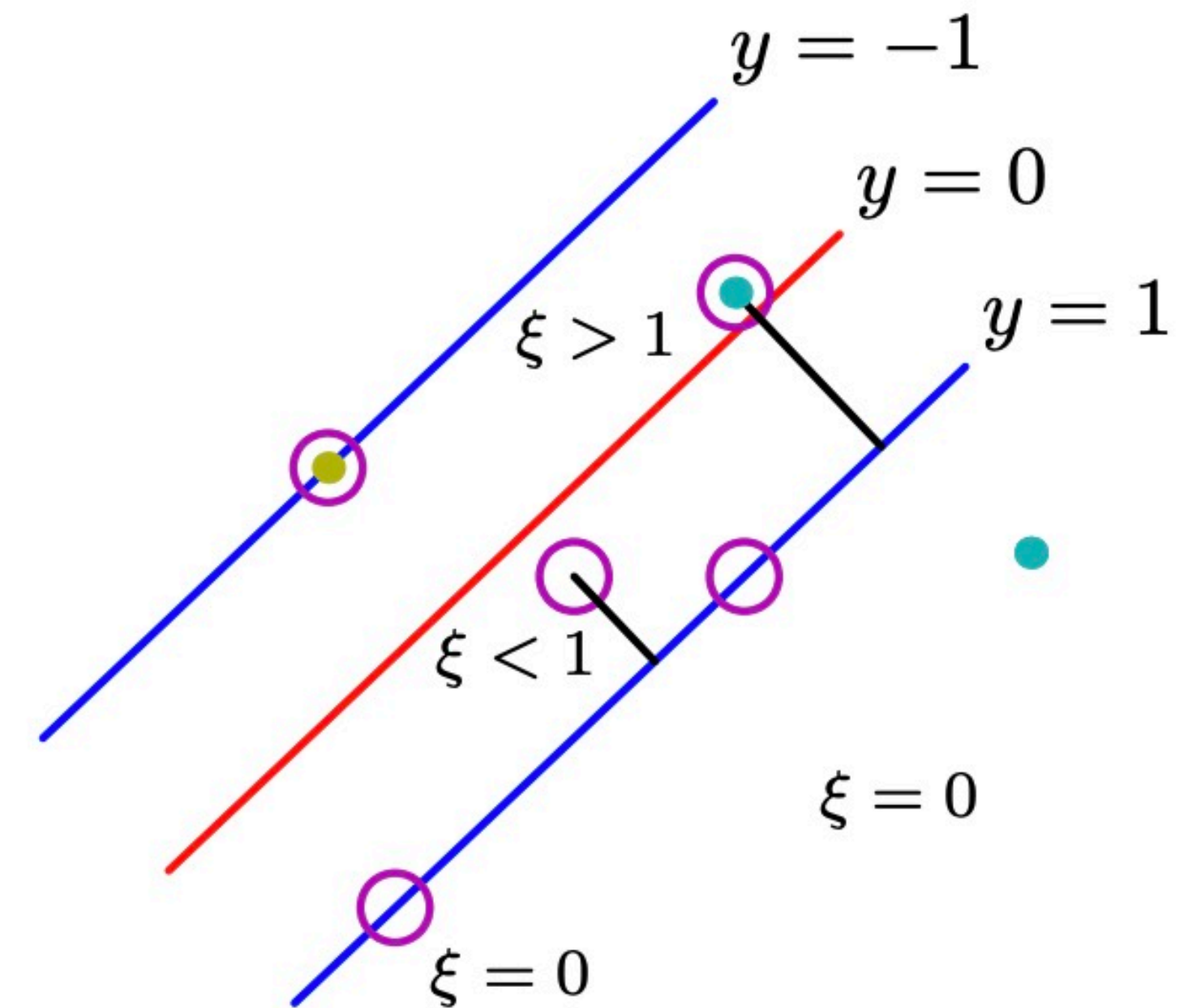
## Objective of soft-margin SVM

$$\min_w \left( \frac{\|w\|^2}{2} + C \sum_{i=1}^n \xi_i \right) \quad \text{s.t.} \quad \begin{aligned} & y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0 \quad (i = 1, \dots, n \text{ data points}) \end{aligned}$$

Use slack variable to 'soft' constraint:  
allow violation of the constraint

$$\xi_i = \begin{cases} 0, & y^{(i)}(w^T x^{(i)} + b) \geq 1, \\ 1 - y^{(i)}(w^T x^{(i)} + b), & \text{otherwise} \end{cases} \quad \text{0 if correctly classified}$$

or  $\xi_i = \max(0, 1 - y^{(i)}(w^T x^{(i)} + b))$  **hinge loss**





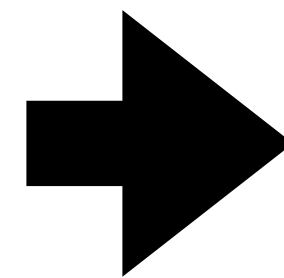
## Objective of soft-margin SVM

$$\min_w \left( \frac{\|w\|^2}{2} + C \sum_{i=1}^n \xi_i \right) \quad \text{s.t.} \quad \begin{aligned} y^{(i)}(w^T x^{(i)} + b) &\geq 1 - \xi_i, \\ \xi_i &\geq 0 \end{aligned} \quad (i = 1, \dots, n \text{ data points})$$

Slack penalty:  $C > 0$

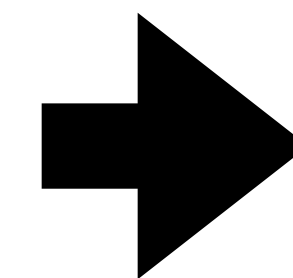
Wrong classified data

If  $C = 0$ : data is ignored



Underfitting

If  $C = \infty$ : data has to be correctly classified



Overfitting

**KKT**

Lagrange function where lambda and beta are both  
Lagrange multiplier

$$L(w, b, \lambda, \beta, \xi) = \frac{\|w\|^2}{2} + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \lambda_i g_i(w, b, \xi) + \sum_{i=1}^n \beta_i (-\xi_i)$$

$$g_i(w, b, \xi) = 1 - \xi_i - y^{(i)}(w^T x^{(i)} + b) \leq 0 \quad -\xi_i \leq 0$$

**Primal feasibility:**  $g_i(w, b, \xi) \leq 0 \quad -\xi_i \leq 0$

**Dual feasibility**  $\lambda_i \geq 0 \quad \beta_i \geq 0$

**Complementary slackness**  $\lambda_i g_i(w, b, \xi) = 0 \quad \beta_i \xi_i = 0$

**Stationarity**  $\frac{\partial L}{\partial w_j} = 0 : w_j = \sum_{i=1}^n \lambda_i y^{(i)} x_j^{(i)} \quad \frac{\partial L}{\partial b} = 0 : \sum_{i=1}^n \lambda_i y^{(i)} = 0$

$$\frac{\partial L}{\partial \xi_i} = 0 : C - \lambda_i - \beta_i = 0$$

**KKT**

**Primal feasibility:**  $g_i(w, b, \xi) = 1 - \xi_i - y^{(i)}(w^T x^{(i)} + b) \leq 0, \quad -\xi_i \leq 0$

**Dual feasibility**  $\lambda_i \geq 0 \quad \beta_i \geq 0$

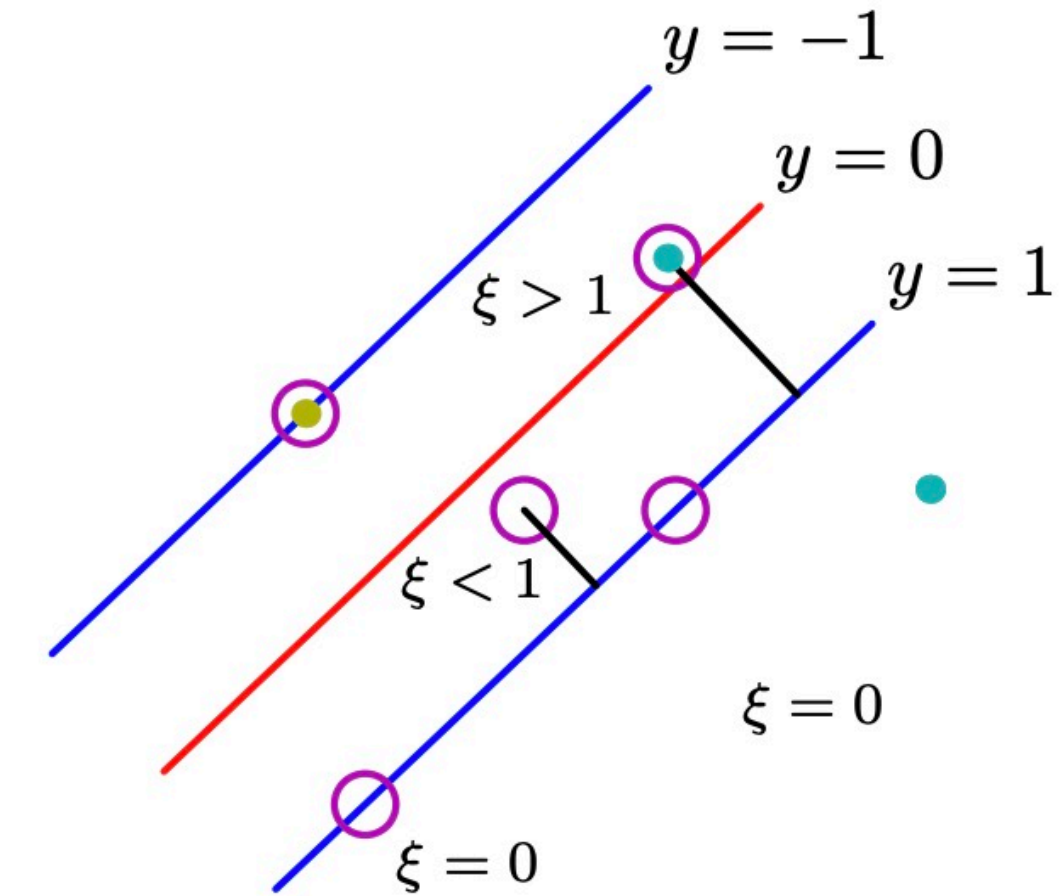
**Complementary slackness**  $\lambda_i g_i(w, b, \xi) = 0 \quad \beta_i \xi_i = 0$

$$C - \lambda_i - \beta_i = 0 : \quad 0 \leq \lambda_i \leq C \quad \text{Since } \beta_i \geq 0, \lambda_i \text{ must } \leq C$$

*if*  $\lambda_i = 0$  :  $\beta = C, \xi_i = 0$  Due to  $\beta_i \xi_i = 0$   $y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i = 1$  These are points we successfully classified

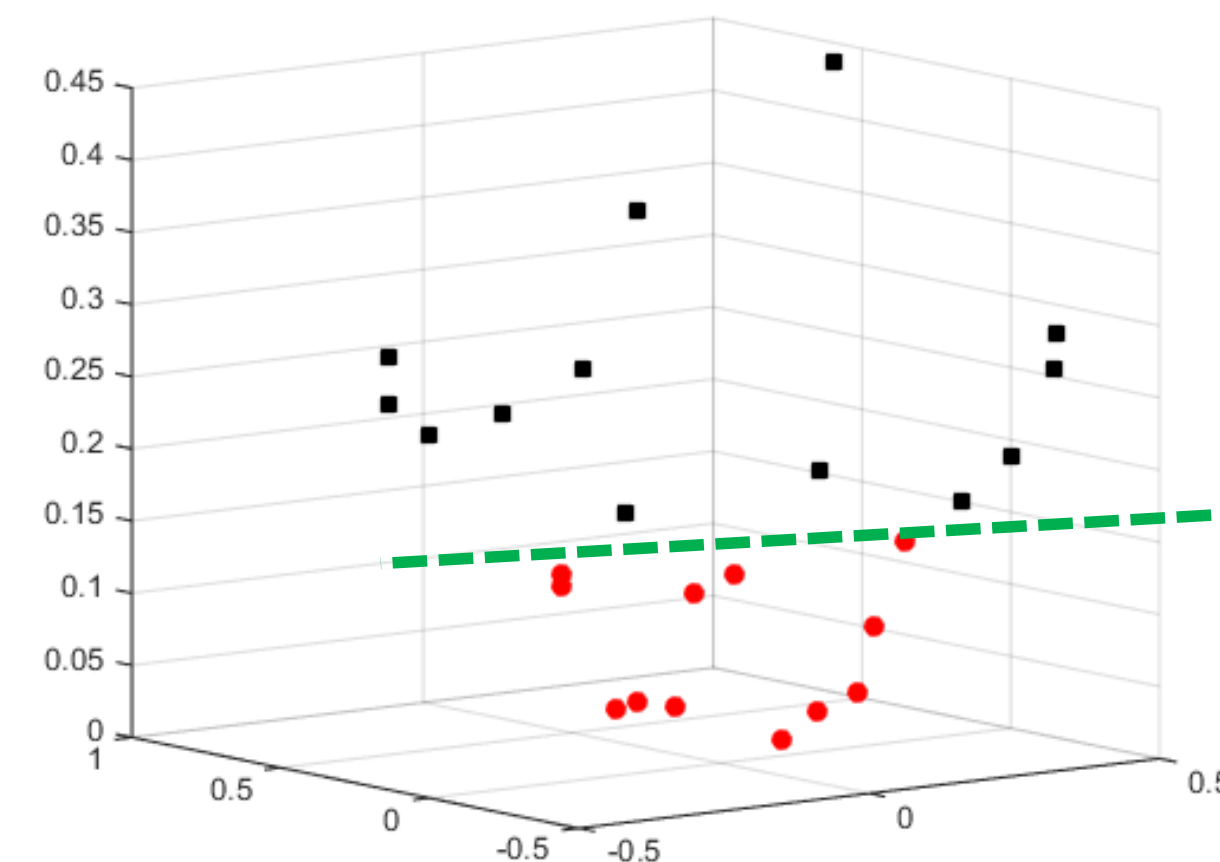
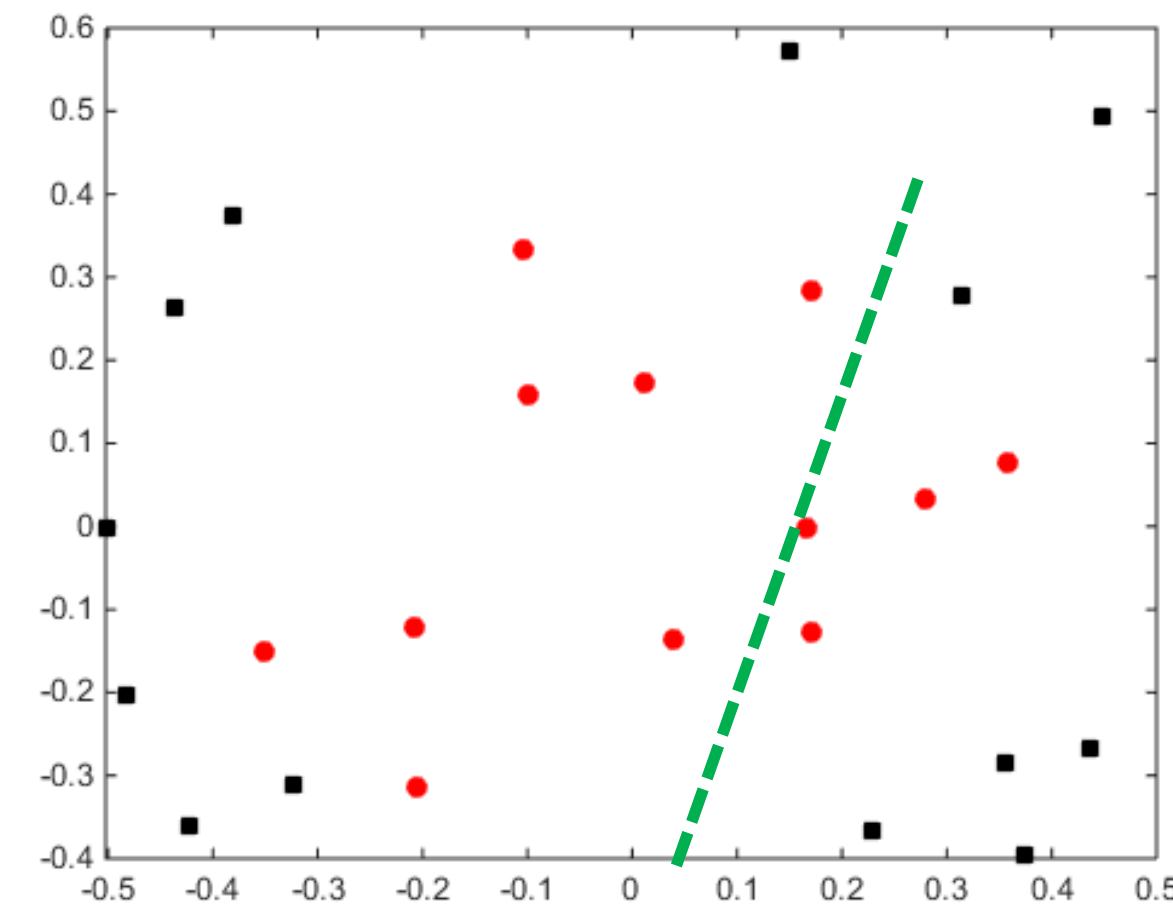
*if*  $\lambda_i = C$  :  $\beta_i = 0, -\xi_i \leq 0$   $y^{(i)}(w^T x^{(i)} + b) = 1 - \xi_i \leq 1$  These are points in the soft margin

*if*  $0 < \lambda_i < C$  :  $\xi_i = 0$   $g_i(w, b, \xi) = 0$   $y^{(i)}(w^T x^{(i)} + b) = 1 - \xi_i = 1$  Sine  $C > \beta_i > 0$  now **The point is a Support vector!**



## Non-linearly separable data

- Consider a binary classification problem
- Each example has features  $[x_1, x_2]$
- Not linearly separable
- Now 'add' a feature  $x_3 = x_1^2 + x_2^2$
- Each point is now  $[x_1, x_2, x_1^2 + x_2^2]$
- Linearly separable!



## Naïve workflow

- Choose/design a linear model
- Choose/design a high-dimensional transformation  $\varphi(\mathbf{x})$ 
  - \* Hoping that after adding a lot of various features some of them will make the data linearly separable
- For each training example, and for each new instance compute  $\varphi(\mathbf{x})$
- Train classifier/Do predictions

## Hard-margin SVM in feature space

Training, solve:

$$\begin{aligned} \max_{\lambda} \quad & \theta_d(\lambda) \\ \text{s.t.} \quad & \lambda_i \geq 0 \text{ and } \sum_{i=1}^n \lambda_i g_i(x) = 0 \end{aligned} \quad \theta_d(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \lambda_i \lambda_k y^{(i)} y^{(k)} \underbrace{(x^{(i)})^T x^{(k)}}$$

$$\text{Prediction: } f(x) = \sum_{i=1}^n \lambda_i y^{(i)} (x^{(i)})^T x + b$$

We just need the dot product!

Training: solve

$$\max_{\lambda} \theta_d(\lambda)$$

$$\theta_d(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \lambda_i \lambda_k y^{(i)} y^{(k)} \underbrace{(\varphi(x^{(i)}))^T \varphi(x^{(k)})}$$

Making predictions:

$$f(x) = w^T x + b = \sum_{i=1}^n \lambda_i y^{(i)} (\varphi(x^{(i)}))^T \varphi(x) + b$$



## Observation: Kernel representation

- Both parameter estimation and computing predictions depend on data only in a form of a **dot product**
  - \* In original space  $\mathbf{u}'\mathbf{v} = \sum_{i=1}^m u_i v_i$
  - \* In transformed space  $\varphi(\mathbf{u})'\varphi(\mathbf{v}) = \sum_{i=1}^l \varphi(\mathbf{u})_i \varphi(\mathbf{v})_i$
- **Kernel** is a function that can be expressed as a dot product in some feature space  $K(\mathbf{u}, \mathbf{v}) = \varphi(\mathbf{u})'\varphi(\mathbf{v})$

Benefits:

- no need to find the mapping function.
- no need to do transformation.
- no need to do dot product.



## Kernel as shortcut

- For *some*  $\varphi(\mathbf{x})$ 's, **kernel is faster to compute** directly than first mapping to feature space then taking dot product.
- For example, consider two vectors  $\mathbf{u} = [u_1]$  and  $\mathbf{v} = [v_1]$  and transformation  $\varphi(\mathbf{x}) = [x_1^2, \sqrt{2c}x_1, c]$ , some  $c$ 
  - \* So  $\varphi(\mathbf{u}) = [u_1^2, \sqrt{2c}u_1, c]'$  and  $\varphi(\mathbf{v}) = [v_1^2, \sqrt{2c}v_1, c]'$
  - \* Then  $\varphi(\mathbf{u})' \varphi(\mathbf{v}) = (u_1^2 v_1^2 + 2cu_1 v_1 + c^2)$
- This can be alternatively **computed directly** as
$$\varphi(\mathbf{u})' \varphi(\mathbf{v}) = (u_1 v_1 + c)^2$$
  - \* Here  $K(\mathbf{u}, \mathbf{v}) = (u_1 v_1 + c)^2$  is the corresponding kernel

**Hard-margin SVM in feature space**Training: solve

$$\max_{\lambda} L(\lambda) \quad L(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \lambda_i \lambda_k y^{(i)} y^{(k)} \overbrace{(\varphi(x^{(i)}))^T \varphi(x^{(k)})}^{K(x^{(i)}, x^{(k)})}$$

Making predictions:

$$f(x) = w^T x + b = \sum_{i=1}^n \lambda_i y^{(i)} \overbrace{(\varphi(x^{(i)}))^T \varphi(x)}^{K(x^{(i)}, x)} + b$$

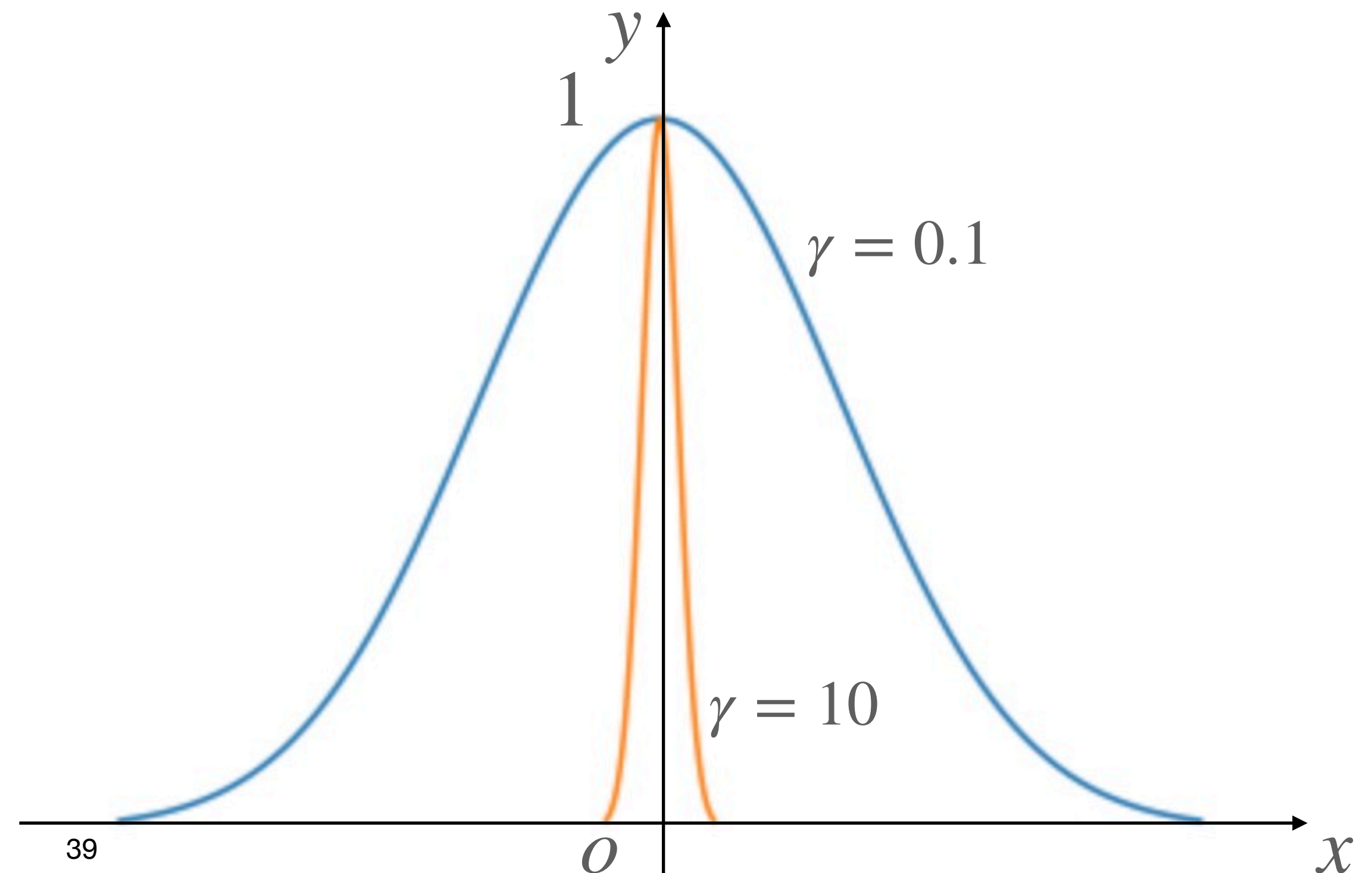
## Radial Basis Function (RBF) kernel

$$K(u, v) = \exp(-\gamma \|u - v\|^2)$$

*$\gamma$  is too small : underfitting*

*$\gamma$  is too large : overfitting*

$$y = \exp(-\gamma x^2) = \exp\left(-\frac{x^2}{2\sigma^2}\right)$$



## Identify new kernels

Mercer's theorem:

Consider a finite sequences of vectors  $x_1, \dots, x_n$

Construct  $n \times n$  matrix  $A$  (Gram matrix) of pairwise values

$K(x_i, x_j)$  is a valid kernel if this matrix is positive semi-definite, and this holds for all possible sequences

$$A = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_n) \\ K(x_2, x_1) & K(x_2, x_2) & \dots & K(x_2, x_n) \\ \vdots & \vdots & \vdots & \vdots \\ K(x_n, x_1) & K(x_n, x_2) & \dots & K(x_n, x_n) \end{bmatrix}$$

## Identify new kernels

Positive semi-definite matrix: a square symmetric matrix satisfies  $v^T A v \geq 0$   
 $v \in \mathbb{R}^{n \times 1}$  any non-zero vector (column),  $A \in \mathbb{R}^{n \times n}$ ,  $A = A^T$

$$A = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_n) \\ K(x_2, x_1) & K(x_2, x_2) & \dots & K(x_2, x_n) \\ \vdots & \vdots & \vdots & \vdots \\ K(x_n, x_1) & K(x_n, x_2) & \dots & K(x_n, x_n) \end{bmatrix}$$

## Identify new kernels

Let  $K_1(u, v)$ ,  $K_2(u, v)$  be kernels,  $c > 0$  be a constant, and  $f(x)$  be a real-valued function.

Then each of the following is also a kernel:

1)  $K(u, v) = K_1(u, v) + K_2(u, v)$

2)  $K(u, v) = c K_1(u, v)$

3)  $K(u, v) = f(u) K_1(u, v) f(v)$

# Summary

- What are the objective and constraints of hard-margin, soft-margin SVM
- What are KKT conditions?
- What are support vectors?
- What are Slack variables & slack penalty of soft-margin SVM?
- What is Kernel?
- How do parameters  $\gamma$ ,  $C$  influence performance of SVM?
- How to identify new kernels?