

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

## OCReP: An Optimally Conditioned Regularization for pseudoinversion based neural training

### **This is the author's manuscript**

*Original Citation:*

*Availability:*

This version is available <http://hdl.handle.net/2318/1542230> since 2017-05-23T18:23:03Z

*Published version:*

DOI:10.1016/j.neunet.2015.07.015

*Terms of use:*

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

# OCReP: An Optimally Conditioned Regularization for Pseudoinversion Based Neural Training

Rossella Cancelliere<sup>a</sup>, Mario Gai<sup>b</sup>, Patrick Gallinari<sup>c</sup>, Luca Rubini<sup>a</sup>

<sup>a</sup>University of Turin, Dep. of Computer Sciences, C.so Svizzera 185, 10149 Torino, Italy

<sup>b</sup>National Institute of Astrophysics, Astrophys. Observ. of Torino, Pino T.se (TO), Italy

<sup>c</sup>Laboratory of Computer Sciences, LIP6, Univ. Pierre et Marie Curie, Paris, France

---

## Abstract

In this paper we consider the training of single hidden layer neural networks by pseudoinversion, which, in spite of its popularity, is sometimes affected by numerical instability issues. Regularization is known to be effective in such cases, so that we introduce, in the framework of Tikhonov regularization, a matricial reformulation of the problem which allows us to use the condition number as a diagnostic tool for identification of instability. By imposing well-conditioning requirements on the relevant matrices, our theoretical analysis allows the identification of an optimal value for the regularization parameter from the standpoint of stability. We compare with the value derived by cross-validation for overfitting control and optimisation of the generalization performance. We test our method for both regression and classification tasks. The proposed method is quite effective in terms of predictivity, often with some improvement on performance with respect to the reference cases considered. This approach, due to analytical determination of the regularization parameter, dramatically reduces the computational load required by many other techniques.

*Keywords:* Regularization parameter, Condition number, Pseudoinversion, Numerical instability

---

## 1. Introduction

In past decades Single Layer Feedforward Neural Networks (SLFN) training was mainly accomplished by iterative algorithms involving the repetition of learning steps aimed at minimising the error functional over the space of

network parameters. These techniques often gave rise to methods slow and computationally expensive.

Researchers therefore have always been motivated to explore alternative algorithms and recently some new techniques based on matrix inversion have been developed. In the literature, they were initially employed to train radial basis function neural networks (Poggio and Girosi, 1990a): the idea of using them also for different neural architectures was suggested for instance in (Cancelliere, 2001).

The work by Huang et al. (see for instance (Huang *et al.*, 2006)) gave rise to a great interest in neural network community: they presented the technique of Extreme Learning Machine (ELM) for which SLFNs with randomly chosen input weights and hidden layer biases can learn sets of observations with a desired precision, provided that activation functions in the hidden layer are infinitely differentiable. Besides, because of the use of linear output neurons, output weights determination can be brought back to linear systems solution, obtained via Moore-Penrose generalised inverse (or pseudoinverse) of the hidden layer output matrix; so doing iterative training is no more required.

Such techniques appear anyway to require more hidden units with respect to conventional neural network training algorithms to achieve comparable accuracy, as discussed in Yu and Deng (Yu and Deng, 2012).

Many application-oriented studies in the last years have been devoted to the use of these single-pass techniques, easy to implement and computationally fast; some are described e.g. in (Nguyen *et al.*, 2010; Kohno *et al.*, 2010; Ajourloo *et al.*, 2007). A yearly conference is currently being held on the subject, the International Conference on Extreme Learning Machines, and the method is currently dealt with in some journal special issue, e.g. Soft Computing (Wang *et al.*, 2012) and the International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems (Wang, 2013).

Because of the possible presence of singular and almost singular matrices, pseudoinversion is known to be a powerful but numerically unstable method: nonetheless in the neural network community it is often used without singularity checks and evaluated through approximated methods.

In this paper we improve on the theoretical framework using singular value analysis to detect the occurrence of instability. Building on Tikhonov regularization, which is known to be effective in this context (Golub *et al.*, 1999), we present a technique, named Optimally Conditioned Regularization for Pseudoinversion (OCReP), that replaces unstable, ill-posed problems with

well-posed ones.

Our approach is based on the formal definition of a new matricial formulation that allows the use of condition number as diagnostic tool. In this context an optimal value for the regularization parameter is analytically derived by imposing well-conditioning requirements on the relevant matrices.

The issue of regularization parameter choice has often been identified as crucial in literature, and dealt with in a number of historical contributions: a conservative guess might put its published estimates at several dozens. Some of the most relevant works are mentioned in section 2, where the related theoretical background is recalled.

Its determination, mainly aimed at overfitting control, has often been done either experimentally via cross-validation, requiring heavy computational training procedures, or analytically under specific conditions on the matrices involved, sometimes hardly applicable to real datasets, as discussed in section 2.

In section 3 we present the basic concepts concerning input and output weights setting, and we recall the main ideas on ill-posedness, regularization and condition number.

In section 4 our matricial framework is introduced, and constraints on condition number are imposed in order to derive the optimal value for the regularization parameter.

In section 5 our diagnosis and control tool is tested on some applications selected from the UCI database and validated by comparison with the framework regularized via cross-validation and with the unregularized one.

The same datasets are used in section 6 to test the technique effectiveness: our performance is compared with those obtained in other regularized frameworks, originated in both statistical and neural domains.

## 2. Recap on ordinary least-square and ridge regression estimators

As stated in the introduction, pseudoinversion based neural training brings back output weights determination to linear systems solution: in this section we recall some general ideas on this issue, that in next sections will be specialized to deal with SLFN training.

The estimate of  $\beta$  through ordinary least-squares (OLS) technique is a classical tool for solving the problem

$$Y = X\beta + \epsilon , \tag{1}$$

where  $Y$  and  $\epsilon$  are column  $n$ -vectors,  $\beta$  is a column  $p$ -vector and  $X$  is an  $n \times p$  matrix;  $\epsilon$  is random, with expectation value zero and variance  $\sigma^2$ .

In (Hoerl, 1962) and (Hoerl and Kennard, 1970) the role of ordinary ridge regression (ORR) estimator  $\hat{\beta}(\lambda)$  as an alternative to the OLS estimator in the presence of multicollinearity is deeply analyzed. In statistics, multicollinearity (also collinearity) is a phenomenon in which two or more predictor variables in a multiple regression model are highly correlated, meaning that one can be linearly predicted from the others with a non-trivial degree of accuracy. In this situation, the coefficient estimates of the multiple regression may change erratically in response to small changes in the model or the data.

It is known in literature that there exist estimates of  $\beta$  with smaller mean square error (MSE) than the unbiased, or Gauss-Markov, estimate (Golub *et al.*, 1979; Berger, 1976)

$$\hat{\beta}(0) = (X^T X)^{-1} X^T Y . \quad (2)$$

Allowing for some bias may result in a significant variance reduction: this is known as the bias-variance dilemma (see e.g. (Tibshirani, 1996; Geman *et al.*, 1992), whose effects on output weights determination will be deepened in section 3.2.

Hereafter we focus on the one parameter family of ridge estimates  $\hat{\beta}(\lambda)$  given by

$$\hat{\beta}(\lambda) = (X^T X + n\lambda I)^{-1} X^T Y . \quad (3)$$

It can be shown that  $\hat{\beta}(\lambda)$  is also the solution to the problem of finding the minimum over  $\beta$  of

$$\frac{1}{n} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 , \quad (4)$$

which is known as the method of regularization in the approximation theory literature (Golub *et al.*, 1979); basing on it we will develop the theoretical framework for our work in the next sections.

There has always been a substantial amount of interest in estimating a good value of  $\lambda$  from the data: in addition to those already cited in this section a non-exhaustive list of well known or more recent papers is e.g (Hoerl and Kennard, 1976; Lawless and Wang, 1976; McDonald and Galarneau, 1975; Nordberg, 1982; Saleh and Kibria, 1993; Kibria, 2003; Khalaf and Shukur, 2005; Mardikyan and Cetin, 2008).

A meaningful review of these formulations is provided in (Dorugade and Kashid, 2010). They first define the matrix  $T$  such that  $T^T X^T X T = \Lambda$  ( $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$  contains the eigen values of the matrix  $X^T X$ ); then they set  $Z = XT$  and  $\alpha = T^T \beta$ , and show that a great amount of different methods require the OLS estimates of  $\alpha$  and  $\sigma$

$$\hat{\alpha} = (Z^T Z)^{-1} Z^T Y, \quad (5)$$

$$\hat{\sigma}^2 = \frac{Y^T Y - \hat{\alpha}^T Z^T Y}{n - p - 1}. \quad (6)$$

to define effective ridge parameter values. It is important to note that often specific conditions on data are needed to evaluate these estimators.

In particular this applies to the expressions of the ridge parameter proposed by (Kibria, 2003) and (Hoerl and Kennard, 1970), that share the characteristic of being functions of the ratio between  $\hat{\sigma}^2$  and a function of  $\hat{\alpha}$ ; they will be used for comparison with our proposed method in section 6.

The alternative technique of generalised cross-validation (GCV) proposed by (Golub *et al.*, 1979) provides a good estimate of  $\lambda$  from the data as the minimizer of

$$V(\lambda) = \frac{1}{n} \frac{\|I - A(\lambda)Y\|_2^2}{\left[\frac{1}{n} \text{Trace}(I - A(\lambda))\right]_2^2}, \quad (7)$$

where

$$A(\lambda) = X(X^T X + n\lambda I)^{-1} X^T. \quad (8)$$

This solution is particularly interesting, since it does not require an estimate of  $\sigma^2$ : because of this, it will be one term of comparison with our experimental results in section 6.

In the next section we will show how the problem of finding a good solution to (1) applies to the context of pseudoinversion based neural training, specializing the involved relevant matrices to deal with this issue.

### 3. Main ideas on regularization and condition number theory

#### 3.1. Generalised inverse matrix for weights setting

We deal with a standard SLFN with  $L$  input neurons,  $M$  hidden neurons and  $Q$  output neurons, non-linear activation functions  $\phi$  in the hidden layer and linear activation functions in the output layer.

Considering a dataset of  $N$  distinct training samples  $(\mathbf{x}_j, \mathbf{t}_j)$ , where  $\mathbf{x}_j \in \mathbb{R}^L$  and  $\mathbf{t}_j \in \mathbb{R}^Q$ , the learning process for a SLFN aims at producing the matrix of desired outputs  $T \in \mathbb{R}^{N \times Q}$  when the matrix of all input instances  $X \in \mathbb{R}^{N \times L}$  is presented as input.

As stated in the introduction, in the pseudoinverse approach the matrix of input weights and hidden layer biases is randomly chosen and no longer modified: we name it  $C$ . After having fixed  $C$ , the hidden layer output matrix  $H = \phi(XC)$  is completely determined; we underline that since  $H \in \mathbb{R}^{N \times M}$ , it is not invertible.

The use of *linear* output neurons allows to determine the output weight matrix  $W^*$  in terms of the OLS solution to the problem  $T = HW + \epsilon$ , in analogy with eq.(1). Therefore from eq.(2), we have

$$W^* = (H^T H)^{-1} H^T T \quad (9)$$

According to (Penrose and Todd, 1956; Bishop, 2006)

$$W^* = H^+ T. \quad (10)$$

$H^+$  is the Moore-Penrose pseudoinverse (or generalized inverse) of matrix  $H$ , and it minimises the cost functional

$$E_D = \|HW - T\|_2^2 \quad (11)$$

Singular value decomposition (SVD) is a computationally simple and accurate way to compute the pseudoinverse (see for instance (Golub and Van Loan, 1996)), as follows.

Every matrix  $H \in \mathbb{R}^{N \times M}$  can be expressed as

$$H = U \Sigma V^T, \quad (12)$$

where  $U \in \mathbb{R}^{N \times N}$  and  $V \in \mathbb{R}^{M \times M}$  are orthogonal matrices and  $\Sigma \in \mathbb{R}^{N \times M}$  is a rectangular diagonal matrix (i.e. a matrix with  $\sigma_{ih} = 0$  if  $i \neq h$ ); its elements  $\sigma_{ii} = \sigma_i$ , called singular values, are non-negative. A common convention is to list the singular values in descending order, i.e.

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p > 0 \quad (13)$$

where  $p = \min \{N, M\}$ , so that  $\Sigma$  is uniquely determined.

The SVD of  $H$  is then used to obtain the pseudoinverse matrix  $H^+$ :

$$H^+ = V\Sigma^+U^T, \quad (14)$$

where  $\Sigma^+ \in \mathbb{R}^{M \times N}$  is again a rectangular diagonal matrix whose elements  $\sigma_i^+$  are obtained by taking the reciprocal of each corresponding element:  $\sigma_i^+ = 1/\sigma_i$  (see also (Rao and Mitra, 1971)). From eq.(9) we than have:

$$W^* = V\Sigma^+U^TT, \quad (15)$$

**Remark**

An interesting case occurs when only  $k < p$  elements in eq.(13) are non-zero, i.e.  $\sigma_{k+1} = \dots = \sigma_p = 0$ ; in this case the rank of matrix  $H$  is  $k$  and  $\Sigma^+$  is defined as:

$$\Sigma^+ = \text{diag}(1/\sigma_1, \dots, 1/\sigma_k, 0, \dots, 0) \in \mathbb{R}^{M \times N}, \quad (16)$$

as shown for instance in (Golub and Van Loan, 1996).

This is also often done in practice, for computational reasons, for elements smaller than a predefined threshold, thus actually computing an approximated version of the pseudoinverse matrix  $H^+$ .

This approach is for example used by default for pseudoinverse evaluation by means of the Matlab **pinv** function <sup>1</sup>, because the tool is widely used by many scientists for example in ELM context, each time that it is applied blindly, i.e. without having decided at what threshold to zero the small  $\sigma_i$ , an approximation *a priori* uncontrolled is introduced in  $H^+$  evaluation.

*3.2. Stability and generalization properties of regularization algorithms*

A key property for any learning algorithm is stability: the learned mapping has to suffer only small changes in presence of small perturbations (for instance the deletion of one example in the training set).

Another important property is generalization: the performance on the training examples (empirical error) must be a good indicator of the performance on future examples (expected error), that is, the difference between the two must be small. An algorithm that guarantees good generalization predicts well if its empirical error is small.

---

<sup>1</sup><http://www.mathworks.com/help/matlab/ref/pinv.html>.



Many studies in literature dealt with the connection between stability and generalization: the notion of stability has been investigated by several authors, e.g. by Devroye and Wagner (Devroye and Wagner, 1979) and Kearns and Ron (Kearns and Ron, 1999).

Poggio et al. in (Mukherjee *et al.*, 2003) introduced a statistical form of leave-one-out stability, named  $CVEE_{loo}$ , building on a cross-validation leave-one-out stability endowed with conditions on stability of both expected and empirical errors; they demonstrated that this condition is necessary and sufficient for generalization and consistency of the class of empirical risk minimization (ERM) learning algorithms, and that it is also a sufficient condition for generalisation for not ERM algorithms (see also (Poggio *et al.*, 2004)).

To turn an original instable, ill-posed problem into a well-posed one, regularization methods of the form (4) are often used (Badeva and Morozov, 1991) and among them, Tikhonov regularization is one of the most common (Tikhonov and Arsenin, 1977; Tikhonov, 1963). It minimises the error functional

$$E \equiv E_D + E_R = \|HW - T\|_2^2 + \|\Gamma W\|_2^2, \quad (17)$$

obtained adding to the cost functional  $E_D$  in eq.(11) a penalty term  $E_R$  that depends on a suitably chosen Tikhonov matrix  $\Gamma$ . This issue has been discussed in its applications to neural networks in (Poggio and Girosi, 1990b), and surveyed in (Girosi *et al.*, 1995; Haykin, 1999).

Besides, Bousquet and Elisseeff (Bousquet and Elisseeff, 2002) proposed the notion of uniform stability to characterize the generalization properties of an algorithm. Their results state that Tikhonov regularization algorithms are uniformly stable and that uniform stability implies good generalization (Mukherjee *et al.*, 2006).

Regularization thus introduces a penalty function that not only improves on stability, making the problem less sensitive to initial conditions, but it is also important to contain model complexity avoiding overfitting.

The idea of penalizing by a square function of weights is also well known in neural literature as weight decay: a wide amount of articles have been devoted to this argument, and more generally to the advantage of regularization for the control of overfitting. Among them we recall (Hastie *et al.*, 2009; Tibshirani, 1996; Bishop, 2006; Girosi *et al.*, 1995; Fu, 1998; Gallinari and Cibas, 1999).

A frequent choice is  $\Gamma = \sqrt{\gamma}I$ , to give preference to solutions with smaller

norm (Bishop, 2006), so eq. (17) can be rewritten as

$$E \equiv E_D + E_R = \|HW - T\|_2^2 + \gamma \|W\|_2^2. \quad (18)$$

We define  $\hat{W} = \min_W(E)$  the regularized solution of (18): it belongs to the family of ridge estimates described by eq.(3) and can be expressed as

$$\hat{W} = (H^T H + \gamma I)^{-1} H^T T \quad (19)$$

or, as show in ((Fuhry and Reichel, 2012)) as

$$\hat{W} = V D U^T T. \quad (20)$$

$V$  and  $U$  are from the singular value decomposition of  $H$  (eq.(12)) and  $D \in \mathbb{R}^{M \times N}$  is a rectangular diagonal matrix whose elements, built using the singular values  $\sigma_i$  of matrix  $\Sigma$ , are:

$$D_i = \frac{\sigma_i}{\sigma_i^2 + \gamma}. \quad (21)$$

We remark on the difference between the minima of the regularized and unregularized error functionals. Increasing values of the regularization parameter  $\gamma$  induce larger and larger departure of the former (eq. (19)) from the latter (eq. (9)). Thus, the regularization process increases the bias of the approximating solution and reduces its variance, as discussed about the bias-variance dilemma in section 2.

A suitable value for the Tikhonov parameter  $\gamma$  has therefore to derive from a compromise between having it sufficiently large to control the approaching to zero of  $\sigma_i$  in eq.(21), while avoiding an excess of the penalty term in eq.(18). Its tuning is therefore crucial.

### 3.3. Condition number as a measure of ill-posedness

The condition number of a matrix  $A \in \mathbb{R}^{N \times M}$  is the number  $\mu(A)$  defined as

$$\mu(A) = \|A\| \|A^+\| \quad (22)$$

where  $\|\cdot\|$  is any matrix norm. If the columns (rows) of  $A$  are linearly independent, e.g. in case of experimental data matrices, then  $A^+$  is a left (right) inverse of  $A$ , i.e.  $A^+ A = I_N$  ( $AA^+ = I_M$ ). The Cauchy-Schwarz inequality in this case then provides  $\mu(A) \geq 1$ ; besides,  $\mu(A) \equiv \mu(A^+)$ .

Matrices are said to be ill-conditioned if  $\mu(A) \gg 1$ .

If  $\|\cdot\|_2$  norm is used, then

$$\mu(A) = \frac{\sigma_1(A)}{\sigma_p(A)}, \quad (23)$$

where  $\sigma_1$  and  $\sigma_p$  are the largest and smallest singular values of  $A$  respectively.

From eq.(23) we can easily understand that large condition numbers  $\mu(A)$  suggest the presence of very small singular values (i.e. of almost singular matrices), whose numerical inversion, required to evaluate  $\Sigma^+$  and the unregularized solution  $W^*$ , is a cause of instability.

From numeric linear algebra we also know that if the condition number is large the problem of finding least-squares solutions to the corresponding system of linear equations is ill-posed, i.e. even a small perturbation in the data can lead to huge perturbations in the entries of solution (see (Golub and Van Loan, 1996)).

According to (Mukherjee *et al.*, 2006) the stability of Tikhonov regularization algorithms can also be characterized using the classical notion of condition number: our proposed regularization method fits within this context. We will see that it specifically aims at analitically determining the value of the  $\gamma$  parameter that minimizes the conditioning of the regularized hidden layer output matrix so that the solution  $\hat{W}$  is stable in the sense of eq.(2.9) of (Mukherjee *et al.*, 2006).

In the next section, we will derive the optimal value of the regularization parameter  $\gamma$  according to this stability criterion (minimum condition number).

The experimental results presented in sections 5 and 6 will evidence that our quest for stable solutions allows us to also achieve good generalization and predictivity. A comparison will be made to this purpose with the performance obtained when  $\gamma$  is determined via the standard cross-validation approach, aimed at overfitting control and generalization performance optimization.

#### 4. Conditioning of the regularized matricial framework

For convenient implementation of our diagnostics, and building on eq.(20), we propose an original matricial framework in which to develop our study tool with the following definition.

**Definition 1.** *We define the matrix*

$$H^{reg} \equiv VDU^T \quad (24)$$

as the **regularized** hidden layer output matrix of the neural network.

This allows us to rewrite eq.(20) as

$$\hat{W} = H^{reg}T, \quad (25)$$

for similarity with eq.(9).

By construction,  $H^{reg}$  is decomposed in three matrices according to the SVD framework, and its singular values are provided by eq.(21) as a function of the singular values  $\sigma_i$  of  $H$ .

This new regularized matricial framework makes easier the comparison of the properties of  $H^{reg}$  with those of the corresponding unregularized matrix  $H^+$ . In fact, when unregularized pseudoinversion is used, nothing prevents the occurrence of very small singular values that make numerically instable the evaluation of  $H^+$  (see eq. 14). On the contrary, even in presence of very small values  $\sigma_i$  of the original unregularized problem, a careful choice of the parameter  $\gamma$  allows to tune the singular values  $D_i$  of the regularized matrix  $H^{reg}$ , preventing numerical instability.

#### 4.1. Condition number definition

According to eq. (23), we define the condition number of  $H^{reg}$  as:

$$\mu(H^{reg}) = \frac{D_{max}}{D_{min}}. \quad (26)$$

where  $D_{max}$  and  $D_{min}$  are the largest and smallest singular values of  $H^{reg}$ .

The shape of the functional relation  $\sigma/(\sigma^2 + \gamma)$  that links regularized and unregularized singular values, defined through eq. (21), is shown in Fig.1 for three different values of  $\gamma$ .

The curves are non-negative, because  $\sigma > 0$  and  $\gamma > 0$ , and have only one maximum, with coordinates  $(\sqrt{\gamma}; \frac{1}{2\sqrt{\gamma}})$ .

A few pairs of corresponding values  $(D_i, \sigma_i)$  are marked by dots on each curve.

For the sake of the determination of  $\mu(H^{reg})$  we are interested in evaluating  $D_{max}$  and  $D_{min}$  of  $H^{reg}$  over the finite, discrete range  $[\sigma_1, \sigma_2, \dots, \sigma_p]$ .

The value  $D_{max}$  is reached in correspondence to a given singular value of  $H$ , a priori not known, that we label  $\sigma_{max}$ , so that:

$$D_{max} = \frac{\sigma_{max}}{\sigma_{max}^2 + \gamma}. \quad (27)$$

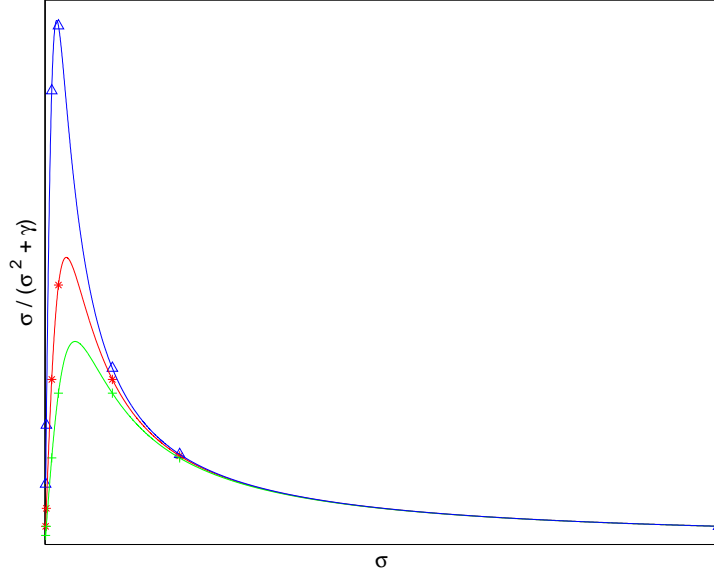


Figure 1: Example of regularized/unregularized singular values relationship via eq. (21)

The variation of  $\gamma$  has the effect of changing the curve and shifting its maximum point within the interval  $[\sigma_1, \sigma_p]$ . Therefore,  $\sigma_{max}$  can coincide with any singular value of  $H$  from eq. (13), including the extreme ones.

Conversely, we now demonstrate that  $D_{min}$  can only be reached in correspondence to  $\sigma_1$  or  $\sigma_p$  (or both when coincident).

### Theorem 3.1

The minimum singular value  $D_{min}$  of matrix  $H^{reg}$  can only be reached in correspondence to the largest singular value  $\sigma_1$  or to the smallest singular value  $\sigma_p$  of the unregularized matrix  $H$  (or both).

*Proof.* Without loss of generality, we can express  $\gamma$  as a function of  $\sigma_1\sigma_p$ , i.e.  $\gamma = \beta\sigma_1\sigma_p$ , where  $\beta$  is a real positive value. By replacement in eq. (21), we get

$$D_1 = \frac{1}{\sigma_1 + \beta\sigma_p}, \quad D_p = \frac{1}{\sigma_p + \beta\sigma_1}$$

To establish their ordering, we evaluate the difference  $\Delta$  of their inverses:

$$\Delta = \frac{1}{D_1} - \frac{1}{D_p} = (\sigma_1 + \beta\sigma_p) - (\sigma_p + \beta\sigma_1) = (1 - \beta)(\sigma_1 - \sigma_p) .$$

Recalling that  $\sigma_1 - \sigma_p > 0$ , we can distinguish three cases:

**Case 1**,  $\beta > 1$  ( $\gamma > \sigma_1\sigma_p$ )  $\rightarrow \Delta < 0 \rightarrow D_1 > D_p$

Because of the  $D_i$  distribution shape,  $D_p$  is also the minimum among all values  $D_i$ , so that  $D_{min} \equiv D_p$ .

**Case 2**,  $\beta < 1$  ( $\gamma < \sigma_1\sigma_p$ )  $\rightarrow \Delta > 0 \rightarrow D_1 < D_p$

Then,  $D_1$  is also the minimum among all values  $D_i$ , so that  $D_{min} \equiv D_1$ .

**Case 3**,  $\beta = 1$  ( $\gamma = \sigma_1\sigma_p$ )  $\rightarrow \Delta = 0 \rightarrow D_1 = D_p$

Thus,  $D_1$  and  $D_p$  are both minima, so that  $D_{min} \equiv D_1 = D_p$ .

□

#### 4.2. Condition number evaluation

The result by Theorem 3.1 allows us to find, according to eq. (26), the following expressions for  $\mu(H^{reg})$  :

**Case 1**,  $\beta > 1$ :

$$\mu(H^{reg}) = \frac{D_{max}}{D_p} = \frac{\sigma_{max}(\sigma_p + \beta\sigma_1)}{\sigma_{max}^2 + \beta\sigma_1\sigma_p}$$

**Case 2**,  $\beta < 1$ :

$$\mu(H^{reg}) = \frac{D_{max}}{D_1} = \frac{\sigma_{max}(\sigma_1 + \beta\sigma_p)}{\sigma_{max}^2 + \beta\sigma_1\sigma_p}$$

**Case 3**,  $\beta = 1$ :

$$\mu(H^{reg}) = \frac{D_{max}}{D_p} = \frac{D_{max}}{D_1} = \frac{\sigma_{max}(\sigma_p + \sigma_1)}{\sigma_{max}^2 + \sigma_1\sigma_p}$$

Bearing in mind that well-conditioned problems are characterized by small condition numbers, we now will look for the  $\beta$  parameter values which, in the three cases above, make the regularized condition number smaller.

In Case 1,  $\mu(H^{reg})$  is an increasing function of  $\beta$ , so that in its domain, i.e.  $(1, \infty)$ , its minimum value is reached when  $\beta \rightarrow 1^+$ . On the contrary, in Case 2,  $\mu(H^{reg})$  is a decreasing function of  $\beta$ , so that in its domain, i.e.  $(0, 1)$ , the minimum is reached when  $\beta \rightarrow 1^-$ .

Fig.2 shows the function behaviour over the whole domain.

Both cases have a common limit:

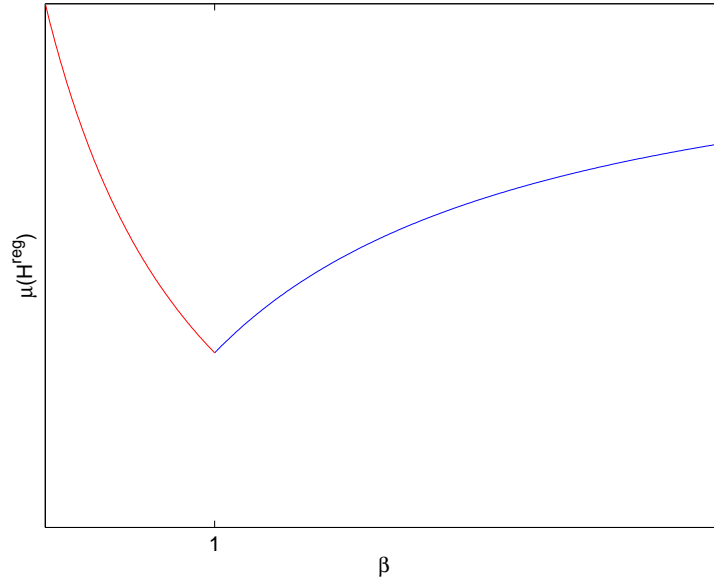


Figure 2: Regularized condition number vs.  $\beta$

$$\lim_{\beta \rightarrow 1^+} \mu(H^{reg}) = \lim_{\beta \rightarrow 1^-} \mu(H^{reg}) = \frac{\sigma_{max}(\sigma_p + \sigma_1)}{\sigma_{max}^2 + \sigma_1 \sigma_p} \quad (28)$$

Such value is just that provided by Case 3, which can therefore be considered the best possible choice to minimize the condition number.

Thus our quest for the best possible conditioning for the matrix  $H^{reg}$  identifies an explicit optimal value for the regularization parameter  $\gamma$ :

$$\gamma = \sigma_1 \sigma_p \quad (29)$$

## 5. Simulation and Discussion

For the numerical experimentation, we use eight benchmark datasets from the UCI repository (Bache and Lichman, 2013) listed in Table 1. All simulations are carried out in Matlab 7.3 environment.

The performance is assessed by statistics over a set of 50 different extractions of input weights, computing either the average RMSE (for regression tasks) or the average percentage of misclassification rate (for classification tasks) on the test set. Either quantity is labeled “Err” in the tables summarising our results. The error standard deviation (labeled “Std”) is also computed to evidence the dispersion of experimental results.

Dataset	Type	N. Instances	N. Attributes	N. Classes
Abalone	Regression	4177	8	-
Machine Cpu	Regression	209	6	-
Delta Ailerons	Regression	7129	5	-
Housing	Regression	506	13	-
Iris	Classification	150	4	3
Diabetes	Classification	768	8	2
Wine	Classification	178	13	3
Segment	Classification	2310	19	7

Table 1: The UCI datasets and their characteristics

Our regularization strategy, labeled Optimally Conditioned Regularization for Pseudoinversion (OCReP), is verified by simulation against the common approach in which cross-validation is used i) to determine the regularization parameter  $\gamma$  at a fixed high number of hidden neurons and ii) to perform also hidden neurons number optimization, respectively in sec. 5.1 and 5.2.

A discussion of the effectiveness of OCReP in terms of minimization of the condition number of the involved matrices is done in sec. 5.3.

#### 5.1. OCReP performance assessment: fixed number of hidden units

In this section we compare OCReP with a regularization approach in which  $\gamma$  is selected by a cross-validation scheme, which is typically used for control of under/overfitting and optimization of the model generalization performance. A 70%/30% split between training and test set is applied; then, a three-fold cross-validation search on the training set identifies the best  $\gamma$  by best performance on the validation set, over the set of 50 values of  $\gamma$  [ $10^{-25}, 10^{-24}, \dots, 10^{25}$ ].

For the sake of comparison, a fixed, high number of hidden units  $M$  is used, selected according to dimension and complexity of the datasets. For the three datasets Machine Cpu, Iris and Wine the simulation is performed for 50 and 100 hidden neurons; for Abalone, Delta Ailerons, Housing and Diabetes, we use 50, 100, 200 and 300 neurons; for Segment, we use 1000 and 1500 units.

Figures 3 and 4 (respectively for regression and classification datasets) show average test errors as a function of the sampled values of  $\gamma$  (red dots);



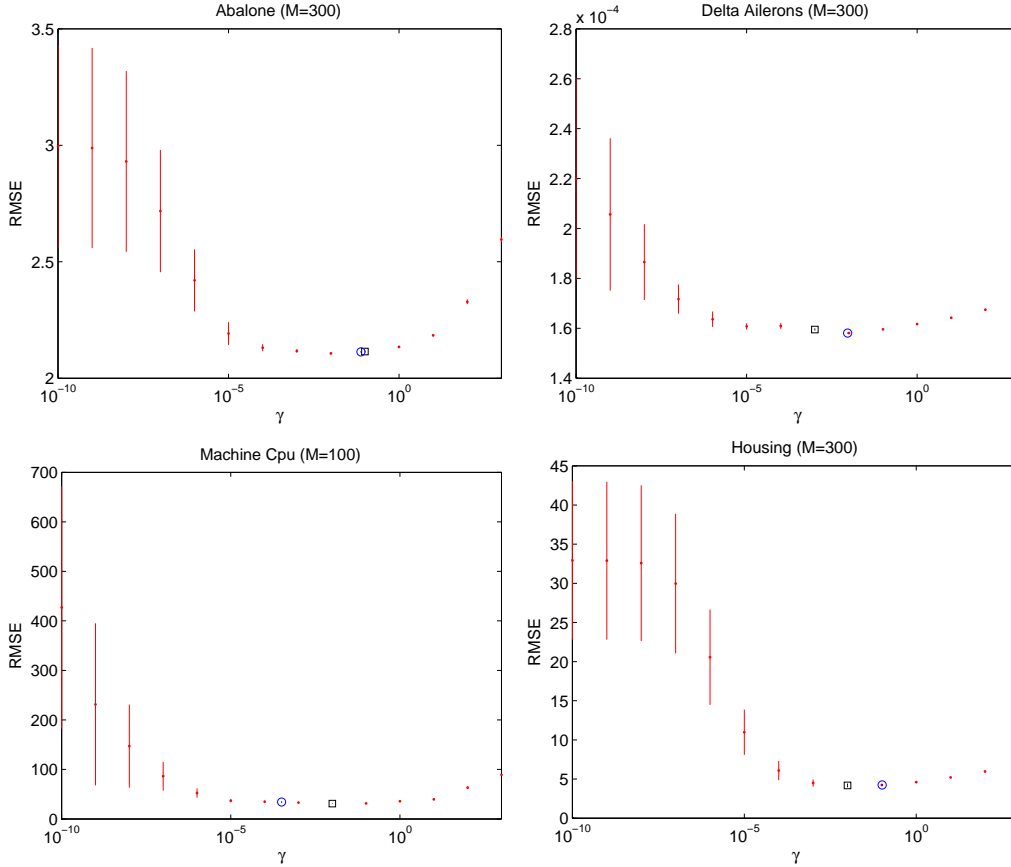


Figure 3: Test error trends for regression datasets as a function of the values of  $\gamma$  over the selected cross-validation range (red dots): the cross-validation selected  $\gamma$  is the black square; the proposed  $\gamma$  from OCRReP is the blue circle.

the standard deviation is shown as an error bar. Our proposed optimal  $\gamma$  is evidenced as a blue circle, whereas the value of  $\gamma$  selected by cross-validation is shown as a black square. The results are in each case related to the highest number of neurons experimented.

The horizontal axis has been zoomed in onto the region of interest, i.e.  $[10^{-10}, 10^5]$ .

It may be noted that the performance from OCRReP and cross-validation are comparable, and also close to the experimental minimum. This may be interpreted as good predictivity for both algorithms.

Also, we remark that the error bars, i.e. experimental result dispersion,

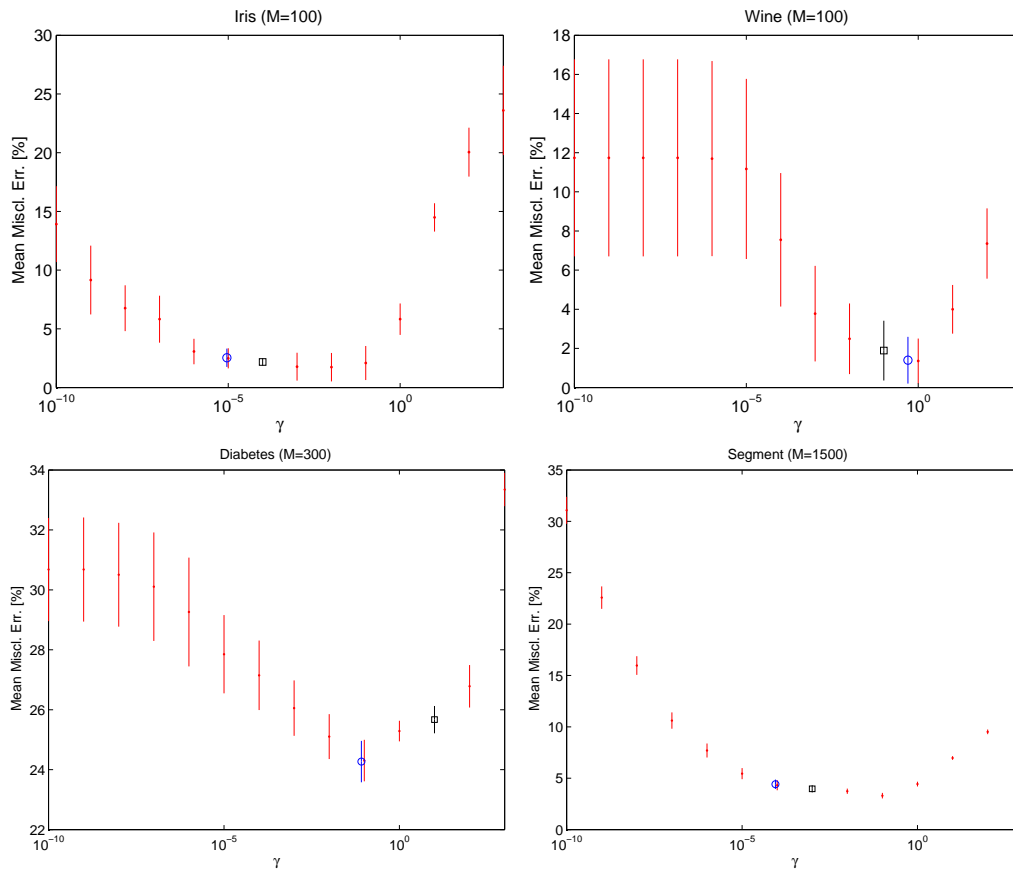


Figure 4: Test error trends for classification datasets as a function of the values of  $\gamma$  over the selected cross-validation range (red dots): the cross-validation selected  $\gamma$  is the black square; the proposed  $\gamma$  from OCRReP is the blue circle.

is large for small values of  $\gamma$ , consistently with expectations on ineffective regularization.

Table 2: Comparison of OCREP vs. cross-validation at fixed number of hidden neurons for small size datasets

				Iris	Wine	Machine Cpu
$M$	50	OCReP	Err.	<b>1.51</b>	2.98	31.21
			Std	1.13	1.75	1.1
	cross-val.	Err.	2.13	<b>3.37</b>	31.1	
		Std	0.77	2.27	1.02	
100	OCReP	Err.	2.53	<b>1.39</b>	34.13	
		Std	0.77	1.19	01.68	
	cross-val.	Err.	<b>2.17</b>	1.88	<b>30.94</b>	
		Std	0.31	1.88	0.69	

Table 3: Comparison of OCREP vs. cross-validation at fixed number of hidden neurons for large size datasets

				Segment
$M$	1000	OCReP	Err.	2.53
			Std	0.77
	cross-val.	Err.	<b>2.17</b>	
		Std	0.31	
1500	OCReP	Err.	4.41	
		Std	0.45	
	cross-val.	Err.	<b>3.97</b>	
		Std	0.35	

The numerical results have been reported in Tab. 2, 3 and 4 according to the grouping based on dimension and complexity of the datasets.

For each dataset and selected number of hidden neurons  $M$ , the best test error is evidenced in bold, whenever the difference is statistically significant<sup>2</sup>.

---

<sup>2</sup>The Student's t-test has been used for assessing the statistical significance through

Table 4: Comparison of OCR<sub>e</sub>P vs. cross-validation at fixed number of hidden neurons for medium size datasets. For Delta Ailerons, average errors and standard deviations have to be multiplied by  $10^{-4}$ .

			Abalone	Delta Ailerons	Housing	Diabetes	
<i>M</i>	50	OCR <sub>e</sub> P	Err.	2.22	1.64	5.54	<b>26.01</b>
			Std	0.16	0.0051	0.12	0.604
		cross-val.	Err.	<b>2.13</b>	<b>1.59</b>	<b>4.79</b>	26.79
			Std	0.017	0.0073	0.37	0.814
	100	OCR <sub>e</sub> P	Err.	2.15	1.62	5.17	25.66
			Std	0.007	0.004	0.08	0.608
		cross-val.	Err.	<b>2.11</b>	<b>1.58</b>	<b>4.49</b>	25.71
			Std	0.006	0.0036	0.28	0.608
	200	OCR <sub>e</sub> P	Err.	2.12	<b>1.59</b>	4.62	<b>25.13</b>
			Std	0.003	0.0031	0.09	0.445
		cross-val.	Err.	<b>2.11</b>	1.61	<b>4.30</b>	25.79
			Std	0.003	0.0096	0.27	0.443
300	OCR <sub>e</sub> P	Err.	2.113	<b>1.58</b>	4.24	<b>24.26</b>	
		Std	0.03	0.0018	0.13	0.689	
	cross-val.	Err.	2.114	1.60	4.18	25.66	
		Std	0.003	0.0042	0.23	0.456	

Thus, for example, on Iris the best performance is achieved using 50 neurons by OCR<sub>e</sub>P, and with 100 neurons by cross-validation. In some cases, e.g. Wine (50 neurons), there is no clear winner from statistical considerations, i.e. the best results are comparable, within the errors.

From the above results it appears that cross-validation has better test error performance on a number of datasets slightly higher, at fixed number of hidden neurons. However, it is important to evidence that the use of OCR<sub>e</sub>P allows to save the hundreds of pseudoinversion steps required by cross-validation, which is a crucial issue for practical implementation.

---

determination of the confidence intervals related to 99% confidence level.

### 5.2. OCREP performance assessment: variable number of hidden units

In order to pursue the double aim of performance and hidden units optimization, a first interesting step is to give a look to the variation as a function of hidden layer dimension of error trends of unregularized models (i.e. models whose output weights are evaluated according to eq.(10)).

A context widely used among researchers using such techniques (see e.g. Helmy and Rasheed (2009); Huang *et al.* (2006)) is to use input weights distributed according to a random uniform distribution in the interval  $(-1, 1)$ , and sigmoidal activation functions for hidden neurons: hereafter we name this framework Sigm-unreg.

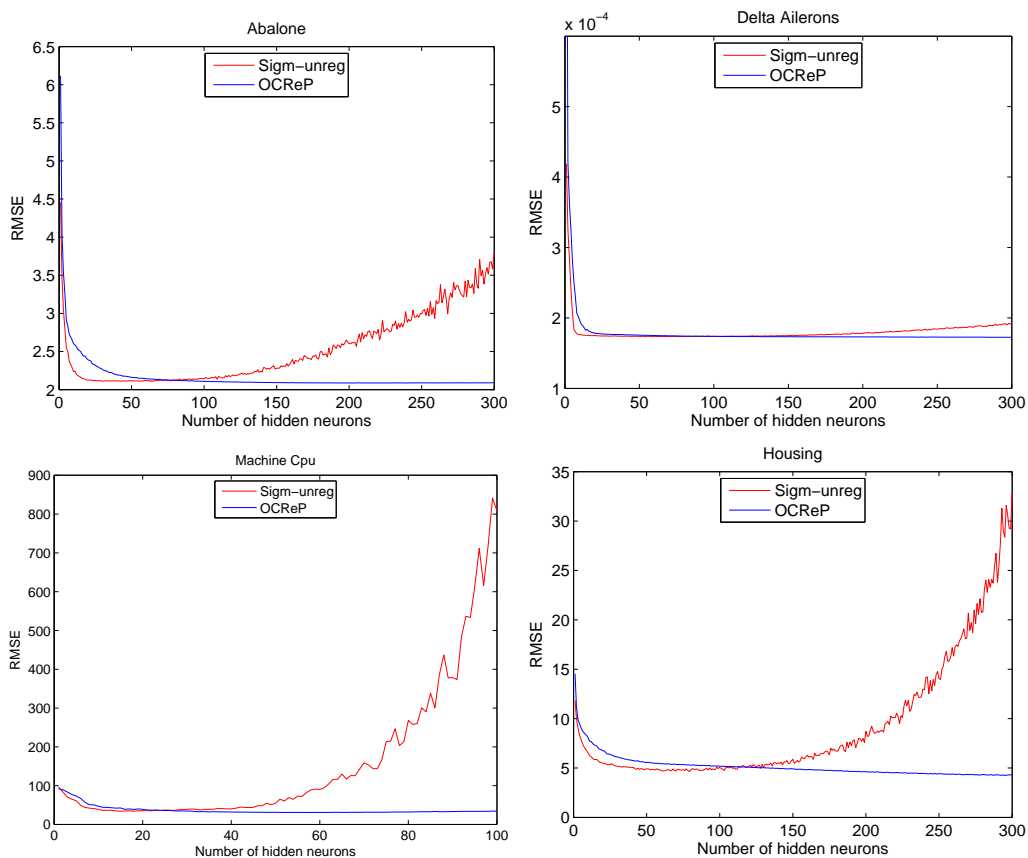


Figure 5: Test error trends for regression datasets: OCREP vs. unregularized pseudoinversion.

Figures 5 and 6 show, respectively for regression and classification datasets,

the average test error values, (over 50 different input weights selections) for both OCREP (blue line) and Sigm-unreg (red line) as a function of the number of hidden nodes, which is gradually increased by unity steps. In all cases, after an initial decrease the Sigm-unreg test error increases significantly.

On the contrary, the OCREP test error curves keep decreasing, albeit at slower and slower rate, thus showing also a good capability of overfitting control of the method.

We aim now at comparing the results obtained when the trade-off value of  $\gamma$  is searched by cross-validation, with the two different frameworks discussed so far, i.e. OCREP and Sigm-unreg.

A 70%/30% split between training and test set is applied; we then perform a three-fold cross-validation for the selection of the number of hidden neurons  $\bar{M}$  at which the minimum error is recorded in all cases. Test errors are again evaluated as the average of 50 different random choices of input weights.

The numerical results of the simulation are presented in Tables 5 and 6, respectively for regression and classification tasks, with their standard deviations (Std) and  $\bar{M}$ .

Best test errors are evidenced in bold, whenever the difference between OCREP and cross-validation is statistically significant.

We see that our proposed regularization technique provides, for regression datasets, performance comparable with the cross-validation option but always a better performance (with statistical significance at 99% level) with respect to the unregularized case.

For classification datasets in three cases out of four OCREP provides a better performance with respect to cross-validation, and always a better performance with respect to the Sigm-unreg case. In all such cases, the statistical significance is at the 99% level.

Also, in almost all cases smaller standard deviations are associated with the OCREP method, suggesting a lower sensitivity to initial input weights conditions.

### 5.3. Additional considerations

The proposed method OCREP presents in our opinion two features of interest: on one side, its computational efficiency, and on the other side its optimal conditioning.

Our goal of optimal analytic determination of the regularization parameter  $\gamma$  results in a dramatic improvement in the computing requirements with respect to experimental tuning by search over a pre-defined large grid

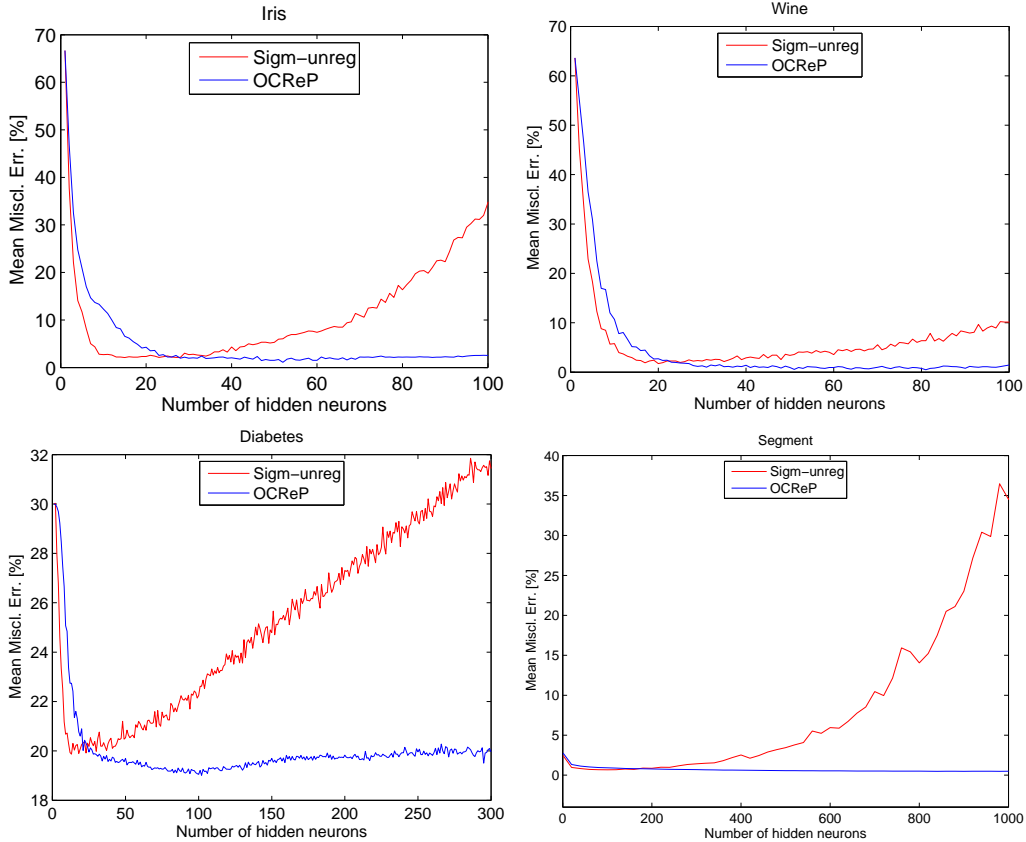


Figure 6: Test error trends for classification datasets: OCREP vs. unregularized pseudoinversion.

of  $N_\gamma$  tentative values. In the latter case, for each choice of  $\gamma$  over the selected range, at least a pseudoinversion is required for every output weight determination, thus increasing the computational load by a factor  $N_\gamma$ .

Besides, our method is designed explicitly for optimal conditioning. In our simulations, we verify that the goal is fulfilled by evaluating average condition numbers of hidden layer output matrices. The statistics is performed over 50 different configurations of input weights and a fixed number of hidden units, namely the largest used in section 5.1 for each dataset. The results are summarised in Tables 7 and 8, respectively for regression and classification datasets. On the first row of each table, we list the ratio of average condition numbers of matrices  $H^{reg}$ , and  $H^+$ , associated respectively to OCREP and Sigm-unreg, i.e. regularized and unregularized approaches. On the second

row, we list the ratio of average condition numbers of matrices  $H^{reg}$  and  $H^{CV}$ , thus comparing our regularization approach with the more conventional one, the latter using cross-validation.

Not surprisingly, our regularization method provides a significant improvement on conditioning with respect to the unregularized approach, as evidenced by ratio values much smaller than unity. Besides, OCREP also provides better conditioned matrices than those derived by selection of  $\gamma$  through cross-validation, since the corresponding condition numbers are systematically smaller in the former case, sometimes up to an order of magnitude.

Table 5: Hidden layer optimization for regression tasks. For Delta Ailerons, average errors and standard deviations have to be multiplied by  $10^{-4}$ .

	Abalone	Housing	Delta Ailerons	Machine Cpu
OCReP				
Err.	2.12	4.25	1.58	<b>31.22</b>
Std.	0.32	0.13	0.0048	0.78
$\bar{M}$	178	255	298	63
Cross-validation				
Err.	<b>2.11</b>	4.19	1.58	31.51
Std.	0.0097	0.25	0.0036	1.25
$\bar{M}$	110	250	93	70
Sigm-unreg				
Err.	2.14	4.73	1.62	34.44
Std.	0.014	0.20	0.57	2.89
$\bar{M}$	31	76	74	15

## 6. Comparison with other approaches

Since the literature provides a host of different recipes for either the choice of the regularization parameter, or the actual regularization algorithm, hereafter we focus on a couple of specific frameworks.

### 6.1. Other choices of regularization parameter

Among the approaches mentioned in section 2, we primary select the technique of generalised cross-validation (GCV) from (Golub *et al.*, 1979),



Table 6: Hidden layer optimization for classification tasks.

	Iris	Wine	Diabetes	Segment
OCReP				
Err.	<b>1.6</b>	<b>1.73</b>	25.53	<b>2.50</b>
Std.	1.10	1.25	0.51	0.32
$\bar{M}$	67	91	291	760
Cross-validation				
Err.	2.12	2.10	25.2	2.65
Std.	1.26	2.27	1.29	0.38
$\bar{M}$	14	137	25	620
Sigm-unreg				
Err.	2.31	3.20	25.92	4.45
Std.	1.48	2.09	1.12	0.47
$\bar{M}$	67	91	291	760

Table 7: Condition number comparison for regression datasets

	Abalone	Housing	Delta Ailerons	Machine Cpu
$\mu(H^{reg})/\mu(H^+)$	0.0002	0.0008	0.00007	0.0001
$\mu(H^{reg})/\mu(H^{CV})$	0.8	0.3	0.3	0.1

described by eqs. (7) and (8), for comparison with our method. The main motivation for our choice is its independence on the estimate of the error variance  $\sigma^2$ , which is a characteristic shared with our case. For each dataset, we select the same fixed numbers of hidden units as in section 5.1: then for each case eq. (7) is minimized over the set of 50 values of  $\gamma$  [ $10^{-25}, 10^{-24} \dots 10^{25}$ ] and for 50 different configurations of input weights.

We evaluate the mean and standard deviation of the corresponding regularized test error, reported in Tables 9, 10 and 11. We also remind that the tabulated error “Err” is either the average RMSE for regression tasks, or the average misclassification rate for classification tasks; “Std” is the corresponding standard deviation. The performance comparison is based on statistical significance at 99% level.

Whenever GCV provides test error values statistically better than OCReP

Table 8: Condition number comparison for classification datasets

	Iris	Wine	Diabetes	Segment
$\mu(H^{reg})/\mu(H^+)$	0.00002	0.005	0.0007	0.000005
$\mu(H^{reg})/\mu(H^{CV})$	0.2	0.4	0.1	0.2

Table 9: GCV results at fixed number of hidden neurons for small datasets

			Iris	Wine	Machine Cpu
$M$	50	Err.	2.47	3.66	33.03
		Std	1.06	2.42	1.27
	100	Err.	3.06	3.77	36.06
		Std	1.08	2.44	1.13

(listed in Tab. 2, 3 and 4), they are marked in bold.

We remark that in all cases listed in Tab. 2 and 3 OCREP provides statistically better results than GCV. The situation of medium size datasets evidences a somewhat mixed behaviour: with 50 hidden neurons, GCV wins; with 100 neurons, for three out of four datasets (i.e. Abalone, Housing and Diabetes) the performance is statistically comparable. In all other cases of Tab. 4 OCREP again provides better statistical results than GCV.

We make two other comparisons, using the ridge estimates described in eq.(13) and eq.(9) of (Dorugade and Kashid, 2010), and proposed respectively by (Kibria, 2003) and (Hoerl and Kennard, 1970):

Table 10: GCV results at fixed number of hidden neurons for large size datasets

			Segment
$M$	1000	Err.	11.39
		Std	0.75
	1500	Err.	14.72
		Std	0.803

Table 11: GCV results at fixed number of hidden neurons for medium size datasets. For Delta Ailerons, average errors and standard deviations have to be multiplied by  $10^{-4}$ .

		Abalone	Housing	Delta Ailerons	Diabetes	
$M$	50	Err.	<b>2.13</b>	<b>4.89</b>	<b>1.60</b>	<b>25.2</b>
		Std	0.017	0.45	0.0103	1.22
	100	Err.	2.15	5.05	1.63	26.66
		Std	0.021	0.70	0.0297	1.39
	200	Err.	2.32	6.78	1.74	27.73
		Std	0.10	2.35	0.0892	1.27
	300	Err.	2.98	8.07	2.20	27.14
		Std	0.42	2.89	0.4054	1.15

$$\gamma_K = \frac{1}{p} \sum_1^p \frac{\hat{\sigma}^2}{\hat{\alpha}_i^2}, \quad (30)$$

$$\gamma_{HK} = \frac{\hat{\sigma}^2}{\hat{\alpha}_{max}^2}. \quad (31)$$

Our experimentation is made only for regression datasets because the theoretical background of (Dorugade and Kashid, 2010), and of most of other works referred in section 2, directly applies to the case in which the quantity  $Y$  in eq.(1) is a one column matrix. In our formulation  $Y$  is the desired target  $T$  and it is a one-column matrix only for regression tasks.

For each dataset we applied both methods described by eq. 30 and 31; we select the same fixed numbers of hidden units as in section 5.1 and perform 50 experiments with different configuration of input weights.

Each step of pseudoinversion is regularized for each method with the corresponding  $\gamma$  value. We evaluate the mean and standard deviation of the regularized test errors, reported respectively in Tables 12 and 13.

Whenever the methods provide test error values statistically better than OCREP (listed in Tab. 2 and 4), they are marked in bold.

We remark that the method by Kibria obtains a better performance in two cases over sixteen, while OCREP in 12 cases over sixteen. Besides, the method by Hoerl and Kennard obtains a better performance in three cases

Table 12: Kibria estimate of ridge parameter: results at fixed number of hidden neurons for regression datasets. For Delta Ailerons, average errors and standard deviations have to be multiplied by  $10^{-4}$ .

			Abalone	Housing	Delta Ailerons	Machine Cpu
$M$	50	Err.	2.32	5.72	<b>1.63</b>	<b>34.28</b>
		Std	0.37	0.84	0.027	4.67
	100	Err.	2.38	5.45	1.64	32.40
		Std	0.90	0.86	0.08	3.72
	200	Err.	2.20	5.31	1.65	
		Std	0.13	0.76	0.15	
	300	Err.	2.34	5.46	1.62	
		Std	1.01	1.60	0.035	

over sixteen, while OCReP in eight cases over sixteen. For both methods, better performance is achieved only for the case of  $M = 50$  neurons.

It may be noted that with respect to processing requirements OCReP has clear advantages, since it requires only a SVD step for each determination of  $\gamma$ , while the above two methods require full spectral decomposition and an additional matrix inversion.

### 6.2. Alternative regularization methods

A first comparison can be done with the work by Huang et al. (Huang *et al.*, 2012), whose technique Extreme Learning Machine (ELM) uses a cost parameter  $C$  that can be considered as related to the inverse of our regularization parameter  $\gamma$ . As authors state, in order to achieve good generalization performance,  $C$  needs to be chosen appropriately. They do this by trying 50 different values of this parameter:  $[2^{-24}, 2^{-23}, \dots, 2^{24}, 2^{25}]$ .

A fair comparison can be done on our classification datasets, using their number of hidden neurons, i.e. 1000. Our optimal choice of  $\gamma$  allows to obtain a better performance on all datasets (with statistical significance assessed at the same confidence level that previous experiments).

Deng et al. (Deng *et al.*, 2009) propose a Regularized Extreme Learning Machine (hereafter, RELM) in which the regularization parameter is selected according to a similar criterion among 100 values:  $[2^{-50}, 2^{-49}, \dots, 2^{50}]$ . Because their performance is optimized with respect to the number of hidden

Table 13: H-K estimate of ridge parameter: results at fixed number of hidden neurons for regression datasets. For Delta Ailerons, average errors and standard deviations have to be multiplied by  $10^{-4}$ .

			Abalone	Housing	Delta Ailerons	Machine Cpu
$M$	50	Err.	<b>2.13</b>	<b>4.87</b>	<b>1.60</b>	34.28
		Std	0.016	0.44	0.01	2.37
	100	Err.	2.14	4.98	1.62	37.39
		Std	0.90	0.67	0.029	3.18
	200	Err.	2.33	8.101	1.73	
		Std	0.10	2.83	0.08	
	300	Err.	2.95	29.06	2.21	
		Std	0.41	9.26	0.41	

Table 14: Comparison between OCREP and ELM

			Iris	Wine	Diabetes	Segment
OCReP	Err.	<b>2.22</b>	<b>1.28</b>	<b>21.06</b>	<b>3.40</b>	
	Std.	0.21	0.88	0.65	0.25	
ELM	Err.	2.4	1.53	22.05	3.93	
	Std	2.29	1.81	2.18	0.69	

neurons, for the sake of comparison we use OCREP values from table 6. We obtain a statistically significant better performance on dataset Segment, while for Diabetes the method RELM performs better (see table 15).

Comparing our results on the common regression datasets with the alternative method TROP-ELM proposed by Miche *et al.* (Miche *et al.*, 2011), we note that OCREP achieves always lower RMSE values <sup>3</sup> (with statistical significance), as can be seen from table 16.

Besides, in our opinion our method is simpler, in the sense that it uses a single step of regularization rather than two.

In (Martinez-Martinez *et al.*, 2011), an algorithm is proposed for pruning

---

<sup>3</sup>In that work, performance and related statistics are expressed in terms of MSE; we only derived the corresponding RMSE for comparison with our results.

Table 15: Comparison between OCREP and RELM

		Diabetes	Segment
OCReP	Err.	25.53	<b>2.50</b>
	Std.	0.51	0.32
	$\bar{M}$	291	760
RELM	Err.	<b>21.81</b>	4.49
	Std.	2.55	0.0074
	$\bar{M}$	15	200

Table 16: Comparison between OCREP and TROP-ELM. For Delta Ailerons, average errors and standard deviations have to be multiplied by  $10^{-4}$ .

		Abalone	Delta Ailerons	Machine Cpu	Housing
OCReP	Err.	<b>2.12</b>	<b>1.58</b>	<b>31.22</b>	<b>4.25</b>
	Std.	0.32	0.0048	0.78	0.13
	$\bar{M}$	178	298	63	255
TROP-ELM	Err.	2.19	1.64	264.03	34.35
	$\bar{M}$	42	80	28	59

ELM networks by using regularized regression methods: the crucial step of regularization parameter determination is solved by creating  $K$  different models, each one based on a different value of this parameter, among which the best one is selected using a Bayesian information criterion. Authors state that a typical value for  $K$  is 100, thus an heavy computational load is required, and the method is focused on regression tasks.

## 7. Conclusions

In the context of regularization techniques for single hidden layer neural networks trained by pseudoinversion, we provide an optimal value of the regularization parameter  $\gamma$  by analytic derivation. This is achieved by defining a convenient regularized matricial formulation in the framework of Singular Value Decomposition, in which the regularization parameter is derived under the constraint of condition number minimization. The OCREP method

has been tested on UCI datasets for both regression and classification tasks. For all cases, regularization implemented using the analytically derived  $\gamma$  is proven to be very effective in terms of predictivity, as evidenced by comparison with implementations of other approaches from the literature, including cross-validation. OCR<sub>e</sub>P avoids hundreds of pseudoinversions usually needed by most other methods, i.e. it is quite computationally attractive.

## Acknowledgements

The activity has been partially carried on in the context of the Visiting Professor Program of the Gruppo Nazionale per il Calcolo Scientifico (GNCS) of the Italian Istituto Nazionale di Alta Matematica (INdAM). This work has been partially supported by ASI contracts (Gaia Mission - The Italian Participation to DPAC) I/058/10/0-1 and 2014-025-R.0.

## References

- Ajorloo, H., Manzuri-Shalmani, M. T., and Lakdashti, A. (2007). Restoration of damaged slices in images using matrix pseudo inversion. In *Proceedings of the 22nd International symposium on computer and information sciences*, pages 1–6.
- Bache, K. and Lichman, M. (2013). UCI machine learning repository.
- Badeva, V. and Morozov, V. (1991). *Problèmes incorrectement posés: Théorie et applications en identification , filtrage optimal, contrôle optimal, analyse et synthèse de systèmes, reconnaissance d'images*. Série Automatique. Masson.
- Berger, J. (1976). Minimax estimation of a multivariate normal mean under arbitrary quadratic loss. *J. Multivariate Analysis*, **6**, 256–264.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- Bousquet, O. and Elisseeff, A. (2002). Stability and generalization. *J. Mach. Learn. Res.*, **2**, 499–526.

- Cancelliere, R. (2001). A high parallel procedure to initialize the output weights of a radial basis function or bp neural network. In *Proceedings of the 5th International Workshop on Applied Parallel Computing, New Paradigms for HPC in Industry and Academia*, PARA '00, pages 384–390, London, UK, UK. Springer-Verlag.
- Deng, W., Zheng, Q., and Chen, L. (2009). Regularised extreme learning machine. In *Proceedings of the IEEE Symposium on Computational Intelligence and Data Mining*.
- Devroye, L. P. and Wagner, T. (1979). Distribution-free performance bounds for potential function rules. *Information Theory, IEEE Transactions on*, **25**(5), 601–604.
- Dorugade, A. and Kashid, D. (2010). Alternative method for choosing ridge parameter for regression. *Applied Mathematical Sciences*, **4**, 447–456.
- Fu, W. (1998). Penalized regressions: the bridge vs. the lasso. *Journal of Computational and Graphical Statistics*, **7**, 397–416.
- Fuhry, M. and Reichel, L. (2012). A new tikhonov regularization method. *Numerical Algorithms*, **59**(3), 433–445.
- Gallinari, P. and Cibas, T. (1999). Practical complexity control in multilayer perceptrons. *Signal Processing*, **74**, 29–46.
- Geman, S., Bienenstock, E., and Doursat, R. (1992). Neural networks and the bias/variance dilemma. *Neural Computation*, **4**, 1–58.
- Girosi, F., Jones, M., and Poggio, T. (1995). Regularization Theory and Neural Networks Architectures. *Neural Computation*, **7**, 219–269.
- Golub, G., Heath, M., and Wahba, G. (1979). Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, **21**, 215–223.
- Golub, G. H. and Van Loan, C. F. (1996). *Matrix computations (3rd ed.)*. Johns Hopkins University Press, Baltimore, MD, USA.
- Golub, G. H., Hansen, P. C., and O’Leary, D. P. (1999). Tikhonov regularization and total least squares. *SIAM J. Matrix Anal. Appl.*, **21**(1), 185–194.



- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of Statistical Learning: Data Mining, Inference, and Prediction (2nd ed.)*. Springer.
- Haykin, S. (1999). *Neural Networks: A Comprehensive Foundation*. International edition. Prentice Hall.
- Helmy, T. and Rasheed, Z. (2009). Multi-category bioinformatics dataset classification using extreme learning machine. In *Proceedings of the Eleventh conference on Congress on Evolutionary Computation, CEC'09*, pages 3234–3240, Piscataway, NJ, USA. IEEE Press.
- Hoerl, A. (1962). Application of ridge analysis to regression problems. *Chemical Engineering Progress*, **58**, 54–59.
- Hoerl, A. and Kennard, R. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, **12**, 55–67.
- Hoerl, A. and Kennard, R. (1976). Ridge regression: iterative estimation of the biasing parameter. *Communications in Statistics*, **A5**, 77–88.
- Huang, G.-B., Zhu, Q.-Y., and Siew, C.-K. (2006). Extreme learning machine: theory and applications. *Neurocomputing*, **70**(1), 489–501.
- Huang, G.-B., Zhu, Q.-Y., and Siew, C.-K. (2012). Extreme learning machine for regression and multiclass classification. *IEEE Transactions on Systems, Man, and Cybernetics–Part B: Cybernetics*, **42**(2), 513–529.
- Kearns, M. and Ron, D. (1999). Algorithmic stability and sanity-check bounds for leave-one-out cross-validation. *Neural Computation*, **11**(6), 1427–1453.
- Khalaf, G. and Shukur, G. (2005). Choosing ridge parameter for regression problem. *Communications in Statistics – Theory and methods*, **34**, 1177–1182.
- Kibria, B. (2003). Performance of some new ridge regression estimators. *Communications in Statistics – Simulation and Computation*, **32**, 419–435.
- Kohno, K., Kawamoto, M., and Inouye, Y. (2010). A matrix pseudoinversion lemma and its application to block-based adaptive blind deconvolution for mimo systems. *Trans. Cir. Sys. Part I*, **57**(7), 1449–1462.

- Lawless, J. and Wang, P. (1976). A simulation study of ridge and other regression estimators. *Communications in Statistics*, **A5**, 307–324.
- Mardikyan, S. and Cetin, E. (2008). Efficient choice of biasing constant for ridge regression. *International Journal of Contemporary Mathematical Sciences*, **3**, 527–547.
- Martinez-Martinez, J., Escandell-Montero, P., Soria-Olivas, E., Martn-Guerrero, J., Magdalena-Benedito, R., and Juan, G.-S. (2011). Regularized extreme learning machine for regression problems. *Neurocomputing*, **74**, 3716–3721.
- McDonald, G. and Galarneau, D. (1975). A monte carlo evaluation of some ridge-type estimators. *J. Amer. Statist. Assoc.*, **70**, 407–416.
- Miche, Y., van Heeswijk, M., Bas, P., Simula, O., and Lendasse, A. (2011). Trop-elm: A double-regularized elm using lars and tikhonov regularization. *Neurocomputing*, **74**(16), 2413 – 2421.
- Mukherjee, S., Niyogi, P., Poggio, T., and Rifkin, R. (2003). Statistical learning: stability is sufficient for generalization and necessary and sufficient for consistency of empirical risk minimization. *CBCL, Paper 223, Massachusetts Institute of Technology*.
- Mukherjee, S., Niyogi, P., Poggio, T., and Rifkin, R. (2006). Learning theory: stability is sufficient for generalization and necessary and sufficient for consistency of empirical risk minimization. *Advances in Computational Mathematics*, **25**(1-3), 161–193.
- Nguyen, T. D., Pham, H. T. B., and Dang, V. H. (2010). An efficient pseudo inverse matrix-based solution for secure auditing. In *Proceedings of the IEEE International Conference on Computing and Communication Technologies, Research, Innovation, and Vision for the Future*, IEEE International Conference.
- Nordberg, L. (1982). A procedure for determination of a good ridge parameter in linear regression. *Communications in Statistics*, **A11**, 285–309.
- Penrose, R. and Todd, J. A. (1956). On best approximate solutions of linear matrix equations. *Mathematical Proceedings of the Cambridge Philosophical Society*, **null**, 17–19.

- Poggio, T. and Girosi, F. (1990a). Networks for approximation and learning. *Proceedings of the IEEE*, **78**(9), 1481–1497.
- Poggio, T. and Girosi, F. (1990b). Regularization algorithms for learning that are equivalent to multilayer networks. *Science*, **247**(4945), 978–982.
- Poggio, T., Rifkin, R., Mukherjee, S., and Niyogi, P. (2004). General conditions for predictivity in learning theory. *Letters to Nature*, **428**, 419/422.
- Rao, C. and Mitra, S. (1971). *Generalized inverse of matrices and its applications*. Wiley series in probability and mathematical statistics: Applied probability and statistics. Wiley.
- Saleh, A. and Kibria, B. (1993). Performances of some new preliminary test ridge regression estimators and their properties. *Communications in Statistics – Theory and methods*, **22**, 2747–2764.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of Royal Statistical Society*, **58**, 267–288.
- Tikhonov, A. and Arsenin, V. (1977). *Solutions of ill-posed problems*. Scripta series in mathematics. Winston, Washington DC.
- Tikhonov, A. N. (1963). Solution of incorrectly formulated problems and the regularization method. *Soviet Math. Dokl.*, **4**, 1035–1038.
- Wang, X. (2013). Special issue on extreme learning machines with uncertainty. *Int. J. Unc. Fuzz. Knowl. Based Syst.*, **21**.
- Wang, X.-Z., D., W., and Huang, G.-B. (2012). Special issue on extreme learning machines. *Soft Computing*, **16**(9), 1461–1463.
- Yu, D. and Deng, L. (2012). Efficient and effective algorithms for training single-hidden-layer neural networks. *Pattern Recogn. Lett.*, **33**(5), 554–558.