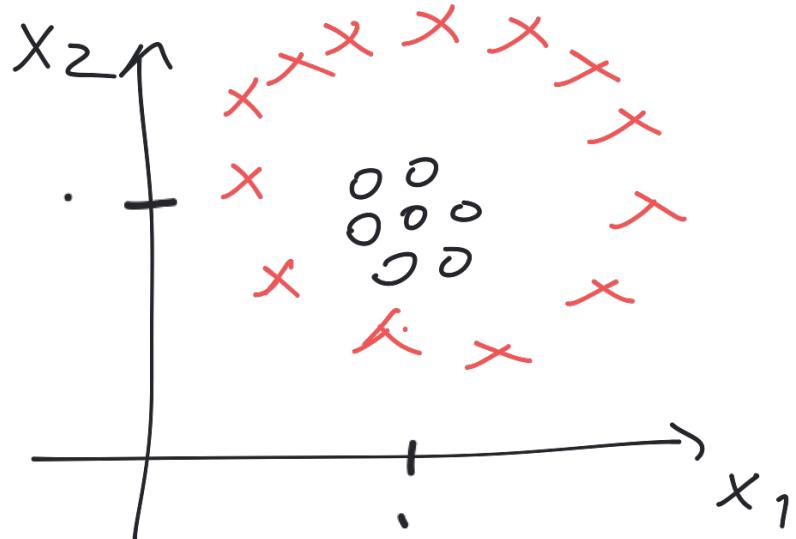
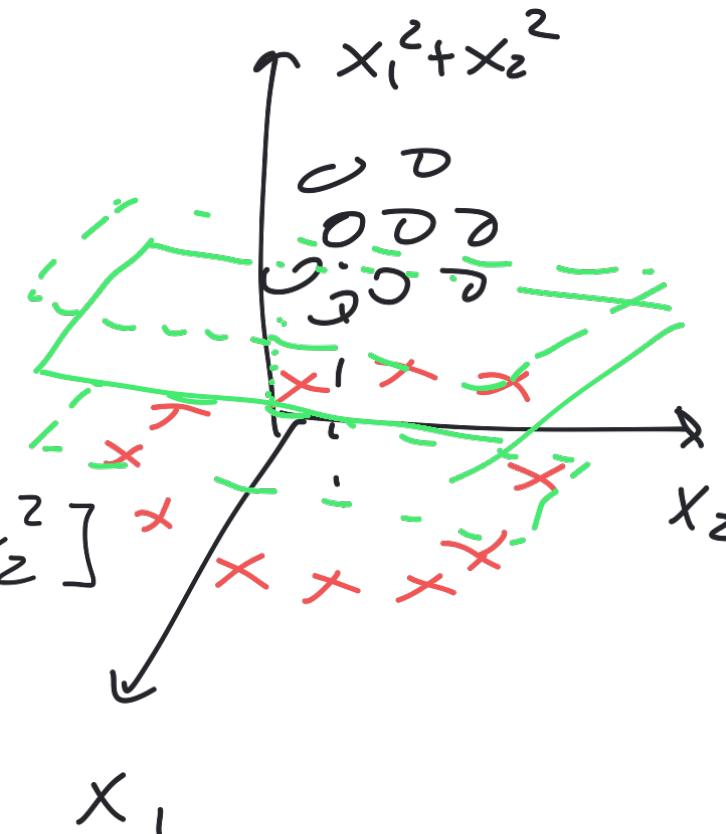


## SVM



Transform  
the data  
in higher-  
dimensional  
space.

$$\phi(x) = [x_1, x_2, x_1^2 + x_2^2]$$



SVM: ① Linear Classifier

② WORKS with a data transformation  
— MAPPING:  $\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$

③  $\phi(x)$  can be any  $x \mapsto \phi(x)$

transformation,  $\phi(x)$  is also has basis functions.

$\phi(x)$  can be infinite dimensional!!

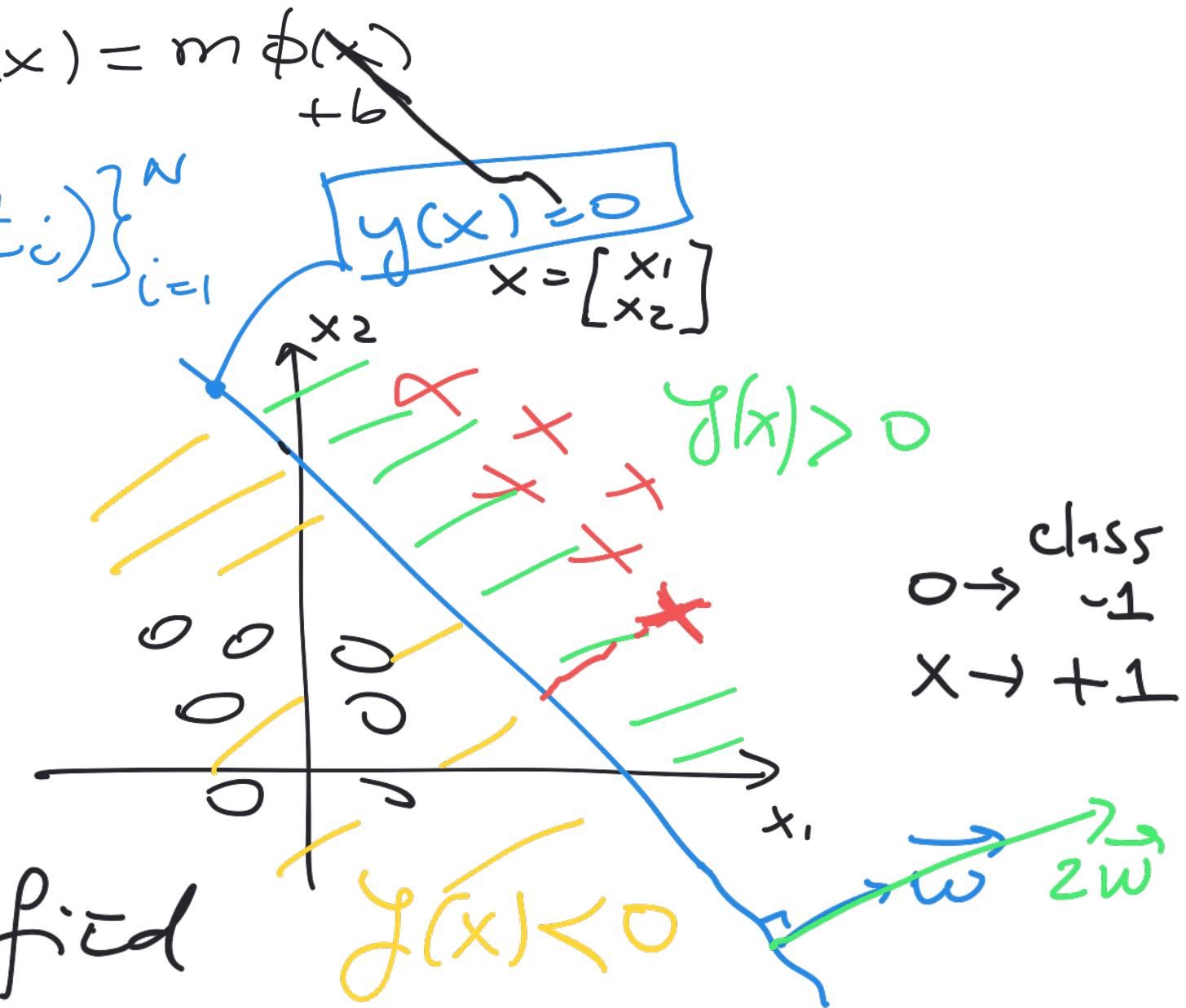
$$y(x) = \underline{w}^T \underline{\phi(x)} + b$$

if  $\phi(x)$  is 1-D then:  $y(x) = m \phi(x) + b$

Suppose we have data  $\{(x_i, t_i)\}_{i=1}^N$   
and  $t_i \in \{-1, 1\}$

$$y(x_n) \cdot t_n > 0$$

$\checkmark x_n$  correctly classified



Distance of a point  $x_n$  to decision surface:

$$\frac{t_n \cdot y(x_n)}{\|w\|} = \frac{t_n(w^T \phi(x_n) + b)}{\|w\|}$$

↳ we want to maximize this

distance  $\forall x_n$ !

↳ But we only  
use values  $x_n$  for  
support vectors!!



support vectors

↳ they have smallest dist. to surface.

$$\text{dist.} = \frac{t_n y(x_n)}{\|w\|}$$

$$\arg \max_{w,b} \left\{ \frac{1}{\|w\|} \min_n \sum_i [t_n(w^T \phi(x_n) + b)] \right\}$$

Scaling both  $\|w \rightarrow K w$  and  
 $b \rightarrow K b$

it does not affect direction  
 nor the distance to the  
 decision surface.

Then, we will consider support vectors to have distance:

$$t_n(\omega^\top \phi(x_n) + b) = 1$$

$\forall x_n \equiv$  support vectors.

Any other  $x_n$  will have distance:

$$t_n(\omega^\top \phi(x_n) + b) > 1$$

For all data points:

$$t_n(\omega^T \phi(x_n) + b) \geq 1, \forall x_n$$

Find  $\omega$  that defines maximum margin such that (constrained to)

$$t_n(\omega^T \phi(x_n) + b) \geq 1.$$

$$\arg \max_{\omega, b} \frac{1}{\|\omega\|} \quad \text{subject to} \\ t_n(\omega^T \phi(x_n) + b) \geq 1$$

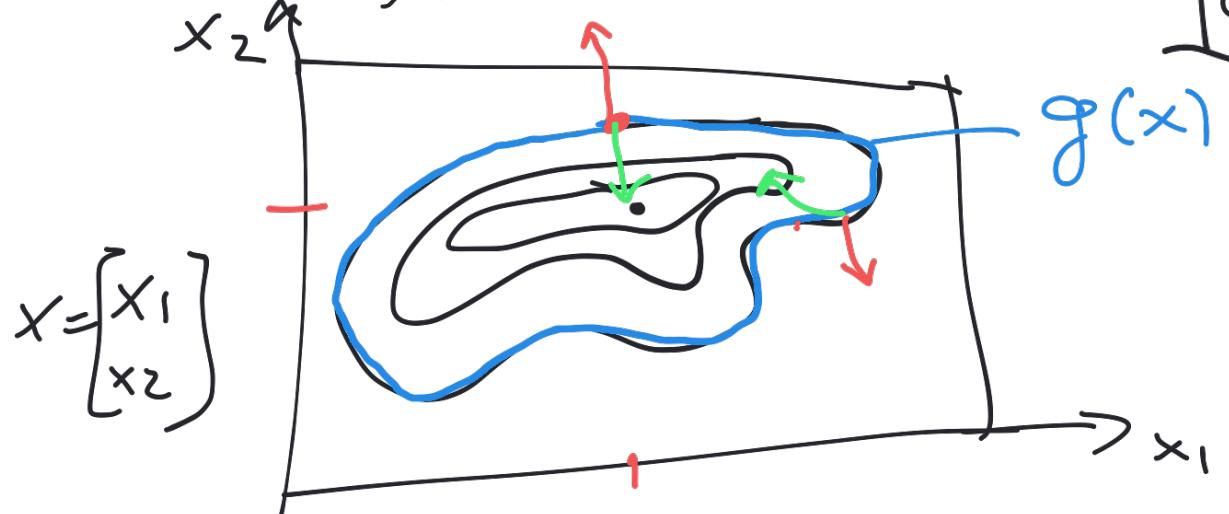
$$\Leftrightarrow \boxed{\arg \min_{\omega, b} \frac{1}{2} \|\omega\|^2 \quad \text{s.t.} \quad \underbrace{t_n(\omega^T \phi(x_n) + b) \geq 1}_{\text{---}}}$$

Optimization //  
of a function  
with equality / inequality  
constraint.

①  $\min_x f(x)$

Necessary condition  
 $f'(x) = 0$

②  $\min_x f(x)$  s.t.  $g(x) = 0$



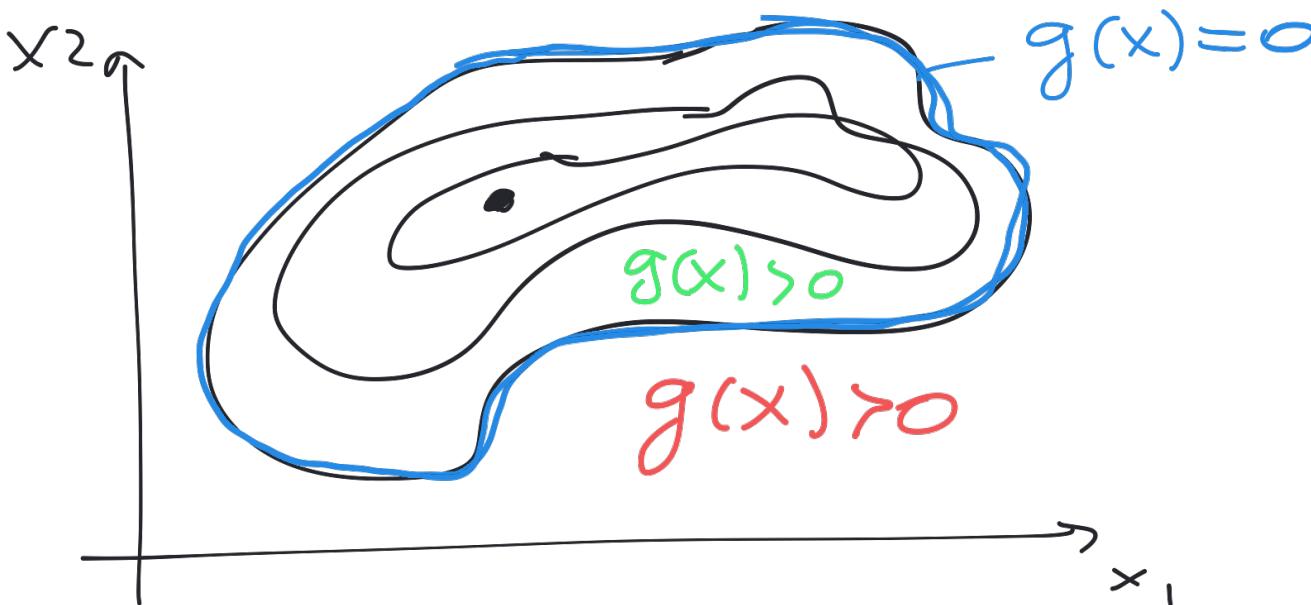
$$L(x, \lambda) = f(x) - \lambda \cdot g(x)$$

Lagrange  
multiplier

Necessary condition:

$$\left| \begin{array}{l} \frac{\partial L}{\partial x} = 0 \Rightarrow f'(x) = \lambda \cdot g'(x) \\ \frac{\partial L}{\partial \lambda} = 0 \Rightarrow +g(x) = 0 \end{array} \right.$$

$$\textcircled{3} \quad \min_x f(x) \quad \text{s.t. } g(x) \geq 0$$



necessary condition:

$$\begin{aligned} L(x, \lambda) &= f(x) - \lambda g(x) \\ \frac{\partial L}{\partial x} &= f'(x) - \lambda g'(x) = 0 \\ \frac{\partial L}{\partial \lambda} &= -g(x) = 0 \end{aligned}$$

Case 1: INACTIVE constraint:  $g(x) > 0, \lambda = 0$

Case 2: Active constraint:  $g(x) = 0, \lambda > 0$

KKT condition :

$$\left\{ \begin{array}{l} g(x) \geq 0 \\ \lambda \geq 0 \\ \lambda \cdot g(x) = 0 \end{array} \right.$$

Lagrangian multiplier for our SVM constrained optimization:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N \alpha_n (t_n (w^T \phi(x_n) + b) - 1)$$

PRIMAL

$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}}{\partial w} = 0 \Leftrightarrow \underline{w} = \sum_{n=1}^N \alpha_n \cdot t_n \cdot \phi(x_n) \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \Leftrightarrow \boxed{\sum_{n=1}^N \alpha_n t_n = 0} \end{array} \right.$$

Let's note that if  $\phi(x)$  is infinite dimensional, we CANNOT compute this.

By plugging in the solution for  $w$ , we are able rewrite PRIMAL  $\mathcal{L}$  into its DUAL representation.

DUAL :

$$L(a) = \sum_{n=1}^N a_n - \sum_{n=1}^N \sum_{m=1}^N a_n \cdot a_m t_n t_m \cdot K(x_n, x_m)$$

such that  $\sum_{n=1}^N a_n t_n = 0$  and

$$a_n > 0, \forall n = 1, 2, \dots, N$$

$K(x_n, x_m)$  known as GRAM MATRIX



kernel operation between points  $x_n$  and  $x_m$ .

The most popular kernel function  
is RADIAL BASIS FUNCTIONS (RBF)

$$\frac{\text{RBF}}{K(x,y)} = \exp\left(-\gamma \frac{\|x-y\|^2}{2}\right)$$



controls the variance

Proximity  
(or similarity)  
between two  
points  $x$  and  $y$ .

$\gamma \rightarrow \infty \rightarrow$  smaller variance

$\gamma \rightarrow 0 \rightarrow$  large variance.