# **KERNEL MACHINES AND KNOWLEDGE REPRESENTATION**

INTRODUCTION

#### **SUPPORT VECTOR MACHINES**

- Support Vector Machines (SVMs) are a machine learning paradigm based on the statistical learning theory [Vapnik, 1995]
  - No need to remember everything, just the discriminating instances (i.e. the support vectors, SV)
  - The classifier corresponds to the linear combination of SVs



# LINEAR CLASSIFIERS AND SEPARABILITY

- In a R<sup>2</sup> space, 3 point can always be separable by a linear classifier
  - but 4 points cannot always be shattered [Vapnik and Chervonenkis(1971)]
- One solution could be a more complex classifier
  - 8 Risk of over-fitting

#### LINEAR CLASSIFIERS AND SEPARABILITY (2)

- ... but things change when projecting instances in a higher dimension feature space through a function  $\phi$
- IDEA: It is better to have a more complex feature space instead a more complex function (i.e. learning algorithm)

#### **THE KERNEL FUNCTION**

- In perceptrons and SVMs the learning algorithm only depends on the scalar product over pairs of example instance vectors
- Basically only the Gram-matrix is involved. In general, we call kernel the following function:

$$K(\vec{x}, \vec{z}) = \Phi(\vec{x}) \cdot \Phi(\vec{z})$$

- The kernel corresponds to a scalar product over the transformed of initial objects x and z
- Notice that the training in most learning machines (such as the perceptron) makes use of instances only through the kernel



## **AN EXAMPLE: A MAPPING FUNCTION**

- Two masses  $m_1$  and  $m_2$ , one is constrained
- A force  $f_a$  is applied to the mass  $m_1$
- Instead of applying an analyitical law we want to experiment
  - The Features of individual experiments are masses  $m_1, m_2$  and the appropriate orce  $f_a$
- It is clear that the Newton law of gravity is involved:

$$f(m_1, m_2, r) = C \frac{m_1 m_2}{r^2}$$

The task corresponds to determine if  $f(m_1, m_2, r) < f_a$ 

97

## **AN EXAMPLE: A MAPPING FUNCTION (2)**

 $\vec{x} = (x_1, \dots, x_n) \to \Phi(\vec{x}) = (\Phi_1(\vec{x}), \dots, \Phi_k(\vec{x}))$ 

This law cannot be expressed linearly. A change of space:

 $(f_a, m_1, m_2, r) \rightarrow (k, x, y, z) = (\ln f_a, \ln m_1, \ln m_2, \ln r)$ 

holds as:

 $\ln f(m_1, m_2, r) = \ln C + \ln m_1 + \ln m_2 - 2 \ln r = c + x + y - 2z$ 

The following hyperplane is the requested function h():

 $\ln f_a - \ln m_1 - \ln m_2 + 2 \ln r - \ln C = 0$ 

 $(1,1,-2,-1) \cdot (\ln m_1, \ln m_2, \ln r, \ln f_a) + \ln C = 0,$ 

We can decide with no error if masses  $m_1, m_2$  get closer or not

#### **FEATURE SPACES AND KERNELS**

- Feature Space
  - The input space is mapped into a new space F with scalar product (called *feature space*) through a (non linear) trasformation  $\phi$

• The kernel function  $\phi = R^N \to F$ 

- The evaluation require the computation of the scalar product over the trasformed vectors  $\phi(x)$  but not the feature vectors themselves
- The scalar product is computed by a specialized function called kernel

 $k(x, y) = (\phi(x) \cdot \phi(y))$ 



## **CLASSIFICATION FUNCTION: THE DUAL FORM**

$$h(x) = sgn(\vec{w} \cdot \vec{x} + b) = sgn(\sum_{J=1}^{l} \alpha_{J} y_{J} \vec{x_{J}} \cdot \vec{x} + b)$$

On the right form, instances only appear in the scalar product
The ony thing that is needed is the Gram matrix,

$$G = \left(\left\langle \mathbf{x}_{i} \cdot \mathbf{x}_{j} \right\rangle\right)_{i,j=1}^{l}$$

i.e. the explicit computation of the scalar product over any pair of training instances  $x_1 \dots x_l$ 

#### **A KERNELIZED PERCEPTRON**

We can rewrite the decision function of a perceptron by taking into account a kernel:

$$h(x) = sgn(\vec{w} \cdot \Phi(\vec{x}) + b) = sgn(\sum_{J=1}^{l} \alpha_j y_j \Phi(\vec{x_j}) \cdot \Phi(\vec{x}) + b)$$
$$= sgn(\sum_{J=1}^{l} \alpha_j y_j k(\vec{x_j}, \vec{x}) + b)$$

and during training the on-line adjustment steps become:

$$y_i(\sum_{J=1}^l \alpha_j y_j \Phi(\overrightarrow{x_j}) \cdot) \Phi(\overrightarrow{x_i}) + b) = \sum_{J=1}^l \alpha_j y_i y_j k(\overrightarrow{x_j}, \overrightarrow{x_i}) + b)$$

#### **KERNELS IN SUPPORT VECTOR MACHINES**

$$\sum_{i=1}^{m} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m} y_i y_j \alpha_i \alpha_j \vec{x_i} \cdot \vec{x_j} + \frac{1}{2C} \vec{\alpha} \cdot \vec{\alpha} - \frac{1}{C} \vec{\alpha} \cdot \vec{\alpha}$$

• By using kernel functions we rewrite the problem as:

In Soft Margin SVMs we need to maxin

$$\begin{aligned} \max & \max \sum_{i=1}^{m} \alpha_{i} - \frac{1}{2} \sum_{i,j=1}^{m} y_{i} y_{j} \alpha_{i} \alpha_{j} \left( k(o_{i}, o_{j}) + \frac{1}{C} \delta_{ij} \right) \\ \alpha_{i} &\geq 0, \quad \forall i = 1, ..., m \\ \sum_{i=1}^{m} y_{i} \alpha_{i} &= 0 \end{aligned}$$

### WHAT MAKES A FUNCTION A KERNEL FUNCTION?

**Def. 2.26** A kernel is a function k, such that  $\forall \vec{x}, \vec{z} \in X$ 

 $k(\vec{x}, \vec{z}) = \phi(\vec{x}) \cdot \phi(\vec{z})$ 

where  $\phi$  is a mapping from X to an (inner product) feature space.

Only such type of functions support implicit mappings such as

 $\vec{x} = (x_1, \dots, x_n) \in \mathbb{R}^n \rightarrow \Phi(\vec{x}) = (\Phi_1(\vec{x}), \dots, \Phi_m(\vec{x})) \in \mathbb{R}^m$ 

## WHAT MAKES A FUNCTION A KERNEL FUNCTION? (2)

Def. B.11 Eigen Values

Given a matrix  $A \in \mathbb{R}^n \times \mathbb{R}^n$ , an egeinvalue  $\lambda$  and an egeinvector  $\vec{x} \in \mathbb{R}^n - {\vec{0}}$  are such that

$$A\vec{x} = \lambda\vec{x}$$

**Def. B.12** Symmetric Matrix A square matrix  $A \in \mathbb{R}^n \times \mathbb{R}^n$  is symmetric iff  $A_{ij} = A_{ji}$  for  $i \neq j$  i = 1, ..., mand j = 1, ..., n, i.e. iff A = A'.

**Def. B.13** Positive (Semi-) definite Matrix A square matrix  $\mathbf{A} \in \mathbb{R}^n \times \mathbb{R}^n$  is said to be positive (semi-) definite if its eigenvalues are all positive (non-negative).

## WHAT MAKES A FUNCTION A KERNEL FUNCTION? (3)

**Proposition 2.27** (Mercer's conditions) Let X be a finite input space with  $K(\vec{x}, \vec{z})$  a symmetric function on X. Then  $K(\vec{x}, \vec{z})$  is a kernel function if and only if the matrix

 $k(\vec{x},\vec{z}) = \underline{\phi}(\vec{x}) \cdot \underline{\phi}(\vec{z})$ 

is positive semi-definite (has non-negative eigenvalues).

#### · .« 1

- IDEA: If the Gram matrix is positive semi-definite then the mapping φ, such that F is an inner-product space whose scalar product corresponds to the kernel k(.,.), exists
- In F the separability should be easier

#### **KERNEL FUNCTIONS: OTHER EXAMPLES**

An example of non linear classification

- $g(\vec{x}) = w_1 x_1^2 + w_{21} x_2^2 + w_0$  (ellipsis in  $\mathbb{R}^2$ )
  - For example circles (of radius *r*)
    - $g(\vec{x}) = x_1^2 + x_2^2 r^2$
- Can it be made linear in some space?



Lesson 1: March 1°, 2023

## LINEAR CLASSIFICATION IN THE KERNEL SPACE

Let try to rewrite the definition for g()

• 
$$g(\vec{x}) = w_1 x_1^2 + w_2 x_2^2 + w_0 = g(\vec{x}) = w_1 \chi_1 + w_2 \chi_2 + w_0 =$$
  
=  $w_1 \phi_1(\vec{x}) + w_2 \phi_2(\vec{x}) + w_0$   
 $(\phi(\vec{x}) = \phi(x_1, x_2) \mapsto (x_1^2, x_2^2))$ 

• G is linear in the  $\phi$ -transformed space!

Lesson 1: March 1°, 2023

#### **FEATURE SPACES AND KERNELS**

# Another example of Kernel

The Polynomial kernel 
$$k(x, y) = (x \cdot y)^d$$
  
 $x, y \in \mathbb{R}^2$   
If  $d=2$  and  $(x \cdot y)^2 = \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right)^2 = \left(\begin{bmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} y_1^2 \\ \sqrt{2} y_1 y_2 \\ y_2^2 \end{bmatrix}\right)$   
 $= (\phi(x) \cdot \phi(y)) = k(x, y)$ 

## **POLYNOMIAL KERNEL**



# **POLYNOMIAL KERNEL (N DIMENSIONS)**

$$(\vec{x} \cdot \vec{z})^2 = \left(\sum_{i=1}^n x_i z_i\right)^2 = \left(\sum_{i=1}^n x_i z_i\right) \left(\sum_{j=1}^n x_i z_i\right) = \sum_{i=1}^n \sum_{j=1}^n x_i x_j z_i z_j = \sum_{i,j \in \{1,\dots,n\}} (x_i x_j) (z_i z_j) = \sum_{k=1}^m X_k Z_k = \vec{X} \cdot \vec{Z}$$

# **GENERAL POLYNOMIAL KERNEL (N DIMENSIONS)**

$$(\vec{x} \cdot \vec{z} + c)^2 = \left(\sum_{i=1}^n x_i z_i + c\right)^2 = \left(\sum_{i=1}^n x_i z_i + c\right) \left(\sum_{j=1}^n x_i z_i + c\right) =$$
$$= \sum_{i=1}^n \sum_{j=1}^n x_i x_j z_i z_j + 2c \sum_{i=1}^n x_i z_i + c^2 =$$
$$= \sum_{i,j \in \{1,..,n\}} (x_i x_j) (z_i z_j) + 2c \sum_{i=1}^n \left(\sqrt{2c} x_i\right) \left(\sqrt{2c} z_i\right) + c^2$$

#### POLYNOMIAL KERNEL AND THE CONJUNCTION OF FEATURES

- The initial vectors can be mapped into a higher dimensional space (c=1)  $\Phi(\langle x_1, x_2 \rangle) \rightarrow (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$
- More expressive, as (x<sub>1</sub>x<sub>2</sub>) encodes original feature pairs, e.g. stock+market vs. downtown+market are contributing (when occurring) togheter
- We can smartly compute the scalar product as  $\Phi(\vec{x}) \cdot \Phi(\vec{z}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1) \quad (z_1^2, z_2^2, \sqrt{2}z_1z_2, \sqrt{2}z_1, \sqrt{2}z_2, 1) = x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 + 2x_1 z_1 + 2x_2 z_2 + 1 = (x_1 z_1 + x_2 z_2 + 1)^2 = (\vec{x} \cdot \vec{z} + 1)^2 = K_{p2} \quad (\vec{x}, \vec{z})$

## THE ARCHITECTURE OF AN SVM

- It is a non linear classifier (based on a kernel)
- Decision function:



#### **OTHER EXAMPLES OF KERNEL FUNCTIONS**

Linear: 
$$k(\vec{x}_i, \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j$$

• Polynomial (degree p):  $k(\vec{x}_i, \vec{x}_j) = (1 + \vec{x}_i \cdot \vec{x}_j)^p$ 

Gaussian (radial-basis function network):

network):  

$$k(\vec{x}_i, \vec{x}_j) = e^{-\frac{\|\vec{x}_i - \vec{x}_j\|^2}{2\sigma^2}}$$

• Perceptron (two stages):  $k(\vec{x}_i, \vec{x}_j) = \tanh(\beta_1 + \beta_0 \vec{x}_i \cdot \vec{x}_j)^p$ 

#### **KERNEL COMBINATION AND NORMALIZATION**

- Kernels can be easily combined so that the evidences captured by several kernel functions can contribute to the learning algorithm
  - Linear operations on kernel functions preserve the kernel properties
  - The sum of kernels is a valid kernel
  - The product of kernels is a valid kernel
- We can also normalize the implicit space operating directly only the kernel function

$$\begin{split} \hat{K}(s,t) &= \left\langle \hat{\phi}(s) \cdot \hat{\phi}(t) \right\rangle = \left\langle \frac{\phi(s)}{\|\phi(s)\|} \cdot \frac{\phi(t)}{\|\phi(t)\|} \right\rangle \\ &= \frac{1}{\|\phi(s)\| \|\phi(t)\|} \left\langle \phi(s) \cdot \phi(t) \right\rangle = \frac{K(s,t)}{\sqrt{K(s,s)K(t,t)}} \end{split}$$

# **STRING KERNEL**

- Given two strings, the number of matches between their substrings is computed
- E.g. *Bank* and *Rank* 
  - B, a, n, k, Ba, Ban, Bank, an, ank, nk
  - R, a, n, k, Ra, Ran, Rank, an, ank, nk
- String kernel over sentences and texts
- Huge space but there are efficient algorithms
  - Lodhi, Huma; Saunders, Craig; Shawe-Taylor, John; Cristianini, Nello; Watkins, Chris (2002). "Text classification using string kernels". Journal of Machine Learning Research: 419–444.

### **STRING KERNEL**

A function that give two strings s and t is able to compute a real number k(s,t) such that

- two vectors exist  $\vec{s}$  and  $\vec{t}$
- $\vec{s}$  and  $\vec{t}$  are unique for s and t
- (the vectors represents strings by embedding their crucial properties!!)
- $k(s,t) = \vec{s} \times \vec{t}$
- We will see how vectors  $\vec{s}$  and  $\vec{t}$  are defined in  $\mathbb{R}^{\infty}$ , as the numer of strings of arbitrary length over an alphabet is infinite

IDEA: Define a space whereas each substring is a dimension

#### **KERNEL TRA BANK E RANK**

B, a, n, k, Ba, Ban, Bank, an, ank, nk, Bn, Bnk, Bk and ak are the substrings of *Bank*.

R, a, n, k, Ra, Ran, Rank, an, ank, nk, Rn, Rnk, Rk and ak are the substrings of *Rank*.

 $\phi$ 

 $\begin{aligned} \phi(\text{Bank}) &= (\lambda \ , \ 0, \ \lambda, \ \lambda, \ \lambda, \ \lambda, \ \lambda^2 \ , \ \lambda^2, \ \lambda^3 \ , \ 0 \ , \ \lambda^4 \ , \ 0 \ , \ \lambda^2 \ , \ \lambda^3 \ , \ \lambda^3 \ , \ \dots \\ \phi(\text{Rank}) &= (0 \ , \ \lambda, \ \lambda, \ \lambda, \ \lambda, \ \lambda, \ 0 \ , \ 0 \ , \ 0 \ , \ \lambda^3 \ , \ 0 \ , \ \lambda^4 \ , \ \lambda^2 \ , \ \lambda^3 \ , \ \lambda^3 \ , \ \dots \\ & \text{B} \ , \ \text{R}, \ \text{a, n} \ , \ \text{k, Ba, Ra, Ban, Ran, Bank, Rank, an, ank \ , ak \ \dots } \end{aligned}$ 

• Common substrings:

*a, n, k, an, ank, nk, ak*Notice how these are the same subsequences as between

Schrianak and Rank

## FORMALLY ...



#### **AN EXAMPLE OF STRING KERNEL COMPUTATION**

- 
$$\phi_{a}(Bank) = \phi_{a}(Rank) = \lambda^{(i_{1}-i_{1}+1)} = \lambda^{(2-2+1)} = \lambda$$
,

- 
$$\phi_n(\text{Bank}) = \phi_n(\text{Rank}) = \lambda^{(i_1 - i_1 + 1)} = \lambda^{(3-3+1)} = \lambda$$
,

- 
$$\phi_k(\text{Bank}) = \phi_k(\text{Rank}) = \lambda^{(i_1-i_1+1)} = \lambda^{(4-4+1)} = \lambda$$
,

- 
$$\phi_{an}(Bank) = \phi_{an}(Rank) = \lambda^{(i_1-i_2+1)} = \lambda^{(3-2+1)} = \lambda^2$$
,

- 
$$\phi_{ank}(Bank) = \phi_{ank}(Rank) = \lambda^{(i_1-i_3+1)} = \lambda^{(4-2+1)} = \lambda^3$$

$$\phi_{\mathrm{nk}}(\mathrm{Bank}) = \phi_{\mathrm{nk}}(\mathrm{Rank}) = \lambda^{(i_1 - i_2 + 1)} = \lambda^{(4 - 3 + 1)} = \lambda^2,$$

$$\phi_{\texttt{ak}}(\texttt{Bank}) = \phi_{\texttt{ak}}(\texttt{Rank}) = \lambda^{(i_1 - i_2 + 1)} = \lambda^{(4-2+1)} = \lambda^3.$$

It follows that  $K(\text{Bank}, \text{Rank}) = (\lambda, \lambda, \lambda, \lambda^2, \lambda^3, \lambda^2, \lambda^3) \cdot (\lambda, \lambda, \lambda, \lambda^2, \lambda^3, \lambda^2, \lambda^3) = 3\lambda^2 + 2\lambda^4 + 2\lambda^6.$