

Imperfect separation

Soft-margin SVM

Definition: Hinge loss: cost of a violation.

$$\xi^{(i)} = \max \left(0, \underbrace{1 - (\vec{x}^{(i)} \cdot \vec{w} + b) y^{(i)}}_{\text{size of violation}} \right)$$

Soft-margin SVM
optimisation:

$$\min_{\vec{w}} \frac{1}{2} \|\vec{w}\|^2 + \lambda \sum_{i=1}^N \xi^{(i)}$$

$$\text{subject to } 1 - y^{(i)} (\vec{w} \cdot \vec{x}^{(i)} + b) \leq \xi^{(i)}$$
$$\xi^{(i)} \geq 0 \quad i=1, \dots, N$$

★ Remember that $y^{(i)} (\vec{x}^{(i)} \cdot \vec{w} + b) \geq 1$
for all points when
there is no violation

Beyond linearity:

Kernelised SVM

based on the 'kernel trick'

Standard SVM (linear)

Linear classifier (\vec{w}^*, b)

$$\vec{w}^* = \sum_{i=1}^N \alpha_i y^{(i)} \vec{x}^{(i)}$$

$\alpha_i = 0$ $\forall i$ not a support vector

Assume

there are only two support vectors (generic case):

Then only $0 \neq \alpha_+$ and $\alpha_- \neq 0$: $\vec{w}^* = \alpha_+ \vec{x}_+ - \alpha_- \vec{x}_-$

$$b = 1 - \vec{x}_+ \cdot [\alpha_+ \vec{x}_+ - \alpha_- \vec{x}_-]$$

$$= 1 - \alpha_+ \boxed{\vec{x}_+ \cdot \vec{x}_+} + \alpha_- \boxed{\vec{x}_+ \cdot \vec{x}_-}$$

These define
our model

Obtained from
maximising of
the dual:

$$L(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} \boxed{(\vec{x}^{(i)} \cdot \vec{x}^{(j)})}$$

$$\text{Given } \vec{x}^{in}, \quad \begin{cases} \boxed{\vec{x}^{in} \cdot \vec{w}^*} + b \geq 0 \Rightarrow \hat{y} = +1 \\ \boxed{\vec{x}^{in} \cdot \vec{w}^*} + b \leq 0 \Rightarrow \hat{y} = -1 \end{cases}$$

Kernel functions

Given $\vec{x}, \vec{y} \in \mathbb{R}^d$, consider $\vec{\phi}(\vec{x}) = \vec{z} \in \mathbb{R}^D$
potentially in $D \gg d$

In some cases, for the right ϕ we have nice properties:

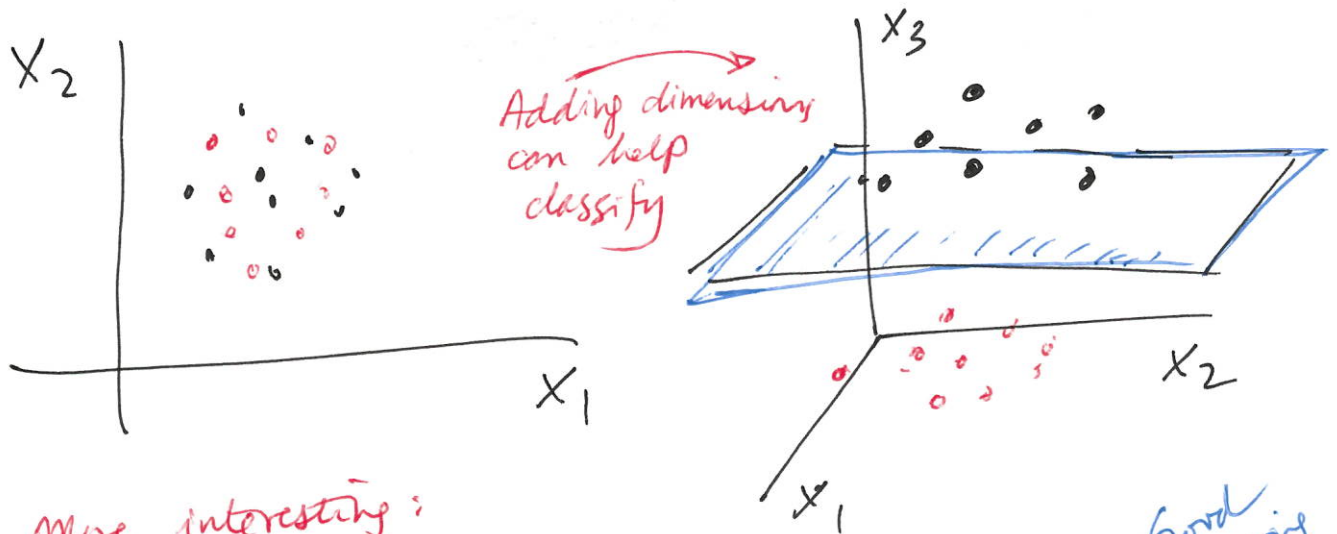
e.g. $\vec{x} = (x_1, x_2) \in \mathbb{R}^2$
 $\vec{\phi}(\vec{x}) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \end{pmatrix} \in \mathbb{R}^3$

$$\vec{\phi}(\vec{x}) \cdot \vec{\phi}(\vec{y}) = x_1^2 y_1^2 + x_2^2 y_2^2 + 2 x_1 x_2 y_1 y_2$$

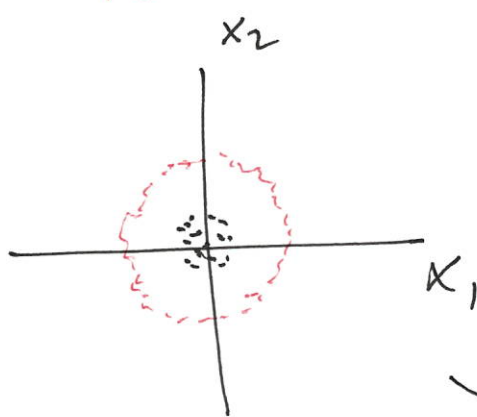
$$\vec{x} \cdot \vec{y} = x_1 y_1 + x_2 y_2$$

$$\begin{aligned} g(\underbrace{\vec{x} \cdot \vec{y}}_{\substack{\uparrow \\ 2D \cdot 2D}}}) &= (\vec{x} \cdot \vec{y})^2 = x_1^2 y_1^2 + x_2^2 y_2^2 + 2 x_1 y_1 x_2 y_2 \\ &= \underbrace{\vec{\phi}(\vec{x})}_{3D} \cdot \underbrace{\vec{\phi}(\vec{y})}_{3D} \end{aligned}$$

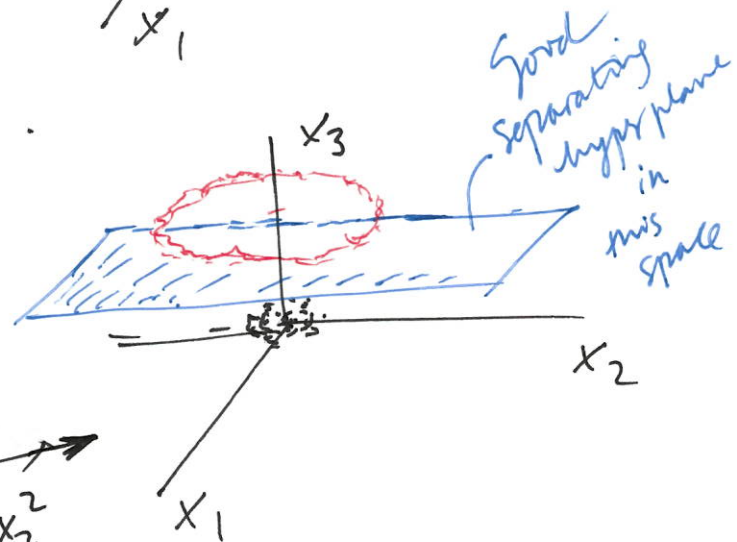
Idea: Nonlinearity and high-dimension can help with classification.



More interesting:



No separating hyperplane



Add nonlinear coordinate

There exists a good separating hyperplane in this space:

$$(x_1, x_2, x_1^2 + x_2^2)$$

In general, a kernel function maps a pair of vectors in the lower dimensional space (\mathbb{R}^d) to give the dot product of transformed (nonlinear) vectors in the high-dimensional space (\mathbb{R}^D)

$$K(\vec{x}, \vec{y}) = \vec{\phi}(\vec{x}) \cdot \vec{\phi}(\vec{y}) \quad \text{where } \vec{x}, \vec{y} \in \mathbb{R}^d \text{ and } \vec{\phi}(\vec{x}), \vec{\phi}(\vec{y}) \in \mathbb{R}^D$$

K is a kernel with associated $\vec{\phi}$ iff K is positive semi-definite

Several important kernels:

$$(1) \quad K(\vec{x}, \vec{y}) = (\vec{x} \cdot \vec{y} + 1)^n \quad \text{Polynomial}$$

$$(2) \quad K(\vec{x}, \vec{y}) = e^{-\|\vec{x} - \vec{y}\|^2 / \sigma} \quad \text{Gaussian Radial basis}$$

$$(3) \quad K(\vec{x}, \vec{y}) = \tanh(\beta(\vec{x} \cdot \vec{y}) + c) \quad \text{Sigmoid.}$$

Choose a kernel and compute the corresponding kernelised SVM.

More formally:

Choose
a kernel
 $k(\bar{x}, \bar{y})$

Kernelised SVM
optimization:

$$\bar{z} = \phi(\bar{x})$$

associated with
 $\underline{k(\bar{x}, \bar{y})}$

$$\min_{\bar{w}_z} \frac{1}{2} \|\bar{w}_z\|^2$$

$$\text{subject to } 1 - y^{(i)} (\bar{w}_z \cdot \bar{z}^{(i)} + b) \leq 0$$

$i=1, \dots, N$

Looking back at the SVM expressions:

$$\bar{w}_z \cdot \bar{z}^{(i)} = \phi(\bar{w}) \cdot \phi(\bar{x}^{(i)}) = k(\bar{w}, \bar{x}^{(i)})$$

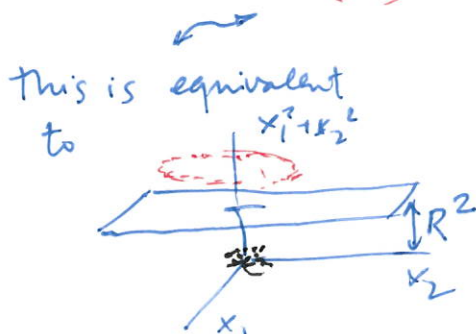
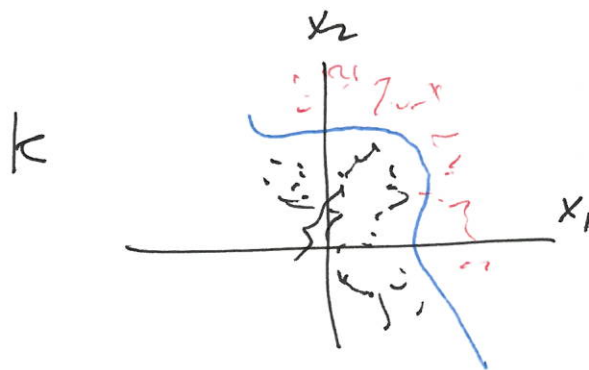
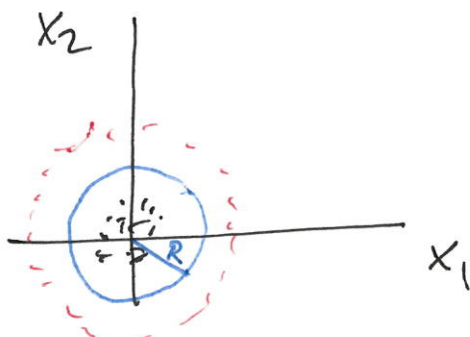
In the Lagrangian:

$$\bar{z}^{(i)} \cdot \bar{z}^{(j)} = k(\bar{x}^{(i)}, \bar{x}^{(j)})$$

Decision:

Given \bar{x}^{in}

$$\begin{cases} k(\bar{x}^{in}, \bar{w}) + b > 0 \Rightarrow \hat{y} = +1 \\ k(\bar{x}^{in}, \bar{w}) + b < 0 \Rightarrow \hat{y} = -1 \end{cases}$$



The decision boundary is
nonlinear in the \bar{x} space