

# A Practical Use of Regularization for Supervised Learning with Kernel Methods

Marco Prato<sup>1</sup>, Luca Zanni<sup>1</sup>

1. Dipartimento di Scienze Fisiche, Informatiche e Matematiche,  
Università di Modena e Reggio Emilia, Italy

Optimization Techniques for Inverse Problems II  
Modena (Italy), September 20-21, 2012

## Abstract

In several supervised learning applications, it happens that reconstruction methods have to be applied repeatedly before being able to achieve the final solution. In these situations, the availability of learning algorithms able to provide effective predictors in a very short time may lead to remarkable improvements in the overall computational requirement. Here we consider the kernel ridge regression problem and we look for predictors given by a linear combination of kernel functions plus a constant term, showing that an effective solution can be obtained very fastly by applying specific regularization algorithms directly to the linear system arising from the Empirical Risk Minimization problem.

## Learning with the quadratic loss

**Regularized least squares** (RLS) for learning: given a training set

$$S = \{(\mathbf{x}_i, y_i) : i = 1, \dots, n\} \subset X \times Y, \quad X \subset \mathbb{R}^d, \quad Y \subset \mathbb{R}$$

find the decision function  $f : X \rightarrow Y$  to predict the label  $y$  of new examples  $\mathbf{x}$  by solving

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \quad (1)$$

where: -  $\lambda$  is a positive regularization parameter,

-  $\mathcal{H}$  is a Reproducing Kernel Hilbert Space with kernel  $K$  [1].

**Representer Theorem**: the solution of (1) in a RKHS assumes the form

$$f_0(\mathbf{x}) = \sum_{i=1}^n c_i K(\mathbf{x}, \mathbf{x}_i), \quad (2)$$

where  $\mathbf{c}$  is the solution of the linear system [2]

$$(\mathbf{K} + n\lambda \mathbf{I})\mathbf{c} = \mathbf{y}, \quad (3)$$

being  $\mathbf{K}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ .

More general prediction function (SVM choice):

$$f(\mathbf{x}) = f_0(\mathbf{x}) + b, \quad b \in \mathbb{R} \quad (4)$$

Two different generalizations of (1) [3,4,5]:

$$(a) \quad \min_{f_0 \in \mathcal{H}, b \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (y_i - f_0(\mathbf{x}_i) - b)^2 + \lambda \|f_0\|_{\mathcal{H}}^2 \quad (5)$$

(constant  $b$  not penalized)

(constant  $b$  penalized)

$$(b) \quad \min_{f_0 \in \mathcal{H}, b \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (y_i - f_0(\mathbf{x}_i) - b)^2 + \lambda (\|f_0\|_{\mathcal{H}}^2 + b^2) \quad (6)$$

leading, respectively, to the linear systems (where  $\mathbf{1} = (1, \dots, 1)^t$ )

$$\begin{cases} (\mathbf{K} + n\lambda \mathbf{I})\mathbf{c} + b\mathbf{1} = \mathbf{y} \\ \mathbf{1}^t \mathbf{c} = 0 \end{cases} \quad (7) \quad \begin{cases} (\mathbf{K} + \mathbf{1}\mathbf{1}^t + n\lambda \mathbf{I})\mathbf{c} + b\mathbf{1} = \mathbf{y} \\ b = \mathbf{1}^t \mathbf{c} \end{cases} \quad (8)$$

CPU time required by these approaches:  $\mathcal{O}(4n^3/3)$ , due to the eigendecomposition of  $\mathbf{K}$  or  $\mathbf{K} + \mathbf{1}\mathbf{1}^t$ .

## Tikhonov and Conjugate Gradient

**From  $\mathcal{H}$  to  $\mathbb{R}^n$** : problem (3) is the Lavrentiev [6] regularized version of the linear system  $\mathbf{K}\mathbf{c} = \mathbf{y}$ , that arises from the minimization of the Least Squares functional when the desired solution is in the form (2).

With solutions in the form (4), such linear system becomes

$$\widetilde{\mathbf{K}}\widetilde{\mathbf{c}} = \mathbf{y}, \quad \widetilde{\mathbf{K}} = (\mathbf{K} \quad \mathbf{1}), \quad \widetilde{\mathbf{c}} = (c_1, \dots, c_n, b)^t \quad (9)$$

Instead of adding a penalty term on  $f$  – as in (5), (6) –, we can directly apply a regularization algorithm to the linear system (9), e.g.:

- **Tikhonov**

$$(\widetilde{\mathbf{K}}^t \widetilde{\mathbf{K}} + n\lambda \mathbf{I})\widetilde{\mathbf{c}} = \widetilde{\mathbf{K}}^t \mathbf{y}$$

- **Conjugate Gradient**

$$\widetilde{\mathbf{c}}^{(0)} \in \mathbb{R}^{n+1}, \widetilde{\mathbf{r}}^{(0)} = \mathbf{y} - \widetilde{\mathbf{K}}\widetilde{\mathbf{c}}^{(0)}, \widetilde{\mathbf{d}}^{(0)} = \widetilde{\mathbf{K}}^t \widetilde{\mathbf{r}}^{(0)}$$

for  $i = 1, \dots, t$

$$\alpha_i = \|\widetilde{\mathbf{K}}^t \widetilde{\mathbf{r}}^{(i-1)}\|_2^2 / \|\widetilde{\mathbf{K}}\widetilde{\mathbf{d}}^{(i-1)}\|_2^2$$

$$\widetilde{\mathbf{c}}^{(i)} = \widetilde{\mathbf{c}}^{(i-1)} + \alpha_i \widetilde{\mathbf{d}}^{(i-1)}$$

$$\widetilde{\mathbf{r}}^{(i)} = \widetilde{\mathbf{r}}^{(i-1)} - \alpha_i \widetilde{\mathbf{K}}\widetilde{\mathbf{d}}^{(i-1)}$$

$$\beta_i = \|\widetilde{\mathbf{K}}^t \widetilde{\mathbf{r}}^{(i)}\|_2^2 / \|\widetilde{\mathbf{K}}^t \widetilde{\mathbf{r}}^{(i-1)}\|_2^2$$

$$\widetilde{\mathbf{d}}^{(i)} = \widetilde{\mathbf{K}}^t \widetilde{\mathbf{r}}^{(i)} + \beta_i \widetilde{\mathbf{d}}^{(i-1)}$$

end

**Why Tikhonov**: the regularized solution can be written in term of the SVD

of  $\widetilde{\mathbf{K}}$ :

$$\widetilde{\mathbf{c}}_{\lambda} = \sum_{\tilde{\sigma}_i \neq 0} \frac{\tilde{\sigma}_i(\mathbf{y}^t \tilde{\mathbf{u}}_i)}{\tilde{\sigma}_i^2 + n\lambda} \tilde{\mathbf{v}}_i$$

For rank-deficient kernels (e.g., the linear kernel when  $d \ll n$ ) many singular values are zero, thus allowing to avoid the calculation of the corresponding singular vectors.

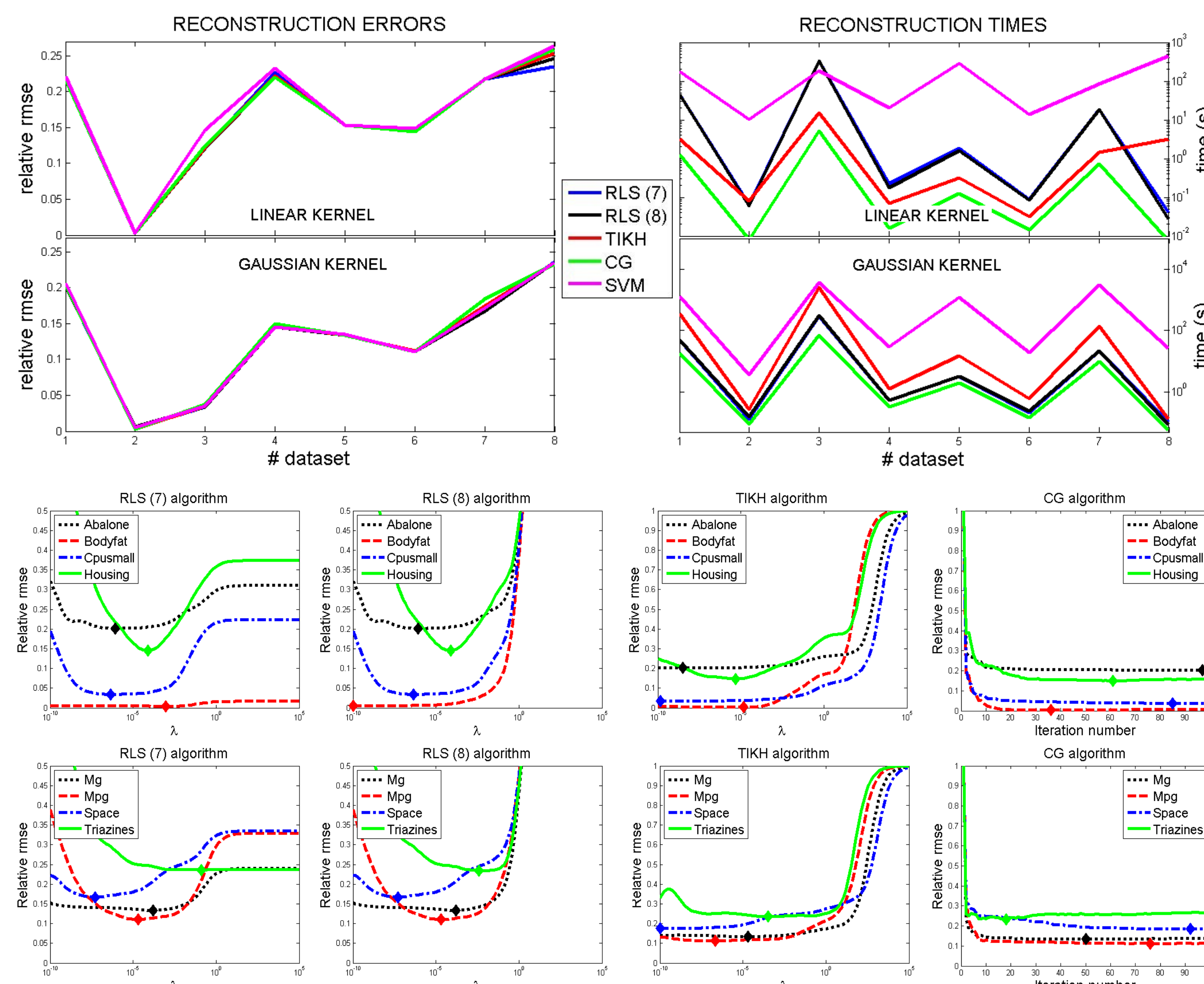
**Why Conjugate Gradient**:  $\mathcal{O}(tn^2)$  instead of the  $\mathcal{O}(n^3)$  required by the “direct” approaches, with  $t$  low due to regularization.

## Numerical experiments

| Dataset     | # examples | # features |
|-------------|------------|------------|
| 1 Abalone   | 4177       | 8          |
| 2 Bodyfat   | 252        | 14         |
| 3 Cpusmall  | 8192       | 12         |
| 4 Housing   | 506        | 13         |
| 5 Mg        | 1385       | 6          |
| 6 Mpg       | 392        | 7          |
| 7 Space     | 3107       | 6          |
| 8 Triazines | 186        | 60         |

### Test settings

- 2/3 training, 1/3 test
- 100 values for  $\lambda$  geometrically distributed in  $[10^{-10}, 10^5]$
- 50 max iterations for CG
- 5x5 grid for the parameters  $(C, \epsilon)$  of SVM<sup>light</sup> [7], within  $[1, 500] \times [10^{-3}, 10^{-1}]$
- RLS, TIKH and CG implemented in Matlab R2010a
- PC 1.60 GHz Intel Core i7, Windows 7 environment



### Figures

Top: reconstruction errors and times for the eight datasets in the case of linear and Gaussian kernels

Bottom: reconstruction errors as functions of  $\lambda$  or  $t$  in the case of Gaussian kernel

### Conclusions

- Regularization approaches are faster than SVM
- TIKH outperforms RLS for rank-deficient kernels, while becomes heavier if the full SVD has to be computed
- CG is always the faster
- TIKH and CG provide flatter error vs  $\lambda$  /  $t$  curves, thus allowing a more stable construction of the solution, e.g. with cross validation.

## References

- [1] Hofmann T, Schölkopf B, Smola AJ 2008. Kernel methods in machine learning. Ann. Stat. 36 (3), 1171–1220
- [2] Schölkopf B, Smola AJ 2002. Learning with Kernels. MIT Press, Cambridge, MA
- [3] De Vito E et al 2004. Some properties of regularized kernel methods. J Mach Learn Res 5, 1363–1390
- [4] Evgeniou T, Pontil M, Poggio T 2000. Regularization Networks and Support Vector Machines. Adv. Comput. Math. 13 (1), 1–50
- [5] Poggio T et al 2002. B. In: Winkler J, Niranjana M Eds, Uncertainty in Geometric Computations. Kluwer Acad Publ, Dordrecht, 131–141
- [6] Lavrentiev MM 1967. Some improperly posed problems of mathematical physics. Springer, New York, NY
- [7] Joachims T 1999. Making large-scale SVM learning practical. In: Schölkopf B, Burges CJC, Smola AJ Eds, Advances in Kernel Methods – Support Vector Learning. MIT Press, Cambridge, MA, 169–184