#### Kernel Methods

#### Madhavan Mukund

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

Data Mining and Machine Learning August-December 2020

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

## Soft margin optimization

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

Subject to

$$\begin{array}{ll} \xi_i \geq 0 \\ \langle w \cdot x \rangle + b > 1 - \xi_i, & \text{if } y_i = 1 \\ \langle w \cdot x \rangle + b < -1 + \xi_i, & \text{if } y_i = -1 \end{array}$$

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



(日) (部) (注) (注)

æ



# 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  $\quad x^2+y^2=1$
- 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

$$\operatorname{sign}\left[\sum_{i \in sv} y_i \alpha_i \langle \underline{x_i \cdot z} \rangle + b\right]$$

• After transformation

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

・ロト ・ 御 ト ・ モト ・ モト

æ

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



## Dot products



•  $K \, {\rm is} \, {\rm a} \, {\it kernel} \,$  for transformation  $\varphi \,$  if

 $K(x,z) = \langle \varphi(x) \cdot \varphi(z) \rangle$ 

- 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

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

$$\mathsf{K}(\mathbf{x}_i, \mathbf{z})$$



-  $K \, {\rm is} \, {\rm a} \, {\it kernel} \,$  for transformation  $\varphi \,$  if

$$K(x,z) = \langle \varphi(x) \cdot \varphi(z) \rangle$$

- 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

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





- 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





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

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

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



æ

## **Known kernels**

K11K2 Kendo -> K1+K2

- Fortunately, there are many known kernels
- Polynomial kernels

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

• Any K(x,z) representing a similarity measure

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

 Gaussian radial basis function – similarity based on inverse exponential distance



SVM + Kernel functions Powerful dessifier Love 1990' ( ) 2010 Identify "good" kernels for domain Manual engeneerig of kernels