

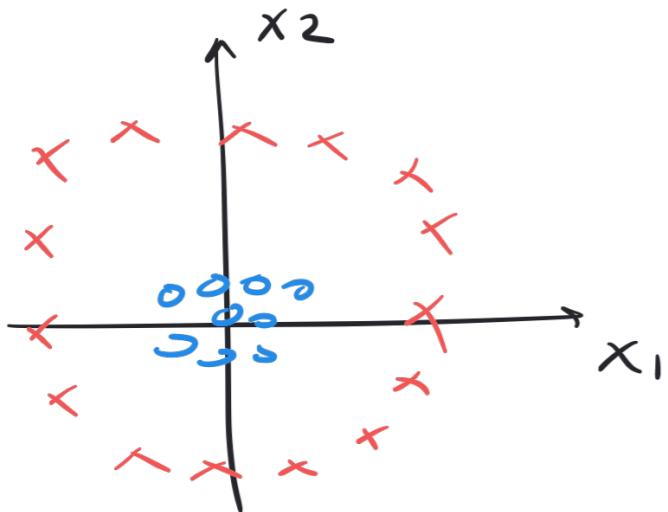
HARD-MARGIN SUPPORT VECTOR MACHINE (SVM)

MAPPER: $y(x) = \omega^T \phi(x) + b$

$\phi(x)$ = transformation space

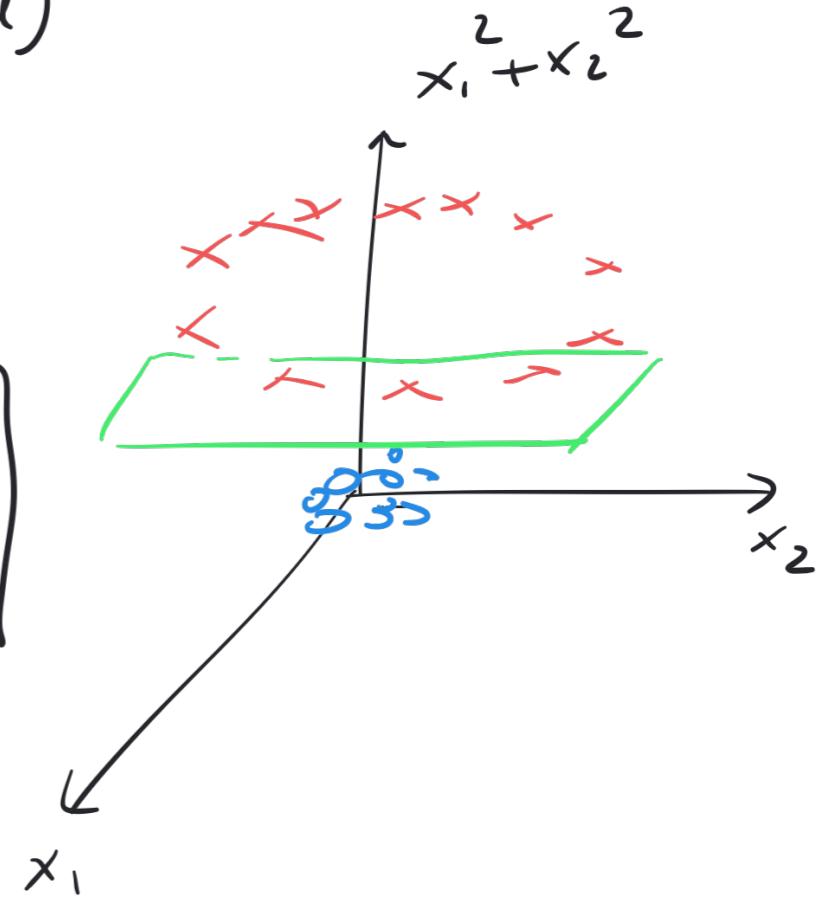
typically $\phi(x)$ maps to
a higher-dimensional space

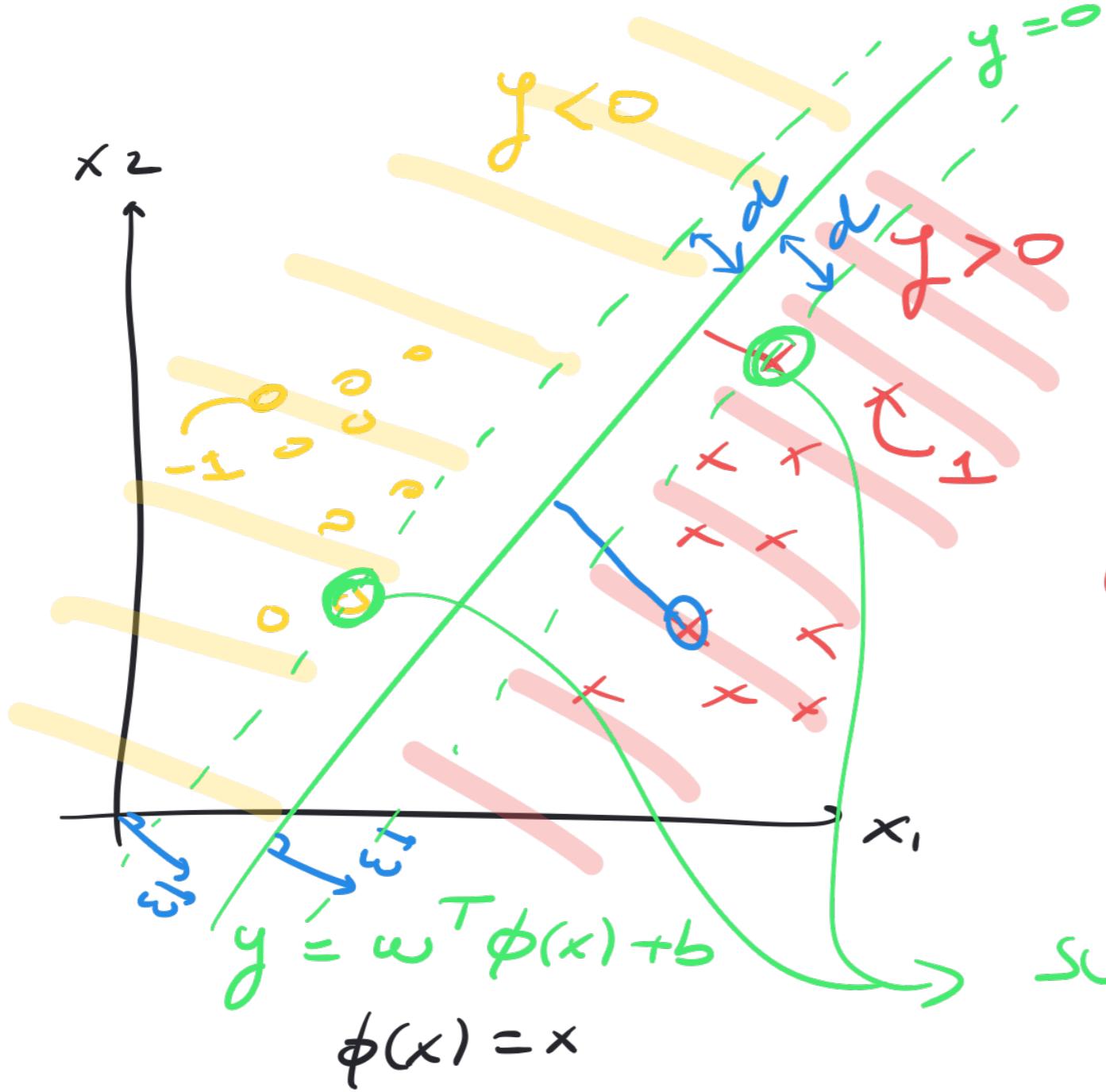
$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D \quad (D > d)$$



$$\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \mapsto \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{bmatrix}$$





x - class 1

o - class -1

$t_n = \text{label}^{\text{class}}$ for point $x_n \in \{1, -1\}$

① $y(x_n) \cdot t_n > 0$

x_n that are
correctly classified

SUPPORT VECTORS
(S.V.s)

Distance of point x_n to discriminant
function $y(x_n)$:

$$\frac{y(x_n) \cdot t_n}{\|\omega\|} \geq 0$$

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

so we want to Maximize this

distance for all points x_n .

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

we are looking to preserve the
smallest amount of points
that are closest to
 $y(x_n)$.

we will scale the parameters w and b

such that the distance of

Support vectors (S.V.s) to $y(x)$

is 1, i.e.,

$$\frac{t_n \cdot y(x_n)}{\|w\|} = 1 \quad \text{for } x_n \text{ that are S.V.s}$$

$$w \leftarrow k \cdot w$$

$$b \leftarrow k \cdot b$$

with this, we have:

$$t_n (\omega^T \phi(x_n) + b) = 1, \text{ if } x_n \text{ that are S.V.s}$$

and

$$t_n (\omega^T \phi(x_n) + b) > 1, \text{ if } x_n \text{ that are not S.V.s}$$

In general,

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

optimization problem w/ inequality

constraints

$$\arg \max_{\{\omega, b\}}$$

$$\frac{1}{\|\omega\|}$$

$$\text{sub. to } t_n \cdot (\omega^T \phi(x_n) + b) \geq 1$$

\Leftrightarrow

$$\arg \min_{\{\omega, b\}}$$

$$\|\omega\|^2$$

subject to

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

"/" (APPENDIX E - Bishop textbook)

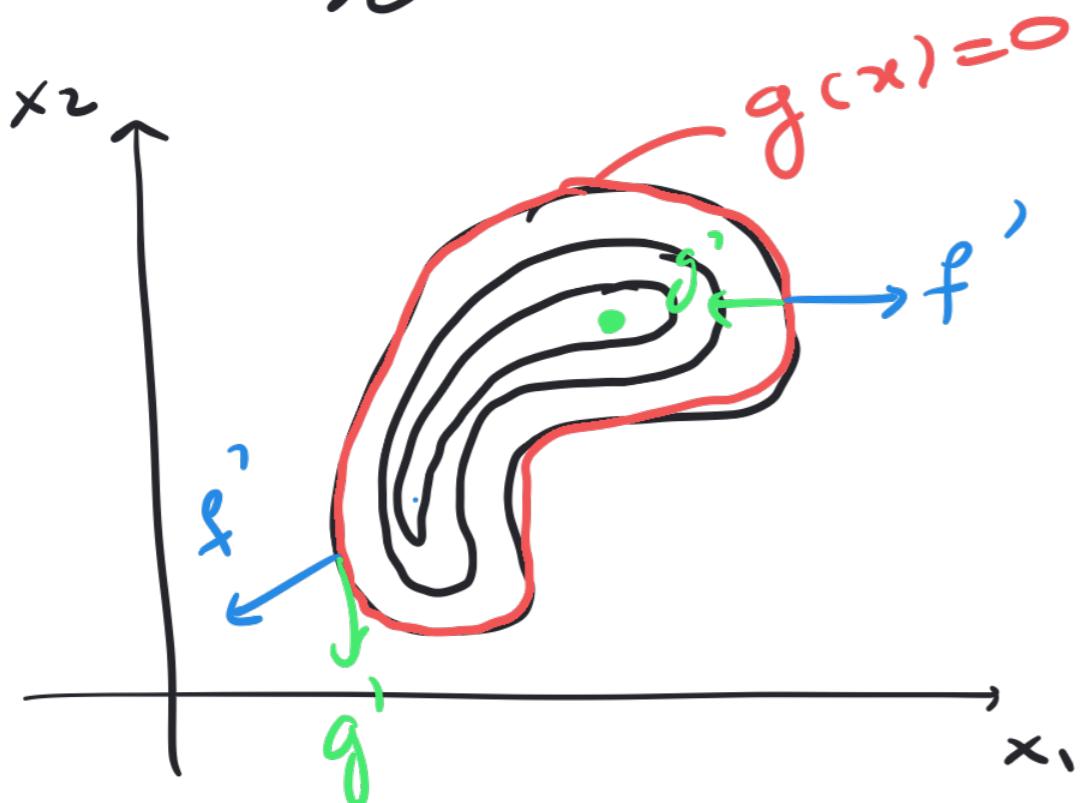
optimization function with
EQUALITY / INEQUALITY
CONSTRAINTS

① $\min_x f(x)$

NECESSARY condition:

$$f'(x) = 0$$

$$\textcircled{2} \quad \arg \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{sub. to} \quad g(\mathbf{x}) = 0$$



Lagrangian function:

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

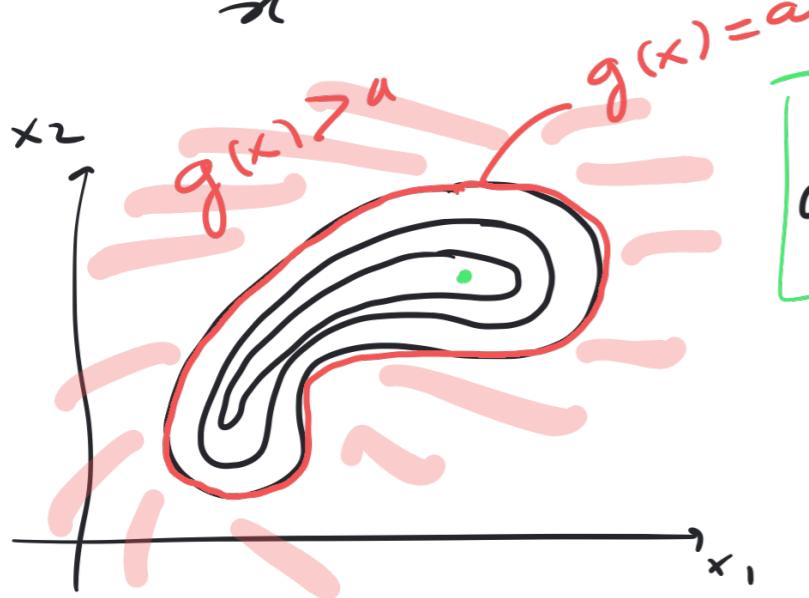
Lagrange
multiplier

Necessary condition:

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

$$③ \arg \min_{\alpha} f(\alpha)$$

sub. to $g(x) \geq a$
 $\Leftrightarrow g(x) - a \geq 0$

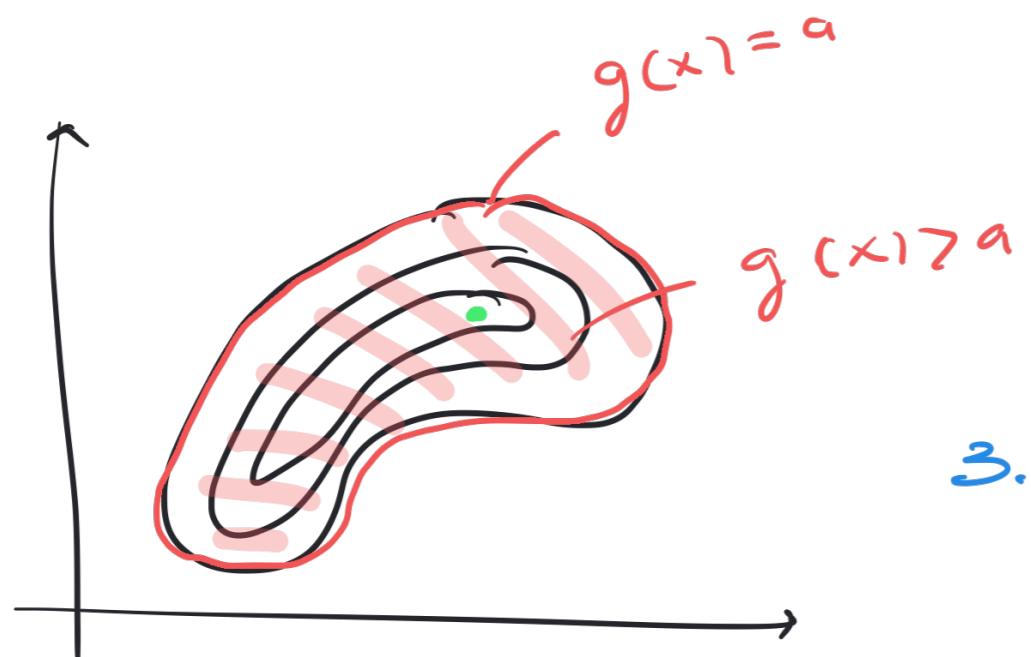


$$\mathcal{L}(x, \lambda) = f(x) - \lambda(g(x) - a)$$

3.1) ACTIVE CONSTRAINT

In this case the closest x that minimizes $f(x)$ is on $g(x) = a$

$$\lambda > 0, g(x) = a$$



3.2) INACTIVE CONSTRAINT

In this case, the minima of $f(x)$ already satisfies $g(x) > a$ thus:

$$\lambda = 0, g(x) > a$$

KARUSH - KUHN - TUCKER (KKT)

conditions

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

← at least
one of them
is zero.

!!

Coming back to the SVM optimization:

$$\arg \min_{\{\omega, b\}} \|\omega\|^2 \text{ sub. to } t_n \cdot y(x_n) \geq 1$$

Lagrangian function:

$$L(\omega, b, a) = \|\omega\|^2 - \sum_{n=1}^N a_n \cdot (t_n y(x_n) - 1)$$

Set of
Lagrange
multipliers

KKT conditions:

$$\begin{cases} a_n \geq 0 \\ t_n \cdot y(x_n) - 1 \geq 0 \\ a_n (t_n \cdot y(x_n) - 1) = 0 \end{cases}$$

These conditions say that for every point x_n :

$$a_n = 0$$

OR

$$a_n > 0$$

$$t_n \cdot y(x_n) > 1$$

inactive constraint

x_n is NOT a SV.

$$t_n \cdot y(x_n) = 1$$

active constraint

x_n is a SV.

$$\mathcal{L}(\omega, b, a) = \frac{1}{2} \|\omega\|^2 - \sum_{n=1}^N a_n (t_n (\omega^\top \phi(x_n) + b) - 1)$$

Optimizing $\mathcal{L}(\omega, b, a)$ w.r.t. ω and b :

$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}}{\partial \omega} = 0 \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \end{array} \right. \quad (\Rightarrow) \quad \left\{ \begin{array}{l} \omega = \sum_{n=1}^N a_n \cdot t_n \cdot \underline{\phi(x_n)} \\ \sum_{n=1}^N a_n \cdot t_n = 0 \end{array} \right.$$

in some cases, $\phi(x)$ is infinite-dimensional. So we cannot compute, nor store it.
 so we will work with the DUAL LAGRANGIAN

Substituting ω in the Lagrangian

function and making sure that

condition $\sum_{n=1}^N a_n \cdot t_n = 0$ is met:

The DUAL LAGRANGIAN is:

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

where $a_n \geq 0$

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

KERNEL function

kernel function

$$K : \mathbb{R}^d \times \mathbb{R}^d \longrightarrow \mathbb{R}$$

$$(x, y) \longmapsto \phi^T(x) \cdot \phi(y)$$

$$\phi : \mathbb{R}^d \longrightarrow \mathbb{R}^D$$

Gram matrix

- ① K is $N \times N$, symmetric

$$K = \begin{bmatrix} \phi^T(x_1) \cdot \phi(x_1) & \phi^T(x_1) \cdot \phi(x_2) & \dots & \phi^T(x_1) \cdot \phi(x_n) \\ \phi^T(x_2) \cdot \phi(x_1) & \phi^T(x_2) \cdot \phi(x_2) & \dots & \phi^T(x_2) \cdot \phi(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \phi^T(x_n) \cdot \phi(x_1) & \phi^T(x_n) \cdot \phi(x_2) & \dots & \phi^T(x_n) \cdot \phi(x_n) \end{bmatrix}$$

- ② K is positive semi-definite

$$z^T K z \geq 0, \text{ for } z \in \mathbb{R}^n \setminus \{\vec{0}\}$$

- ③ All eigenvalues are ≥ 0 . which means that some columns / rows are collinear.

$$\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \mapsto \begin{bmatrix} x_1^2 \\ \sqrt{2} \cdot x_1 \cdot x_2 \\ x_2^2 \end{bmatrix}$$

KERNEL
TRICK

$$K(x, y) = \phi^T(x) \cdot \phi(y)$$

$$= [x_1^2, \sqrt{2} \cdot x_1 \cdot x_2, x_2^2] \begin{bmatrix} y_1^2 \\ \sqrt{2} \cdot y_1 \cdot y_2 \\ y_2^2 \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$= x_1^2 \cdot y_1^2 + 2 \cdot x_1 \cdot x_2 \cdot y_1 \cdot y_2 + x_2^2 \cdot y_2^2$$

$$= (x_1 y_1 + x_2 y_2)^2$$

$$= \left(\begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right)^2 = (x^T \cdot y)^2$$

RADIAL BASIS FUNCTION (RBF)

it maps the input space to an

infinite-dimensional

space.

$$\phi_{\text{RBF}} : \mathbb{R}^d \rightarrow \mathbb{R}^\infty$$

$$K_{\text{RBF}} : \mathbb{R}^d \times \mathbb{R}^d \longrightarrow \mathbb{R}$$
$$(x, y) \longmapsto e^{-\gamma \cdot \|x-y\|^2}$$

equivalent
 $e^{-\frac{\|x-y\|^2}{\sigma^2}}, r = \frac{1}{\sigma^2}$

Proof: $\gamma = \frac{1}{2}$ (without loss of generality)

$$\begin{aligned} K_{RBF}(x, y) &= \exp\left(-\frac{1}{2} \|x - y\|^2\right) \\ &= \exp\left(-\frac{1}{2} (x - y)^T (x - y)\right) \\ &= \exp\left(-\frac{1}{2} (x^T x - x^T y - y^T x + y^T y)\right) \\ &= \exp\left(-\frac{1}{2} x^T x - \frac{1}{2} y^T y + x^T y\right) \\ &= \underbrace{\exp\left(-\frac{1}{2} x^T x - \frac{1}{2} y^T y\right)}_{= c \text{ constant}} \cdot \exp(x^T y) \end{aligned}$$

$$= c \cdot \exp(x^T y) = c \cdot \sum_{n=0}^{\infty} \frac{(x^T y)^n}{n!} = c \cdot \sum_{n=0}^{\infty} \frac{K_{poly}^{(n)}(x, y)}{n!}$$

Taylor Series
Expansion

Q.E.D.

- The dual Lagrangian takes a form of a **Quadratic programming** optimization. We can use known algorithms to find the solution for a.
- The solution to a quadratic optimization problem on N variables is $\Theta(N^3)$.
- Going to the dual Lagrangian, we have $D \rightarrow N$ variables ($N \equiv \# \text{ samples}$ and $D \equiv \# \text{ of features in } \phi$), $N \gg D$.

$$\omega = \sum_{n=1}^N a_n \cdot t_n \cdot \phi(x_n)$$

and

$$y(x) = \omega^T \phi(x) + b$$

During test, we only use the S.V.s

to make a decision:

$$y(x) = \sum_{n=1}^N a_n \cdot t_n \cdot \phi^T(x_n) \cdot \phi(x) + b$$

$$\text{test point} = \sum_{m \in S} a_m \cdot t_m \cdot \phi^T(x_m) \cdot \phi(x) + b$$

$$= \sum_{m \in S} a_m \cdot t_m \cdot K(x_m, x) + b$$

$a_n > 0$ and $t_n \cdot y(x_n) = 1$ for x_n that
are S.V.s

where $y(x) = \sum_{n=1}^N a_n \cdot t_n \cdot \phi^T(x_n) \cdot \phi(x) + b$

$$= \sum_{m \in S} a_m \cdot t_m \cdot \phi^T(x_m) \cdot \phi(x) + b$$

Substituting,

$$t_n \left(\sum_{m \in S} a_m \cdot t_m \cdot \phi^T(x_m) \cdot \phi(x_n) + b \right) = 1$$

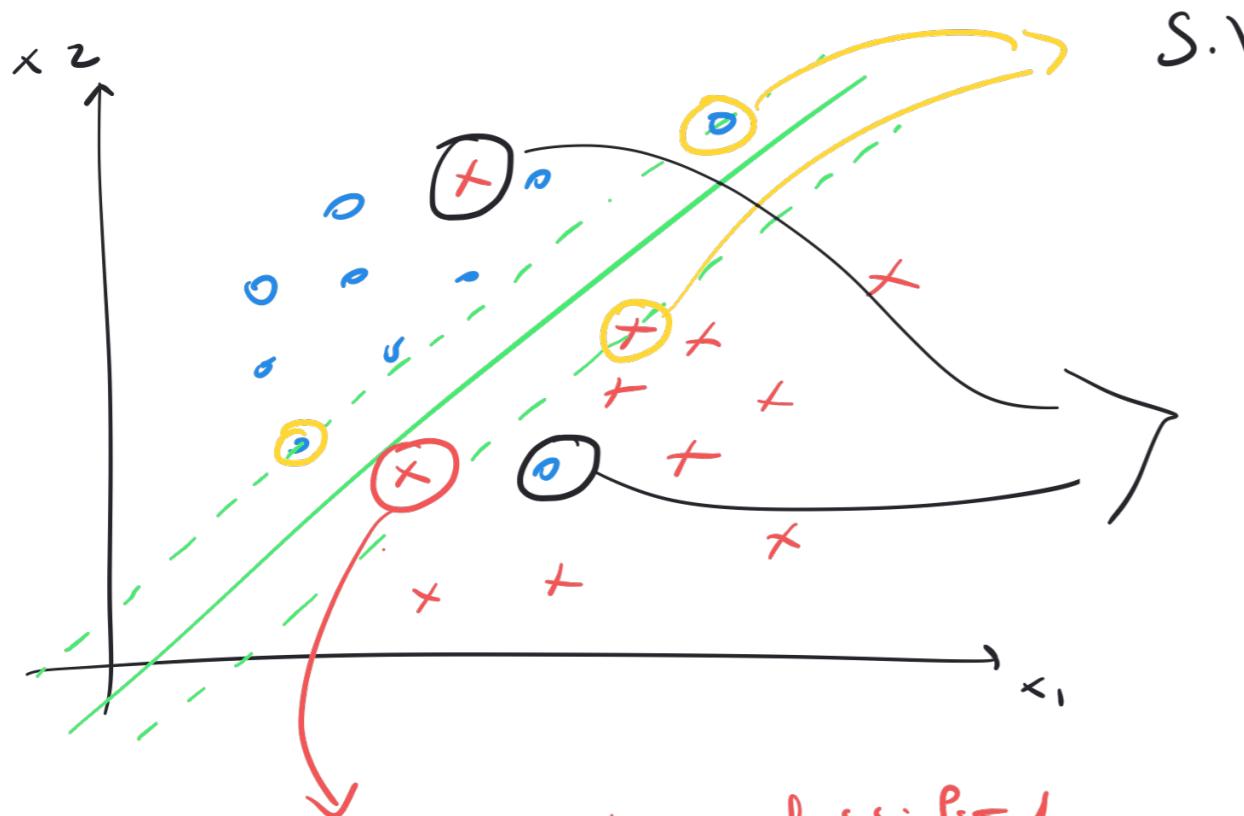
$$b = \frac{1}{N_S} \sum_{m \in S} \left[t_n - \left(\sum_{n \in S} a_n \cdot t_n \cdot K(x_m, x_n) \right) \right]$$

$N_S \equiv \# \text{ of S.V.s}$

$S \equiv \text{set of all samples that are S.V.s}$

Soft - Margin

SVM



S.V.s

incorrectly
classified

correctly classified
but inside the margin

$$\epsilon_n = |t_n - y(x_n)| \geq 0$$

Slack variable

$\epsilon_n = 0$ for all
 x_n that
correctly
classified
and on or
outside of
margin