

# Assignment 4: Kernels and Regularization

Submission: Thursday April 10th  
Maximum of 2 students per group

Prof. Fabio A. González  
Machine Learning - 2008-I  
Maestría en Ing. de Sistemas y Computación

---

## 1. Regression on strings

- (a) Implement a function that calculates a kernel over fixed-length strings,

$$k : \Sigma^d \times \Sigma^d \rightarrow \mathbb{R},$$

which counts the number of **coincidences** between two strings.

- (b) Implement the kernel ridge-regression (KRR) algorithm.
- (c) Use the KRR implementation and the kernel  $k$  to train a model using the training data set in <http://dis.unal.edu.co/~fgonza/courses/2008-I/ml/assign4-train.txt>. Evaluate the error of the model on the training data set. Plot the output of the model on the training data along with the real output values.
- (d) Evaluate the trained model on the test data set <http://dis.unal.edu.co/~fgonza/courses/2008-I/ml/assign4-test.txt>. Plot the results and discuss them.
- (e) Build a new kernel,  $k'$ , composing the kernel  $k$  with a polynomial kernel. Repeat steps (c) and (d).
2. Let  $\mathbf{x} = \{x_1, \dots, x_n\}$  be a subset of a input data set  $X$ . Consider a kernel function  $k : X \times X \rightarrow \mathbb{R}$ , which induces a feature space  $\phi(X)$ :

- (a) Deduce an expression, that allows to calculate the average distance to the center of mass of the image of set  $\mathbf{x}$  in the feature space:

$$\frac{1}{n} \sum_{i=1}^n \|\phi(x_i) - \phi_S(\mathbf{x})\|_{\phi(X)},$$

where the center of mass is defined as

$$\phi_S(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \phi(x_i).$$

- (b) Use previous expression to calculate the average distance to the center of mass of the following point set in  $\mathbb{R}^2$ ,  $\mathbf{x} = \{(0, 1), (-1, 3), (2, 4), (3, -1), (-1, -2)\}$ , in the feature spaces induced by the following kernels:
- $k(x, y) = \langle x, y \rangle$
  - $k(x, y) = \langle x, y \rangle^2$
  - $k(x, y) = (\langle x, y \rangle + 1)^5$

- iv. Gaussian kernel
- 3. Controlling the model complexity
  - (a) Download the Wisconsin Breast Cancer data set from [http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](http://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)) and divide it in a training set and a test set (50/50).
  - (b) Train a SVM using a linear kernel. Find an optimal complexity parameter,  $C$ , plotting the training and test error vs. the complexity parameter. Discuss the results.
  - (c) Repeat item (b) using a Gaussian kernel.
- 4. [Alp04] Exercise 9 (chap. 4, page 83)

## References

- [Alp04] Alpaydin, E. 2004 Introduction to Machine Learning (Adaptive Computation and Machine Learning). The MIT Press.