

Lecture 17: 24 March, 2026

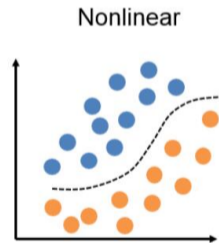
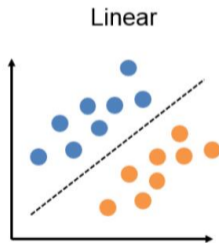
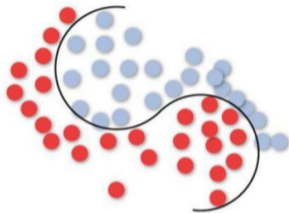
Madhavan Mukund

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

Data Mining and Machine Learning
January–April 2026

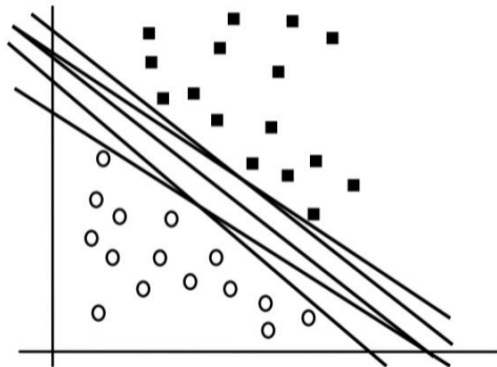
A geometric view of supervised learning

- Think of data as points in space
- Find a separating curve (surface)
- Separable case
 - Each class is a connected region
 - A single curve can separate them
- Simplest case — linearly separable data
- Dual of linear regression
 - Find a line that passes close to a set of points
 - Find a line that separates the two sets of points



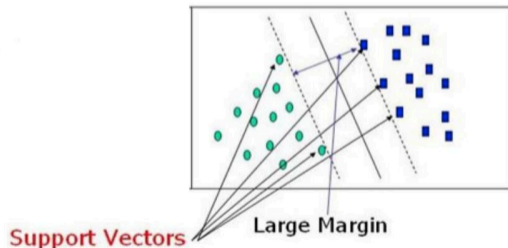
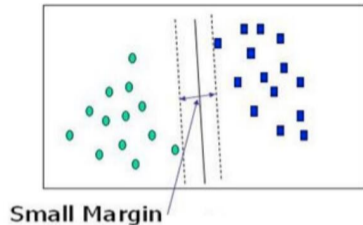
Linear separators

- Perceptron algorithm is a simple procedure to find a linear separator, if one exists
- Many lines are possible
 - Does the Perceptron algorithm find the best one?
 - What is a good notion of “cost” to optimize?



Margin

- Each separator defines a margin
 - Empty corridor separating the points
 - Separator is the centre line of the margin
- Wider margin makes for a more robust classifier
 - More gap between the classes
- Optimum classifier is one that maximizes the width of its margin
- Margin is defined by the training data points on the boundary
 - Support vectors



Finding a maximum margin classifier

Minimize $\frac{|w|}{2}$

Subject to

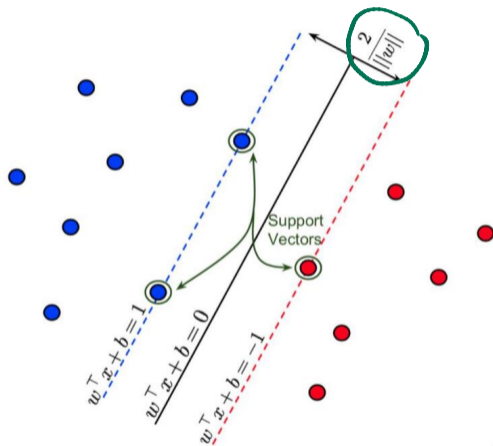
$$w_1 x_1^i + w_2 x_2^i + \dots + w_n x_n^i + b > 1, \text{ if } y_i = 1$$

$$w_1 x_1^i + w_2 x_2^i + \dots + w_n x_n^i + b < -1, \text{ if } y_i = -1$$

- The constraints are linear
- The objective function is not linear

$$|w| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$$

- This is a **quadratic optimization problem**, not linear programming



Solution to optimization problem

- Rewrite as dual

- Maximize

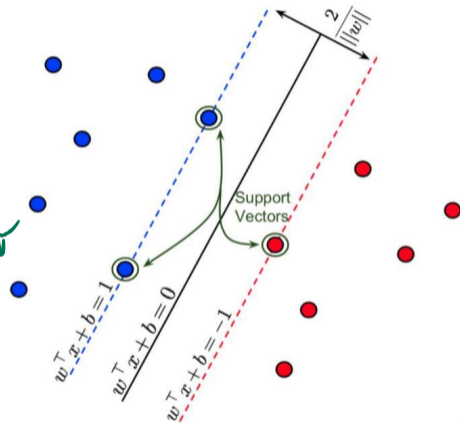
$$\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j x_i \cdot x_j$$

Subject to

$$\alpha_i \geq 0, i \in 1, 2, \dots, N$$

- Lagrange multipliers $\alpha_1, \alpha_2, \dots, \alpha_N$, one multiplier per training input

One per constraint



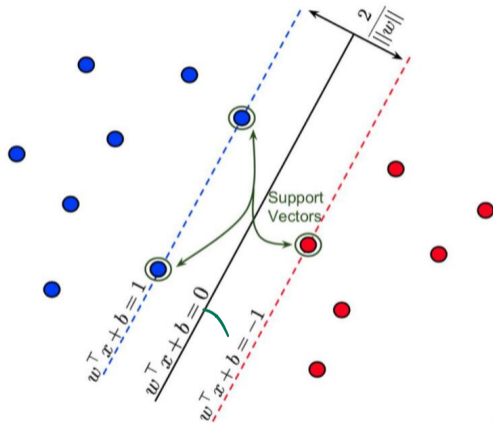
Support Vector Machine (SVM)

Final classifier is of the form

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

Note: In the original image, green wavy lines underline 'sign' and 'i ∈ sv', and a yellow wavy line underlines '(x_i · z)'. A green arrow points from the top of the sum to the dot product term.

- Solution depends only on support vectors
 - If we add more training data away from support vectors, separator does not change
- Solution uses dot product of support vectors with new point
 - Will be used later, in the non-linear case



Soft margin optimization

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

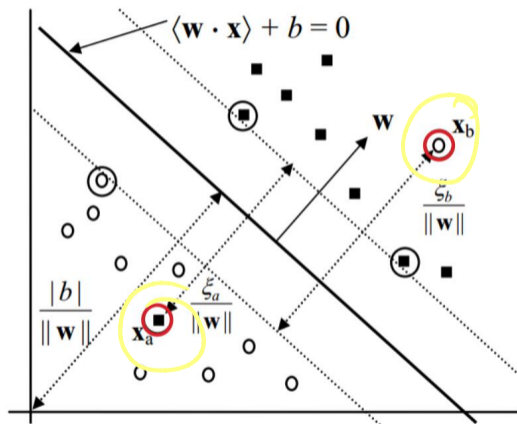
Subject to

$$\xi_i \geq 0$$

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

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

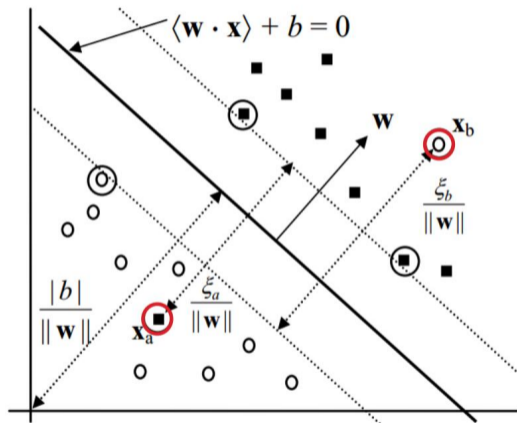
- Constraints include requirement that error terms are non-negative
- Again the objective function is quadratic



Soft margin optimization

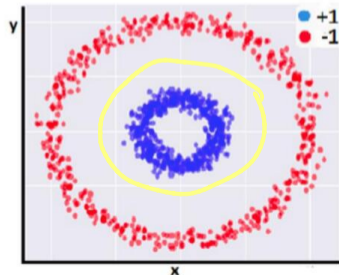
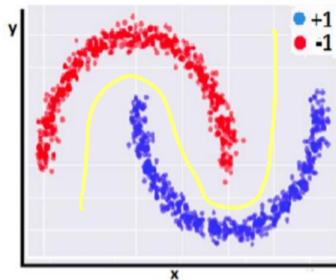
- Can again be solved using the dual
- Form of the solution turns out to be the same as the hard margin case
 - Expression in terms of Lagrange multipliers α_j
 - Only terms corresponding to support vectors are actively used

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



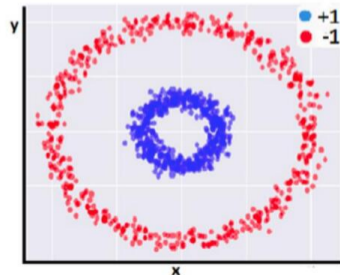
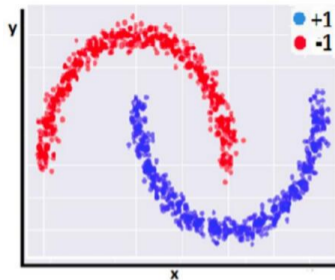
The non-linear case

- How do we deal with datasets where the separator is a complex shape?



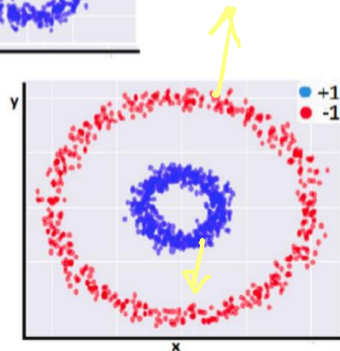
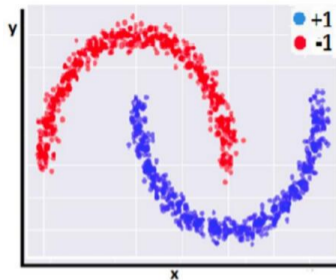
The non-linear case

- How do we deal with datasets where the separator is a complex shape?
- Geometrically transform the data
 - Typically, add dimensions



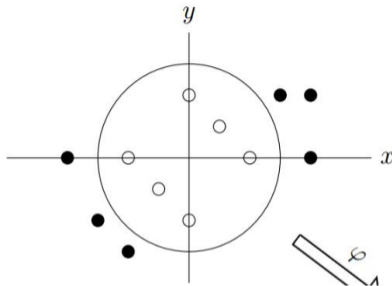
The non-linear case

- 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



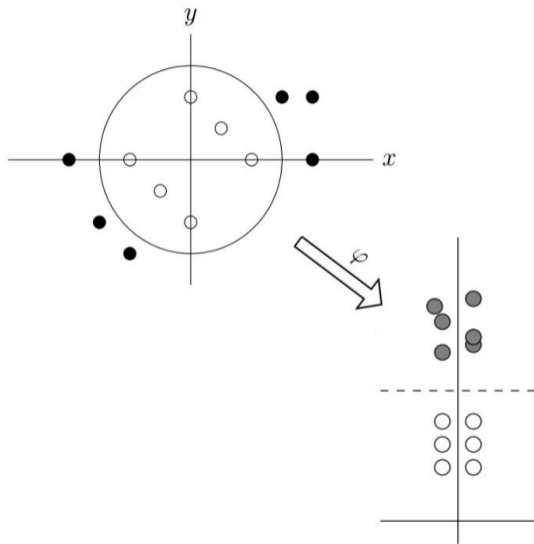
Geometric transformation

- Consider two sets of points separated by a circle of radius 1
- Equation of circle is $x^2 + y^2 = 1$



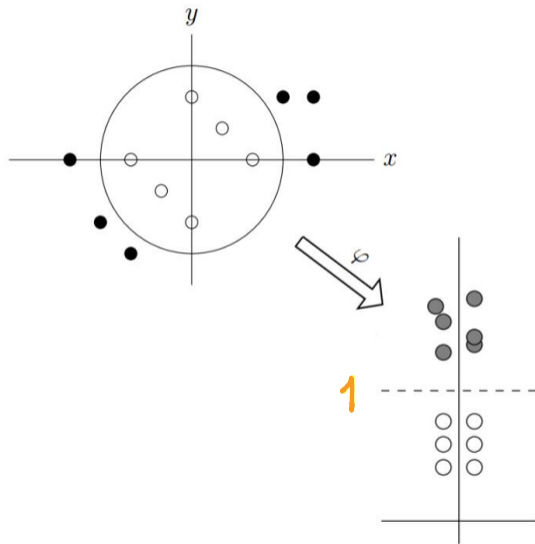
Geometric transformation

- 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, $x^2 + y^2 > 1$



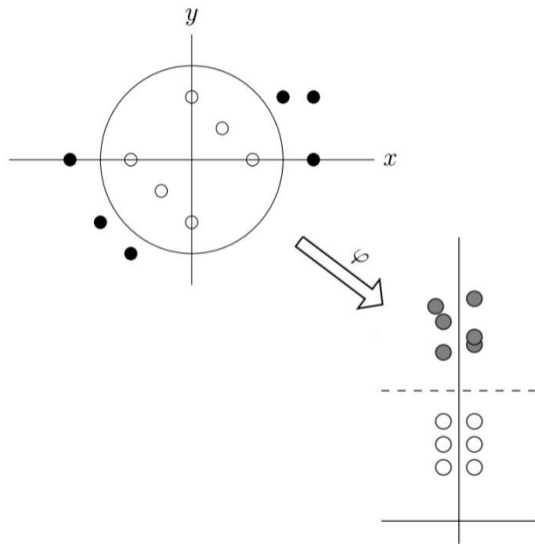
Geometric transformation

- 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
$$\varphi : (x, y) \mapsto (x, y, x^2 + y^2)$$



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 $z = 1$

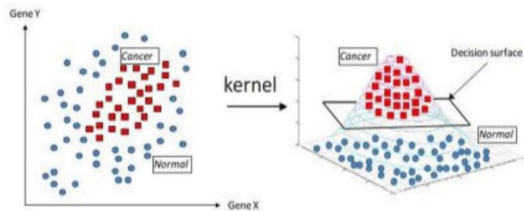


SVM after transformation

- SVM in original space

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

$$x \mapsto \varphi(x)$$



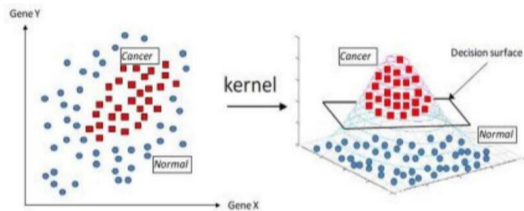
SVM after transformation

- SVM in original space

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

- After transformation

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



SVM after transformation

- SVM in original space

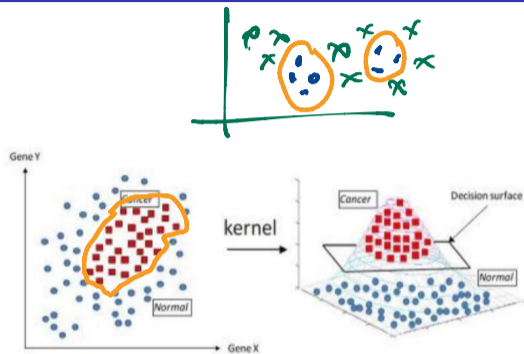
$$\text{sign} \left[\sum_{i \in \text{sv}} y_i \alpha_i (x_i \cdot z) + b \right]$$

- After transformation

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

- Training: maximize

$$\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \varphi(x_i) \cdot \varphi(x_j)$$



What is φ ?
Don't care

SVM after transformation

- SVM in original space

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

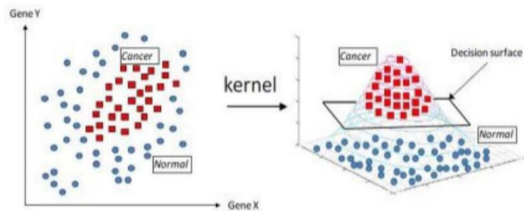
- After transformation

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

- Training: maximize

$$\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \varphi(x_i) \cdot \varphi(x_j)$$

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

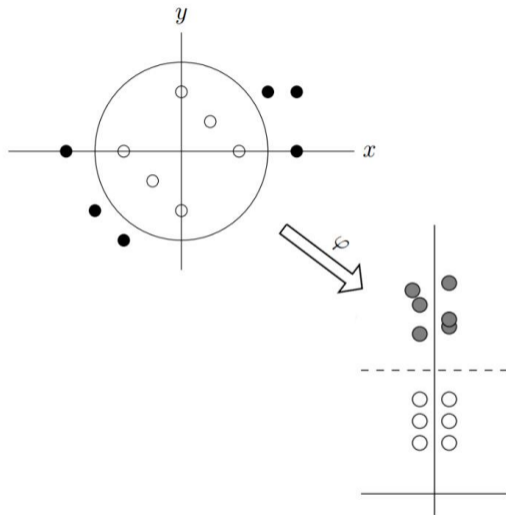


Dot products

- Consider the transformation

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

2D \longrightarrow 6D



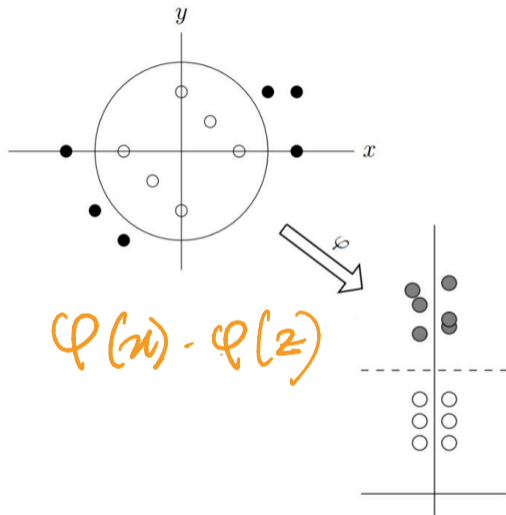
Dot products

- 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

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



Dot products

- Consider the transformation

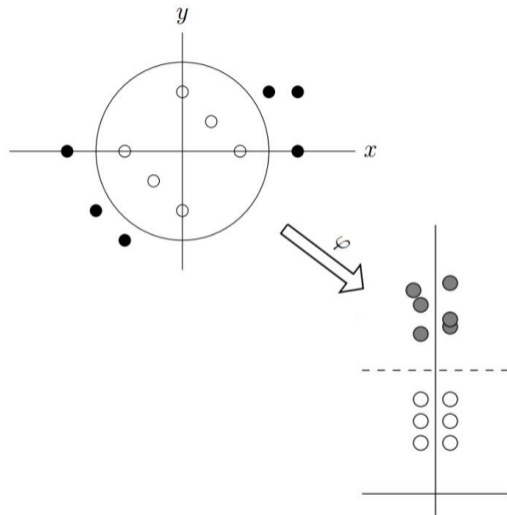
$$\varphi : (x_1, x_2) \mapsto (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

- Dot product in transformed space

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

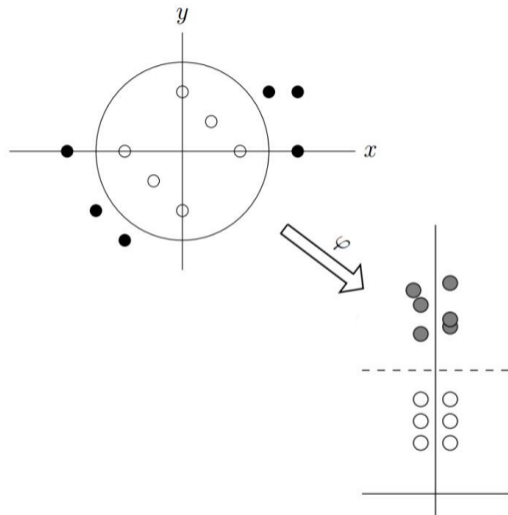
- Transformed dot product can be expressed in terms of original inputs

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



Kernels

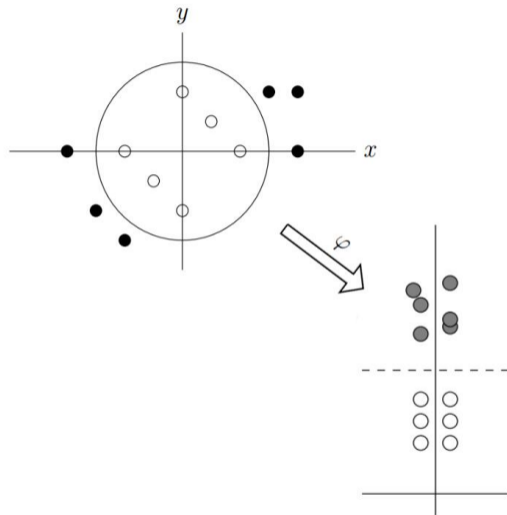
- K is a **kernel** for transformation φ if
$$K(x, z) = \varphi(x) \cdot \varphi(z)$$



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

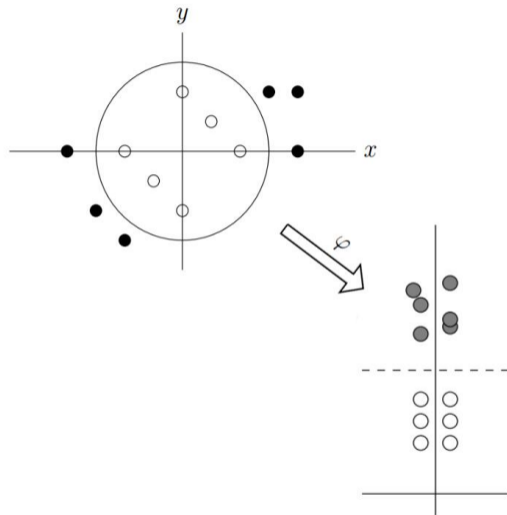
$$\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \underbrace{\varphi(x_i) \cdot \varphi(x_j)}$$
$$\text{sign} \left[\sum_{i \in S^+} y_i \alpha_i \underbrace{(\varphi(x_i) \cdot \varphi(z))} + b \right]$$



Kernels

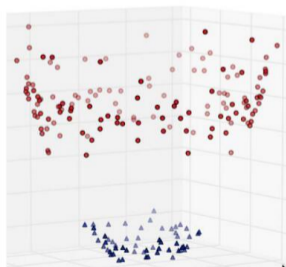
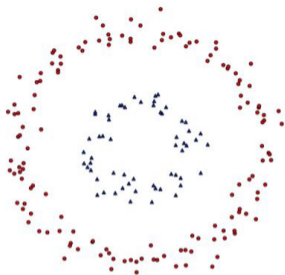
- K is a **kernel** for transformation φ if $K(x, z) = \varphi(x) \cdot \varphi(z)$
- If we have a kernel, we don't need to explicitly compute transformed points
- All dot products can be computed implicitly using the kernel on original data points

$$\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \underline{K(x_i, x_j)}$$
$$\text{sign} \left[\sum_{i \in \text{sv}} y_i \alpha_i \underline{K(x_i, z)} + b \right]$$



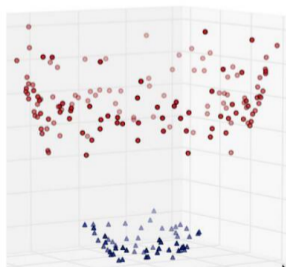
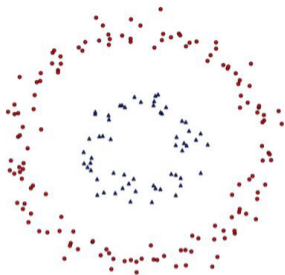
Kernels

- If we know K is a kernel for some transformation φ , we can blindly use K without even knowing what φ looks like!

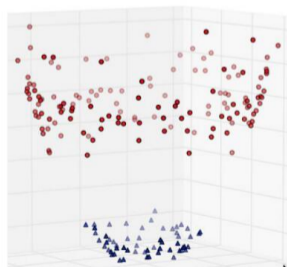
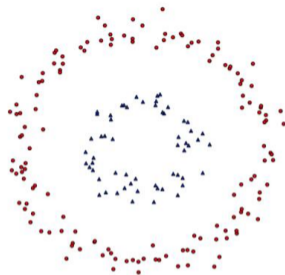


Kernels

- 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?

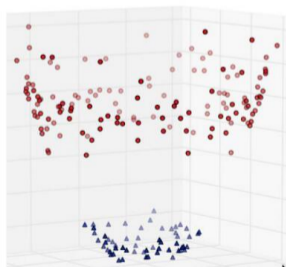
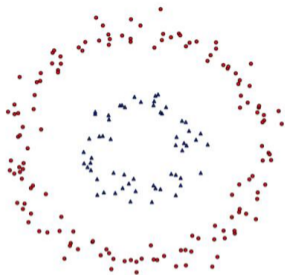


- If we know K is a kernel for some transformation φ , we can blindly use K without even knowing what φ looks like!
- When is a function a valid kernel?
- Has been studied in mathematics — **Mercer's Theorem**
 - Criteria are non-constructive



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

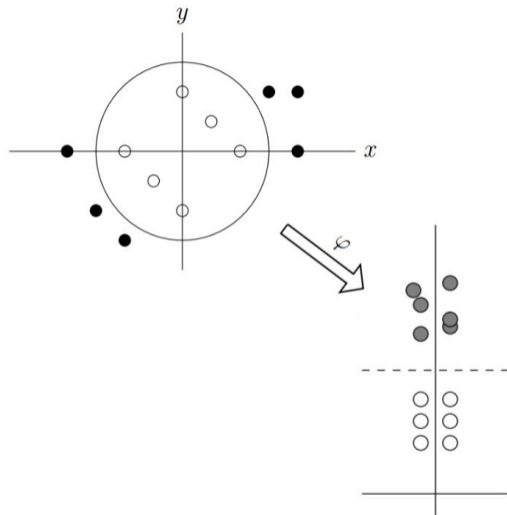


Kernels

- 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 $K(x_i, x_j)$

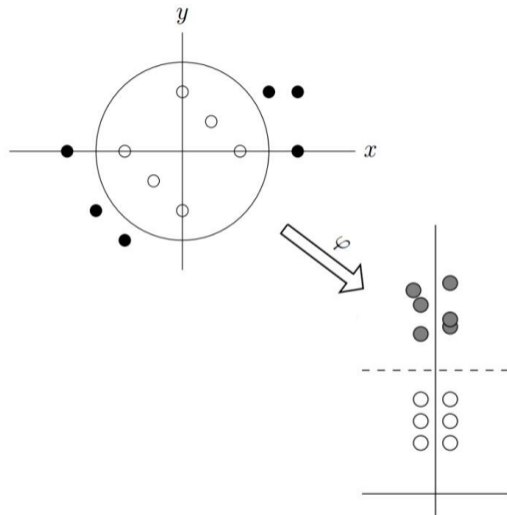


Kernels

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

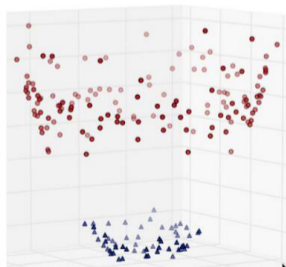
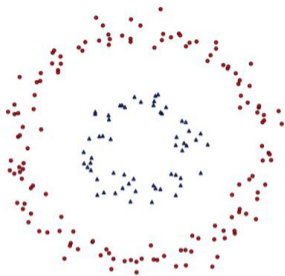
$$K = \begin{matrix} & x_1 & x_2 & \cdots & x_N \\ \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{matrix} & \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \end{matrix}$$

- 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

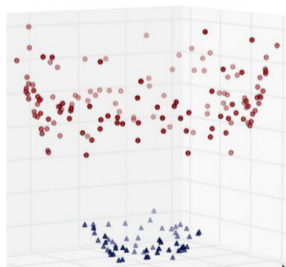


Known kernels

- Fortunately, there are many known kernels
- Polynomial kernels

$$K(x, z) = (1 + x \cdot z)^k$$

$$(1 + (x \cdot z))^2 k$$



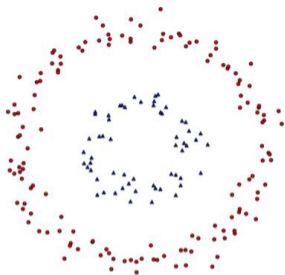
Known kernels

- Fortunately, there are many known kernels

- Polynomial kernels

$$K(x, z) = (1 + x \cdot z)^k$$

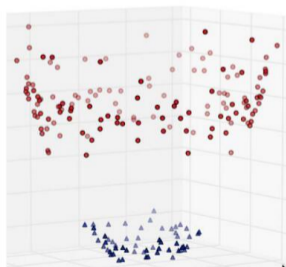
- Any $K(x, z)$ representing a similarity measure



$$d(x, z)$$

similarity

$$\frac{1}{d(x, z)}$$

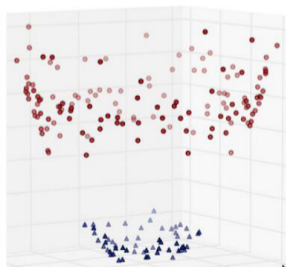
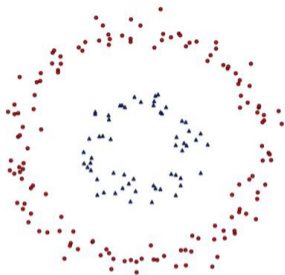


Known kernels

- 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

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

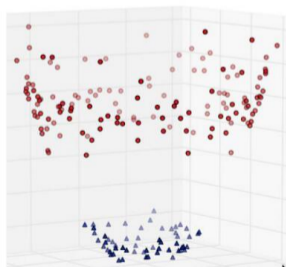
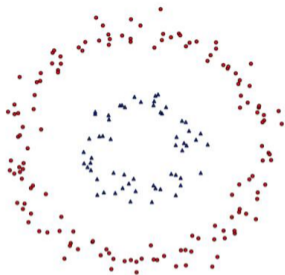
|
distance



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

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



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

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

