

Lecture 18: 21 March, 2023

Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Data Mining and Machine Learning
January–April 2023

Soft margin optimization

$$\text{Minimize } \frac{\|w\|}{2} + \sum_{i=1}^N \xi_i^2$$

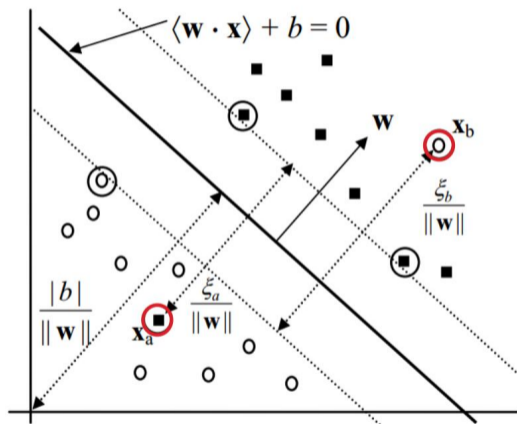
Subject to

$$\xi_i \geq 0$$

$$w \cdot x_i + b > 1 - \xi_i, \text{ if } y_i = 1$$

$$w \cdot x_i + b < -1 + \xi_i, \text{ if } y_i = -1$$

- Constraints include requirement that error terms are non-negative
- Again the objective function is quadratic



Dualization

Wolfe dual

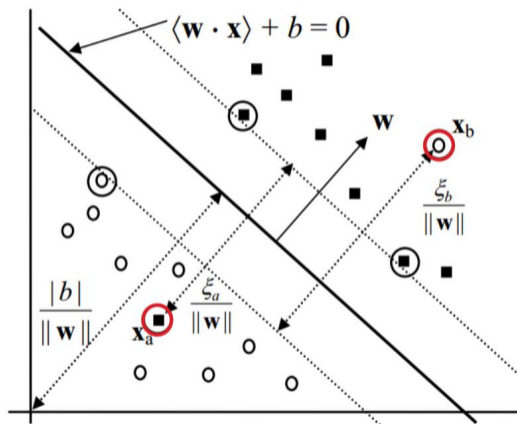
$$\text{Maximize } \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (x_i \cdot x_j)$$

Subject to

$$0 \leq \alpha_i \leq 1$$

$$\sum_i \alpha_i y_i = 0$$

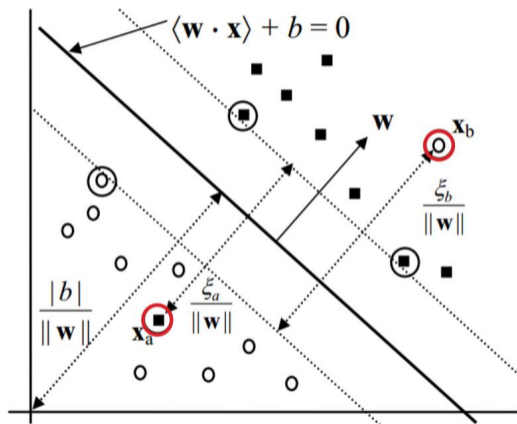
- α_i are Lagrange multipliers



Soft margin optimization

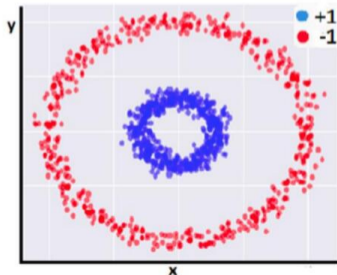
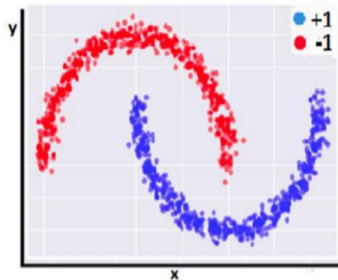
- Can again be solved using convex optimization theory
- Form of the solution turns out to be the same as the hard margin case
 - Expression in terms of Lagrange multipliers α_j
 - Only terms corresponding to support vectors are actively used

$$\text{sign} \left[\sum_{i \in \text{sv}} y_i \alpha_i (x_i \cdot z) + b \right]$$



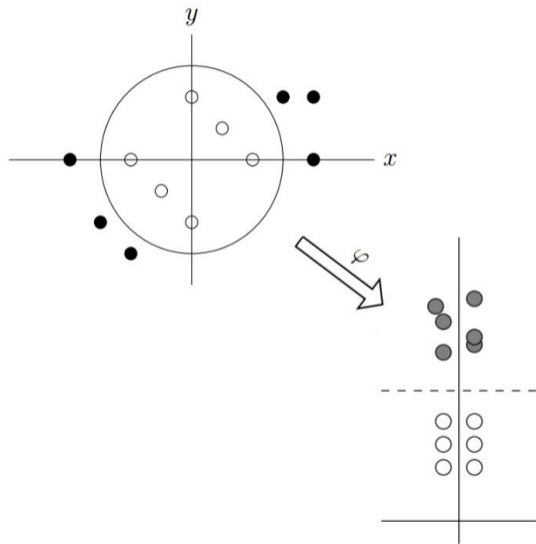
The non-linear case

- How do we deal with datasets where the separator is a complex shape?
- Geometrically transform the data
 - Typically, add dimensions
- For instance, if we can “lift” one class, we can find a planar separator between levels



Geometric transformation

- Consider two sets of points separated by a circle of radius 1
- Equation of circle is $x^2 + y^2 = 1$
- Points inside the circle, $x^2 + y^2 < 1$
- Points outside circle, $x^2 + y^2 > 1$
- Transformation
$$\varphi : (x, y) \mapsto (x, y, x^2 + y^2)$$
- Points inside circle lie below $z = 1$
- Point outside circle lifted above $z = 1$



SVM after transformation

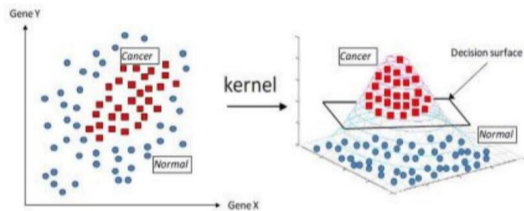
- SVM in original space

$$\text{sign} \left[\sum_{i \in Sv} y_i \alpha_i (x_i \cdot z) + b \right]$$

- After transformation

$$\text{sign} \left[\sum_{i \in Sv} y_i \alpha_i (\varphi(x_i) \cdot \varphi(z)) + b \right]$$

- All we need to know is how to compute dot products in transformed space



Dot products

- Consider the transformation

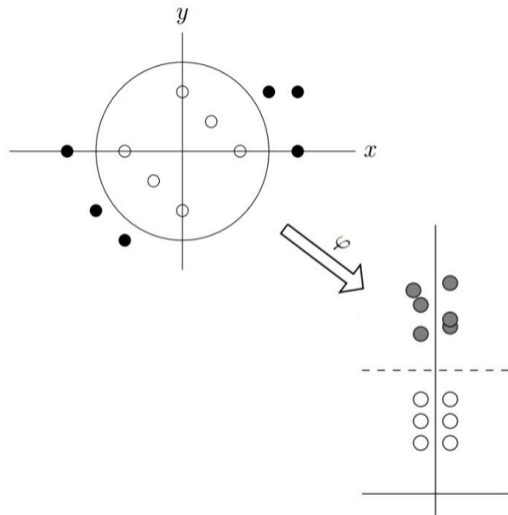
$$\varphi : (x_1, x_2) \mapsto (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

- Dot product in transformed space

$$\begin{aligned}\varphi(x) \cdot \varphi(z) &= 1 + 2x_1z_1 + 2x_2z_2 + x_1^2z_1^2 \\ &\quad + 2x_1x_2z_1z_2 + x_2^2z_2^2 \\ &= (1 + x_1z_1 + x_2z_2)^2\end{aligned}$$

- Transformed dot product can be expressed in terms of original inputs

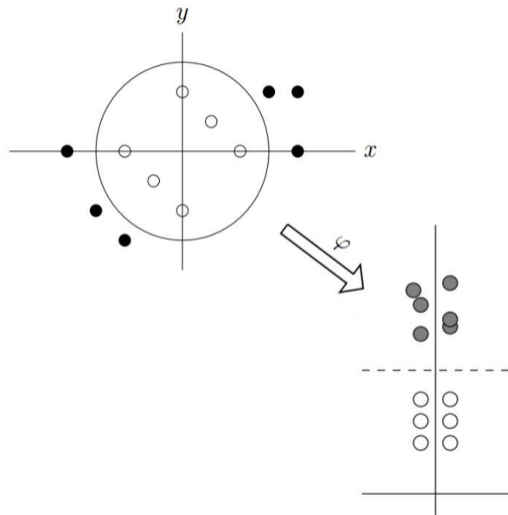
$$\varphi(x) \cdot \varphi(z) = K(x, z) = (1 + x_1z_1 + x_2z_2)^2$$



Kernels

- K is a **kernel** for transformation φ if $K(x, z) = \varphi(x) \cdot \varphi(z)$
- If we have a kernel, we don't need to explicitly compute transformed points
- All dot products can be computed implicitly using the kernel on original data points

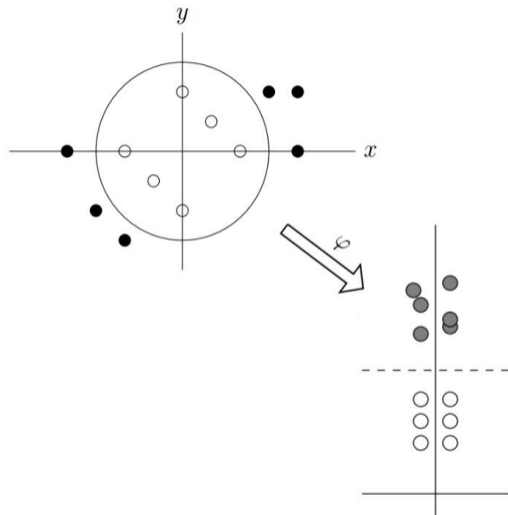
$$\text{sign} \left[\sum_{i \in sv} y_i \alpha_i (\varphi(x_i) \cdot \varphi(z)) + b \right]$$



Kernels

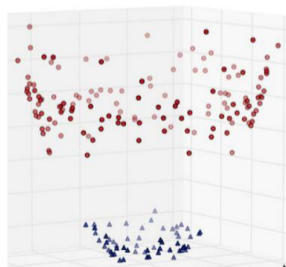
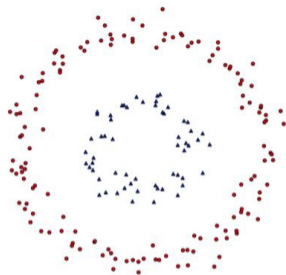
- K is a **kernel** for transformation φ if $K(x, z) = \varphi(x) \cdot \varphi(z)$
- If we have a kernel, we don't need to explicitly compute transformed points
- All dot products can be computed implicitly using the kernel on original data points

$$\text{sign} \left[\sum_{i \in S^+} y_i \alpha_i K(x_i, z) + b \right]$$



Kernels

- If we know K is a kernel for some transformation φ , we can blindly use K without even knowing what φ looks like!
- When is a function a valid kernel?
- Has been studied in mathematics — **Mercer's Theorem**
 - Criteria are non-constructive
- Can define sufficient conditions from linear algebra

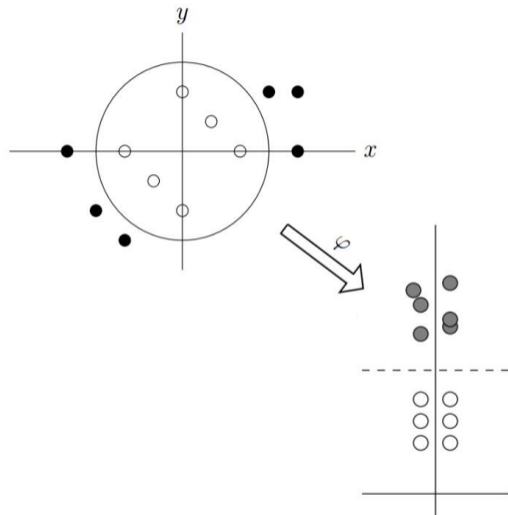


Kernels

- Kernel over training data x_1, x_2, \dots, x_N can be represented as a **gram matrix**

$$K = \begin{matrix} & x_1 & x_2 & \cdots & x_n \\ \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} & \left[\begin{array}{cccc} & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{array} \right] \end{matrix}$$

- Entries are values $K(x_i, x_j)$
- Gram matrix should be **positive semi-definite** for all x_1, x_2, \dots, x_N



Known kernels

- Fortunately, there are many known kernels
- Polynomial kernels
$$K(x, z) = (1 + x \cdot z)^k$$
- Any $K(x, z)$ representing a similarity measure
- Gaussian radial basis function — similarity based on inverse exponential distance

$$K(x, z) = e^{-c|x-z|^2}$$

