

kernlab: Machine Learning with kernels in R

Alexandros Karatzoglou¹
joint work with Alex Smola and Kurt Hornik

¹INSA de Rouen
LITIS
France

December 12, 2008

Outline

1 kernlab Introduction

2 Demo

kernlab



kernlab R package:

- Contains a wide range of kernel-based Machine Learning methods
- Uses modern Open Source R environment
- Extensible by exploiting the inherent modularity of kernel methods
- GPL 2

R

The screenshot displays the RStudio environment with several open windows:

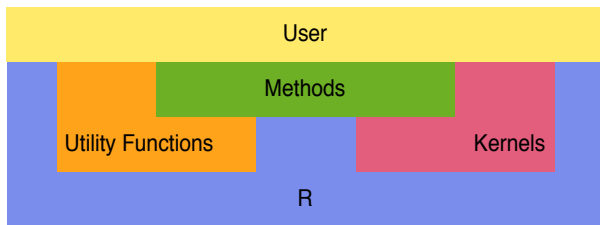
- R Console:** Contains R code for plotting a 3D surface. The code includes:


```
rgl.srs$ylen <- ylen[2] - ylen[1] + 1
rgl.srs$colorlut <- terrain.colors(ylen)
rgl.srs$col <- colorlut[y - ylen[1] + 1]
rgl.srs$rgl.clear()
rgl.srs$rgl.surface(x, z, y, color = col)
```
- R Data Editor:** A window showing a table of data with columns 'height' and 'weight'. The data ranges from 58 to 72 for height and 115 to 164 for weight.
- Workspace Browser:** Displays the objects in the workspace, including 'data.frame', 'factor', 'numeric', 'list', and 'data.frame' objects with their dimensions and levels.
- R Package Manager:** Shows the status of installed packages, including 'graphics', 'grid', 'lattice', 'methods', and 'rgl'.
- 3D Surface Plot:** A window titled 'RGL device 1 (active)' showing a 3D surface plot of the data, with axes labeled 'x', 'y', and 'z'.

R

- Environment for statistical data analysis, inference and visualization.
- Ports for Unix, Windows and MacOSX
- Highly extensible through user-defined functions
- Generic functions and conventions for standard operations like plot, predict etc.
- ~ 1200 add-on packages contributed by developers from all over the world
- e.g. Multivariate Statistics, Machine Learning, Natural Language Processing, Bioinformatics (Bioconductor).
- Interfaces to C, C++, Fortran, Java

kernlab Package Content



- 9 kernel functions, 4 kernel expression functions
- 18 kernel-based ML methods:
 - Regression, Quantile regression, Classification and Novelty detection (SVM, RVM, Gaussian Processes)
 - Clustering (kernel k -means, Spectral Clustering)
 - Ranking and dimensionality reduction (kernel PCA, kernel CCA, etc.) kernel-based two sample test (KMMD)
- Utility Functions:
 - Quadratic Problem solver
 - Incomplete Cholesky Decomposition

Some kernel functions in `kernlab`

- Linear kernel
- Gaussian radial basis kernel
- Hyperbolic tangent kernel
- String kernel
- Polynomial kernel
- ...

Runtime

Inputs

- `kernlab` methods take kernel function and data or kernel matrix as input
- Function parameters are given as options

Outputs

- Return objects that contain model parameters.
- Objects can be used with generic functions such as `predict`, `plot` etc.

Performance

- Nystrom Method for eigendecomposition based methods (e.g. spectral clustering). Eigenvectors can be calculated using only part of the kernel matrix yielding speedups and better scaling $O(n^3) + O(nN)$.
- Build in kernel function computations are vectorized
- C++ on performance bottlenecks
- Use of scalable algorithms
- ...

Features

Most function implement additional features such as e.g.

- Cross-validation (svm, GP's)
- Model-selection in particular hyper-parameter estimation
- Data Scaling
- Extensive documentation of every function and options

Demo