

Kernel regularized least squares (KRLS)

Zalán Bodó

1 Least squares classification

In least squares classification (or regression) we are searching for a hyperplane separating positive and negative examples such that the labels output by the classifier and the true labels to be as close to each other as possible, in terms of squared difference. That is

$$\min_{\mathbf{w}} \frac{1}{\ell} \|\mathbf{X}'\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

where ℓ denotes the number of training examples (i.e. number of points placed in the columns of \mathbf{X}) $\lambda \|\mathbf{w}\|_2^2$ is a regularization tag. The solution is given by

$$\mathbf{w}^* = (\mathbf{X}\mathbf{X}' + \lambda\ell\mathbf{I})^{-1} \mathbf{X}\mathbf{y}$$

2 Kernelization

In this section we show a method to kernelize, different from the one shown usually in other manuals.

To kernelize RLS we need to form dot products of the examples. Our classifier is given by

$$f(\mathbf{x}) = \mathbf{w}'\mathbf{x} = \mathbf{y}'\mathbf{X} (\mathbf{X}\mathbf{X}' + \lambda\ell\mathbf{I})^{-1} \mathbf{x}$$

and the problem is that the input variables do not form dot products, neither on the left and right hand side of the expression, nor inside the parenthesis. Therefore, we will use a formula that rewrites $(\mathbf{X}\mathbf{X}' + \lambda\ell\mathbf{I})^{-1} \mathbf{X}$ in another form. Consider the Searle formula [1]:

$$(\mathbf{A} + \mathbf{B}\mathbf{B}')^{-1} \mathbf{B} = \mathbf{A}^{-1} \mathbf{B} (\mathbf{I} + \mathbf{B}'\mathbf{A}^{-1} \mathbf{B})^{-1}$$

Applying it to our case we obtain:

$$\begin{aligned} (\mathbf{X}\mathbf{X}' + \lambda\ell\mathbf{I})^{-1} \mathbf{X} &= \frac{1}{\lambda\ell} \mathbf{X} \left(\mathbf{I} + \frac{1}{\lambda\ell} \mathbf{X}'\mathbf{X} \right)^{-1} \\ &= \mathbf{X} (\lambda\ell\mathbf{I} + \mathbf{X}'\mathbf{X})^{-1} \end{aligned}$$

Substituting this back into our classifier we get

$$f(\mathbf{x}) = \mathbf{w}'\mathbf{x} = \mathbf{y}'(\lambda\ell\mathbf{I} + \mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{x}$$

which now can be kernelized. If \mathbf{K} denotes the kernel or Gram matrix of the training data, and \mathbf{k}_x is the vector $[k(\mathbf{x}_1, \mathbf{x}), k(\mathbf{x}_2, \mathbf{x}), \dots, k(\mathbf{x}_N, \mathbf{x})]'$, by the *kernel trick*¹ we can write

$$f(\mathbf{x}) = \mathbf{y}'(\lambda\ell\mathbf{I} + \mathbf{K})^{-1}\mathbf{k}_x$$

3 Kernelization using the representer theorem

The optimization problem can be written as

$$\min_f \frac{1}{\ell} \sum_{i=1}^{\ell} (f(\mathbf{x}_i) - y_i)^2 + \lambda \|f\|_{\mathbf{K}}^2$$

where $\|f\|_{\mathbf{K}}^2$ denotes the RKHS norm of the decision function. The representer theorem [2] says that the solution can be written as

$$f(\cdot) = \sum_{i=1}^{\ell} \alpha_i k(\mathbf{x}_i, \cdot)$$

Therefore we can rewrite the optimization problem as

$$\min_{\alpha} \frac{1}{\ell} \|\mathbf{K}\alpha - \mathbf{y}\|_2^2 + \lambda \|f\|_{\mathbf{K}}^2$$

In order to solve the optimization problem we need to express the regularization tag, that is

$$\begin{aligned} \|f\|_{\mathbf{K}}^2 &= \langle f, f \rangle_{\mathbf{K}} \\ &= \left\langle \sum_{i=1}^{\ell} \alpha_i k(\mathbf{x}_i, \cdot), \sum_{j=1}^{\ell} \alpha_j k(\mathbf{x}_j, \cdot) \right\rangle_{\mathbf{K}} \\ &= \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \alpha_i \alpha_j \langle k(\mathbf{x}_i, \cdot), k(\mathbf{x}_j, \cdot) \rangle_{\mathbf{K}} \\ &= \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \\ &= \alpha' \mathbf{K} \alpha \end{aligned}$$

¹The kernel trick actually refers to the application of a feature mapping ϕ to the points without specifying the mapping. The mapping is given by the kernel function, $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)' \phi(\mathbf{x}_j)$, where the kernel function can be any semi-positive definite, symmetric and continuous function.

Thus, if we rewrite the optimization problem,

$$\min_{\boldsymbol{\alpha}} \quad \frac{1}{\ell} \|\mathbf{K}\boldsymbol{\alpha} - \mathbf{y}\|_2^2 + \lambda \boldsymbol{\alpha}' \mathbf{K} \boldsymbol{\alpha}$$

by setting the derivative of the objective function to zero and expressing $\boldsymbol{\alpha}$ we obtain

$$\mathbf{K}\boldsymbol{\alpha} - \mathbf{y} + \lambda \ell \boldsymbol{\alpha} = 0$$

from which

$$\boldsymbol{\alpha} = (\mathbf{K} + \lambda \ell \mathbf{I})^{-1} \mathbf{y}$$

Substituting back the result into the decision function we obtain the same formula as earlier:

$$f(\mathbf{x}) = \mathbf{y}' (\mathbf{K} + \lambda \ell \mathbf{I})^{-1} \mathbf{k}_{\mathbf{x}}$$

References

- [1] K. B. Petersen and M. S. Pedersen. *The matrix cookbook*, November 2012.
- [2] B. Schölkopf and A. J. Smola. *Learning with Kernels*. The MIT Press, Cambridge, MA, 2002.