

Introduction to Machine Learning

Lecture 8 - SVM and The Kernel Trick

1 Linear Separable SVM (Support Vector Machine)

1.1 Introduction

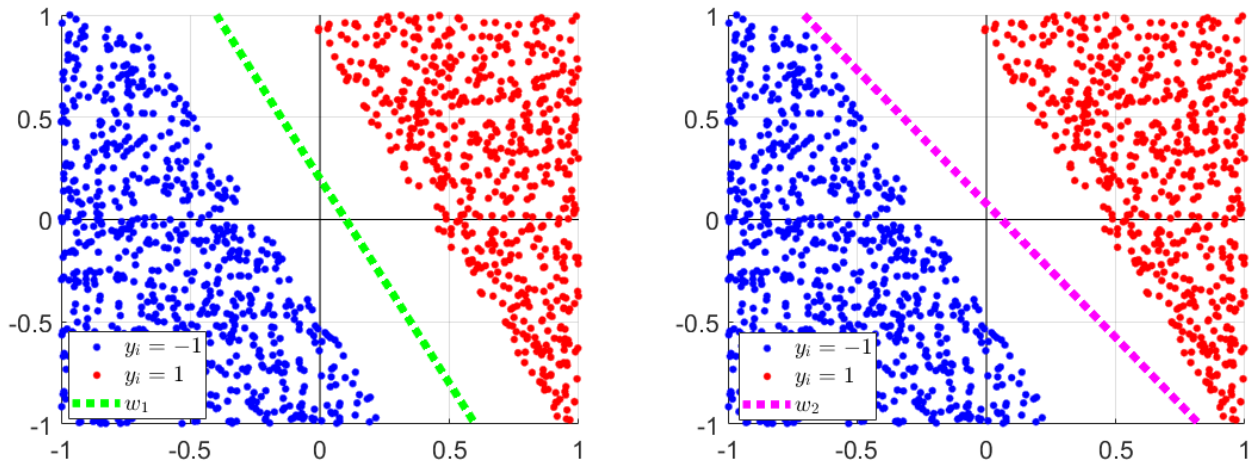
Consider some **linear separable** training set $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ such that:

$$y_i \in \{-1, 1\}$$

Since that data is linear separable, there may be several linear classifiers $\{\mathbf{w}, b\}$ such that:

$$y_i = \text{sign}(\mathbf{w}^T \mathbf{x}_i - b), \quad \forall i$$

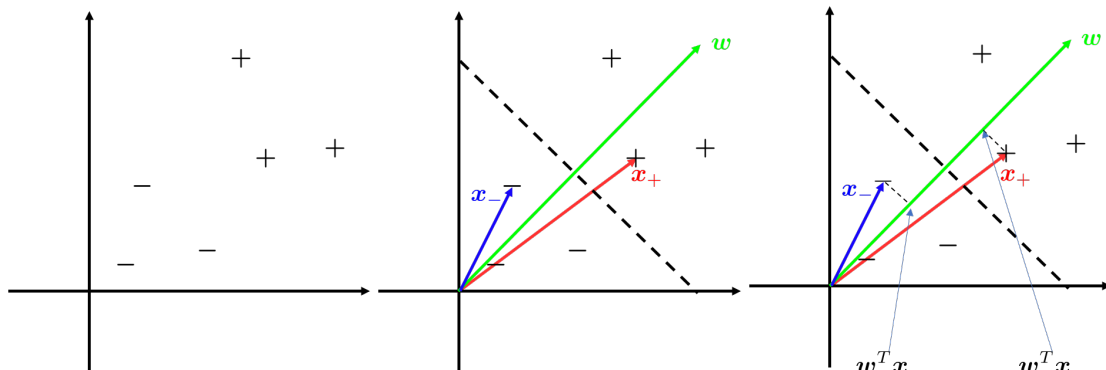
For example, both \mathbf{w}_1 and \mathbf{w}_2 provide perfect classification (so which one is better? w_1 or w_2):



The Support Vector Machine (SVM) classifier search for the largest margin between the boundary and the training set. In this figure, w_1 provides a larger margin than w_2 .

1.2 Derivation

Consider some \mathbf{w} and b :



Since we want large separation we require:

$$\begin{cases} \mathbf{w}^T \mathbf{x}_+ - b \geq 1 \\ \mathbf{w}^T \mathbf{x}_- - b \leq -1 \end{cases}$$

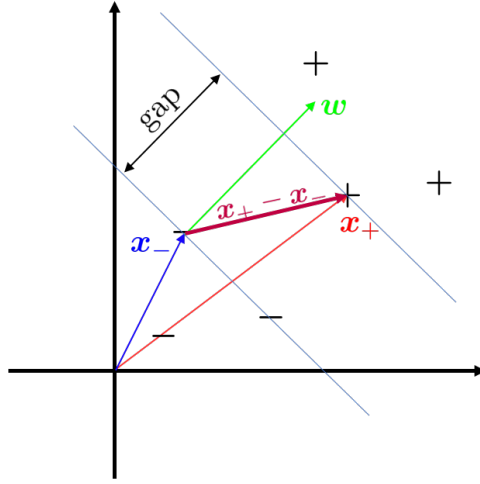
Alternatively, since $y_i = \pm 1$ (1 for the positive class and -1 for the negative class), we can write both equation as:

$$\boxed{y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1}$$

Notes:

1. The value 1 was chosen arbitrarily since the length of \mathbf{w} is not given (yet).
2. The points which are in the gutter (closest to the boundary) result with:

$$\begin{aligned} y_i (\mathbf{w}^T \mathbf{x}_i - b) &= 1, \quad \text{for } \mathbf{x}_i \text{ in the gutter} \\ \Rightarrow \mathbf{w}^T \mathbf{x}_i &= b + y_i \end{aligned}$$



The gap is given by (projection of $(\mathbf{x}_+ - \mathbf{x}_-)$ to the normalized \mathbf{w}):

$$\begin{aligned} \text{gap} &= \left\langle \frac{\mathbf{w}}{\|\mathbf{w}\|}, \mathbf{x}_+ - \mathbf{x}_- \right\rangle, \quad \mathbf{x}_+ \text{ and } \mathbf{x}_- \text{ are in the gutter} \\ &= \frac{\mathbf{w}^T}{\|\mathbf{w}\|} (\mathbf{x}_+ - \mathbf{x}_-) \\ &= \frac{1}{\|\mathbf{w}\|} (\mathbf{w}^T \mathbf{x}_+ - \mathbf{w}^T \mathbf{x}_-) \\ &= \frac{1}{\|\mathbf{w}\|} (b + 1 - b + 1) \\ &= \frac{2}{\|\mathbf{w}\|} \end{aligned}$$

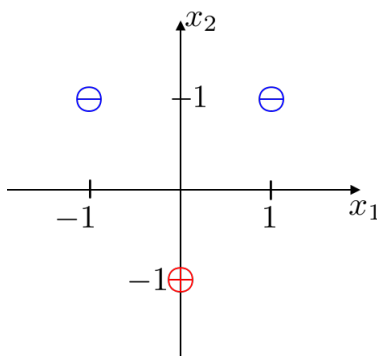
Thus, to maximize the gap one should minimize $\|\mathbf{w}\|$.

Equivalently, one can minimize $\frac{1}{2} \|\mathbf{w}\|^2$:

$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \\ y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 \quad \forall i \end{cases}$$

1.3 Example I

Consider the following training set:

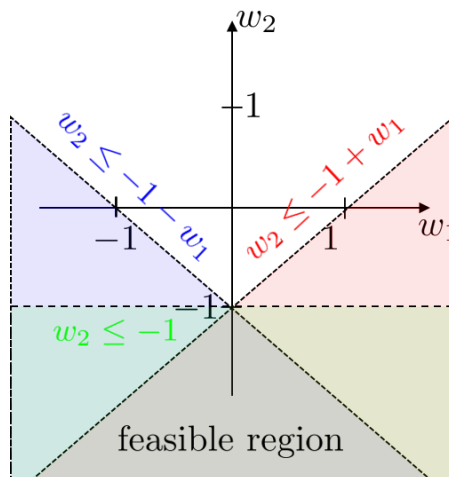


For simplicity, let us assume $b = 0$.

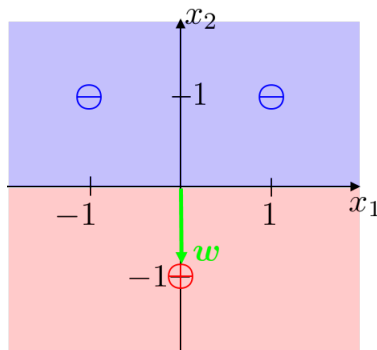
The SVM problem is given by:

$$\begin{cases} \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \\ (-1) \cdot \begin{bmatrix} 1 & 1 \end{bmatrix} \mathbf{w} \geq 1 \\ (-1) \cdot \begin{bmatrix} -1 & 1 \end{bmatrix} \mathbf{w} \geq 1 \\ (+1) \cdot \begin{bmatrix} 0 & -1 \end{bmatrix} \mathbf{w} \geq 1 \end{cases} = \begin{cases} \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \\ \begin{bmatrix} 1 & 1 \end{bmatrix} \mathbf{w} \leq -1 \\ \begin{bmatrix} -1 & 1 \end{bmatrix} \mathbf{w} \leq -1 \\ \begin{bmatrix} 0 & -1 \end{bmatrix} \mathbf{w} \geq 1 \end{cases} = \begin{cases} \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \\ w_2 \leq -1 - w_1 \\ w_2 \leq -1 + w_1 \\ w_2 \leq -1 \end{cases}$$

The feasible region is given by:



Notice that $\mathbf{w} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ is the closet point (in the feasible region) to the origin:



1.4 Example II

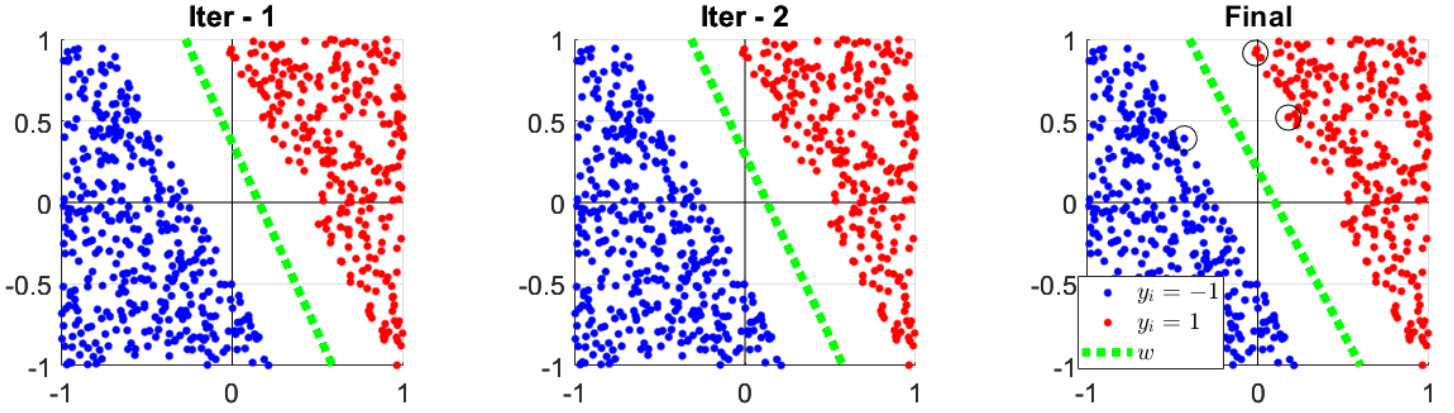
1.4.1 The linear separable case

The Lagrangian is given by:

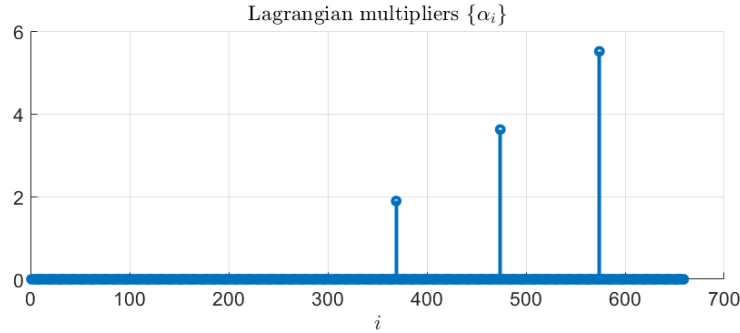
$$\mathcal{L}(\mathbf{w}, b) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1)$$

where $\{\alpha_i \geq 0\}_i$ are the Lagrangian multipliers.

The Augmented Lagrangian algorithm can solve the given (constrained) optimization problem:



The obtained Lagrangian multipliers $\{\alpha_i\}_i$ are:



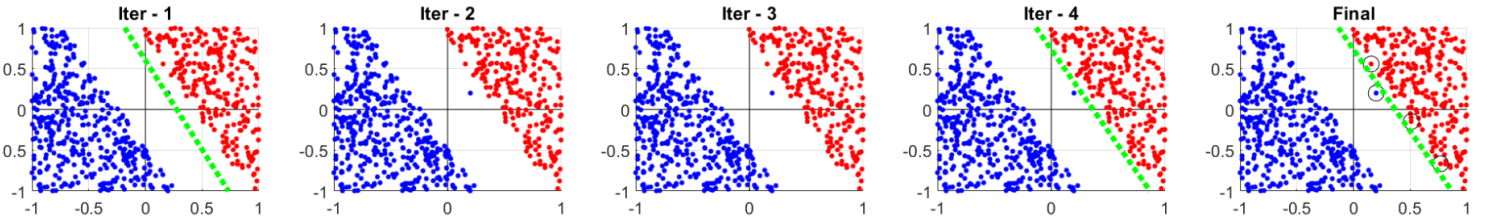
The three points correspond to the three $\alpha_i \neq 0$ are marked with black circles in the “Final” figure.

These three points are exactly in the gutter and they are called the support vectors.

The other points do not affect the optimization result (the value of \mathbf{w} and b).

1.4.2 The linear barely separable case

Notice how a single point can change the classifier (one blue point close to the red group):



For non-separable problems, the constraints cannot be satisfied.

Thus, a soft (relaxed) version of SVM should be used.

2 Soft SVM

Reminder: The SVM optimization problem is given by:

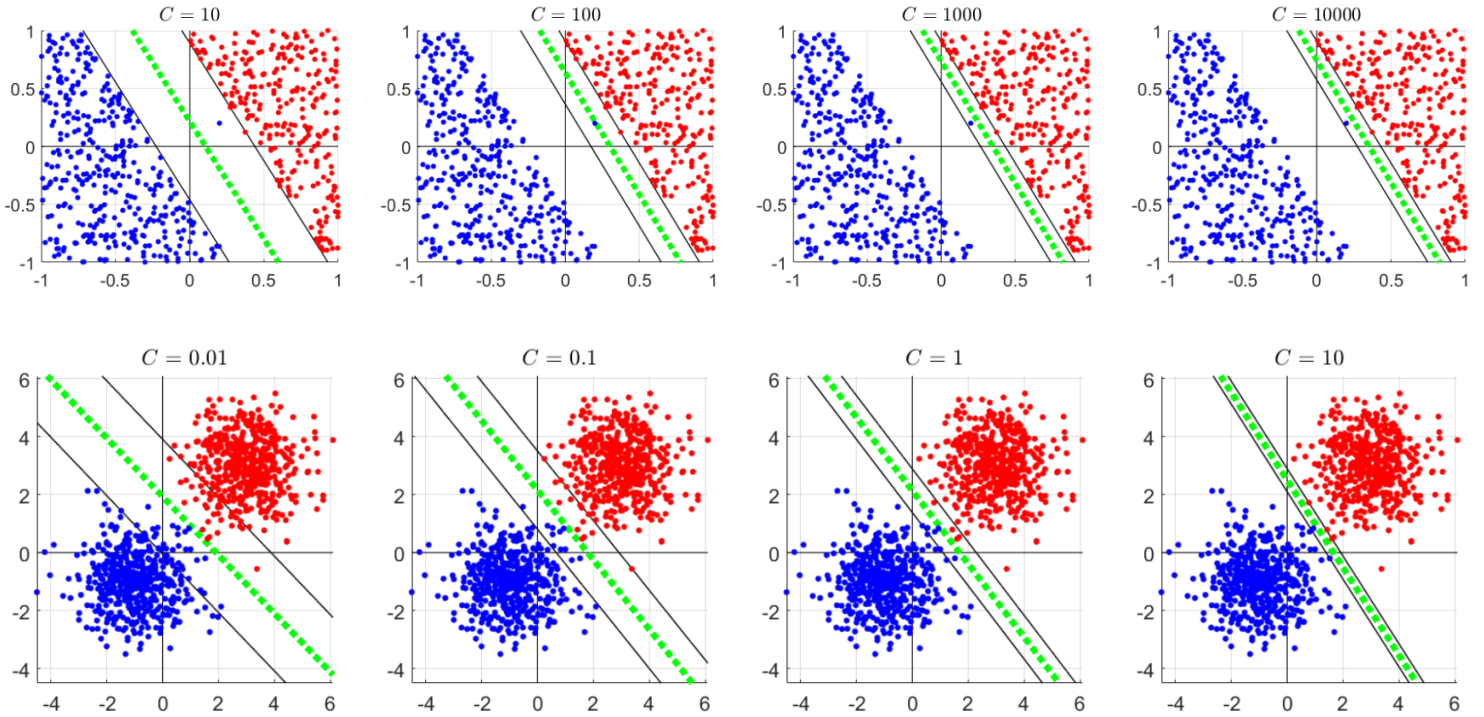
$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \\ y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1 \quad \forall i \end{cases}$$

Instead of strict constraints we can write a relaxed unconstrained version of SVM:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \max \{0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i - b)\}$$

1. The new term is a cost penalty for each point \mathbf{x}_i which does not satisfy the original constrain: $y_i (\mathbf{w}^T \mathbf{x}_i - b) \geq 1$.
2. The $C \geq 0$ is the cost factor (regularization factor).
3. For $C = 0$, there is no cost for violating the original constraints.
Hence, the solution will be $\mathbf{w} = 0$
4. For $C \rightarrow \infty$, the solution must not violate the constraints, and we obtained the original constrained problem.

Examples:



3 Dual Problem

3.1 Derivation

The Lagrangian is given by:

$$\mathcal{L}(\mathbf{w}, b, \{\alpha_i\}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1)$$

Comparing the gradient to zero:

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L} &= 0 \\ \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i &= 0 \\ \Rightarrow \boxed{\mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i} \end{aligned}$$

$$\begin{aligned} \nabla_b \mathcal{L} &= 0 \\ \sum_i \alpha_i y_i &= 0 \\ \Rightarrow \boxed{\sum_i \alpha_i y_i = 0} \end{aligned}$$

Plugging these relations back to \mathcal{L} results in:

$$\begin{aligned} \mathcal{L}(\mathbf{w}^*, b, \{\alpha_i\}) &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i (y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1) \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_i \alpha_i y_i \mathbf{w}^T \mathbf{x}_i + b \underbrace{\sum_i \alpha_i y_i}_{=0} + \sum_i \alpha_i \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \underbrace{\left(\sum_i \alpha_i y_i \mathbf{x}_i^T \right)}_{=\mathbf{w}^T} \mathbf{w} + \sum_i \alpha_i \\ &= \sum_i \alpha_i - \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ &= \sum_i \alpha_i - \frac{1}{2} \left(\sum_i \alpha_i y_i \mathbf{x}_i \right)^T \left(\sum_j \alpha_j y_j \mathbf{x}_j \right) \\ &= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ &= \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \end{aligned}$$

3.2 Formulation

Overall, the dual problem can be written as:

$$\begin{cases} \max_{\{\alpha_i\}} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{s.t.} \\ \alpha_i \geq 0 \\ \sum_i \alpha_i y_i = 0 \end{cases} \quad \forall i$$

Remarks

1. The maximum value of the dual problem is the minimal value of the primal problem (not always true for general constrained problems).
2. The dual problem dependent only on the inner products $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ and not the vector $\{\mathbf{x}_i\}_i$ themselves.
3. The vector \mathbf{w} is obtained by:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

4. b can be obtained by using a point in the gutter:

$$y_i (\mathbf{w}^T \mathbf{x}_i - b) = 1, \quad \text{for } \mathbf{x}_i \text{ in the gutter}$$

$$\Rightarrow b = \mathbf{w}^T \mathbf{x}_i - y_i$$

5. The decision rule (of a new point \mathbf{x}_0) is given by:

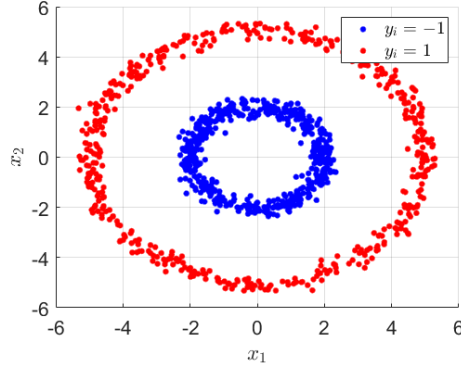
$$\begin{aligned} \hat{y}_0 &= \text{sign}(\mathbf{w}^T \mathbf{x}_0 - b) \\ &= \text{sign}\left(\sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_0 - b\right) \\ &= \text{sign}\left(\sum_i \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x}_0 \rangle - b\right) \end{aligned}$$

6. The decision rule dependent only on the inner product.

4 The Kernel Trick (Non-linear SVM)

4.1 Feature transform

Consider the following problem:



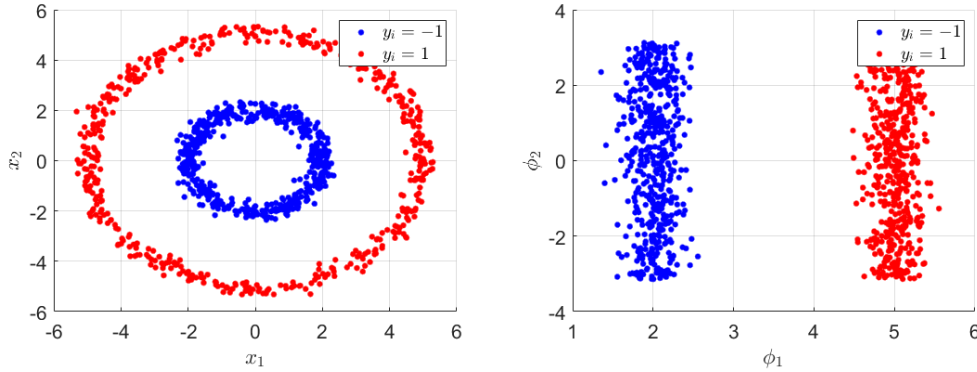
A linear SVM classifier is not able to separate these two clusters.

Reminder: The dual problem is given by:

$$\begin{cases} \max_{\{\alpha_i\}} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{s.t.} \\ \alpha_i \geq 0 \\ \sum_i \alpha_i y_i = 0 \end{cases} \quad \forall i$$

To solve the problem above, one can use the following feature transform (polar coordinates):

$$\phi(\mathbf{x}_i) = \begin{bmatrix} \|\mathbf{x}_i\| \\ \angle \mathbf{x}_i \end{bmatrix}$$



The dual optimization problem using the new features is given by:

$$\begin{cases} \max_{\{\alpha_i\}} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \\ \text{s.t.} \\ \alpha_i \geq 0 \\ \sum_i \alpha_i y_i = 0 \end{cases} \quad \forall i$$

with the decision rule:

$$\hat{y}_0 = \text{sign} \left(\sum_i \alpha_i y_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle - b \right)$$

Notice again that the algorithm needs to compute only the inner products $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$.
 Lets define the following function (**kernel**):

$$K(\mathbf{x}_i, \mathbf{x}_j) \triangleq \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

Hence, we can write the optimization problem as follows:

$$\begin{cases} \max_{\{\alpha_i\}} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t.} \\ \alpha_i \geq 0 \\ \sum_i \alpha_i y_i = 0 \end{cases} \quad \forall i$$

Definition 1. A **Kernel** function $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies:

1. Symmetry: $K(\mathbf{x}_i, \mathbf{x}_j) = K(\mathbf{x}_j, \mathbf{x}_i)$
2. K is positive definite, that is, the matrix $\mathbf{K}[i, j] = K(\mathbf{x}_i, \mathbf{x}_j)$ is SPD for any set $\{\mathbf{x}_i\}_i$

Theorem 2. Under some technical conditions, the kernel function K can be written as:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^M \phi_m(\mathbf{x}_i) \phi_m(\mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

where M might be infinity and $\{\phi_m\}_m$ are the basis of K .

4.1.1 Common Kernels

Polynomial Kernel Consider for example the following kernel:

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$$

where $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$.

This kernel can also be written as:

$$\begin{aligned} K(\mathbf{x}_i, \mathbf{x}_j) &= (1 + \mathbf{x}_i^T \mathbf{x}_j)^2 \\ &= 1 + 2\mathbf{x}_i^T \mathbf{x}_j + (\mathbf{x}_i^T \mathbf{x}_j)^2 \\ &= 1 + 2 \sum_{k=1}^d \mathbf{x}_i[k] \mathbf{x}_j[k] + \left(\sum_{k=1}^d \mathbf{x}_i[k] \mathbf{x}_j[k] \right) \left(\sum_{m=1}^d \mathbf{x}_i[m] \mathbf{x}_j[m] \right) \\ &= 1 + \sum_{k=1}^d \sqrt{2} \mathbf{x}_i[k] \cdot \sqrt{2} \mathbf{x}_j[k] + \sum_{k=1}^d \sum_{m=1}^d (\mathbf{x}_i[k] \mathbf{x}_i[m]) (\mathbf{x}_j[k] \mathbf{x}_j[m]) \\ &= \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \end{aligned}$$

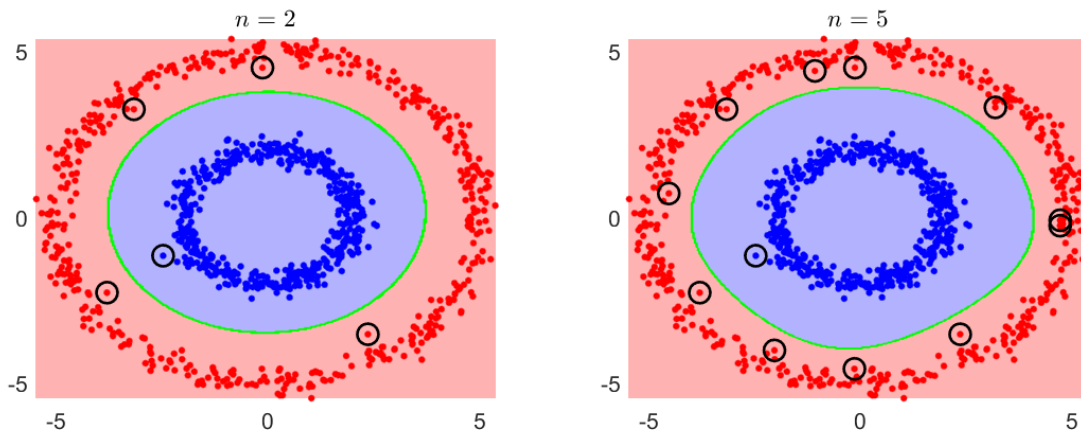
where:

$$\phi(\mathbf{x}) \triangleq \begin{bmatrix} 1 \\ \sqrt{2}\mathbf{x}[1] \\ \vdots \\ \sqrt{2}\mathbf{x}[d] \\ \mathbf{x}[1]\mathbf{x}[1] \\ \vdots \\ \mathbf{x}[1]\mathbf{x}[d] \\ \mathbf{x}[2]\mathbf{x}[1] \\ \vdots \\ \mathbf{x}[d]\mathbf{x}[d] \end{bmatrix}$$

One can also use higher order polynomials:

$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^n, \quad n \geq 1$$

Examples:



The black circles indicate the support vectors.

Gauss Kernel For $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right), \quad \sigma > 0$$

For simplicity let's assume $x_i, x_j \in \mathbb{R}$ (scalars).

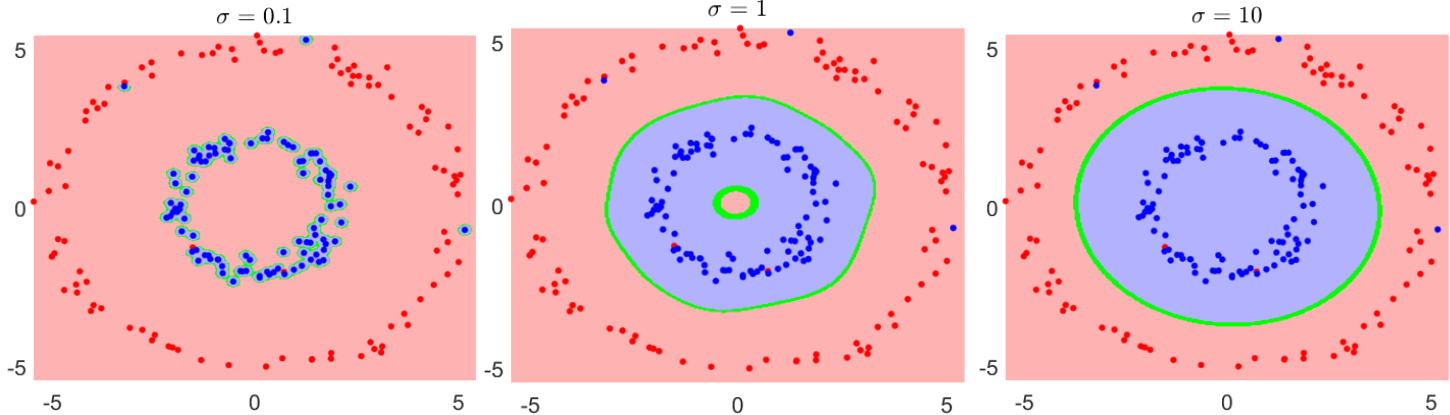
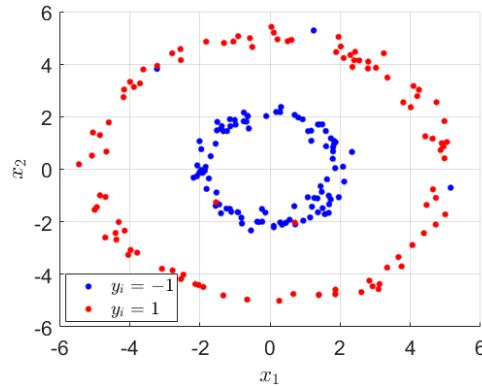
Using Taylor series:

$$\begin{aligned} K(x_i, x_j) &= \exp\left(-\frac{(x_i - x_j)^2}{2\sigma^2}\right) = \exp\left(-\frac{x_i^2 - 2x_i x_j + x_j^2}{2\sigma^2}\right) \\ &= \exp\left(-\frac{x_i^2}{2\sigma^2}\right) \exp\left(-\frac{x_j^2}{2\sigma^2}\right) \exp\left(\frac{x_i x_j}{\sigma^2}\right) \\ &= \exp\left(-\frac{x_i^2}{2\sigma^2}\right) \exp\left(-\frac{x_j^2}{2\sigma^2}\right) \sum_{n=0}^{\infty} \frac{\left(\frac{x_i x_j}{\sigma^2}\right)^n}{n!} \\ &= \exp\left(-\frac{x_i^2}{2\sigma^2}\right) \exp\left(-\frac{x_j^2}{2\sigma^2}\right) \sum_{n=0}^{\infty} \frac{x_i^n}{\sigma^n \sqrt{n!}} \cdot \frac{x_j^n}{\sigma^n \sqrt{n!}} \\ &= \langle \phi(x_i), \phi(x_j) \rangle \end{aligned}$$

where:

$$\phi(x) \triangleq \exp\left(-\frac{x^2}{2\sigma^2}\right) \cdot \begin{bmatrix} 1 \\ \frac{x}{\sigma} \\ \frac{x^2}{\sigma^2 \sqrt{2!}} \\ \frac{x^3}{\sigma^3 \sqrt{3!}} \\ \vdots \end{bmatrix} \in \mathbb{R}^{\infty}$$

Similar training set with a few mixed label points:



5 Extra

5.1 A simple example for the dual problem

Consider the following convex constrained optimization:

$$\begin{cases} \min f(x) \\ \text{s.t.} \\ g(x) \leq 0 \end{cases} = \begin{cases} \min x^2 \\ \text{s.t.} \\ x + 1 \leq 0 \end{cases}$$

One can easily notice that the optimal solution is $x^* = -1$ and $f(x^*) = 1$.

The Lagrangian is given by:

$$\mathcal{L}(x, \lambda) = x^2 + \lambda(x + 1)$$

$$\begin{aligned} \frac{\partial}{\partial x} \mathcal{L} &= 0 \\ 2x + \lambda &= 0 \\ \Rightarrow x^* &= -\frac{\lambda}{2} \end{aligned}$$

$$\begin{aligned} \Rightarrow \mathcal{L}(x^*, \lambda) &= \left(-\frac{\lambda}{2}\right)^2 + \lambda\left(-\frac{\lambda}{2} + 1\right) \\ &= \frac{\lambda^2}{4} + -\frac{\lambda^2}{2} + \lambda \\ &= -\frac{\lambda^2}{4} + \lambda \end{aligned}$$

The dual problem is given by:

$$\begin{cases} \max \mathcal{L}(x^*, \lambda) \\ \text{s.t.} \\ \lambda \geq 0 \end{cases}$$

The solution is given by:

$$\begin{aligned} \frac{\partial}{\partial \lambda} \mathcal{L} &= 0 \\ -\frac{\lambda}{2} + 1 &= 0 \\ \Rightarrow \lambda^* &= 2 \end{aligned}$$

$$\Rightarrow \mathcal{L}(x^*, \lambda^*) = -1 + 2 = 1$$

$$\Rightarrow x^*(\lambda^*) = -\frac{2}{2} = -1$$

Note that the solution of the dual problem is exactly the same solution of the primal problem:

$$\begin{cases} \mathcal{L}(x^*, \lambda^*) = f(x^*) = 1 \\ x^* = -1 \end{cases}$$

Thus, one can solve the dual problem instead of the primal problem.