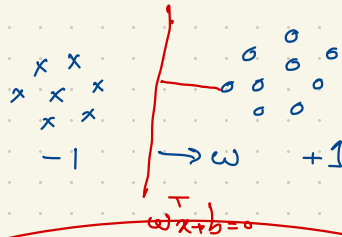$$\gamma \triangleq \min_{i=1,\ldots,N} \text{dist}(x^i, \omega^T x + b = 0)$$

$$\begin{cases} \max \ \gamma \\ \text{s.t.} \ y^i(\omega^T x^i + b) \geq \gamma \ , \ i=1,\ldots,N \end{cases}$$

$$\equiv$$

$$\begin{cases} \min_{\omega, b} \ \frac{1}{2} \|\omega\|_2^2 \qquad \underline{\text{Vanilla SVM}} \\ \text{s.t.} \ y^i(\omega^T x^i + b) \geq 1 \quad \forall i=1,\ldots,N \end{cases}$$

Convex opt! $\implies$ Solve using convex Solvers!

$\overline{\text{Training}} \longrightarrow$

Testing : $x^{new}$ $\qquad y(x^{new}) = \text{sgn}\left(\omega^{*T} x^{new} + b^*\right)$

**Challenge** : could be costly to solve when $d$ ($x^i \in \mathbb{R}^d$) is large!

$$\omega \in \mathbb{R}^d, b \in \mathbb{R} \implies \# \text{pars} = O(d)$$

$$\text{Complexity of solver} \quad O(d^3) \quad \ddot\frown$$

How to overcome the computational bottleneck?

Build the Lagrangian function:

$$\min_{\omega, b} \frac{1}{2} \|\omega\|_2^2$$

$$\text{s.t. } g^i(\omega^T x^i + b) \geq 1, \quad \forall i = 1, \cdots, N$$

$\Rightarrow$

$$\min_{\omega, b} \frac{1}{2} \|\omega\|_2^2$$

$$\mathcal{Z} = \begin{pmatrix} \omega \\ b \end{pmatrix}$$

$$\text{s.t. } 1 - y^i(\omega^T x^i + b) \leq 0 \quad \forall i = 1, -, N$$

$\alpha_i \geq 0$  Lagrange multiplier

$$L(\omega, b, \alpha_1, \cdots, \alpha_N) = \frac{1}{2} \|\omega\|_2^2 + \sum_{i=1}^{N} \alpha_i \left( 1 - y^i(\omega^T x^i + b) \right)$$

$\underbrace{\frac{1}{2} \omega^T \omega}$

$\in \mathbb{R}$, $\pm 1$

$$\boxed{- \alpha_i y^i x^{iT} \omega}$$

$$\frac{\partial \alpha^T \omega}{\partial \omega} = a$$

$$\frac{\partial L}{\partial \omega} = 0 \Rightarrow \omega + \sum_{i=1}^{N} - \alpha_i y^i x^i = 0$$

$$\Rightarrow \boxed{\omega^* = \sum_{i=1}^{N} \alpha_i^* y^i x^i}$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^{N} - \alpha_i y^i = 0 \longrightarrow$$

$$\boxed{\sum_{i=1}^{N} \alpha_i^* y^i = 0}$$

[ Skipping mathematical derivations of Lagrangian wrt $\alpha_i$'s ]

plug back $\omega^* \longrightarrow L(\overset{*}{\omega}, b, \alpha_1, \ldots, \alpha_N) = -\frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} y^i y^j (z^{iT} x^j) \alpha_i \alpha_j + \sum_{i=1}^{N} \alpha_i$

Dual SVM

$$\begin{cases} \max\limits_{\alpha_1, \ldots, \alpha_N} -\frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} y^i y^j \underline{(z^{iT} x^j)} \alpha_i \alpha_j + \sum_{i=1}^{N} \alpha_i \\ \\ \text{s.t.} \quad \sum_{i=1}^{N} \alpha_i y^i = 0 \quad , \quad \alpha_i \geq 0 \quad \forall i = 1, \ldots, N \end{cases}$$

$\Longrightarrow$

unknowns

$$\begin{cases} \min\limits_{\omega, b} \frac{1}{2} \|\omega\|_2^2 \longrightarrow d+1 \\ \\ \text{s.t.} \quad y^i (\omega^T x_i + b) \geq 1 \end{cases}$$

$d \gg N$

only $N$ unknowns

Convex optimization!

e.g. MRI scans ( $d$ : # pixels )
( $N$ : # patients )

$\alpha_1^*, \ldots, \alpha_N^*$

$O(N^3)$ (Naively) $\longrightarrow$ $O(N^2)$ intelligent implementations!

To classify, we need $\omega^*, b^*$ to compute $\text{sgn}(\underline{\omega^{*T}} x^{new} + \underline{b^*})$ !

But from dual SVM we only have $\alpha_1^*, \ldots, \alpha_N^*$.

$$\omega^* = \sum_{i=1}^{N} \alpha_i^* y^i x^i$$

$$b^* = \frac{1}{2} \left[ \min_{i: y^i = +1} \omega^{*T} x^i - \max_{i: y^i = -1} \omega^{*T} x^i \right]$$

# Kernel SVM :

Basis function expansion :

$$\varphi(x) = \begin{pmatrix} \varphi_1(x) \\ \varphi_2(x) \end{pmatrix} = \begin{pmatrix} x_1^2 \\ x_2^2 \end{pmatrix}$$
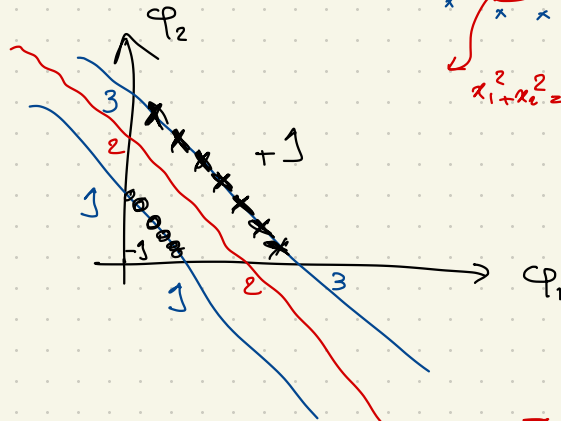
+1 $x$ : $1 x_1^2 + 1 x_2^2 = 3$ $\qquad \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix}$

$\qquad 1\varphi_1 + 1\varphi_2 = 3$

-1 o : $1 x_1^2 + 1 x_2^2 = 1$

$\qquad 1\varphi_1 + 1\varphi_2 = 1$



$\omega^T \varphi(x) + b = 0$

$1\varphi_1 + 1\varphi_2 = 2$

$1 x_1^2 + 1 x_2^2 = 2$

$x_1^2 + x_2^2 = 1$
$x_1^2 + x_2^2 = 3$
$x_1^2 + x_2^2 = 2$

## Vanilla SVM

$$\begin{cases} \min_{w,b} \frac{1}{2} \|w\|_2^2 \\ s.t. \quad y^i(\omega^T \varphi(x^i) + b) \geqslant 0 \quad \forall i = 1, \ldots, N \end{cases}$$

## Dual SVM

$$\begin{cases} \max_{\alpha_1, \ldots, \alpha_N} -\frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} y^i y^j (\varphi(x^i)^T \varphi(x^j)) \alpha_i \alpha_j + \sum_i \alpha_i \\ s.t. \quad \sum_{i=1}^{N} \alpha_i y^i = 0 \quad, \quad \alpha_i \geqslant 0 \quad \forall i = 1, \ldots, N \end{cases}$$

**Challenge:**

$$\varphi : x \longmapsto \varphi(x)$$
$$\mathbb{R}^d \longrightarrow \mathbb{R}^{d'}$$

$d' \gg d$

Vanilla SVM $O(d'^3)$ :(

Dual SVM ; have to compute $\varphi(x^i)^\top \varphi(x^j)$

$N^2$ computations of $\nearrow$

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \quad , \quad \varphi(x) = \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_1 x_3 \\ \vdots \\ x_1 x_d \\ x_2^2 \\ x_2 x_3 \\ \vdots \\ x_2 x_d \\ \vdots \\ x_d^2 \end{pmatrix} = \begin{pmatrix} \vdots \\ x_\ell x_{\ell'} \\ \vdots \end{pmatrix}$$

$x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$ is $d \times 1$

$\varphi(x) \simeq d^2 \times 1$

$\ell, \ell' = 1, \ldots, d$

$d \times d$

$x_1 \times x_1 = x_1^2$

$x_1 \times x_2 = x_1 x_2$

$\vdots$

$$\hookrightarrow \overline{\varphi(x^i)^\top \varphi(x^j)}$$

$$\underset{d^2 \times 1}{\begin{pmatrix} \\ \\ \end{pmatrix}}^\top \underset{d^2 \times 1}{\begin{pmatrix} \\ \\ \end{pmatrix}} = \langle \underset{1 \times d^2}{\begin{pmatrix} \\ \end{pmatrix}} \underset{d^2 \times 1}{\begin{pmatrix} \\ \end{pmatrix}} \rangle$$

$$O(1 \times d^2 \times 1) = O(d^2)$$

Computation of $\left\{ \varphi(x^i)^\top \varphi(x^j) \right\}_{i,j=1,\ldots,N}$ : $O(N^2 d^2)$ cost :(

$$\varphi(x) = \begin{pmatrix} x_1^3 \\ x_1^2 x_2 \\ x_1^2 x_3 \\ \vdots \\ x_1^2 x_d \\ x_1 x_2 x_3 \\ x_1 x_2 x_4 \\ \vdots \\ x_d^3 \end{pmatrix} = \begin{pmatrix} \vdots \\ x_\ell \, x_{\ell'} \, x_{\ell''} \\ \end{pmatrix}$$

$$\ell, \ell', \ell'' = 1, \ldots, d$$

$$x_1 \times x_1 \times x_1 = x_1^3 \qquad \ell = \ell' = \ell'' = 1$$

$$x_1 \times x_2 \times x_5 \longrightarrow 1+1=2 \qquad \ell = 1, \; \ell' = 2, \; \ell'' = 5$$

$$\underbrace{\phantom{x}}_{d} \quad \underbrace{\phantom{x}}_{d} \quad \underbrace{\phantom{x}}_{d}$$

$$\simeq d^3 \times 1$$

$$\left\{ \varphi(x^i)^T \varphi(x^j) \right\} \underset{i,j=1,\ldots,N}{\longrightarrow} \quad O(N^2 d^3) \;\; \ddot{\frown}$$

$$\varphi(x) = \begin{pmatrix} x_1^n \\ x_1^{n-1} x_2 \\ \vdots \\ x_1^{n-1} x_d \\ x_1^{n-2} x_2 x_3 \\ \}_n \\ x_d^n \end{pmatrix} \simeq d^n \times 1$$

$$\longrightarrow \left\{ \varphi(x^i)^T \varphi(x^j) \right\} \underset{i,j=1,\ldots,N}{\longrightarrow} \quad O(N^2 d^n) \;\; \ddot{\frown}$$

Key question: Can we compute $\varphi(x^i)^\top \varphi(x^j)$ implicitly, without explicitly computing $\varphi(x^i)$ ti and then taking their inner product? **Yes!**

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}, \qquad \varphi(x) = \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ \vdots \\ x_1 x_d \\ \vdots \\ x_d^2 \end{pmatrix} = \begin{pmatrix} \vdots \\ x_\ell x_{\ell'} \\ \vdots \end{pmatrix} \qquad \ell, \ell' = 1, 2, \ldots, d$$

$$\varphi(x^i)^\top \varphi(x^j) \quad \Rightarrow \quad \varphi(x)^\top \varphi(z)$$
$$\underbrace{\phantom{\varphi(x^i)}}_{\text{call it } x} \quad \underbrace{\phantom{\varphi(x^j)}}_{\text{call it } z}$$

$$\varphi(x)^\top \varphi(z) = \begin{pmatrix} \vdots \\ x_\ell x_{\ell'} \\ \vdots \end{pmatrix}^\top \begin{pmatrix} \vdots \\ z_\ell z_{\ell'} \\ \vdots \end{pmatrix} = \sum_{\ell=1}^{d} \sum_{\ell=1}^{d} x_\ell x_{\ell'} z_\ell z_{\ell'} = \sum_{\ell=1}^{d} \sum_{\ell'=1}^{d} (x_\ell z_\ell)(x_{\ell'} z_{\ell'})$$

$$\boxed{O(d^2)}$$

$$x_\ell x_{\ell'} \longleftrightarrow z_\ell z_{\ell'}$$

$$= \underbrace{\sum_{\ell=1}^{d} x_\ell z_\ell}_{x^\top z} \underbrace{\sum_{\ell'=1}^{d} x_{\ell'} z_{\ell'}}_{x^\top z} = x^\top z \times x^\top z = \left( x^\top z \right)^2$$

$$\underbrace{(\;)^\top_{d \times 1} (\;)_{d \times 1}}_{\phantom{x}} = \boxed{O(d)}$$

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \implies \varphi(x) = \begin{pmatrix} x_1^n \\ x_1^{n-1} x_2 \\ \vdots \\ x_d^n \end{pmatrix}_{d^n \times 1}$$

$$\boxed{\begin{array}{c} \varphi(x)^T \varphi(z) \longrightarrow O(d^n) \\ \| \\ (x^T z)^n \longrightarrow O(d) \end{array}}$$

Kernel SVM

$$\begin{cases} \underset{\alpha_1, \dots, \alpha_N}{\max} \; -\frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} y^i y^j \left( \overbrace{\cancel{\varphi(x^i)^T \varphi(x^j)}}^{K(x^i, x^j)} \right) \alpha_i \alpha_j + \sum_{i=1}^{N} \alpha_i \\ \text{s.t.} \quad \sum_{i=1}^{N} \alpha_i y^i = 0 \;,\quad \alpha_i \geq 0 \quad \forall i = 1, \dots, N \end{cases}$$

$$K(x, z) = \varphi(x)^T \varphi(z) \quad \begin{array}{c} \nearrow (2) \; (x^T z)^2 \\ \searrow (n) \; (x^T z)^n \end{array}$$

↙ kernel between $x$ & $z$

$$K(x^i, x^j) = x^{i^T} x^j$$
linear kernel

**Choice of Kernel:**      use cross-validation or hold out data to pick the best kernel!

✳   $K(x, z) = (x^T z)^n$    n-th degree monomial    $\longrightarrow$    $\varphi(x) = \begin{pmatrix} x_1^n \\ x_1^{n-1} x_2 \\ \vdots \\ x_d^n \end{pmatrix}$    associated BFE for this kernel

✳   $K(x, z) = \left(x^T z + c\right)^n$    monomials upto degree $n$    $\longrightarrow$    $\varphi(x) = \begin{pmatrix} x_1^n \\ \vdots \\ x_d^n \end{pmatrix}$ } degree n $\\ \begin{pmatrix} x_1^{n-1} \\ \vdots \\ x_d^{n-1} \end{pmatrix}$ } degree n-1 $\\ \vdots \\ \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$ } degree 1 $\\ 1$ } degree 0

$c \geq 0$

$\varphi(x)^T \varphi(z)$

✳   $K(x, z) = e^{-\|x - z\|_2^2 / 2\sigma^2}$   $\longrightarrow$ hyperpar.

Gaussian / RBF Kernel