

Aprenentatge Automàtic 2

GCED

Lluís A. Belanche
belanche@cs.upc.edu



Soft Computing Research Group
Dept. de Ciències de la Computació (Computer Science)
Universitat Politècnica de Catalunya

2021-2022

LECTURE 5: Kernels redux: definitions, properties, examples in \mathbb{R}^p . Kernel design for data objects not in \mathbb{R}^p .

Kernel design: practical issues

Euclidean space \mathbb{R}^p , but not only ...

- Kernels on real vectors (whole families)
- Kernels on binary vectors (bitstrings = sets)
- General structured kernels:
 - All-subsets kernel
 - Convolution kernels
- Kernels on discrete structures:
 - Tree kernels
 - Graph kernels
- Kernels on distributions (generative kernels):
 - P-kernels
 - Marginalized kernels
- String kernels (text)

... and many others (functional data, categorical data, permutations, ...)

Kernel design: practical issues

All-subsets kernel

Consider a feature space with one feature for every subset $A \subseteq \{1, \dots, d\}$ of the input variables:

For $\mathbf{x} \in \mathbb{R}^p$, feature A is given by $\phi_A(\mathbf{x}) := \prod_{i \in A} x_i$ (note $\phi_\emptyset(\mathbf{x}) = 1$)

The kernel is defined by the mapping $\phi : \mathbf{x} \rightarrow (\phi_A(\mathbf{x}))_{A \subseteq \{1, \dots, d\}}$

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \sum_{A \subseteq \{1, \dots, d\}} \phi_A(\mathbf{x}) \phi_A(\mathbf{x}') \\ &= \sum_{A \subseteq \{1, \dots, d\}} \prod_{i \in A} x_i x'_i = \prod_{i=1}^d (1 + x_i x'_i) \end{aligned}$$

The last step is obtained by expanding $(1 + x_1 x'_1)(1 + x_2 x'_2) \dots (1 + x_d x'_d)$

Kernel design: practical issues

All-subsets kernel

We have the freedom to downplay some features (and thus emphasize others) by introducing weighting factors $w_i \geq 0$ for each feature i :

$$\phi_A(\mathbf{x}) := \prod_{i \in A} \sqrt{w_i} x_i$$

therefore

$$k_w(\mathbf{x}, \mathbf{x}') = \prod_{i=1}^d (1 + w_i x_i x'_i)$$

Kernel design: practical issues

Bitstring/Binary variables/Sets

Let $x, x' \in \{0, 1\}^p$, representing absence/presence of a binary trait:

1. The Simple Matching Coefficient (SMC) is the fraction of 1 – 1 matches, and it is a kernel on $\{0, 1\}^p$. ¶
2. The Jaccard Coefficient is the fraction of 1 – 1 matches among the traits present in either data vector, and it is a kernel on $\{0, 1\}^p$. ¶

Kernel design: practical issues

Bitstring/Binary variables/Sets

- Given two sets $A, B \subset U$, where U is finite, consider

$$k(A, B) := \frac{1}{|U|} \sum_{a \in A} \sum_{b \in B} k_{\text{base}}(a, b)$$

- If k_{base} is the overlap kernel $k(a, b) = \begin{cases} 1 & \text{if } a = b; \\ 0 & \text{otherwise.} \end{cases}$

we get $k(A, B) = \frac{|A \cap B|}{|U|}$, the equivalent to the SMC.

- The equivalent to the Jaccard kernel would be $k(A, B) = \frac{|A \cap B|}{|A \cup B|}$.

Kernel design: practical issues

Generative kernels (I)

Given a conditional probability distribution on $\mathcal{X}|\mathcal{Z}$, we can compare data points by assigning a high value if both have high (conditional) probability:

$$k(\mathbf{x}, \mathbf{x}') := \sum_{z \in \mathcal{Z}} p(\mathbf{x}|z)p(\mathbf{x}'|z)P(z) \quad \text{discrete case}$$

$$k(\mathbf{x}, \mathbf{x}') := \int_{\mathcal{Z}} p(\mathbf{x}|z)p(\mathbf{x}'|z)p(z) dz \quad \text{continuous case}$$

The feature maps are $(\phi(\mathbf{x}))_z = p(\mathbf{x}|z)\sqrt{p(z)}$

Hint: $p(\mathbf{x}, \mathbf{x}', z) = p(\mathbf{x}, \mathbf{x}'|z)p(z) = p(\mathbf{x}|z)p(\mathbf{x}'|z)p(z)$

Kernel design: practical issues

Generative kernels (II)

Given a probability distribution on $\mathcal{X} \times \mathcal{Z}$, and a kernel on $\mathcal{X} \times \mathcal{Z}$ pairs, we can define:

$$k(\mathbf{x}, \mathbf{x}') := \sum_z \sum_{z'} k((\mathbf{x}, z), (\mathbf{x}', z')) p(\mathbf{x}|z) p(\mathbf{x}'|z')$$

Typical applications of generative kernels are found in **graphical models**:

- \mathcal{X} are the **observed** variables and \mathcal{Z} are the **hidden** (latent) variables
- A kernel for the observed ones is obtained by taking the expectation w.r.t. the hidden ones (*marginalizing* them away)
- Applications: HMMs for DNA sequences, or stochastic context-free grammars for RNA sequences –see e.g. “Kernel methods in genomics and computational biology” by J.P. Vert

Kernel design: practical issues

The Spectrum (aka n -Gram) kernel

- Let Σ be a finite alphabet: an n -Gram is a block of n adjacent characters in Σ

$$\text{Define } k(x, x') := \sum_{s \in \Sigma^n} |s \in x| \cdot |s \in x'|$$

Example: Word aababc in alphabet $\Sigma = \{a, b, c\}$, $n = 2$:

aa	ab	ac	ba	bb	bc	...
1	2	0	1	0	1	...

-
- The feature vectors are sparse and the kernel can be computed in $O(|x| + |x'|)$ time and memory

Kernel design: practical issues

Kernels from graphs

- Consider a graph $G = (V, E)$, where the set of vertices (nodes) V are the data points and E is the set of edges. Call $N = |V|$, the number of nodes
- The idea is to compute a (base) matrix $S_{N \times N}$ whose entries are the weights of the edges and consider $S^2 = SS$ (S is only assumed to be symmetric)
- Typical use: **connectivity matrix** of G : the (i, j) element of S^2 is the number of paths of length exactly 2 between i and j

Examples:

1. protein-protein interactions
2. people-to-people interactions

In 2, the (i, j) element of S^2 is the number of common friends between data points i and j (it can be thought of as a measure of their similarity)

Kernel design: practical issues

Kernels from graphs

Notes:

- The entries of S may be real-valued numbers (e.g., symmetric bounded similarities)
- Higher powers of S measure higher-order similarities
- Only the even powers are guaranteed to be PSD

Consider, for a given $\lambda \in (0, 1)$:

$$\sum_{k=0}^{\infty} \frac{1}{k!} \lambda^k S^k =: \exp(\lambda S)$$

1. If S is symmetric, then $S = U\Lambda U^T$ (spectral decomposition), so $S^2 = (U\Lambda U^T)(U\Lambda U^T) = U\Lambda^2 U^T$.
2. In general, we have $S^k = U\Lambda^k U^T$ and therefore:

$$K := \exp(\lambda S) = U \exp(\lambda \Lambda) U^T$$

is an example of a **diffusion** kernel.

Example in a real application domain

Handling missing values in microbiology

- Modern modelling problems are difficult for a number of reasons, including the challenge of dealing with a significant amount of missing information
- **Missing values** almost always represent a serious problem because they force to preprocess the dataset and a good deal of effort is normally put in this part of the modelling
- In order to process such datasets with kernel methods, an imputation procedure is then deemed a necessary but demanding step

Example in a real application domain

Handling missing values in microbiology

- The study of fecal source pollution in waterbodies is a major problem in ensuring the welfare of human populations
- **Microbial source tracking** (MST) methods attempt to identify the source of contamination, allowing for improved risk analysis and better water management
- The available dataset includes 148 observations about 10 chemical, microbial, and eukaryotic markers of fecal pollution in water
- All variables (except the class variable) are binary, i.e., they signal the presence or absence of a particular marker

Example in a real application domain

Handling missing values in microbiology

Origin	HF183	HF134	CF128	Humito	Pomito	Bomito	ADO	DEN
Human :50	0 :68	0 :81	0 :104	0 :35	0 :83	0 :78	0 :56	0 :80
Cow :26	1 :40	1 :26	1 : 5	1 :79	1 :32	1 :32	1 :59	1 :34
Poultry:31	? :31	? :32	? :30	? :25	? :24	? :29	? :24	? :25
Pig :32								

Summary (counts) table for the full dataset. The first column is the target class. The symbol ? denotes a missing value.

The percentage of missing values is around 19.8 %, and all the predictive variables have percentages between 17 % and 23 %

Example in a real application domain

Handling missing values in microbiology

Theorem. Let the symbol $?$ denote a missing element, for which only equality is defined, and \mathcal{X} a finite discrete set. Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a kernel in \mathcal{X} and P a probability mass function in \mathcal{X} . Then the function $k^?(x, y)$ given by

$$k^?(x, y) := \begin{cases} k(x, y), & \text{if } x, y \neq ? ; \\ g(x) := \sum_{y' \in X} P(y')k(x, y'), & \text{if } x \neq ? \text{ and } y = ?; \\ g(y) := \sum_{x' \in X} P(x')k(x', y), & \text{if } x = ? \text{ and } y \neq ?; \\ G := \sum_{x' \in X} P(x') \sum_{y' \in X} P(y')k(x', y'), & \text{if } x = y = ? \end{cases}$$

is a kernel in $\mathcal{X} \cup \{?\}$.

Example in a real application domain

Handling missing values in microbiology

For the particular case of binary variables $x, y \in \{v_1, v_2\}$, a convenient approach is to define the kernel:

$$k_{0/1}(x, y) := \mathbb{I}_{\{x=y\}}$$

where

$$\mathbb{I}_{\{z\}} = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{if } z \text{ is false} \end{cases}$$

Example in a real application domain

Handling missing values in microbiology

Consider now $\mathbf{x}, \mathbf{y} \in \{0, 1, ?\}^p$. When we apply the Theorem to this kernel, we obtain an extended multivariate kernel:

$$\mathcal{K}_1(\mathbf{x}, \mathbf{y}) := \frac{1}{p} \sum_{i=1}^p \begin{cases} 1 & \text{if } x_i = y_i = 1 ; \\ P_i(x_i), & \text{if } x_i \neq ? \text{ and } y_i = ?; \\ P_i(y_i), & \text{if } x_i = ? \text{ and } y_i \neq ?; \\ (P_i(0))^2 + (P_i(1))^2, & \text{if } x_i = y_i = ?; \\ 0, & \text{otherwise} \end{cases}$$

This kernel is a generalization of the classical *Simple Matching Coefficient*, proposed by Sokal and Michener for numerical taxonomy

Example in a real application domain

Handling missing values in microbiology

Alternatives???

Given $\mathbf{x}, \mathbf{y} \in$.

Let $c(\mathbf{x})$ be the set of completions of \mathbf{x} . Given two vectors $\mathbf{x}, \mathbf{y} \in \{0, 1, ?\}^p$, the function

$$\mathcal{K}_2(\mathbf{x}, \mathbf{y}) := \frac{1}{|c(\mathbf{x})||c(\mathbf{y})|} \sum_{\mathbf{x}' \in c(\mathbf{x})} \sum_{\mathbf{y}' \in c(\mathbf{y})} k(\mathbf{x}', \mathbf{y}')$$

is a kernel in $\{0, 1, ?\}^p$.

Example in a real application domain

Handling missing values in microbiology

Approach	C	10x10cv	10x10cv for each class			
			Human	Cow	Poultry	Swine
\mathcal{K}_1	2.0	79.3	95.4	64.5	75.2	69.4
\mathcal{K}_2	1.6	78.2	92.6	62.8	71.8	74.2
MI-1	1.0	79.9	92.7	66.4	69.4	80.2
MI-2	1.0	79.0	94.5	57.5	70.8	78.8

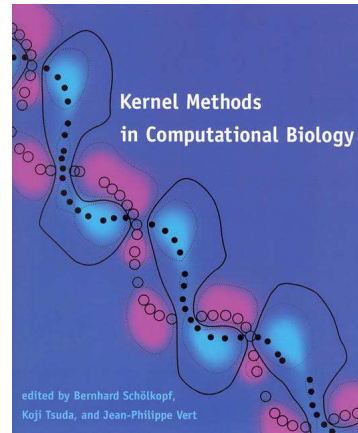
Mean 10x10cv accuracies for the four approaches to handle missing values. Also shown are best cost parameter C and detailed class performance.

(joint work with G. Nebot, T. Aluja and V. Kobayashi)

Kernel design: practical issues

More Kernels!

Kernels abound in computational biology and computational chemistry (e.g., phylogenetic profiles, protein 3D structures)



Example: the prediction of **interacting proteins** to reconstruct an interaction network can be posed as a binary classification problem: given a pair of proteins, do they interact or not?

→ we need kernel between *pairs* of proteins!

Kernel design: practical issues

More Kernels!

The available data is about each single protein; it is then natural to derive kernels for **pairs** of proteins k_{pair} from any kernel k for **single** proteins:

$$k_{\text{pair}}((A, B), (C, D)) := k(A, C)k(B, D) + k(A, D)k(B, C)$$

(there is usually no order in a protein pair, so we try both matches)

-
- Using Product Kernels to Predict Protein Interactions. *Advances in Biochemical Engineering/Biotechnology* (110), pp 215-245 (2007)
 - Kernel methods for predicting protein-protein interactions. *Bioinformatics*. 2005

Kernel design: practical issues

Conclusions

- The power of kernel methods partly relies in the ability to process virtually any sort of data as soon as a valid kernel is defined
- Importance of designing kernels that do not constitute explicit inner products between objects, and therefore fully exploit the kernel trick
- Possibility of learning the kernel function (or the kernel matrix) from the training data
- Theoretical analyses are needed on the implications of the kernel choice for the success of specific kernel-based methods