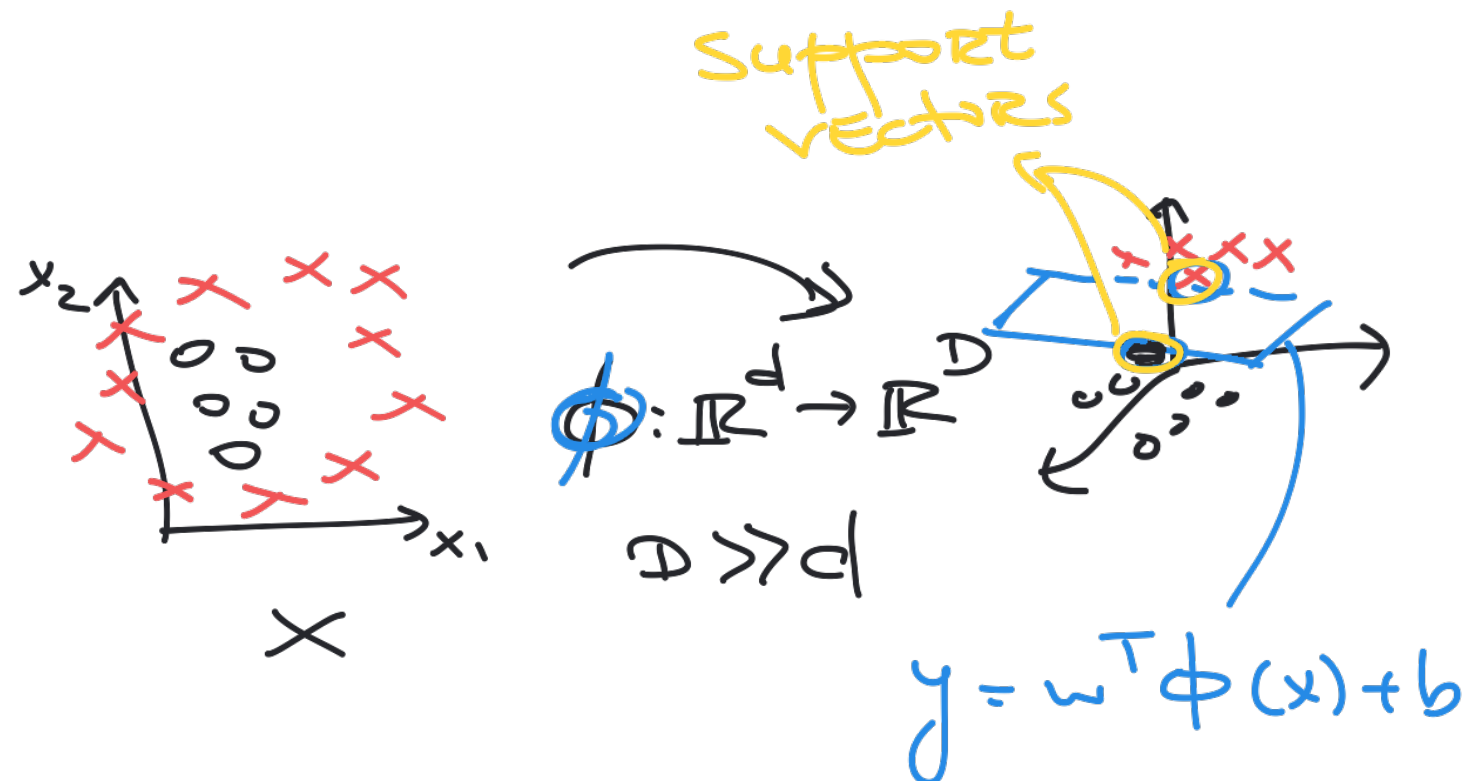


SVM

- Kernel Machine
- Maximizes Margin



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

- Support vectors satisfy: $t_n (w^T \phi(x_n) + b) = 1$, then

$$\parallel \arg \min_{w, b} \frac{1}{2} \|w\|^2 \quad \text{Sub. to } \underline{t_n (w^T \phi(x_n) + b) \geq 1}$$

- Primal Lagrangian:

$$// \mathcal{L}(w, b, a) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N \underline{a_n} [\underline{t_n (w^T \phi(x_n) + b) - 1}]$$

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

$$// \begin{cases} \frac{\partial \mathcal{L}}{\partial w} = 0 \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \end{cases} \Leftrightarrow \begin{cases} \underline{w} = \sum_{n=1}^N a_n t_n \phi(x_n) \\ 0 = \sum_{n=1}^N a_n t_n \end{cases}$$

- DUAL Lagrangian: (plugging in w in $\mathcal{L}(w, b, a)$)

$$// \tilde{\mathcal{L}}(a) = \sum_{n=1}^N a_n - \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \underline{k(x_n, x_m)} = \phi(x_n)^T \phi(x_m)$$

GRAM
MATRIX

such that $a_n \geq 0, n=1, 2, \dots, N$
 $\sum_{n=1}^N a_n t_n = 0$

The solution for

$\tilde{\mathcal{L}}(a)$ is the same as the solution for \mathcal{L} .

So in the dual representation, we need to compute the GRAM MATRIX

GRAM MATRIX

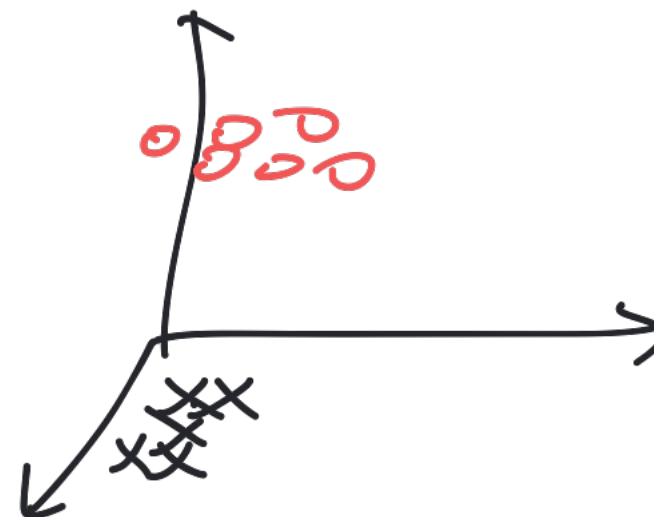
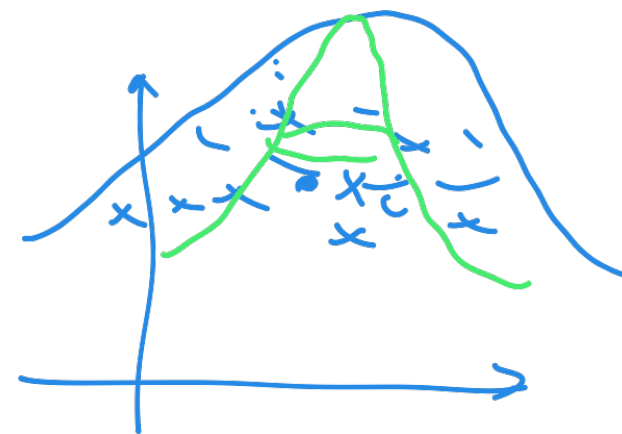
$$K = \begin{bmatrix} \phi^T(x_1)\phi(x_1) & \phi^T(x_1)\phi(x_2) & \dots & \phi^T(x_1)\phi(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \phi^T(x_n)\phi(x_1) & \phi^T(x_n)\phi(x_2) & \dots & \phi^T(x_n)\phi(x_n) \end{bmatrix}$$

$$\phi^T(x_i)\phi(x_j) = K(x_i, x_j) = \exp(-\underline{\sigma} \|x_i - x_j\|^2)$$

e.g. we can
USE RBF

If σ is small \rightarrow variance large

σ is large \rightarrow variance small





So, γ , in the kernel
function controls
the neighborhood of

Each data sample.

→ It needs to be selected using

Cross-validation.

In dual rep., we use quadratic
dynamic programming to solve
for a.

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

Using the solution for w :

$$y(x) = \sum_{n=1}^N a_n t_n \boxed{K(x, x_n)} + b$$

In training, we can define the
GRAM MATRIX K for all samples

$$y(\underset{\uparrow}{x}) = \sum_{n=1}^N a_n t_n \underbrace{K(x, x_n)} + \underline{\underline{b}}$$

If $x \equiv x_T$, then we compute its
similarity (in the kernel fct. sense) as:

$$\left[\phi^T(x_T) \phi(x_1), \phi^T(x_T) \phi(x_2), \dots, \phi^T(x_T) \phi(x_N) \right]^T$$

Linear kernel: $K(x, y) = x^T y$

CASE OVERLAPPING CLASSES:

$$S_n = |t_n - y(x_n)|$$

$S_n = 0$: x_n is correctly classified and far away from margin

$$\arg \min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N S_n$$

$$\text{sub. to } \|t_n y(x_n)\| \geq 1 - S_n$$

$$\|S_n \geq 0, n = 1, \dots, N$$

$0 < S_n \leq 1$: x_n lies inside the (correct side) margin

$S_n > 1$: x_n lies on wrong side of decision boundary

$C \rightarrow \infty$: RECOVERS the HARD MARGIN w/out misclassification

$C \rightarrow 0$: allows for a lot misclassified samples.

(PRIMAL)
① Define Lagrangian:

$$\mathcal{L}(w, b, a) = \frac{1}{2} \|w\|^2 + c \sum_{n=1}^N s_n - \sum_{n=1}^N a_n (\tan y(x_n) - 1 + s_n) - \sum_{n=1}^N \mu_n s_n$$

Lagrange
multiplier for
each sample

② Define the KKT conditions.

③ Take derivatives of \mathcal{L} w.r.t. to w, b, s_n .

④ Construct the DUAL Lagrangian by plugging in solutions in ③.