

# TRADING APPROXIMATION QUALITY VERSUS SPARSITY WITHIN INCREMENTAL AUTOMATIC RELEVANCE DETERMINATION FRAMEWORKS

*Dmitriy Shutin*

Institute for Communications and Navigation  
German Aerospace Center (DLR)  
Münchner Str. 20, 82234 Wessling Germany

*Thomas Buchgraber\**

Signal Processing and Speech Comm. Lab  
Graz University of Technology  
Inffeldgasse 16c, 8010 Graz, Austria

## ABSTRACT

In this paper a trade-off between sparsity and approximation quality of models learned with incremental automatic relevance determination (IARD) is addressed. An IARD algorithm is a class of sparse Bayesian learning (SBL) schemes. It permits an intuitive and simple adjustment of estimation expressions, with the adjustment having a simple interpretation in terms of signal-to-noise ratio (SNR). This adjustment allows for implementing a trade-off between sparsity of the estimated model versus its accuracy in terms of residual mean-square error (MSE). It is found that this adjustment has a different impact on the IARD performance, depending on whether the measurement model coincides with the used estimation model or not. Specifically, in the former case the value of the adjustment parameter set to the true SNR leads to an optimum performance of the IARD with the smallest MSE and estimated signal sparsity; moreover, the estimated sparsity then coincides with the true signal sparsity. In contrast, when there is a model mismatch, the lower MSE can be achieved only at the expense of less sparser models. In this case the adjustment parameter simply trades the estimated signal sparsity versus the accuracy of the model.

## 1. INTRODUCTION

The problem of recovering a sparse representations from a noisy measurement has received considerable attention in recent years [1–4]. Typically, the goal is to find an  $L$ -dimensional vector  $\mathbf{w} = [w_1, \dots, w_L]^T$ , which is assumed to contain only  $K$  nonzero elements,  $K \ll L$ , from a noise-perturbed measurement vector  $\mathbf{t}$  generated from the following linear measurement model:

$$\mathbf{t} = \Phi \mathbf{w} + \boldsymbol{\xi}. \quad (1)$$

Here  $\Phi = [\phi_1, \dots, \phi_L]$  is a dictionary matrix with  $L$  columns corresponding to some fixed signal waveforms

$\phi_l \in \mathbb{R}^N$ ,  $l = 1, \dots, L$ . The additive perturbation  $\boldsymbol{\xi}$  is typically assumed to be a white Gaussian random vector with zero mean and covariance matrix  $\tau^{-1} \mathbf{I}$ , where  $\tau$  is a noise precision parameter. Despite the linearity of (1), the additional sparsity constraints imposed on the weight vector  $\mathbf{w}$  make the estimation problem nontrivial. Specifically, learning the sparse vector  $\mathbf{w}$  requires imposing formal constraints on the objective function for the model parameters  $\mathbf{w}$  to ensure the resulting solution is sparse [2–6].

Sparse Bayesian learning (SBL) [7, 8] is a class of such learning schemes where these constraints are introduced through the use of a hierarchical prior. Specifically, the prior on  $\mathbf{w}$  is introduced as  $p(\mathbf{w}|\boldsymbol{\alpha})p(\boldsymbol{\alpha})$ , where  $p(\mathbf{w}|\boldsymbol{\alpha})$  is a product of Gaussian probability density functions (pdfs) with zero mean and precision parameters  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_L]^T$  – also called sparsity parameters – that regulate the width of these pdfs.<sup>1</sup> The weights  $\mathbf{w}$  and the sparsity parameters  $\boldsymbol{\alpha}$  are then found from the mode of the corresponding posterior pdfs.

Recently, it has been shown that SBL can be realized very efficiently when the hyperprior  $p(\boldsymbol{\alpha})$  is assumed to be flat [10, 11] or non-informative [12]<sup>2</sup>. These special cases of SBL, commonly referred to as automatic relevance determination (ARD), have several important theoretical as well as practical advantages. First, the ARD formulation of the hyperprior  $p(\boldsymbol{\alpha})$  is related to weighted versions of minimum  $\ell_1$ -norm regression and basis pursuit denoising (see [2, 6, 11, 14]) – more traditional “non-Bayesian” methods for learning sparse representations. Second, the ARD formulation can be exploited to construct very fast and efficient inference schemes [10–12] – incremental ARD (IARD).

Typically, the weights  $\mathbf{w}$  are determined from the mode of the posterior pdf  $p(\mathbf{w}|\mathbf{t}, \boldsymbol{\alpha}, \tau)$  that can be computed analytically; the sparsity parameters  $\boldsymbol{\alpha}$  and the noise precision  $\tau$  are then estimated by maximizing the posterior pdf  $p(\boldsymbol{\alpha}, \tau|\mathbf{t})$ . The latter optimization is generally analytically intractable,

<sup>1</sup>It is also possible to extend the SBL prior formulation to priors involving three layers of hierarchy [9].

<sup>2</sup>Note that while these priors are improper, i.e., they do not integrate to one, in our case the resulting weight posterior under such hyperprior is, nonetheless, a proper pdfs (see [13] for more details).

\*This work was supported in part by the Austrian Science Fund (FWF) under Award S10610-N13 within the national research network SISE

yet using the IARD scheme it can be solved very efficiently. Specifically, it can be shown [10, 12] that the optimum of the corresponding objective function with respect to the sparsity parameter of the  $l$ th component can be computed in closed form provided the other sparsity parameters  $\alpha_k$ ,  $k \neq l$ , and the noise precision parameter  $\tau$  are fixed. The posterior pdf  $p(\alpha, \tau | \mathbf{t})$  is then optimized incrementally with respect to  $\alpha$  by cycling through  $L$  sparsity parameters in a round-robin fashion [10, 12]. An important consequence of this algorithm is that the analytical analysis of the corresponding inference expressions allows for deriving exact conditions, termed pruning conditions, that determine if the sparsity parameter that optimizes the posterior  $p(\alpha, \tau | \mathbf{t})$  is finite or infinite<sup>3</sup>.

The pruning conditions are the key features of the IARD scheme that differentiate it from the other sparse learning techniques. In contrast to traditional minimum  $\ell_1$ -norm methods, where the recovery of a sparse solution is established by guaranteeing (with some high probability) an upper bound on the norm of the error between the estimated vector  $\hat{\mathbf{w}}$  and a true sparse vector  $\mathbf{w}_0$  [2, 15, 16], the pruning conditions obtained with IARD are “non-asymptotic”; they determine which elements in the weight vector  $\mathbf{w}$  must be set to zero. More importantly, the structure of these conditions can be used to “explain” why a specific weight has been set to zero. In particular, it has been shown [12] that the pruning conditions depend on an estimate of the per-component signal-to-noise ratio (SNR). This interpretation allows for an empirical tuning of the pruning conditions and, as a result, of the estimated signal sparsity.

The goal of this paper is to analyze the performance of IARD schemes as a function of this adjustment in two practically relevant scenarios: (i) when the measurement model behind  $\mathbf{t}$  coincides with the estimation model (1), and thus the measurement  $\mathbf{t}$  can be optimally represented with exactly  $K$  components of  $\Phi$ , and (ii) when (1) is used as an approximation to the actual measurement model. In the latter case  $\mathbf{t}$  is only “approximately” sparse in  $\Phi$ ; some elements of  $\mathbf{w}$  thus have small, yet larger than zero magnitudes even in noise-free cases. Therefore, it makes sense to inquire if in the second scenario it is possible to trade the sparsity of the resulting model versus the approximation quality and how this trading is to be realized. In this work we demonstrate that the IARD scheme with adjusted pruning condition effectively realizes this trade-off. Specifically, we show that in the model mismatch case the tuning of the pruning conditions effectively trades the sparsity of the estimated signal versus the quality of the resulting models in terms of mean-square error (MSE). In contrast, in the first scenario, there is no trade-off between sparsity and residual MSE; on the contrary, an optimal performance is achieved by setting the adjustment level to the actual SNR. This effectively minimizes the residual MSE, while at the same time estimates the correct signal sparsity.

<sup>3</sup>An infinite value of the hyperparameter forces the posterior value of the component weight to zero.

The rest of the paper is organized as follows. In Section 2 we give an outline of the SBL signal model. In Section 3 we explain the incremental ARD solution to the SBL problem and discuss the adjustment of the pruning conditions. Simulation results that demonstrate the impact of the adjusted pruning conditions on the performance of IARD schemes are discussed in Section 4.

## 2. SBL SIGNAL MODEL

A standard solution to SBL with ARD is a two step procedure that alternates between estimating the weight vector  $\mathbf{w}$  and estimating the corresponding sparsity parameters  $\alpha$  and noise precision  $\tau$ . Given an estimate of the noise precision parameter  $\hat{\tau}$  and sparsity parameters  $\hat{\alpha}$ , an estimate of the weight vector  $\hat{\mathbf{w}}$  is obtained as a mode of the posterior pdf  $p(\mathbf{w} | \mathbf{t}, \hat{\alpha}, \hat{\tau}) = p(\mathbf{t} | \mathbf{w}, \hat{\tau}) p(\mathbf{w} | \hat{\alpha})$ , which can be shown to be a Gaussian pdf with the mean  $\hat{\mathbf{w}}$  and covariance matrix  $\hat{\mathbf{S}}$  given by

$$\hat{\mathbf{S}} = \left( \hat{\tau} \Phi^T \Phi + \hat{\mathbf{A}} \right)^{-1} \quad \text{and} \quad \hat{\mathbf{w}} = \hat{\tau} \hat{\mathbf{S}} \Phi^T \mathbf{t}, \quad (2)$$

where  $\hat{\mathbf{A}} = \text{diag}(\hat{\alpha})$  is a diagonal matrix with the elements of  $\hat{\alpha}$  on the main diagonal. The estimates of  $\alpha$  and  $\tau$  are computed by maximizing  $p(\alpha, \tau | \mathbf{t}) \propto (\mathbf{t} | \alpha, \tau) p(\alpha) p(\tau)$ , where

$$\begin{aligned} p(\mathbf{t} | \alpha, \tau) &= \int p(\mathbf{t} | \mathbf{w}, \tau) p(\mathbf{w} | \alpha) d\mathbf{w} \\ &= \mathcal{N}(\mathbf{t} | \mathbf{0}, \tau^{-1} \mathbf{I} + \Phi \mathbf{A}^{-1} \Phi^T) \end{aligned} \quad (3)$$

is also known as a marginal likelihood function [7]. The maximizers  $\hat{\alpha}$  and  $\hat{\tau}$  of  $p(\alpha, \tau | \mathbf{t})$  can be obtained using the EM algorithm where the weights  $\mathbf{w}$  are used as complete data (see [7] for more details). For the ARD case the corresponding estimation expressions are then given as

$$\hat{\alpha}_l = (|\hat{w}_l|^2 + \hat{S}_{ll})^{-1} \quad (4)$$

$$\text{and} \quad \hat{\tau} = \frac{N}{\|\mathbf{t} - \Phi \hat{\mathbf{w}}\|^2 + \text{Trace}(\hat{\mathbf{S}} \Phi^T \Phi)}, \quad (5)$$

where  $\hat{w}_l$  is the  $l$ th element of the vector  $\hat{\mathbf{w}}$ , and  $\hat{S}_{ll}$  is the  $l$ th element on the main diagonal of the matrix  $\hat{\mathbf{S}}$ . The update expressions (2) and (4)–(5) are then repeatedly evaluated until convergence [7].

In cases when  $\mathbf{t}$  allows for a sparse representation in terms of  $\Phi$ , the sparsity parameters  $\hat{\alpha}$  of some of the components will diverge, forcing the posterior value of the corresponding weights in (2) to converge toward zero and, thus, encouraging a sparse solution.

## 3. INCREMENTAL ARD SCHEME

Unfortunately, due to the EM-based maximization of (3) the rate at which sparsity parameters diverge is very low;

many iterations are needed for the hyperparameters to reach a threshold at which they can be treated as “numerically” infinite.<sup>4</sup> This has motivated the use of alternative, more efficient schemes to maximize (3) with respect to  $\alpha$  [10–12, 17].

A more efficient algorithm can be obtained for a flat hyperprior  $p(\alpha)$ .<sup>5</sup> Surprisingly, the optimum of  $p(\mathbf{t}|\alpha, \tau)$  with respect to a single sparsity parameter  $\alpha_l$ , assuming that the other sparsity parameters  $\hat{\alpha}_{\bar{l}} = [\hat{\alpha}_1, \dots, \hat{\alpha}_{l-1}, \hat{\alpha}_{l+1}, \dots, \hat{\alpha}_L]^T$  and the noise precision parameter  $\hat{\tau}$  are fixed, can be computed in closed form. Specifically, the logarithm of  $p(\mathbf{t}|\alpha_l, \hat{\alpha}_{\bar{l}}, \hat{\tau})$  can be decomposed as [10]

$$\log p(\mathbf{t}|\alpha_l, \hat{\alpha}_{\bar{l}}, \hat{\tau}) = \mathcal{L}(\alpha_l; \hat{\alpha}_{\bar{l}}, \hat{\tau}) = \mathcal{L}(\hat{\alpha}_{\bar{l}}, \hat{\tau}) + \frac{1}{2} \left( \log(\alpha_l) - \log(\alpha_l - s_l) + \frac{q_l^2}{\alpha_l + s_l} \right), \quad (6)$$

where

$$s_l = \phi_l^T \hat{\mathbf{C}}_{\bar{l}}^{-1} \phi_l, \quad q_l = \phi_l^T \hat{\mathbf{C}}_{\bar{l}}^{-1} \mathbf{t}, \quad (7)$$

$$\hat{\mathbf{C}}_{\bar{l}} = \hat{\tau}^{-1} \mathbf{I} + \sum_{k \neq l} \hat{\alpha}_k^{-1} \phi_k \phi_k^T, \quad (8)$$

and  $\mathcal{L}(\hat{\alpha}_{\bar{l}}, \hat{\tau})$  is a part of the marginal log-likelihood that is independent of  $\alpha_l$ . Then, the maximum of (6) with respect to  $\alpha_l$  is obtained at [10, 12]

$$\hat{\alpha}_l = \begin{cases} s_l^2 (q_l^2 - s_l)^{-1}, & q_l^2 > s_l, \\ \infty, & q_l^2 \leq s_l. \end{cases} \quad (9)$$

Using (9) the marginal likelihood  $p(\mathbf{t}|\alpha_l, \hat{\alpha}_{\bar{l}}, \hat{\tau})$  can be maximized incrementally with respect to the sparsity parameter of one component at a time. Note that the result (9) not only tells if a particular element in  $\hat{\mathbf{w}}$  is zero, but also dramatically accelerates the rate of SBL convergence since the optimum of the marginal likelihood with respect to a single sparsity parameters can be computed analytically [10, 12, 17].

### 3.1. Pruning condition of IARD schemes

The key to sparsity inducing properties of the incremental ARD schemes lies in the pruning condition (9) that determines which elements in the vector  $\mathbf{w}$  are set to zero.<sup>6</sup> Here we study these conditions in greater detail.

Let us define  $\mu_l = q_l/s_l$  and  $\varsigma_l = s_l^{-1}$ . In [12] it has been demonstrated that  $\mu_l$  and  $\varsigma_l$  correspond, respectively, to the  $l$ th element of the vector  $\hat{\mathbf{w}}$  and the  $l$ th diagonal element of the matrix  $\hat{\mathbf{S}}$  in (2) computed when the sparsity parameter  $\hat{\alpha}_l$  is set to zero, i.e., when the basis vector  $\phi_l$  is not regularized<sup>7</sup>.

<sup>4</sup>For instance, this threshold can be set to  $10^{15}$  or  $10^{16}$ .

<sup>5</sup>Using variational Bayesian analysis a similar scheme can be derived for non-informative hyperpriors  $p(\alpha_l) \propto \alpha_l^{-1}$ ,  $l = 1, \dots, L$ , [12].

<sup>6</sup>Strictly speaking, the order in which basis vectors are processed also plays a role, in particular when  $\Phi$  has a high mutual coherence.

<sup>7</sup>From (2) it follows that  $\hat{\alpha}$  is a vector of regularization parameters for computing  $\hat{\mathbf{w}}$ .

Naturally, for  $\varsigma_l > 0$ , the pruning condition  $q_l^2 > s_l$  in (9) is equivalent to

$$\frac{\mu_l^2}{\varsigma_l} > 1, \quad (10)$$

where  $\mu_l^2/\varsigma_l$  can be recognized as an estimate of the  $l$ th component SNR. In other words, a finite estimate of  $\hat{\alpha}_l$  is obtained provided an estimate of the  $l$ th component SNR, given by  $\mu_l^2/\varsigma_l$ , exceeds a 0dB threshold. This is a very simple and intuitive result: it demonstrates that the condition (10) determines if the quality of the  $l$ th component (as measured by its SNR) is above the level of the residual noise after the impact of the other  $L - 1$  components is taken into account. Quite naturally, the condition (10) can be adjusted (or tuned) by requiring that

$$\frac{\mu_l^2}{\varsigma_l} > \eta_l \quad (11)$$

where  $\eta_l$  is some predefined SNR level; in other words, only components with a certain SNR level  $\eta_l \geq 1$  have to be retained in the model. Note that (11) is an empirical adjustment since for  $\eta_l > 1$  the resulting estimate  $\hat{\alpha}$  no longer maximizes  $p(\mathbf{t}|\alpha, \hat{\tau})$ . Nonetheless, this adjustment can be motivated from the perspective of statistical hypothesis testing (see [18] for more details).

The condition (11) is a key to understanding how sparsity of the resulting model can be traded for the model approximation quality. Further in the text we will term the IARD scheme that makes use of pruning condition (11) as  $\eta$ -IARD algorithm.

In what follows we investigate the impact of the adjustment parameter  $\eta_l$  on the performance of  $\eta$ -IARD algorithm.

## 4. EXPERIMENTAL ANALYSIS

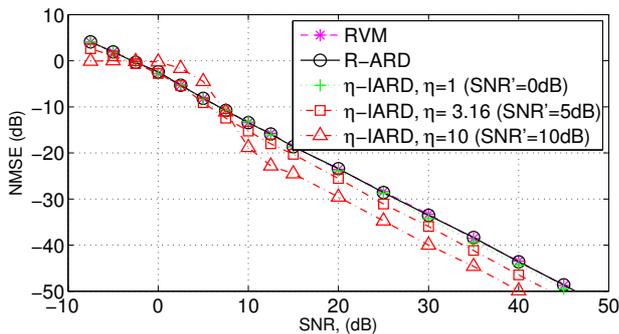
In order to demonstrate the impact of the parameter  $\eta_l$  on the performance of  $\eta$ -IARD scheme, three experimental setups are considered. In the first setting we consider a standard compressive sampling toy problem, where the data  $\mathbf{t}$  is generated using (1); in other words, there is no model mismatch. In the second setting we consider an application of  $\eta$ -IARD scheme to two non-parametric regression problems: (a) a sinc regressions problem as used in [7], and (b) a sparse regression of real data taken from UCI Machine Learning Repository. The last two examples are aimed to demonstrate the performance of the  $\eta$ -IARD scheme when there is a model mismatch and parameter  $\eta_l$  is used to trade model sparsity versus residual MSE.

### 4.1. Compressive sampling toy problem

In this experiment we study the impact of  $\eta_l$  using a compressive sampling toy problem. Here the basis functions  $\phi_l \in \mathbb{R}^N$ ,  $l = 1, \dots, L$ , are generated by drawing  $N = 100$  samples from a Gaussian distribution with zero mean and variance 1. A sparse vector  $\mathbf{w}$  is constructed by setting  $K = 10$

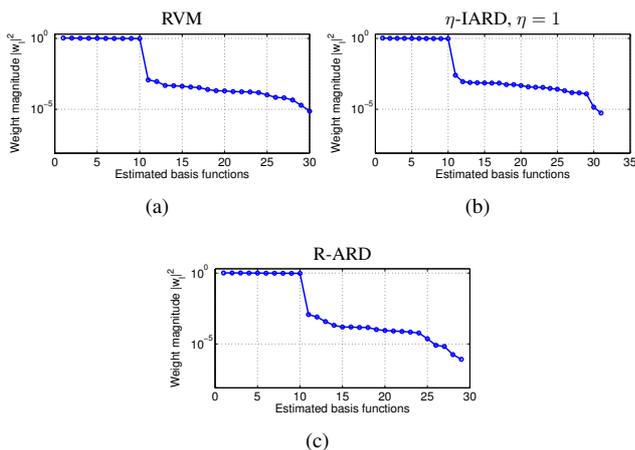
elements of  $\mathbf{w}$  to 1 at random locations. Note that in this setting each component has the same SNR; thus, we assume that  $\eta_1 = \dots = \eta_L = \eta$ .<sup>8</sup> The target vector  $\mathbf{t}$  is then generated according to (1).

We begin by computing the normalized mean-square error (NMSE) of the estimation schemes versus the SNR. For comparison purposes we also include the simulation results for a Relevance Vector Machine (RVM) [7] and reformulated ARD (R-ARD) [11].  $\eta$ -IARD schemes are simulated with  $\eta = 1, 3.16$  and  $10$ , which corresponds to  $0\text{dB}, 5\text{dB}$ , and  $10\text{dB}$  adjustment SNR', respectively. For simplicity we will refer to the values of adjustment parameter  $\eta_l$  in dB-scale. The corresponding simulation results are summarized in Fig. 1. Note



**Fig. 1.** Normalized mean-square error (NMSE) versus SNR.

that RVM, R-ARD and  $\eta$ -IARD with  $\eta = 1$  perform comparably in terms of NMSE. Yet increasing the value of  $\eta$  improves the algorithm performance. The reason for that is an underestimated model sparsity due to the presence of noise. In fact, many of the estimated components have very small, yet non-zero weights (see Fig.2). The introduction of the ad-



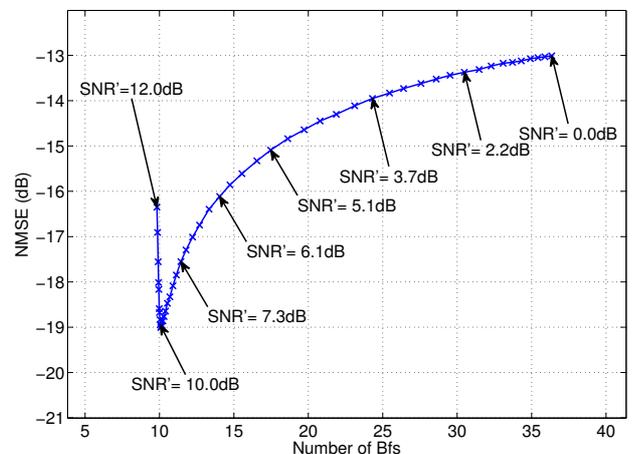
**Fig. 2.** The magnitude of the non-zero component weights for SNR fixed at  $30\text{ dB}$ . a) RVM algorithm; b)  $\eta$ -IARD for  $\eta = 1$  algorithm; c) R-ARD algorithm.

<sup>8</sup>In further experiments we will also make use of this assumption.

justment  $\eta$  effectively prunes the noisy components, thus improving the sparsity estimation and, as the result, the NMSE.

Let us also point to a distinctive threshold-like behavior of the  $\eta$ -IARD schemes: their performance in terms of NMSE is largely insensitive to the level of the noise when  $\eta$  is set above the actual SNR. However, when the actual SNR level approaches the selected value of the adjustment  $\eta$ , the NMSE improves rapidly and then continues to drop linearly<sup>9</sup> with the increasing SNR level.

Let us now consider the performance of the  $\eta$ -IARD scheme in terms of NMSE value and estimated sparsity as a function of the adjustment level  $\eta$ . For that we fix the signal SNR at  $10\text{ dB}$  and consider the  $\eta$ -IARD algorithm with different settings of the adjustment parameter  $\eta$ . The corresponding results, averaged over 10000 Monte Carlo runs, are summarized in Fig. 3. Observe that for small  $\eta$  the number of



**Fig. 3.** NMSE and the estimated number of non-zero components as a function of the adjustment  $\eta$  for true SNR fixed at  $10\text{ dB}$ . SNR' corresponds to the adjustment  $\eta$  in dB-scale.

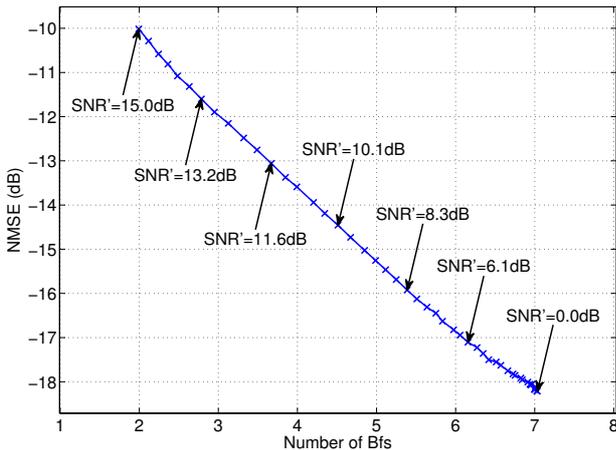
estimated components is high. As  $\eta$  grows, the NMSE value as well as the number of the estimated components drops. Note that the performance of the RVM and R-ARD schemes (not shown in Fig.3) is very similar to the performance of the  $\eta$ -IARD scheme with  $\eta = 1$ , i.e., the sparsity is significantly underestimated. In case of the  $\eta$ -IARD scheme, setting the value of  $\eta$  to the actual SNR of  $10\text{ dB}$  leads to the minimum of the NMSE; also, the estimated number of components coincides with the true signal sparsity. Obviously, there is a value of the parameter  $\eta$  that optimizes the performance of the  $\eta$ -IARD scheme. Specifically, in noisy scenarios the true sparsity can be recovered provided the value of  $\eta$  is selected to correspond to the actual signal SNR; at the same time the residual NMSE is minimized. This observation is also supported by our previous investigations [12, 17].

<sup>9</sup>Over a logarithmic scale.

## 4.2. Non-parameteric sinc regression problem

In this experiment we consider a reconstruction of a noise-perturbed  $\text{sinc}(x) = \sin(x)/x$  function based on  $N = 100$  randomly-spaced samples in the interval  $[-10, 10]$ . The measurements  $\mathbf{t}$  are generated by evaluating  $\text{sinc}(x)$  at sampled locations and perturbing them using a zero-mean normally distributed white noise with variance  $\sigma^2$ . The variance of the noise is chosen such as to ensure a desired SNR value. Reconstruction is implemented using a set of  $L = 100$  Gaussian kernels  $\phi_l$ ,  $l = 1, \dots, L$  centered at measured samples and having a fixed variance of 2.3. Note that the measurement  $\mathbf{t}$  generated with the sinc model is only “approximately” sparse in  $\Phi$ . The noise precision  $\tau$  has been automatically estimated for each run of the algorithm. This example demonstrates the performance of the  $\eta$ -IARD scheme in this case.

Due to space limitations we will not evaluate the NMSE performance as a function of SNR for the following examples. Instead, we go directly to the dependency of NMSE and estimated sparsity on the value of the adjustment parameter  $\eta$  for the SNR level fixed at 10dB. The corresponding results, averaged over 10000 Monte Carlo runs, are summarized in Fig. 4. Note that in contrast to the previous scenario,



**Fig. 4.** NMSE and the estimated number of non-zero components as a function of the adjustment  $\eta$  for sinc regression problem.

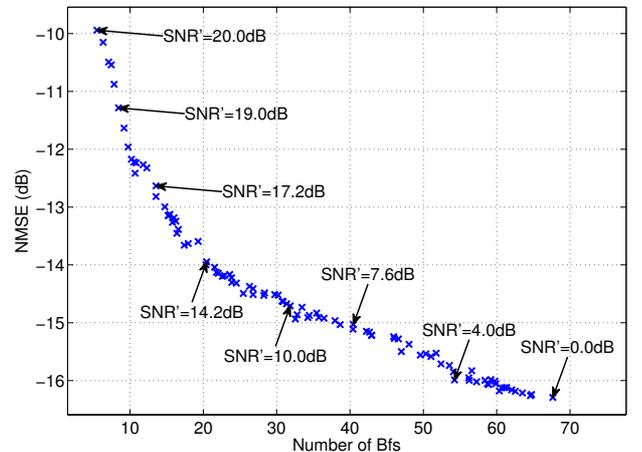
the resulting curve does not have a distinct optimum. Instead, higher values of  $\eta$  increase the sparsity of the model at the expense of higher NMSE; in other words, the adjustment parameter  $\eta$  monotonically trades the approximation quality of the learned model versus estimated signal sparsity.

Let us stress that in practice the model (11) is often used as a mere approximation to the true measurement model. As a result the sparsity is not exact. The result (11) thus provides a mechanism to adjust the IARD scheme so as to account for the model mismatches. More specifically, the parameter  $\eta$  specifies the minimum quality of the component, as measured

by its SNR, that is to be retained in the model. Quite naturally, a specific value of  $\eta$  is typically problem dependent.

## 4.3. Concrete Compressive Strength data

In this example we evaluate the algorithm performance on a real data set – the Concrete Compressive Strength (CCS) data from the UCI Machine Learning Repository. It is a multivariate regression data set with 8 attributes and 1030 instances. To evaluate the model performance a 20-fold cross-validation has been used. Both training and test data outcomes are normalized to have zero mean and unit variance. The dictionary matrix  $\Phi$  consists of Gaussian kernels  $\phi_l$  centered at the measurement samples. The variance of the kernels is fixed at 4.3. The noise precision  $\tau$  has been automatically estimated for each run of the algorithm. In Fig. 5 we plot the dependency of NMSE and estimated signal sparsity on the value of  $\eta$ . Observe that the behavior of the NMSE-sparsity curve as a



**Fig. 5.** NMSE and the estimated number of non-zero components as a function of the adjustment  $\eta$  for Concrete Compressive Strength data.

function of  $\eta$  is very similar to the sinc-regression case. Indeed, due to the model mismatch, the parameter  $\eta$  can be used to trade the sparsity of the resulting model versus the value of the residual NMSE.

## 5. CONCLUSION

In this paper the performance of the adjusted incremental automatic relevance determination (IARD) algorithm has been analyzed. The considered adjustment allows for a trade-off between the sparsity and the approximation quality of the learned models. The latter is realized by modifying the pruning condition of the IARD scheme. The modification has a simple interpretation in terms of the per-component signal-to-noise ratio (SNR).

The impact of the modification has been studied in two scenarios. In the first scenario the estimation model coincides with the measurement model; in the second one the estimation model is used as an approximation to the true model behind the measurement data. It has been determined that in the first case, the performance of the IARD scheme has a distinct optimum, which was achieved by setting the adjustment level to the true signal SNR. In this way the mean-square error (MSE) is minimized and the true signal sparsity is recovered. In the second case, the MSE is minimized only at the expense of less sparse models; fewer nonzero components inevitably leads to a higher MSE. In this case the adjustment of the pruning condition can be seen as a trade-off between signal sparsity and the achieved MSE.

## 6. REFERENCES

- [1] R. Baraniuk, "Compressive sensing," *IEEE Signal Processing Magazine*, vol. 24, no. July, pp. 118–121, 4 2007.
- [2] M. Duarte and Y. Eldar, "Structured compressed sensing: From theory to applications," *IEEE Trans. on Sig. Proc.*, vol. 59, no. 9, pp. 4053–4085, sept. 2011.
- [3] D. Donoho, "For most large underdetermined systems of linear equations, the minimal  $\ell_1$  norm solution is also the sparsest solution," *Comm. on Pure and Appl. Math.*, vol. 59, no. 6, pp. 797–829, June 2006.
- [4] M. Wakin, "An introduction to compressive sampling," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 21–30, Mar. 2008.
- [5] S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM J. Scientific Computing*, vol. 20, no. 1, pp. 33–61, 1998.
- [6] E. J. Cands, J. K. Romberg, and T. Tao, "Stable signal recovery from incomplete and inaccurate measurements," *Communic. on Pure and Applied Math.*, vol. 59, no. 8, pp. 1207–1223, 2006.
- [7] M. Tipping, "Sparse Bayesian learning and the relevance vector machine," *J. Machine Learning Res.*, vol. 1, pp. 211–244, June 2001.
- [8] M. Figueiredo, "Adaptive sparseness for supervised learning," *IEEE Trans. on Pattern Analysis and Machine Intel.*, vol. 25, no. 9, pp. 1150–1159, 2003.
- [9] N. L. Pedersen, D. Shutin, C. N. Manchon, and B. H. Fleury, "Sparse estimation using Bayesian hierarchical prior modeling for real and complex models," *submitted to IEEE Trans. on Sig. Proc.*, 2011. [Online]. Available: arXiv:1108.4324v1
- [10] M. E. Tipping and A. C. Faul, "Fast marginal likelihood maximisation for sparse Bayesian models," in *Proc. 9th International Workshop on Artificial Intelligence and Statistics*, Key West, FL, January 2003.
- [11] D. Wipf and S. Nagarajan, "A new view of automatic relevance determination," in *Proc. 21 Annual Conf. on Neural Inform. Process. Systems*. Vancouver, British Columbia, Canada: MIT Press, Dec. 2007.
- [12] D. Shutin, T. Buchgraber, S. R. Kulkarni, and H. V. Poor, "Fast variational sparse Bayesian learning with automatic relevance determination for superimposed signals," *IEEE Transactions on Signal Processing*, vol. 59, no. 12, pp. 6257–6261, Dec. 2011.
- [13] J. O. Berger, J. M. Bernardo, and D. Sun, "The formal definition of reference priors," *Ann. Statist.*, vol. 37, no. 2, pp. 905–938, 2009.
- [14] E. J. Candes, M. B. Wakin, and S. Boyd, "Enhancing sparsity by reweighted  $\ell_1$  minimization," *J. of Fourier Analysis and Applic.*, vol. 14, no. 5, pp. 877–905, December 2008.
- [15] E. Cands and J. Romberg, "Sparsity and incoherence in compressive sampling," *Inverse Prob.*, vol. 23, no. 3, pp. 969–985, 2007.
- [16] M. J. Wainwright, "Sharp thresholds for high-dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming (Lasso)," *IEEE Trans. on Inform. Theory*, vol. 55, no. 5, pp. 2183–2202, May 2009.
- [17] D. Shutin, S. R. Kulkarni, and H. V. Poor, "Incremental reformulated automatic relevance determination," *IEEE Trans. on Sig. Proc.*, 2012, to appear.
- [18] —, "A novel view of incremental sparse Bayesian learning with automatic relevance determination," submitted to *IEEE Transactions on Signal Processing*.