

Kernel methods

Shikui Tu

Department of Computer Science and
Engineering, Shanghai Jiao Tong University

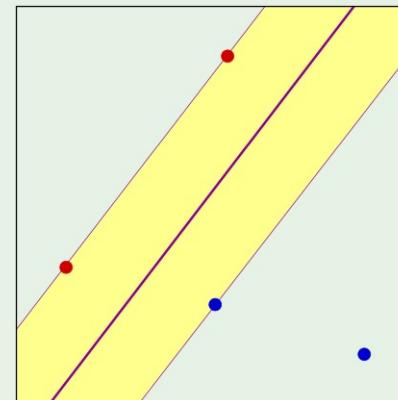
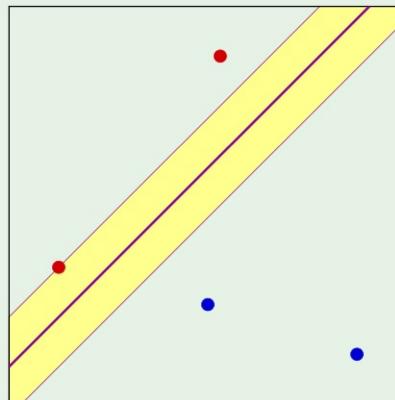
2021-05-28

Outline

- **Review of SVM in the previous lecture**
- A bit of history of kernel methods and SVM
- Kernel methods
- Designing kernels

Maximize margin

- The margin

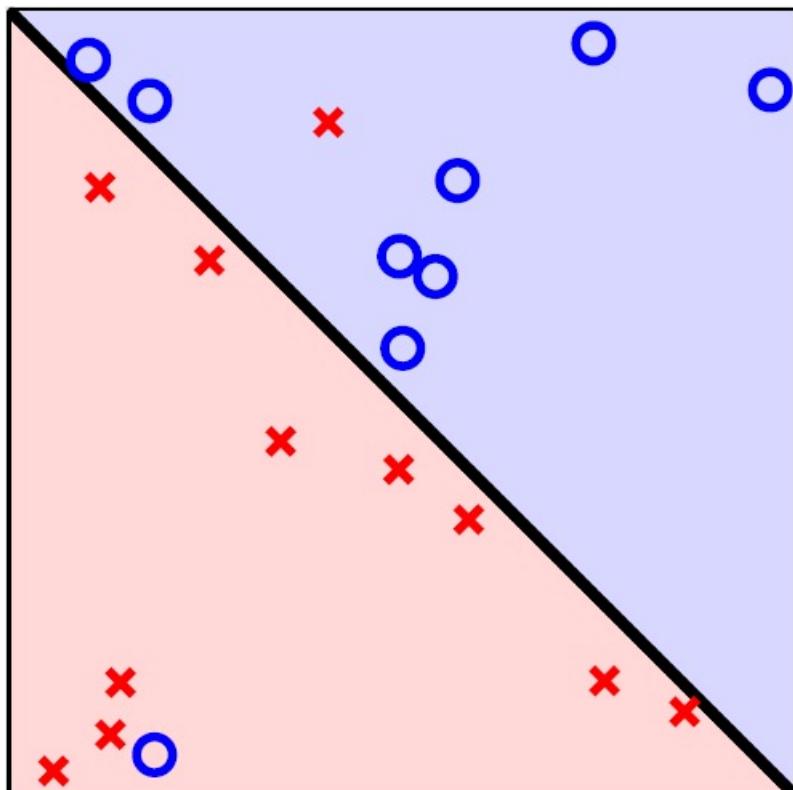


Maximizing the margin \implies dual problem:

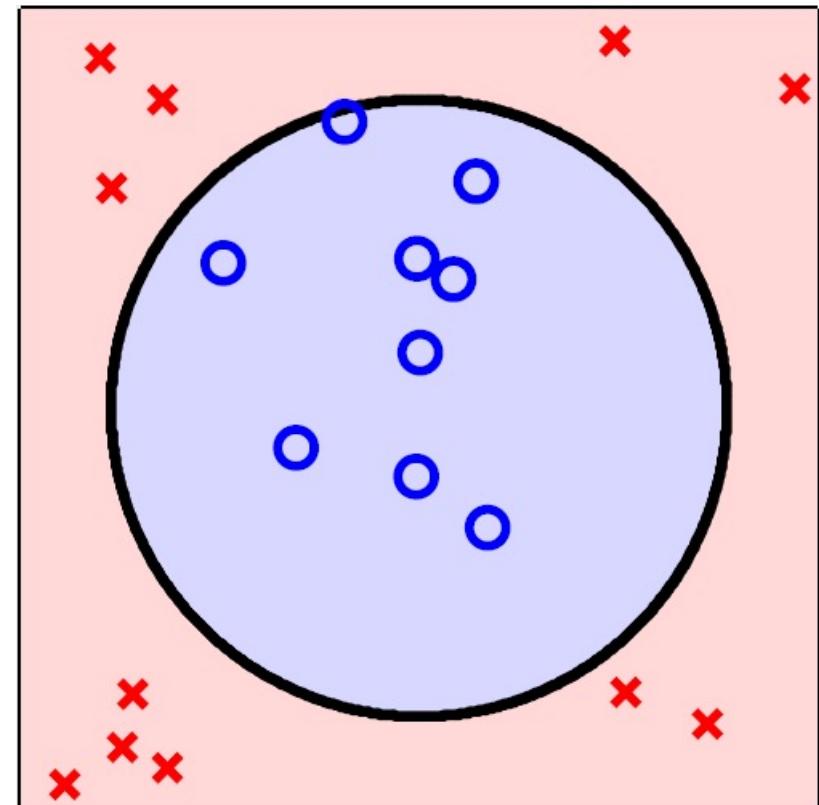
$$\mathcal{L}(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{x}_n^\top \mathbf{x}_m$$

quadratic programming

Non-separable Data



Slightly
non-separable

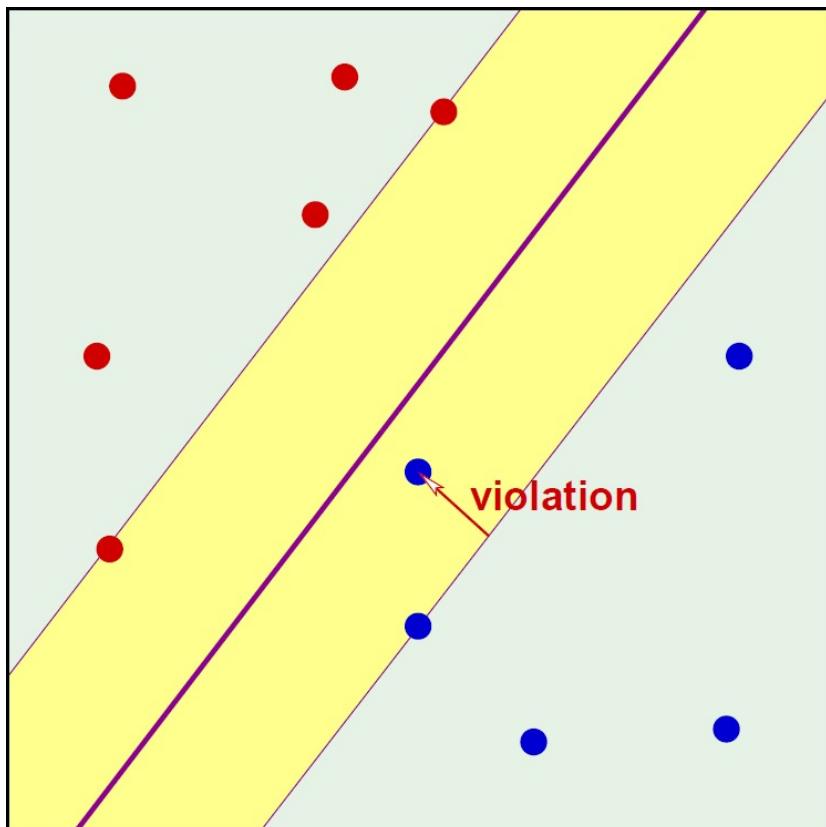


Seriously
non-separable

Soft-margin SVM

Margin violation: $y_n (\mathbf{w}^\top \mathbf{x}_n + b) \geq 1$ fails

Quantify: $y_n (\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi_n$



$$\text{Total violation} = \sum_{n=1}^N \xi_n$$

ξ_n = ksi = “amount” of margin violation

$$\boxed{\xi_n \geq 0}$$

Soft-margin SVM optimization

- Similar to hard-margin optimization, but get compromise between maximizing the margin and allowing for violations

Still get large margin after term minimization

Allow for small violations by minimizing

Minimize

$$\frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{n=1}^N \xi_n$$

subject to

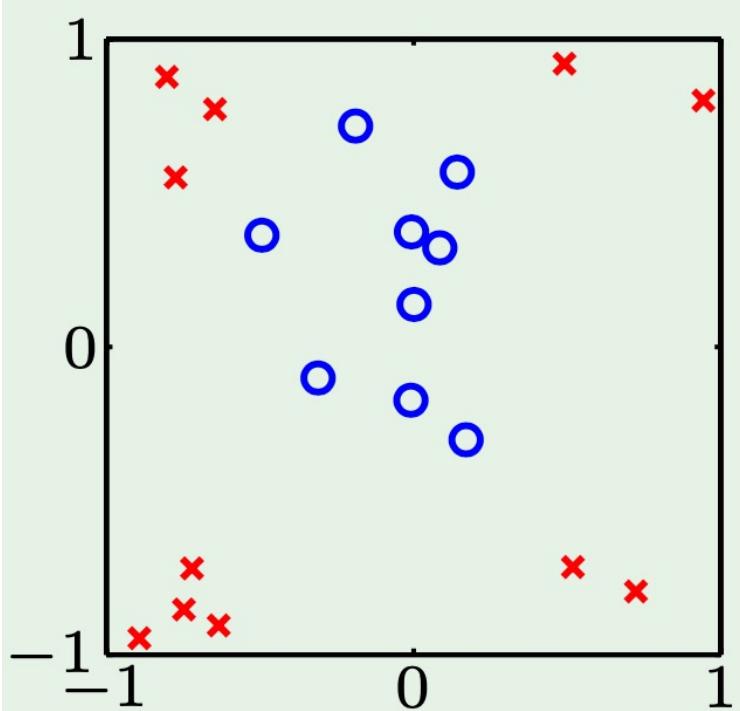
$$y_n (\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi_n$$

for $n = 1, \dots, N$

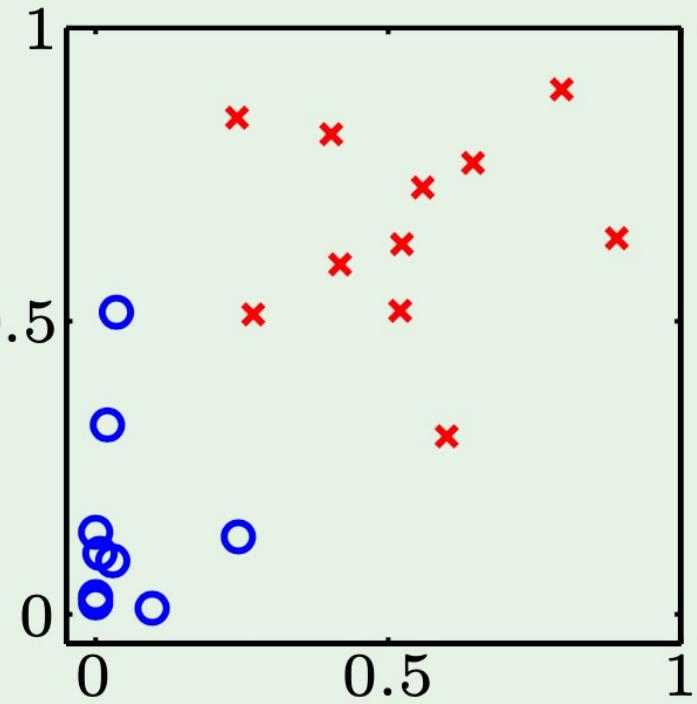
and $\xi_n \geq 0$ for $n = 1, \dots, N$

Kernel transform

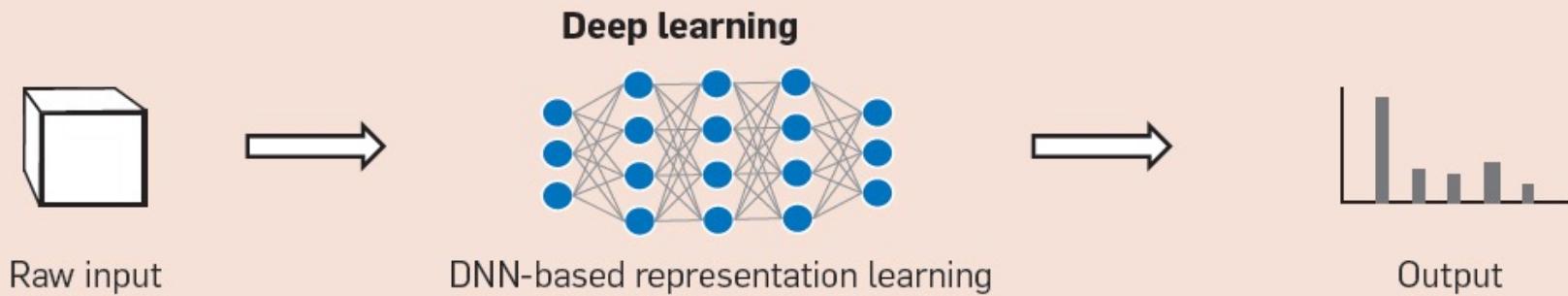
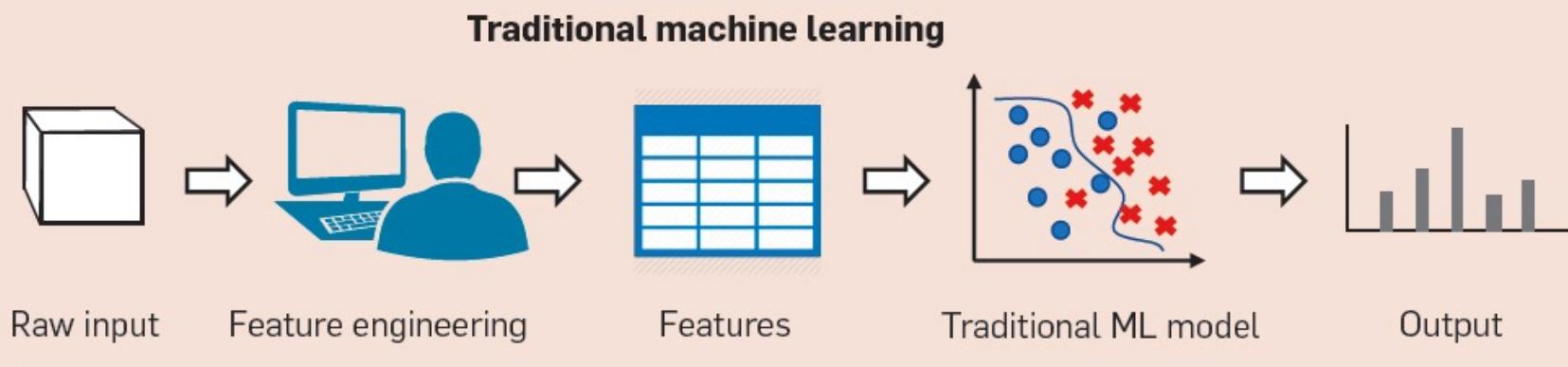
$$\mathcal{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{z}_n^\top \mathbf{z}_m$$



$\mathcal{X} \longrightarrow \mathcal{Z}$



Traditional machine learning



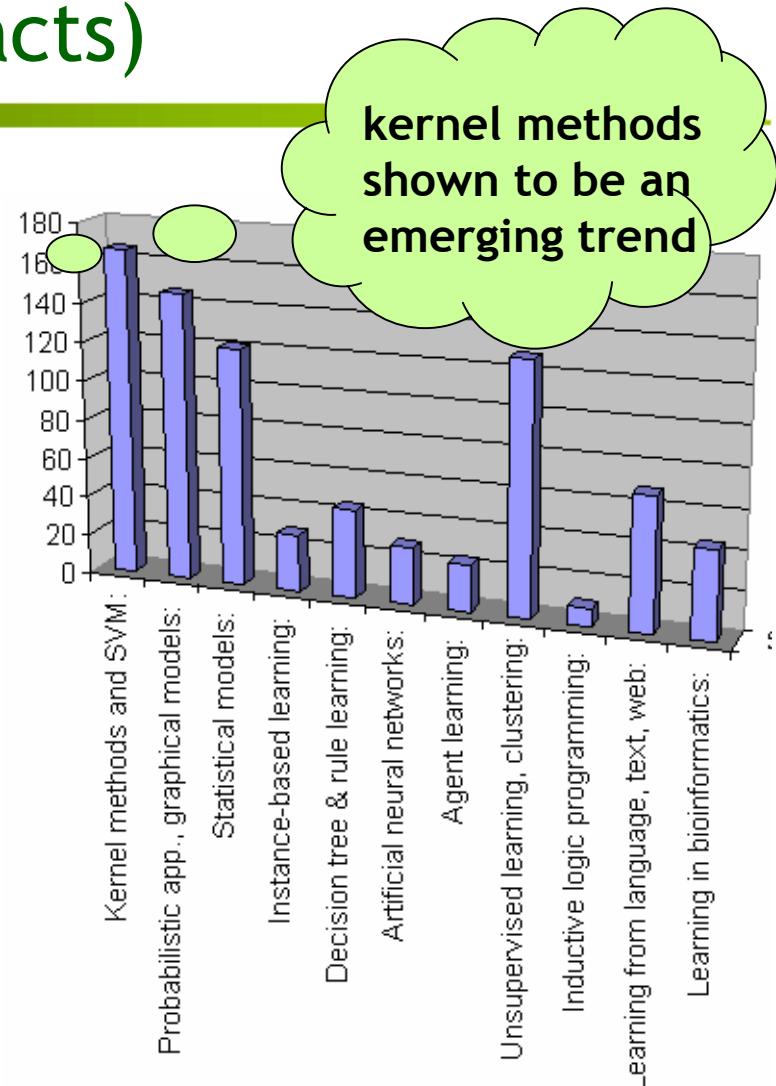
Outline

- Review of SVM in the previous lecture
- **A bit of history of kernel methods and SVM**
- Kernel methods
- Designing kernels

Very popular around 2006

ICML 2006 (720 abstracts)

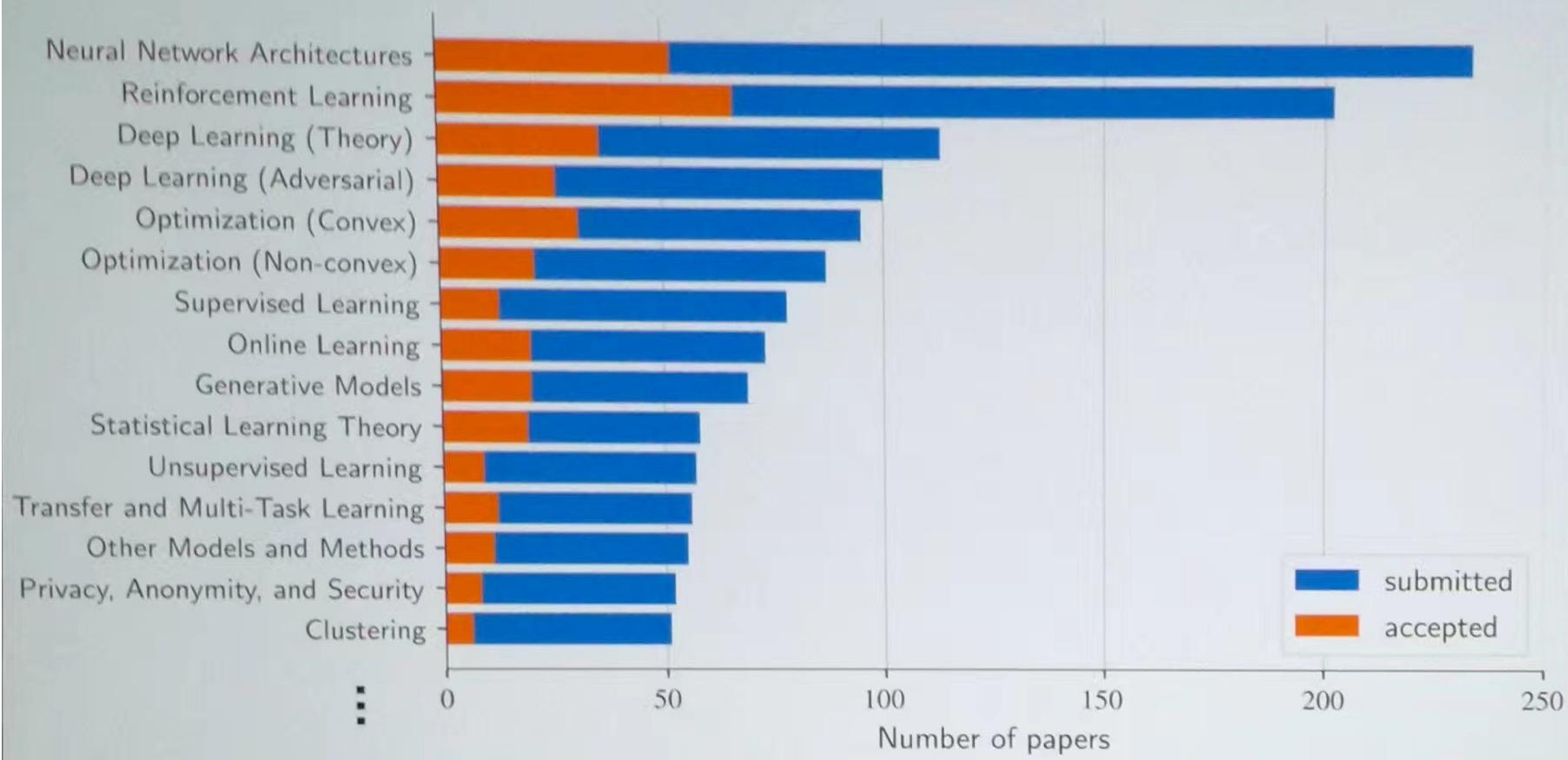
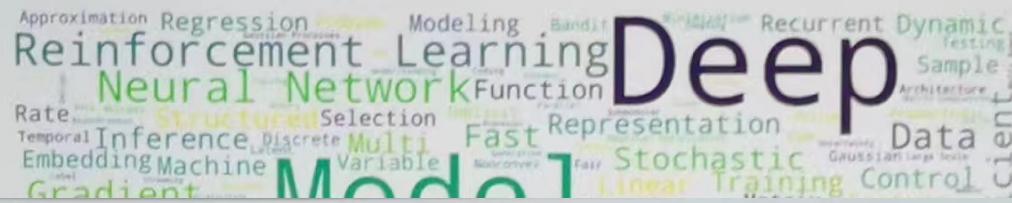
Kernel methods & SVM	166
Probabilistic, graphical models	146
Unsupervised learning, clustering	128
Statistical models	121
Language, Text & web	68
Learning in bioinformatics	45
ANN	29
ILP	9
CRF	13



ICML 2018

Program Committee

160 Area Chairs



A bit of history

- Aronszajn (1950) and Parzen (1962) were some of the first to employ positive definite **kernels** in statistics.
- Aizerman et al. (1964) used positive definite kernels in a way closer to the **kernel trick**.
- Boser et al. (1992) constructed the SVMs, a generalization of optimal hyperplane algorithm worked for **vectorial data**.
- Scholkopf (1997): kernels can work with **nonvectorial data** by providing a vectorial representation of the data in the feature space.
- Scholkopf et al. (1998): kernels can be used to build generalizations of any algorithm that can be carried out in terms of **dot products**.
- A large number of “**kernelizations**” of various algorithms since last 5 years.

Computation from the Z space

$$\mathcal{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{z}_n^\top \mathbf{z}_m$$

Constraints: $\alpha_n \geq 0$ for $n = 1, \dots, N$ and $\sum_{n=1}^N \alpha_n y_n = 0$



$$g(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{z} + b) \quad \text{need } \mathbf{z}_n^\top \mathbf{z}$$

where $\mathbf{w} = \sum_{\mathbf{z}_n \text{ is SV}} \alpha_n y_n \mathbf{z}_n$

and b : $y_m (\mathbf{w}^\top \mathbf{z}_m + b) = 1$ need $\mathbf{z}_n^\top \mathbf{z}_m$

Generalized inner product

Given two points \mathbf{x} and $\mathbf{x}' \in \mathcal{X}$, we need $\mathbf{z}^\top \mathbf{z}'$

Let $\mathbf{z}^\top \mathbf{z}' = K(\mathbf{x}, \mathbf{x}')$ (the kernel) “inner product” of \mathbf{x} and \mathbf{x}'

Second order transform:

Example: $\mathbf{x} = (x_1, x_2) \longrightarrow$ 2nd-order Φ

$$\mathbf{z} = \Phi(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$$

$$K(\mathbf{x}, \mathbf{x}') = \mathbf{z}^\top \mathbf{z}' = 1 + x_1 x'_1 + x_2 x'_2 +$$

$$x_1^2 x'^2_1 + x_2^2 x'^2_2 + x_1 x'_1 x_2 x'_2$$

Without transforming explicitly

- Can we compute $K(\mathbf{x}, \mathbf{x}')$ without transforming \mathbf{x} and \mathbf{x}' ?

Example:

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= (1 + \mathbf{x}^\top \mathbf{x}')^2 = (1 + x_1 x'_1 + x_2 x'_2)^2 \\ &= 1 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x'_1 x_2 x'_2 \end{aligned}$$



Inner product

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

$$(1, x'^2_1, x'^2_2, \sqrt{2}x'_1, \sqrt{2}x'_2, \sqrt{2}x'_1x'_2)$$

The polynomial kernel

$\mathcal{X} = \mathbb{R}^d$ and $\Phi : \mathcal{X} \rightarrow \mathcal{Z}$ is polynomial of order Q

The “equivalent” kernel $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^Q$

$$= (1 + x_1 x'_1 + x_2 x'_2 + \cdots + x_d x'_d)^Q$$

Adjust the scales:

$$K(\mathbf{x}, \mathbf{x}') = (a \mathbf{x}^\top \mathbf{x}' + b)^Q$$

We only need Z to exist

- If $K(x, x')$ is an inner product in some space Z, that is ok.
- For example,

$$K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$$

Infinite-dimensional \mathcal{Z} : take simple case

$$K(x, x') = \exp(-(x - x')^2)$$

$$= \exp(-x^2) \exp(-x'^2) \underbrace{\sum_{k=0}^{\infty} \frac{2^k (x)^k (x')^k}{k!}}_{\exp(2xx')}$$

- How do we know whether Z exists or not?
 1. By construction
 2. Math properties (Mercer's condition)
 3. Who cares?

Kernel SVM

$$\min_{\alpha} \frac{1}{2} \alpha^\top \underbrace{\begin{bmatrix} y_1 y_1 \mathbf{x}_1^\top \mathbf{x}_1 & y_1 y_2 \mathbf{x}_1^\top \mathbf{x}_2 & \dots & y_1 y_N \mathbf{x}_1^\top \mathbf{x}_N \\ y_2 y_1 \mathbf{x}_2^\top \mathbf{x}_1 & y_2 y_2 \mathbf{x}_2^\top \mathbf{x}_2 & \dots & y_2 y_N \mathbf{x}_2^\top \mathbf{x}_N \\ \dots & \dots & \dots & \dots \\ y_N y_1 \mathbf{x}_N^\top \mathbf{x}_1 & y_N y_2 \mathbf{x}_N^\top \mathbf{x}_2 & \dots & y_N y_N \mathbf{x}_N^\top \mathbf{x}_N \end{bmatrix}}_{\text{quadratic coefficients}} \alpha + \underbrace{(-1^\top) \alpha}_{\text{linear}}$$

subject to

$$\underbrace{\mathbf{y}^\top \alpha = 0}_{\text{linear constraint}}$$

$$\underbrace{0}_{\text{lower bounds}} \leq \alpha \leq \underbrace{\infty}_{\text{upper bounds}}$$

Kernel formulation

$$\left[\begin{array}{cccc} y_1 y_1 K(\mathbf{x}_1, \mathbf{x}_1) & y_1 y_2 K(\mathbf{x}_1, \mathbf{x}_2) & \dots & y_1 y_N K(\mathbf{x}_1, \mathbf{x}_N) \\ y_2 y_1 K(\mathbf{x}_2, \mathbf{x}_1) & y_2 y_2 K(\mathbf{x}_2, \mathbf{x}_2) & \dots & y_2 y_N K(\mathbf{x}_2, \mathbf{x}_N) \\ \dots & \dots & \dots & \dots \\ y_N y_1 K(\mathbf{x}_N, \mathbf{x}_1) & y_N y_2 K(\mathbf{x}_N, \mathbf{x}_2) & \dots & y_N y_N K(\mathbf{x}_N, \mathbf{x}_N) \end{array} \right]$$

quadratic coefficients

Solution to kernel SVM

Express $g(\mathbf{x}) = \text{sign} (\mathbf{w}^\top \mathbf{z} + b)$ in terms of $K(-, -)$

$$\mathbf{w} = \sum_{\mathbf{z}_n \text{ is SV}} \alpha_n y_n \mathbf{z}_n \implies g(\mathbf{x}) = \text{sign} \left(\sum_{\alpha_n > 0} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}) + b \right)$$

$$\text{where } b = y_m - \sum_{\alpha_n > 0} \alpha_n y_n K(\mathbf{x}_n, \mathbf{x}_m)$$

for any support vector ($\alpha_m > 0$)

Mercer's condition

$K(\mathbf{x}, \mathbf{x}')$ is a valid kernel iff

1. It is symmetric
2. The matrix:

$$\begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \dots & K(\mathbf{x}_1, \mathbf{x}_N) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & \dots & K(\mathbf{x}_2, \mathbf{x}_N) \\ \dots & \dots & \dots & \dots \\ K(\mathbf{x}_N, \mathbf{x}_1) & K(\mathbf{x}_N, \mathbf{x}_2) & \dots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

is **positive semi-definite**

for any $\mathbf{x}_1, \dots, \mathbf{x}_N$ (Mercer's condition)

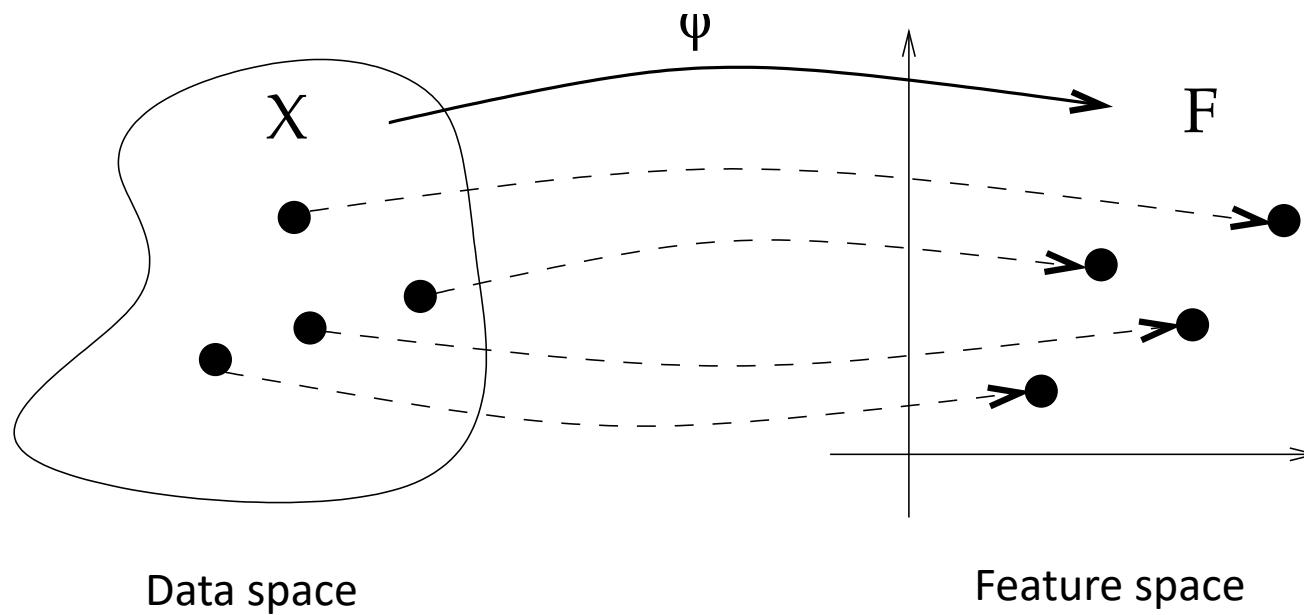
Then, you can design your own kernel as above.

Outline

- Review of SVM in the previous lecture
- A bit of history of kernel methods and SVM
- **Kernel methods**
- Designing kernels

Kernels

- a general framework to represent data and must satisfy some math conditions that give them properties useful to understand kernel methods and kernel design



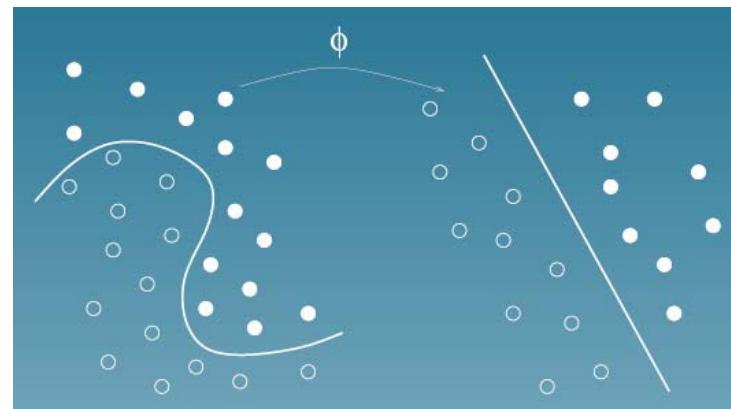
The issue of data representation

- Let $\mathcal{S} = (x_1, \dots, x_n)$ a set of n objects to be analyzed.
- Suppose that each object x_i is an element of a set \mathcal{X} , which may be images, molecules, texts, etc.
- Majority of data analysis methods represent \mathcal{S} by:
 - defining a representation $\phi(x) \in \mathcal{F}$ for each object $x \in \mathcal{X}$, where $\phi(x)$ can be a real-valued vector ($\mathcal{F} = \mathbb{R}^p$) or a finite-length string, or more complex representation.
 - representing \mathcal{S} by a set of representations of the objects

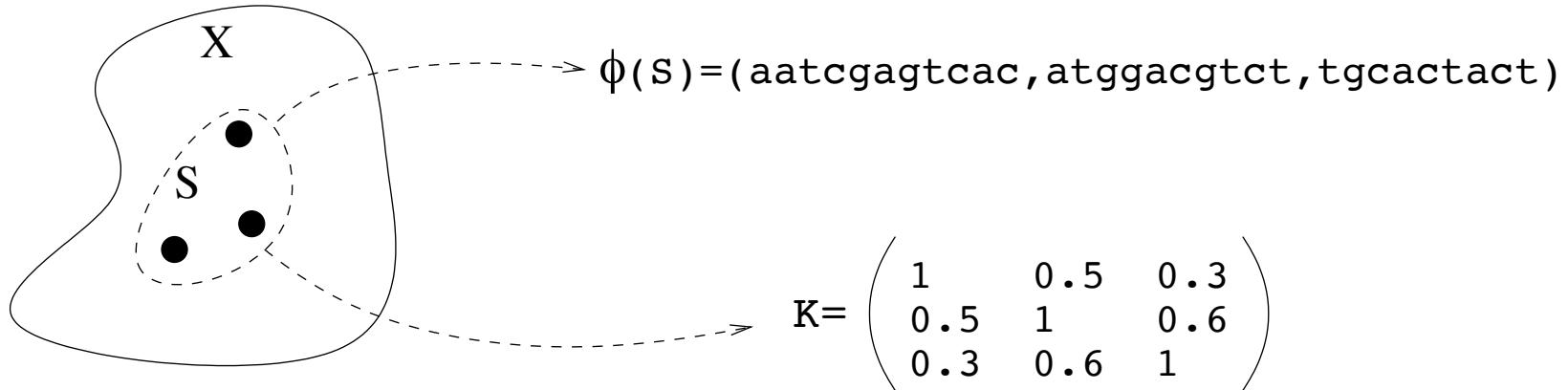
$$\phi(\mathcal{S}) = (\phi(x_1), \dots, \phi(x_n))$$

The idea of kernel representation

- Data are not represented individually anymore, but only through a set of pairwise comparisons.
- Instead of using a mapping $\phi: \mathcal{X} \rightarrow \mathcal{F}$ to represent each object $\mathbf{x} \in \mathcal{X}$ by $\phi(\mathbf{x}) \in \mathcal{F}$, a real-valued “comparison function” (called **kernel**) $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is used, and the data set \mathcal{S} is represented by the $n \times n$ matrix (Gram matrix) of pairwise comparisons $k_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$.
- A main question is how to find a kernel k such that, in the new space, problem solving is easier (e.g. linear).
- All **kernel methods** have two parts:
 - (1) find such a kernel k , and
 - (2) process such Gram matrices.

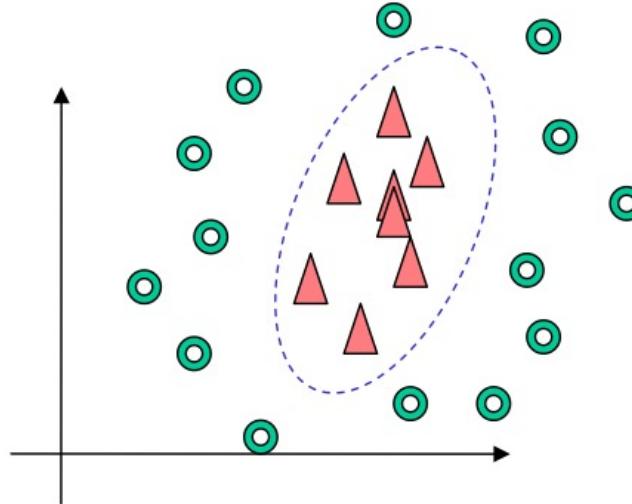


Kernel representation: example



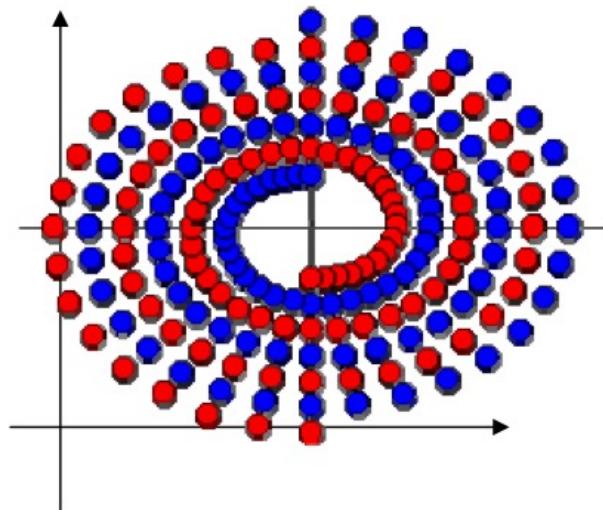
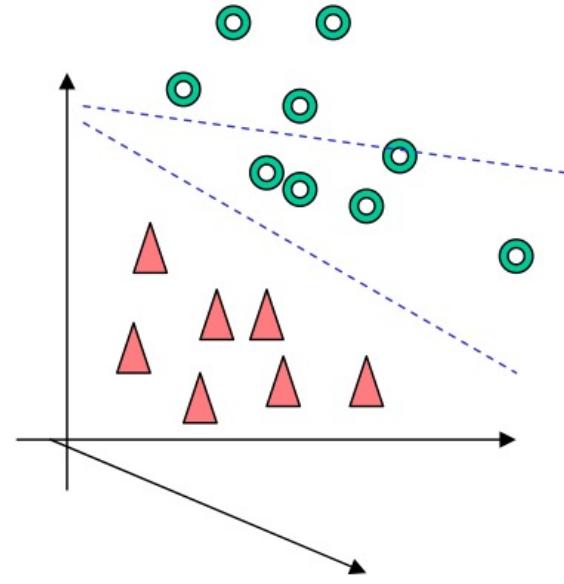
- X is the set of all oligonucleotides, S consists of three oligonucleotides.
- Traditionally, each oligonucleotide is represented by a sequence of letters.
- In kernel methods, S is represented as a matrix of pairwise similarity between its elements.

Examples of kernels



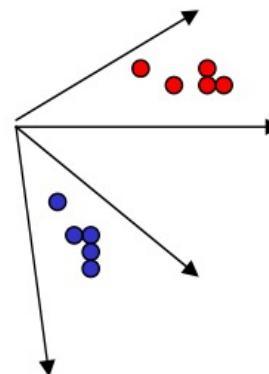
by polynomial
kernel ($n=2$)

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



by RBF kernel
($n=2$)

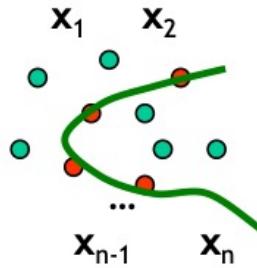
$$\phi: (x_1, x_2) \mapsto \exp\left(-\frac{d(x_1, x_2)^2}{2\sigma^2}\right)$$



(Jean-Michel Renders, ICFT'04)

General scheme of kernel methods

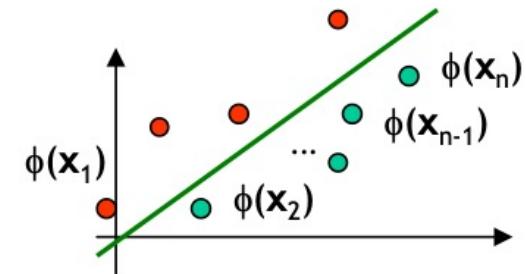
Input space \mathcal{X}



inverse map ϕ^{-1}

$\phi(\mathbf{x})$

Feature space \mathcal{F}



$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

Gram matrix $K_{n \times n}$

kernel function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

kernel-based algorithm on K

- Map the data from \mathcal{X} into a (high-dimensional) vector space, the feature space \mathcal{F} , by applying the **feature map** ϕ on the data points \mathbf{x} .
- Find a **linear** (or other easy) pattern in \mathcal{F} using a well-known algorithm (that works on the Gram matrix).
- By applying the **inverse map**, the linear pattern in \mathcal{F} can be found to correspond to a complex pattern in \mathcal{X} .
- This implicitly by only making use of inner products in \mathcal{F} (**kernel trick**).

Outline

- Review of SVM in the previous lecture
- A bit of history of kernel methods and SVM
- Kernel methods
- **Designing kernels**

How to chose kernels

- There is no absolute rules for choosing the right kernel, adapted to a particular problem
- Kernel design can start from the desired feature space, from combination or from data
- Some considerations are important:
 - Use kernel to introduce a priori (domain) knowledge
 - Be sure to keep some information structure in the feature space
 - Experimentally, there is some “robustness” in the choice, if the chosen kernels provide an acceptable trade-off between
 - simpler and more efficient structure (e.g. linear separability), which requires some “explosion”
 - Information structure preserving, which requires that the “explosion” is not too strong.

Some operations on kernels

- The class of kernel functions on \mathcal{X} has several useful closure properties. It is a convex cone, which means that if k_1 and k_2 are two kernels, then any **linear combination**, $\lambda_1 k_1 + \lambda_2 k_2$, with $\lambda_1, \lambda_2 \geq 0$ is a kernel.
- If k_1 and k_2 are two kernels, then $k(\mathbf{x}, \mathbf{x}') := k_1(\mathbf{x}, \mathbf{x}') k_2(\mathbf{x}, \mathbf{x}')$ is also a kernel. Several other operations are possible.
- Some other operations are forbidden: if k is a kernel, then $\log(k)$ is not positive definite in general, and neither is k^β for $0 < \beta < 1$. However, these two operations can be linked.

Combining kernels

- Multiple kernel matrices k_1, k_2, \dots, k_c for the same set of objects are available.
- We may design a single kernel k from several basic kernels k_1, k_2, \dots, k_c . A simple way to achieve this is to take the **sum of the kernels**:

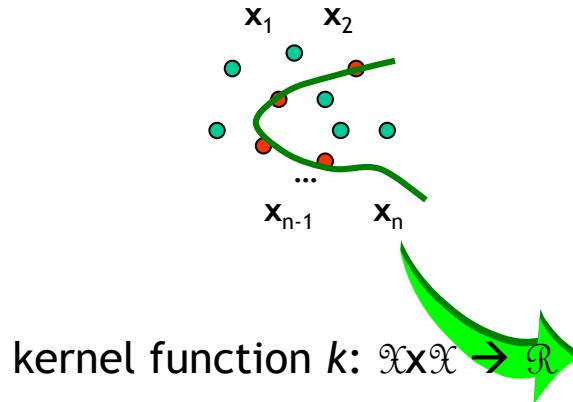
$$k = \sum_{i=1}^c k_i$$

- Multiple kernel matrices k_1, k_2, \dots, k_c for the same set of objects are available.

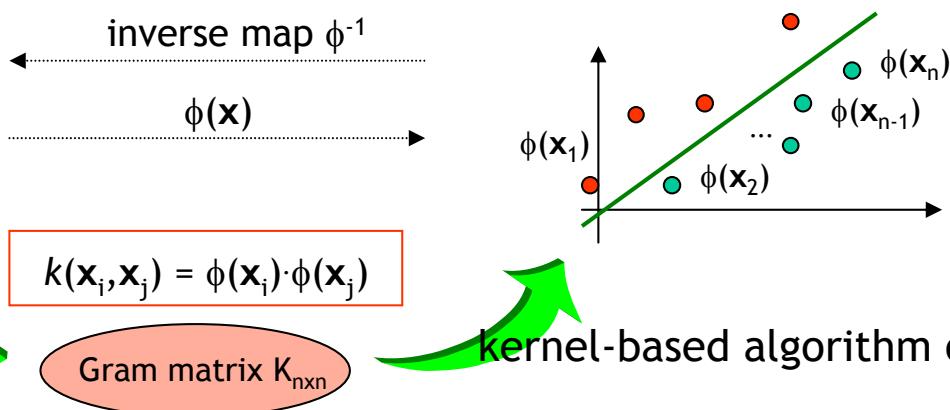
$$k = \sum_{i=1}^c \mu_i k_i$$

Summary

Input space \mathcal{X}



Feature space \mathcal{F}



- Map the data from \mathcal{X} into a (high-dimensional) vector space, the feature space \mathcal{F} , by applying the **feature map** ϕ on the data points x .
- Find a **linear** (or other easy) **pattern in \mathcal{F}** using a well-known algorithm (that works on the Gram matrix).
- By applying the **inverse map**, the linear pattern in \mathcal{F} can be found to correspond to a complex pattern in \mathcal{X} .
- This implicitly by only making use of inner products in \mathcal{F} (**kernel trick**)

Thank you!