# Lecture 18: 21 March, 2023

Madhavan Mukund

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

Data Mining and Machine Learning
January–April 2023

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

*Add* $C \cdot \sum \xi_i^2$

*Use higher power than 2*

Subject to

$\xi_i \geq 0$

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

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

- Constraints include requirement that error terms are non-negative

- Again the objective function is quadratic



$\langle \mathbf{w} \cdot \mathbf{x} \rangle + b = 0$

$\mathbf{w}$

$\mathbf{x}_b$

$\frac{\xi_b}{\|\mathbf{w}\|}$

$\frac{|b|}{\|\mathbf{w}\|}$

$\frac{\xi_a}{\|\mathbf{w}\|}$

$\mathbf{x}_a$

## Wolfe dual

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

only dot products of $x_i$'s

$\varphi(x_i) \cdot \varphi(x_j)$

Subject to

$0 \leq \alpha_i \geq 1$

$\sum_i \alpha_i y_i = 0$

Lagrange Multiplier

No $x_i$'s here

$\langle \mathbf{w} \cdot \mathbf{x} \rangle + b = 0$

$\mathbf{w}$

$\mathbf{x}_b$

$\frac{\xi_b}{\|\mathbf{w}\|}$

$\frac{|b|}{\|\mathbf{w}\|}$

$\frac{\xi_a}{\|\mathbf{w}\|}$
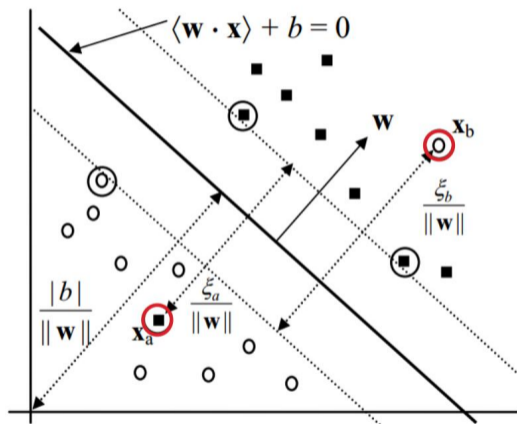
$\mathbf{x}_a$

# 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 $\alpha_i$

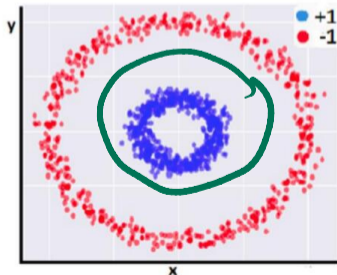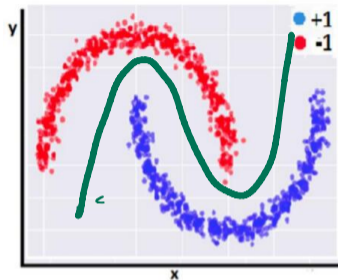  - Only terms corresponding to support vectors are actively used

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

*only dot prod.*



$\langle \mathbf{w} \cdot \mathbf{x} \rangle + b = 0$

$\mathbf{w}$

$\mathbf{x}_b$

$\dfrac{\xi_b}{\|\mathbf{w}\|}$

$\dfrac{|b|}{\|\mathbf{w}\|}$

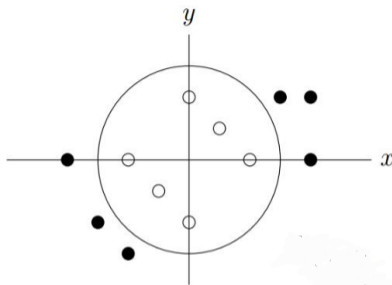$\dfrac{\xi_a}{\|\mathbf{w}\|}$

$\mathbf{x}_a$

# 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
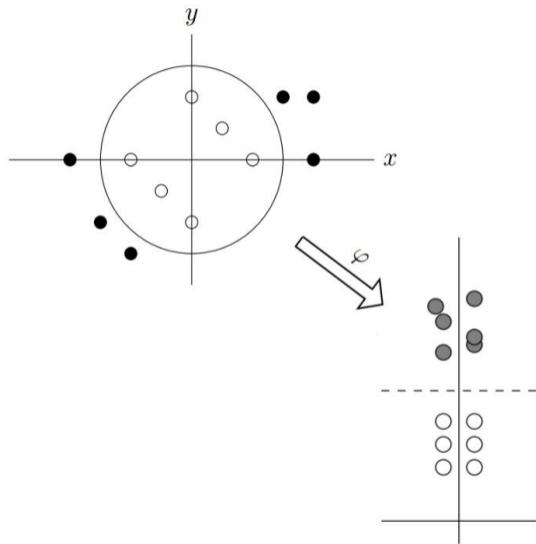
# Geometric tranformation

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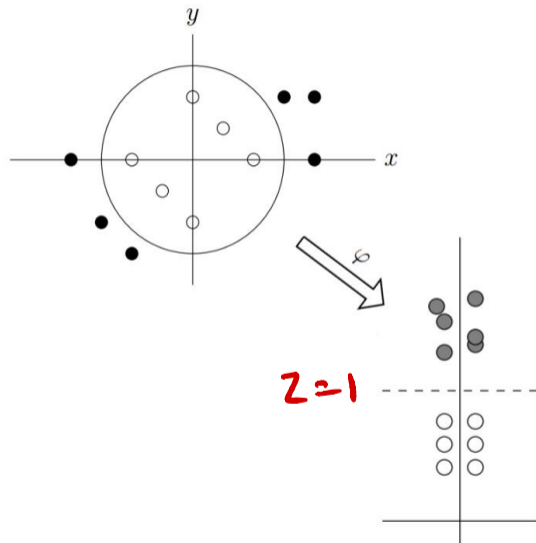
- Equation of circle is $x^2 + y^2 = 1$

# Geometric tranformation
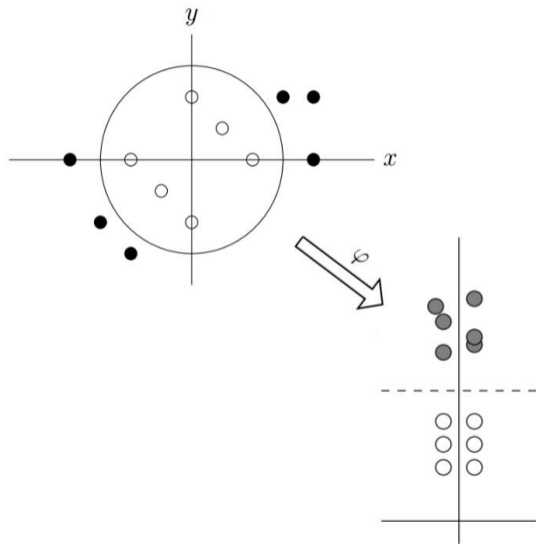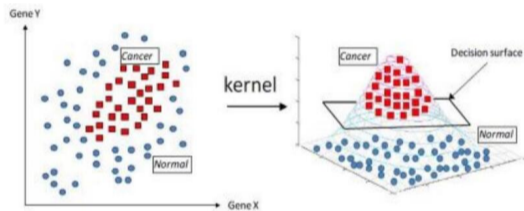
- Consider two sets of points separated by a circle of radius $1$

- Equation of circle is $x^2 + y^2 = 1$

- Points inside the circle, $x^2 + y^2 < 1$

- Points outside circle, $x^2 + y^2 > 1$

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



$z = 1$

# Geometric tranformation

- 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

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

# 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

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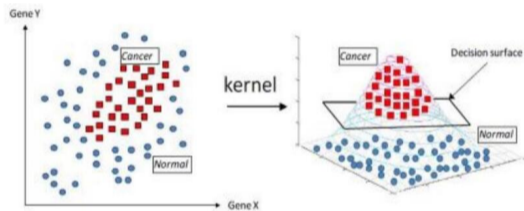
# SVM after transformation

- 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

- Consider the transformation

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



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

- Consider the transformation

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

  $(1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, \sqrt{2}z_1z_2, z_2^2)$

- Dot product in transformed space

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

  $(x_1, x_2) \quad (z_1, z_2)$

# Dot products

- Consider the transformation

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

- Dot product in transformed space

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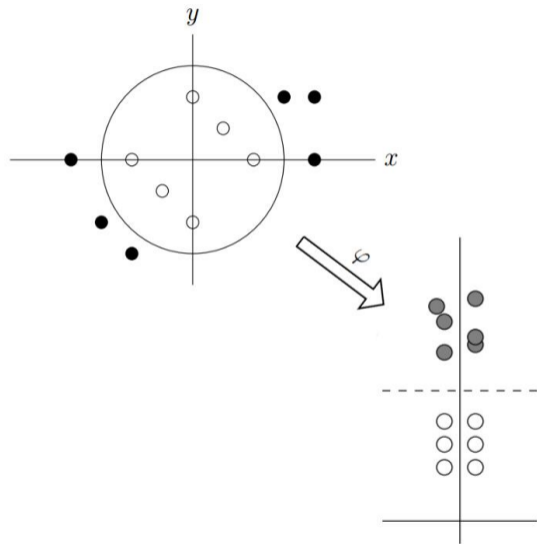&= (1 + x_1z_1 + x_2z_2)^2
\end{aligned}$$

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

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

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

# Kernels

- $K$ is a kernel for transformation $\varphi$ 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
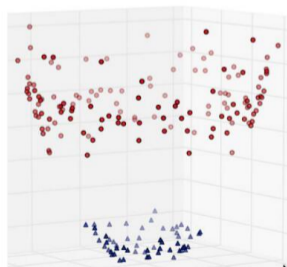
$$\text{sign}\left[\sum_{i \in sv} y_i \alpha_i \left(\varphi(x_i) \cdot \varphi(z)\right) + b\right]$$

- $K$ is a kernel for transformation $\varphi$ 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

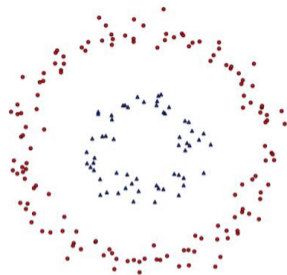$$\text{sign} \left[ \sum_{i \in sv} y_i \alpha_i K(x_i, z) + b \right]$$



Also in $\varphi(x_i) \cdot \varphi(x_j)$ of dual $\rightarrow K(x_i, x_j)$

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

$$K(x_i, x_j) \text{ in training}$$
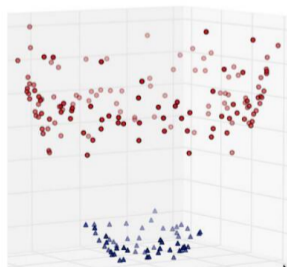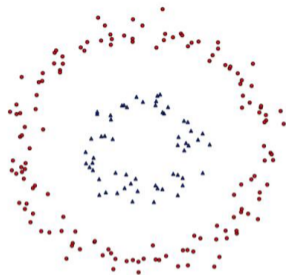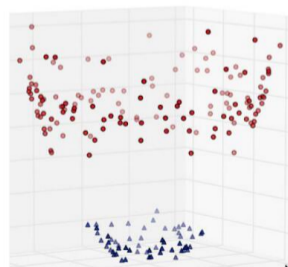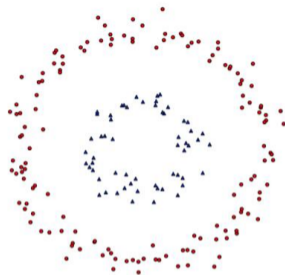
$$K(x_i, z) \text{ for classification}$$

# Kernels

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

- When is a function a valid kernel?

# Kernels

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

- When is a function a valid kernel?

- Has been studied in mathematics — Mercer's Theorem

    - Criteria are non-constructive
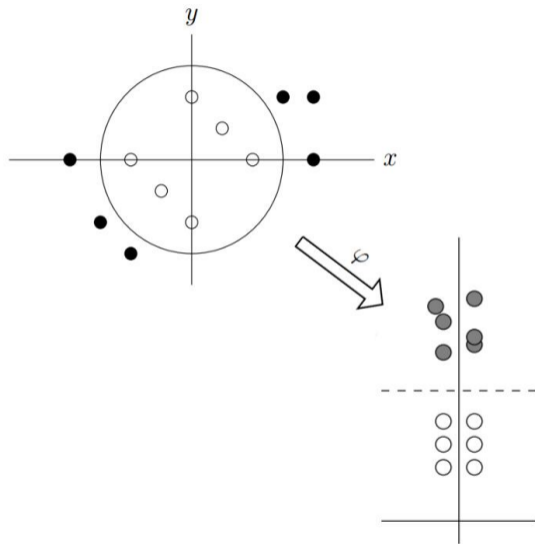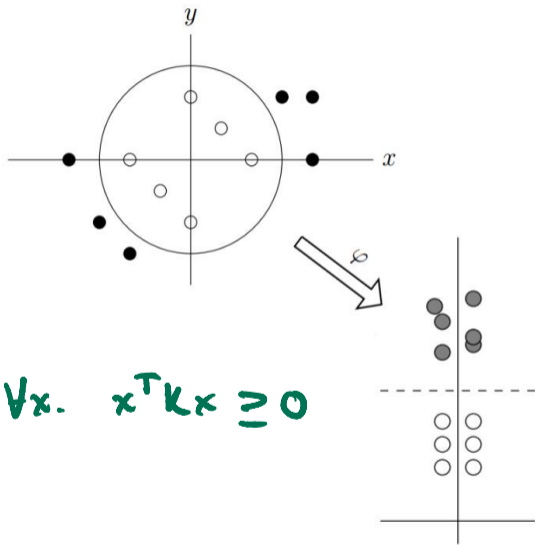
- Can define sufficient conditions from linear algebra

# Kernels

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

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

- Entries are values $K(x_i, x_j)$

# Kernels

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

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

- Entries are values $K(x_i, x_j)$

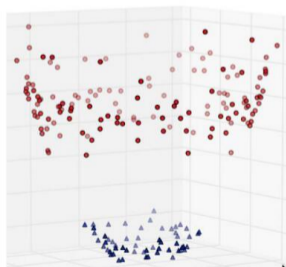- Gram matrix should be positive semi-definite for all $x_1, x_2, \ldots, x_N$
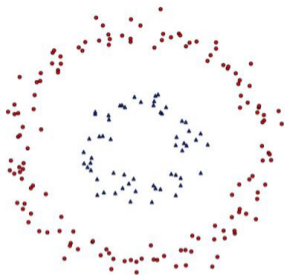
$$\longrightarrow \quad \forall x. \quad x^T k x \geq 0$$

# Known kernels

- Fortunately, there are many known kernels
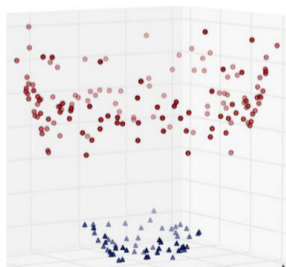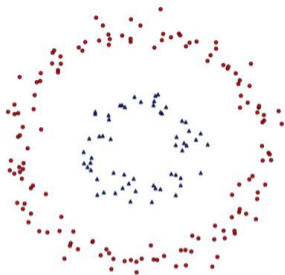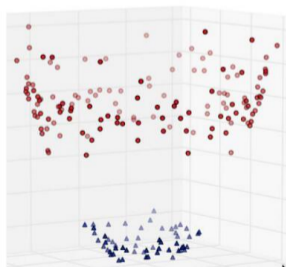
# Known kernels

- Fortunately, there are many known kernels

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



$$\left(1 + x_1 z_1 + x_2 z_2\right)^2$$
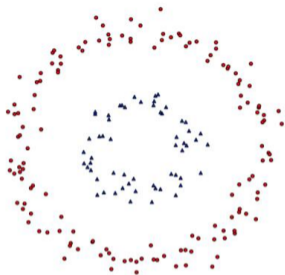
$$\left(1 + x \cdot z\right)^2$$

# Known kernels

- Fortunately, there are many known kernels

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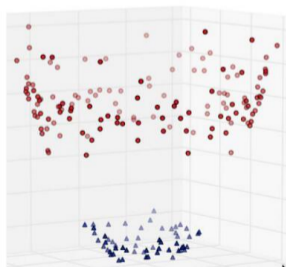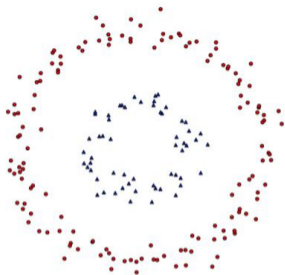$$K(x, z) = (1 + x \cdot z)^k$$

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

# 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

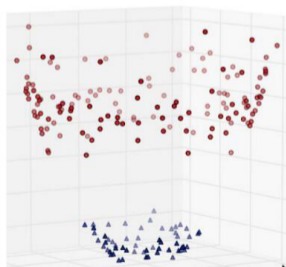- Gaussian radial basis function — similarity based on inverse exponential distance
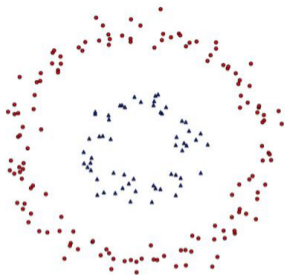$$K(x, z) = e^{-c|x-z|^2}$$

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

Till about 2010

SVM + manually constructed kernels
were "best" classifiers

Then came neural networks

This happens "automatically"