

The Variational Garrote

Hilbert J. Kappen · Vicenç Gómez

Received: 7 January 2012 / Accepted: 18 November 2013 / Published online: 6 December 2013
© The Author(s) 2013. This article is published with open access at Springerlink.com

Abstract We analyze the variational method for sparse regression using ℓ_0 regularization. The variational approximation results in a model that is similar to Breiman’s Garrote model. We refer to this method as the Variational Garrote (VG). The VG has the effect of making the problem effectively of maximal rank even when the number of samples is small compared to the number of variables. We propose a naive mean field approximation combined with a maximum a posteriori (MAP) approach to estimate the model parameters and use an annealing and reheating schedule of the sparsity hyper-parameter to avoid local minima. The hyper-parameter is set by cross-validation. We compare the VG with the lasso, ridge regression and the recently introduced Bayesian paired mean field method (PMF) (Titsias and Lázaro-Gredilla in *Advances in neural information processing systems*, vol. 24, pp. 2339–2347, 2011). For fair comparison, we implemented a similar annealing-reheating schedule for the PMF sparsity parameter. Numerical results show that the VG and PMF yield more accurate predictions and more accurately reconstruct the true model than the other methods. The VG finds correct solutions when the lasso solution is inconsistent due to large input correlations. In the experiments that we consider we find that the VG, although based on a simpler approximation than the PMF, yields qualitatively similar or better results and is computationally more efficient. The naive implementation of the VG scales cubic with the number of features. By introducing Lagrange multipliers we obtain a dual formulation of the problem that scales cubic in the number of samples, but close to linear in the number of features.

Keywords Sparse regression · Variational approximation · Spike-and-slab · Mean-field · Nonnegative garrote

Editor: Yee-Whye Teh.

H.J. Kappen · V. Gómez (✉)
Donders Institute for Brain Cognition and Behaviour, Radboud University Nijmegen, 6525 EZ,
Nijmegen, The Netherlands
e-mail: v.gomez@science.ru.nl

H.J. Kappen
e-mail: b.kappen@science.ru.nl

1 Introduction

One of the most common problems in statistics is linear regression. Given p samples of n -dimensional input data x_i^μ , $i = 1, \dots, n$ and 1-dimensional output data y^μ , with $\mu = 1, \dots, p$, find weights w_i, w_0 that best describe the relation

$$y^\mu = \sum_{i=1}^n w_i x_i^\mu + w_0 + \xi^\mu \quad (1)$$

for all μ . ξ^μ is zero-mean noise with inverse variance β .

The ordinary least square (OLS) solution is given by $\mathbf{w} = \chi^{-1} \mathbf{b}$ and $w_0 = \bar{y} - \sum_i w_i \bar{x}_i$, where χ is the input covariance matrix \mathbf{b} is the vector of input-output covariances and \bar{x}_i, \bar{y} are the mean values. There are several problems with the OLS approach. When p is small, it typically has a low prediction accuracy due to over fitting. In particular, when $p < n$, χ is not of maximal rank and so its inverse is not uniquely defined. In addition, the OLS solution is not sparse: it will find a solution $w_i \neq 0$ for all i . Therefore, the interpretation of the OLS solution is often difficult.

These problems are well-known, and there exist a number of approaches to overcome these problems. The simplest approach is called ridge regression. It adds a regularization term $\frac{1}{2} \lambda \sum_i w_i^2$ with $\lambda > 0$ to the OLS criterion. This has the effect that the input covariance matrix χ gets replaced by $\chi + \lambda I$ which is of maximal rank for all p . One optimizes λ by cross validation. Ridge regression improves the prediction accuracy but not the interpretability of the solution.

Another approach is lasso (Tibshirani 1996). It solves the OLS problem under the linear constraint $\sum_i |w_i| \leq t$. This problem is equivalent to adding an ℓ_1 regularization's term $\lambda \sum_i |w_i|$ to the OLS criterion. The optimizations of the quadratic error under linear constraints can be solved efficiently. See Friedman et al. (2010) for a recent account. Again, λ or t may be found through cross validation. The advantage of the ℓ_1 regularization is that the solution tends to be sparse. This improves both the prediction accuracy and the interpretability of the solution.

The ℓ_1 or ℓ_2 regularization terms are known as shrinkage priors because their effect is to shrink the size of w_i . The idea of shrinkage prior has been generalized by Frank and Friedman (1993) to the form $\lambda \sum_i |w_i|^q$ with $q > 0$ and $q = 1, 2$ corresponding to the lasso and ridge case, respectively. Better solutions can be obtained for $q < 1$, however the resulting optimization problem is no longer convex and therefore more difficult to solve.

An alternative Bayesian approach to obtain a sparse solution using an ℓ_0 penalty was proposed by George and McCulloch (1993), Mitchell and Beauchamp (1988) under the “spike and slab” formulation. There are n variational selector variables s_i such that the prior distribution over w_i is a mixture of a narrow (spike) and wide (slab) Gaussian distribution, both centered on zero. The posterior distribution over s_i indicates whether the input feature i is included in the model or not. Since the number of subsets of features is exponential in n , for large n one cannot compute the solution exactly. In addition, the posterior is a complex high dimensional distribution of the w_i and the other (hyper) parameters of the model. The computation of the posterior requires thus the use of Markov chain Monte Carlo (MCMC) sampling (George and McCulloch 1993; Brown et al. 1998; Clyde and George 2004; Ishwaran and Rao 2005) or variational Bayesian approximations (Carbonetto and Stephens 2012; Titsias and Lázaro-Gredilla 2011; Logsdon et al. 2010).

Although Bayesian approaches tend to over fit less than a maximum likelihood or maximum a posteriori (MAP) estimators, they also tend to be relatively slow. Here we propose a partial Bayesian approach, where we apply a variational approximation to integrate out the binary (selector) variables in combination with a MAP approach for the remaining parameters. For clarity, we analyze this idea in its most simple form, in the absence of (hierarchical) priors. Instead, we infer the sparsity prior through cross validation. As we will motivate below, we call the method the Variational Garrote (VG).

The paper is organized as follows. In Sect. 2 we introduce the model and derive the variational approximation. Related work is described in Sect. 3. In Sect. 4 we study the case when the design matrix is orthogonal. In this case the solution can be computed exactly in closed form with no need to resort to approximations. In Sect. 5 we compare numerically the VG with a number of other MAP methods, such as lasso and ridge regression and with the paired mean field method (PMF) (Titsias and Lázaro-Gredilla 2011), a recently proposed variational Bayesian method. We conclude with discussion in Sect. 6.

2 The variational approximation

Consider the regression model of the form¹

$$y^\mu = \sum_{i=1}^n w_i s_i x_i^\mu + \xi^\mu \quad \sum_{i=1}^n s_i \leq t \tag{2}$$

with $s_i = 0, 1$. The bits $s_i = 1$ will identify the predictive inputs i . Using a Bayesian description, and denoting the data by $D : \{\mathbf{x}^\mu, y^\mu\}, \mu = 1, \dots, p$, the likelihood term is given by

$$\begin{aligned} p(y|\mathbf{x}, \mathbf{s}, \mathbf{w}, \beta) &= \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2} \left(y - \sum_{i=1}^n w_i s_i x_i\right)^2\right) \\ p(D|\mathbf{s}, \mathbf{w}, \beta) &= \prod_{\mu} p(y^\mu|\mathbf{x}^\mu, \mathbf{s}, \mathbf{w}, \beta) \\ &= \left(\frac{\beta}{2\pi}\right)^{p/2} \exp\left(-\frac{\beta p}{2} \left(\sum_{i,j=1}^n s_i s_j w_i w_j \chi_{ij} - 2 \sum_{i=1}^n w_i s_i b_i + \sigma_y^2\right)\right) \end{aligned} \tag{3}$$

with $b_i = \frac{1}{p} \sum_{\mu} x_i^\mu y^\mu, \sigma_y^2 = \frac{1}{p} \sum_{\mu} (y^\mu)^2, \chi_{ij} = \frac{1}{p} \sum_{\mu} x_i^\mu x_j^\mu$.

We should also specify prior distributions over $\mathbf{s}, \mathbf{w}, \beta$. For concreteness, we assume that the prior over \mathbf{s} is factorized over the individual s_i , each with identical prior probability:

$$p(\mathbf{s}|\gamma) = \prod_{i=1}^n p(s_i|\gamma) \quad p(s_i|\gamma) = \frac{\exp(\gamma s_i)}{1 + \exp(\gamma)} \tag{4}$$

with γ given which specifies the sparsity of the solution. We denote by $p(\mathbf{w}, \beta)$ the prior over the inverse noise variance β and the feature weights \mathbf{w} . We will leave this prior unspecified since its choice does not affect the variational approximation.

¹We assume from here on without loss of generality that $\frac{1}{p} \sum_{\mu=1}^p x_i^\mu = \frac{1}{p} \sum_{\mu=1}^p y^\mu = 0$.

The posterior becomes

$$p(\mathbf{s}, \mathbf{w}, \beta|D, \gamma) = \frac{p(\mathbf{w}, \beta)p(\mathbf{s}|\gamma)p(D|\mathbf{s}, \mathbf{w}, \beta)}{p(D|\gamma)} \tag{5}$$

Computing the MAP estimate or computing statistics from the posterior is complex in particular due to the discrete nature of \mathbf{s} . We propose to compute a variational approximation to the marginal posterior $p(\mathbf{w}, \beta|D, \gamma) = \sum_{\mathbf{s}} p(\mathbf{s}, \mathbf{w}, \beta|D, \gamma)$ and computing the MAP solution with respect to \mathbf{w}, β . Since $p(D|\gamma)$ does not depend on \mathbf{w}, β we can ignore it.

The posterior distribution equation (5) for given \mathbf{w}, β is a typical Boltzmann distribution involving terms linear and quadratic in s_i . It is well-known that when the effective couplings $w_i w_j \chi_{ij}$ are small, one can obtain good approximations using methods that originated in the statistical physics community and where s_i denote binary spins. Most prominently, one can use the mean field or variational approximation (Jordan et al. 1999), the TAP approximation (Kappen and Spanjers 2000) or belief propagation (BP) (Murphy et al. 1999). For introductions into these methods also see Oppen and Saad (2001), Wainwright and Jordan (2008). Here, we will develop a solution based on the simplest possible variational approximation and leave the possible improvements using BP or structured mean field approximations to the future.

We approximate the sum by the variational bound (by Jensen’s inequality)

$$\begin{aligned} \log \sum_{\mathbf{s}} p(\mathbf{s}|\gamma)p(D|\mathbf{s}, \mathbf{w}, \beta) &\geq - \sum_{\mathbf{s}} q(\mathbf{s}) \log \frac{q(\mathbf{s})}{p(\mathbf{s}|\gamma)p(D|\mathbf{s}, \mathbf{w}, \beta)} \\ &= -F(q, \mathbf{w}, \beta). \end{aligned} \tag{6}$$

The probability distribution $q(\mathbf{s})$ is called the variational approximation and can be any positive probability distribution on \mathbf{s} and $F(q, \mathbf{w}, \beta)$ is called the variational free energy. The optimal $q(\mathbf{s})$ is found by minimizing $F(q, \mathbf{w}, \beta)$ with respect to $q(\mathbf{s})$ so that the tightest bound (best approximation) is obtained.

In order to be able to compute the variational free energy efficiently, $q(\mathbf{s})$ must be a tractable probability distribution, such as a chain or a tree with limited tree-width (Barber and Wiergerinck 1999). Here we consider the simplest case where $q(\mathbf{s})$ is a fully factorized distribution: $q(\mathbf{s}) = \prod_{i=1}^n q_i(s_i)$ with $q_i(s_i) = m_i s_i + (1 - m_i)(1 - s_i)$, so that q is fully specified by the expected values $m_i = q_i(s_i = 1)$, which we collectively denote by \mathbf{m} .

The expectation values with respect to q can now be easily evaluated and the result is

$$\begin{aligned} F &= \frac{\beta p}{2} \left(\sum_{i,j} m_i m_j w_i w_j \chi_{ij} + \sum_i m_i (1 - m_i) w_i^2 \chi_{ii} - 2 \sum_{i=1}^n m_i w_i b_i + \sigma_y^2 \right) \\ &\quad - \gamma \sum_{i=1}^n m_i + \sum_{i=1}^n (m_i \log m_i + (1 - m_i) \log(1 - m_i)) - \frac{p}{2} \log \frac{\beta}{2\pi}, \end{aligned} \tag{7}$$

where we have omitted terms independent of $\mathbf{m}, \beta, \mathbf{w}$. The first line is due to the likelihood term, the second line is due to the prior on \mathbf{s} and the entropy of $q(\mathbf{s})$. The approximate marginal posterior is then

$$\begin{aligned} p(\mathbf{w}, \beta|D, \gamma) &\propto p(\mathbf{w}, \beta) \sum_{\mathbf{s}} p(\mathbf{s}|\gamma)p(D|\mathbf{s}, \mathbf{w}, \beta) \\ &\approx p(\mathbf{w}, \beta) \exp(-F(\mathbf{m}, \mathbf{w}, \beta, \gamma)). \end{aligned}$$

We can compute the variational approximation \mathbf{m} for given $\mathbf{w}, \beta, \gamma$ by minimizing F with respect to \mathbf{m} . In addition, $p(\mathbf{w}, \beta|D, \gamma)$ needs to be maximized with respect to \mathbf{w}, β . Note, that the variational approximation only depends on the likelihood term and the prior on γ , since these are the only terms that depend on \mathbf{s} . Thus, for given \mathbf{w} , the variational approximation does not depend on the particular choices for the prior $p(\mathbf{w}, \beta)$. For concreteness, we assume a flat prior $p(\mathbf{w}, \beta) \propto 1$. We set the derivatives of F with respect $\mathbf{m}, \mathbf{w}, \beta$ equal to zero. This gives the following set of fixed point equations:

$$m_i = \sigma \left(\gamma + \frac{\beta p}{2} w_i^2 \chi_{ii} \right) \tag{8}$$

$$\mathbf{w} = (\chi')^{-1} \mathbf{b} \quad \chi'_{ij} = \chi_{ij} m_j + (1 - m_j) \chi_{jj} \delta_{ij} \tag{9}$$

$$\frac{1}{\beta} = \sigma_y^2 - \sum_{i=1}^n m_i w_i b_i \tag{10}$$

with $\sigma(x) = (1 + \exp(-x))^{-1}$ and where in Eq. (10) we have used Eq. (9). Equations (8)–(10) provide the final solution. They can be solved by fixed point iteration as outlined in Algorithm 1.

```

input   : Data  $D : \{\mathbf{x}^\mu, y^\mu\}, \mu = 1, \dots, p; \epsilon$  and step-size  $\Delta\gamma$ 
output :  $\mathbf{w}, \mathbf{m}, \beta, \gamma$  solution with minimal cross validation error
1 Preprocess data such that  $\sum_\mu x_i^\mu = \sum_\mu y^\mu = 0$  and partition  $D$  in  $D^{\text{train}}, D^{\text{val}}$ 
2 Compute  $b_i = \frac{1}{p} \sum_\mu x_i^\mu y^\mu$  and if  $n < p$  compute  $\chi_{ij} = \frac{1}{p} \sum_\mu x_i^\mu x_j^\mu$ 
3 Compute  $\gamma_{\min}$  from  $\epsilon$  and  $\gamma_{\max}$  from  $\gamma_{\min}$  and  $\Delta\gamma$ 
4 Initialize  $\mathbf{m} \leftarrow \mathbf{0}$ 
5 for  $\gamma = \gamma_{\min} : \Delta\gamma : \gamma_{\max}$  do // FORWARD PASS
6    $\eta \leftarrow 1$ 
7   while not converged do
8     Compute  $\mathbf{w}, \beta$  from Eqs. (9)–(10) ( $n < p$ ) or Eqs. (22), (25)–(28) ( $n > p$ )
9     Compute  $\mathbf{m}'$  using a smoothed version of Eq. (8):  $m'_i \leftarrow (1 - \eta)m_i + \eta\sigma(\dots)$ 
10    if  $\max_j |m'_j - m_j| > 0.1$  then
11       $\eta \leftarrow \eta/2$ 
12     $\mathbf{m} \leftarrow \mathbf{m}'$ 
13  Store solution  $(\mathbf{w}_1, \mathbf{m}_1, \beta_1)_\gamma$  and  $F_1(\gamma) \leftarrow F((\mathbf{w}_1, \mathbf{m}_1, \beta_1)_\gamma)$  from Eq. (7)
14 for  $\gamma = \gamma_{\max} : -\Delta\gamma : \gamma_{\min}$  do // BACKWARD PASS
15   As 5 – 11
16  Store solution  $(\mathbf{w}_2, \mathbf{m}_2, \beta_2)_\gamma$  and  $F_2(\gamma) \leftarrow F((\mathbf{w}_2, \mathbf{m}_2, \beta_2)_\gamma)$  from Eq. (7)
17 for  $\gamma = \gamma_{\min} : \Delta\gamma : \gamma_{\max}$  do
18   Choose solution  $(\mathbf{w}, \mathbf{m}, \beta)_\gamma$  that has minimal  $F_{1,2}(\gamma)$ 
19   Compute cross validation error on  $D^{\text{val}}$  using Eq. (11)
20 Select  $\mathbf{w}, \mathbf{m}, \beta, \gamma$  with minimal cross validation error
    
```

Algorithm 1: The Variational Garrote algorithm

Within the variational/MAP approximation the predictive model is

$$y = \sum_i m_i w_i x_i + \xi \tag{11}$$

with $\langle \xi^2 \rangle = 1/\beta$ and $\mathbf{m}, \mathbf{w}, \beta$ as estimated by the above procedure.

Equation (11) has some similarity with Breiman’s non-negative Garrote method (Breiman 1995). It computes the solution in a two step approach: it computes first w_i using OLS and then finds m_i by minimizing

$$\sum_{\mu} \left(y^{\mu} - \sum_{i=1}^n x_i^{\mu} w_i m_i \right)^2 \quad \text{subject to} \quad m_i \geq 0 \quad \sum_i m_i \leq t.$$

Because of this similarity, we refer to our method as the Variational Garrote (VG). Note, that because of the OLS step the non-negative garrote requires that $p \geq n$. Instead, the variational solution of Eqs. (8)–(10) computes the entire solution in one step (and as we will see does not require $p \geq n$).

The model in Eqs. (2), (4) is also equivalent to a “spike and slab” prior on the weights parametrized as a product of a Gaussian random variable w_i and a Bernoulli random variable s_i

$$p(w_i, s_i) = \mathcal{N}(w_i|0, \sigma_w^2) \pi^{s_i} (1 - \pi)^{1-s_i} \quad \forall i, \tag{12}$$

under the identification that the VG assumes a constant (improper) prior on w_i ($\sigma_w^2 = \infty$) and the relation between the sparsity γ and π is given by $\gamma = \log(\pi/(1 - \pi))$.

Let us pause to make some observations about the VG solution. One might naively expect that the variational approximation would simply consist of replacing $w_i s_i$ in Eq. (2) by its variational expectation $w_i m_i$. If this were the case, \mathbf{m} would disappear entirely from the equations and one would expect in Eq. (9) the OLS solution with the normal input covariance matrix χ instead of the new matrix χ' (note, that in the special case that $m_i = 1$ for all i , $\chi' = \chi$ and Eq. (9) does reduce to the OLS solution). Instead, \mathbf{m} and \mathbf{w} are both to be optimized, giving in general a different solution than the OLS solution.²

When $m_i < 1$, χ' differs from χ by rescaling with m_i and adding a positive diagonal to it, a ‘variational ridge’. This is similar to the mechanism of ridge regression, but with the important difference that the diagonal term depends on i and is dynamically adjusted depending on the solution for \mathbf{m} . Thus, the sparsity prior together with variational approximation provides a mechanism that solves the rank problem. When all $m_i < 1$, χ' is of maximal rank. Each m_i that approaches 1, reduces the rank by one. Thus, if χ has rank $p < n$, χ' can be still of rank n when no more than p of the $m_i = 1$, the remaining $n - p$ of the $m_i < 1$ making up for the rank deficiency. Note, that the size of m_i (and thus the rank of χ') is controlled by γ through Eq. (8).

In the above procedure, we compute the VG solution for fixed γ and choose its optimal value through cross validation on independent data (Mitchell and Beauchamp 1988). This has the advantage that our result is independent of our (possibly incorrect) prior belief.

²The technical reason that this does not occur is that in the computation of the expectation with respect to the distribution q that results in Eq. (7) one has $\langle s_i s_j \rangle = m_i m_j$ for $i \neq j$, but $\langle s_i^2 \rangle = \langle s_i \rangle = m_i$.

Another important advantage of varying γ manually is that it helps to avoid local minima. When we increase γ from a negative value γ_{\min} to a maximal value γ_{\max} in small steps, we obtain a sequence of solutions with decreasing sparseness. These solutions will better fit the data and as a result β increases with γ . Thus, increasing γ implements an annealing mechanism where we sequentially obtain solutions at lower noise levels. We found empirically that this approach is effective to reduce the problem of local minima. To further deal with the effect of hysteresis (see Sect. 4) we recompute the solution from γ_{\max} down to γ_{\min} and choose the solution with lowest free energy.

The minimal value of γ is chosen as the largest value such that $m_i = \epsilon$, with ϵ small. We find from Eqs. (8)–(10) that

$$\gamma_{\min} = -\frac{pb_i^2 \chi_{ii}}{2\sigma_y^2} + \sigma^{-1}(\epsilon) + \mathcal{O}(\epsilon) \quad (13)$$

with $\sigma^{-1}(x) = \log(x/(1-x))$. We heuristically set the maximal value of γ as well as the step size.

In Appendix B we provide an alternative fixed point iteration scheme that is more efficient in the large n small p limit. Whereas Eqs. (8)–(10) require the repeated solution of a n -dimensional linear system, the dual formulation, Eqs. (8), (22), (25)–(28), requires the repeated solution of a p dimensional linear system. Algorithm 1 summarizes the VG method.

3 Related work

The “spike and slab” model is one of the most widely approaches to sparse Bayesian variable selection. Inference in this model has been performed usually by MCMC sampling. These methods address the combinatorial problem of searching all possible 2^n combinations of predictors by sampling from the posterior distribution. There is an extensive literature on MCMC methods for this model, e.g. George and McCulloch (1993), Brown et al. (1998), Clyde and George (2004), Ishwaran and Rao (2005), O’Hara and Sillanpää (2009). However, their applicability is limited on large-scale problems, since designing a Markov chain that explores the parameter space efficiently is a difficult task. In this paper, we focus on the alternative Bayesian variational approach.

A mean field variational approximation for the spike and slab prior was proposed initially in Logsdon et al. (2010) in the context of genetic association studies. Their model differs from the VG in the sense that they use separate and different priors for positive and negative effects. They also use truncated normal distributions for the feature weights and place hyperpriors on γ .

More recently, an alternative variational approximation called paired mean field (PMF) has been proposed in Titsias and Lázaro-Gredilla (2011). It is defined on a model for multiple outputs and considers a linear combination of an input layer of basis functions governed by a Gaussian process, thereby unifying several sparse linear models such as sparse factor analysis or sparse matrix factorization. To relate the PMF model to the VG, we consider the uni-variate response case without the extra input layer. Instead of assuming a fully factorized variational approximation, PMF places each weight w_i and bit s_i in the same factor, i.e. $q(\mathbf{w}, \mathbf{s}) = \prod_{i=1}^n q_i(w_i, s_i)$.

An important difference between the VG and the two previous methods is the algorithm used for parameter optimization. The VG method computes expectation of \mathbf{s} (called \mathbf{m}) but finds MAP solution for \mathbf{w} and β . Hyper-parameter γ is optimized using an annealing-reheating schedule and a validation dataset. In contrast, Logsdon et al. (2010) and Titsias and Lázaro-Gredilla (2011) rely exclusively on the expectation-maximization algorithm with random restarts. As we will show later, this can have important consequences in terms of sub-optimality in cases where inputs are highly correlated.

Around the time of publication of this paper, we became aware of the work of Carbonetto and Stephens (2012). Their approach also considers the fully factorized case but assumes a joint prior for the hyper-parameters and uses importance sampling to compute their posterior distribution. Similarly to the VG, their algorithm considers an inner-loop of coordinate ascent updates for m_i and w_i .³ The difference is that it considers β as a hyper-parameter, together with σ_w^2 and π , and they are jointly integrated using importance sampling. The sampling step is in practice performed using a three-dimensional grid with resolution selected heuristically. For each setting of the hyper-parameters, they compute the largest marginal likelihood solution ($\mathbf{m}^{(\text{init})}$, $\mathbf{w}^{(\text{init})}$) using random initializations and, instead of annealing, the coordinate ascent updates are run separately again for each setting of the hyper-parameters with ($\mathbf{m}^{(\text{init})}$, $\mathbf{w}^{(\text{init})}$) as initialization.

The fully factorized approximation considered here is also closely related to the one proposed for independent factor analysis (Attias 1999). Combined with a more complex form of annealing for MAP search has been proposed in Yoshida and West (2010) in the context of sparse latent factor analysis. They have shown that this type of optimization strategy can be useful to address the local minima problem and lead to robust estimation.

An alternative to the aforementioned variational approaches is the work of Hernández-Lobato et al. (2010), in which the expectation propagation (EP) algorithm is used in a multi-task setting where the latent variables indicate whether the corresponding features are used for classification in all the tasks or in none of them. EP also considers a factorized approximate distribution.

The problem of inconsistency of the lasso's penalty has been addressed by many authors and lead to several generalizations (see Tibshirani 2011 and references therein). Two popular approaches that, similarly to the VG, consider non convex penalties, are the Smoothed Clipped Absolute Deviation penalty (known as SCAD) (Fan and Li 2001) and the SparseNet (Mazumder et al. 2011). The SCAD (Fan and Li 2001) replaces the lasso penalty with a continuous differentiable function that reduces the amount of shrinkage for larger values of w_i , with eventually no shrinkage for $w_i \rightarrow \infty$. The SparseNet (Mazumder et al. 2011) performs a coordinate-wise optimization of λ and q , covering the bridge of possible solution surfaces between lasso $q = 1$ and variable selection $q = 0$.

From a Bayesian point of view, the lasso estimator can be viewed as solving a MAP estimation problem when the feature weights have independent double exponential (Laplace) priors. A complete Bayesian analysis for the lasso prior is developed in Park and Casella (2008). Fully Bayesian approaches compute posterior mean and median estimates using MCMC sampling and may lead to solutions that are not necessarily sparse. Recently, a Bayesian model that extends the double exponential prior with a normal-exponential-gamma distribution (NEG) and uses MAP estimation has been proposed in Griffin and Brown (2011). The NEG prior has a finite spike at zero and heavy tails, thus preventing over-shrinkage of weights with large absolute values. The authors propose an EM method

³The notation in Carbonetto and Stephens (2012) uses α_k for m_i and μ_k for w_i .

that alternates between estimation of the prior variances of the weights (E-Step) and the weight values conditioned on the variances (M-Step). Other hyper-parameters are chosen using cross validation.

4 Orthogonal and uni-variate case

In this section we show for the uni-variate case that the solution is either unique or has two solutions, depending on the input-output correlations, the number of samples p and on the sparsity prior γ . We derive a phase plot and show that the solution is unique, when the sparsity prior is not too strong *or* when the input-output correlation is not too large. The input-output behavior of the VG is shown to be close to optimal as a smoothed version of hard feature selection. We argue that this behavior also holds in the multi-variate case.

Consider the case in which the inputs are uncorrelated: $\chi_{ij} = \delta_{ij}$. In this case, we can derive the MAP solution of Eq. (5) exactly, without the need to resort to the variational approximation. Equation (5) reduces to a distribution that factorizes over i with log probability proportional to

$$L = \frac{p}{2} \log \beta - \frac{\beta p}{2} \left(\sum_{i=1}^n s_i (w_i^2 - 2w_i b_i) + \sigma_y^2 \right) + \gamma \sum_{i=1}^n s_i$$

Maximizing wrt w_i, β yields $w_i = b_i, \beta^{-1} = \sigma_y^2 - \sum_{i=1}^n s_i b_i^2$ and

$$L = \frac{p}{2} \log \beta + \sum_{i=1}^n s_i \left(\frac{\beta p}{2} b_i^2 + \gamma \right) - \frac{\beta p}{2} \sigma_y^2$$

Assume without loss of generality that b_i^2 are sorted in decreasing order. L is maximized by setting $s_i = 1$ when $\frac{\beta p}{2} b_i^2 + \gamma > 0$ and $s_i = 0$ otherwise. Thus, the optimal solution is $s_{1:k} = 1, s_{k+1:n} = 0, \beta^{-1} = \sigma_y^2 - \sum_{i=1}^k b_i^2$ with k the smallest integer such that

$$\frac{\beta p}{2} b_{k+1}^2 + \gamma < 0 \tag{14}$$

By varying γ from small to large, we find a sequence of solutions with decreasing sparsity.

In the variational approximation the solution is very similar but not identical. Equation (9) gives the same solution $w_i = b_i$. Equations (8) and (10) become

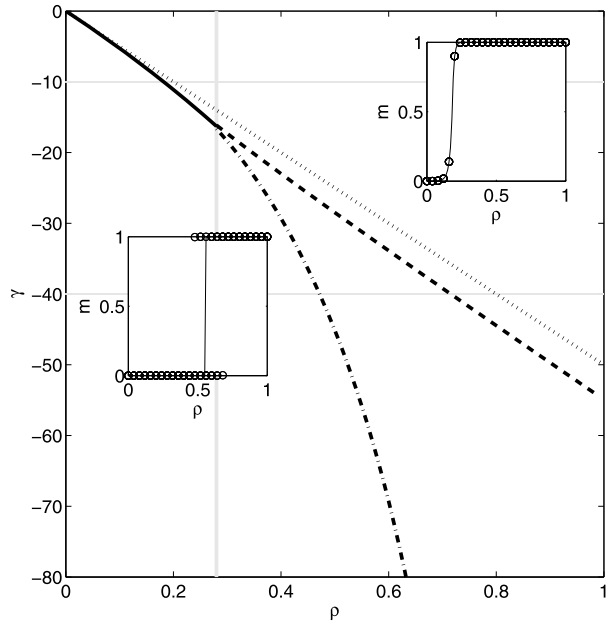
$$m_i = \sigma \left(\gamma + \frac{\beta p}{2} b_i^2 \right)$$

$$\frac{1}{\beta} = \sigma_y^2 - \sum_i b_i^2 m_i$$

which we can interpret as the variational approximations of Eq. (14), with $m_{1:k} \approx 1$ and $m_{k+1:n} \approx 0$. The term $\sum_i b_i^2 m_i$ is the explained variance and is subtracted from the total output variance to give an estimate of the noise variance $1/\beta$.

Note that the posterior is factorized in s_i , the variational approximation is not identical to the exact map solution equation (14), although the results are very similar. The relation is $s_i = 0 \Leftrightarrow m_i < 0.5$ and $s_i = 1 \Leftrightarrow m_i > 0.5$.

Fig. 1 Phase plot ρ, γ for $p = 100$ giving the different solutions for m . *Dashed and dot-dashed lines* for $\rho > \rho^* = 0.28$ are from Eq. (19) where two solutions for m exist. *Solid line* for $\rho < \rho^*$ is the solution for γ when $m = 1/2$, to indicate the transition from the unique solution $m \approx 0$ to the unique solution $m \approx 1$. *Dotted line* is the exact transition from $s = 0$ to $s = 1$ from Eq. (14). *Insets* indicate solutions for m versus ρ for $\gamma = -10, p = 100$ (*top-right*) and for $\gamma = -40, p = 100$ (*bottom-left*). In the *lower left corner* of the insets, the unique solution $m \approx 0$ is found. In the *top right corner*, the unique solution $m \approx 1$ is found. Between the *dot-dashed* and the *dashed line*, the two variational solutions $m \approx 0$ and $m \approx 1$ co-exist



In order to further analyze the variational solution, we consider the 1-dimensional case. The variational equations become

$$m = \sigma \left(\gamma + \frac{p}{2} \frac{\rho}{1 - \rho m} \right) = f(m) \tag{15}$$

$$\frac{1}{\beta} = \sigma_y^2 (1 - m\rho) \tag{16}$$

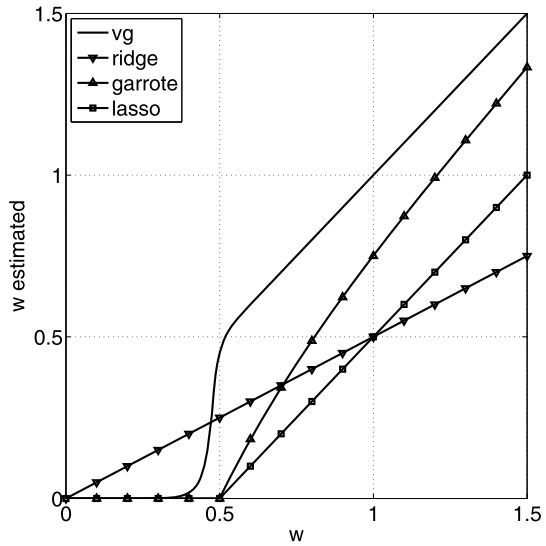
with $\rho = b^2/\sigma_y^2$ the squared correlation coefficient.

In Eq. (15), we have eliminated β and we must find a solution for m for this non-linear equation. We see that it depends on the input-output correlation ρ , the number of samples p and the sparsity γ . For $p = 100$, the solution for different ρ, γ is illustrated in Fig. 12 (see Appendix A). Equation (15) has one or three solutions for m , depending on the values of γ, ρ, p . The three solutions correspond to two local minima and one local maximum of the free energy F . For $\gamma = -40$ and $\gamma = -10$, we plot the stable solution(s) for different values of ρ in the insets in Fig. 1. The best variational solution for m is given by the solution with the lowest free energy, indicated by the solid lines in the insets in Fig. 1.

Figure 1 further shows the phase plot of γ, ρ that indicates that the variational solution is unique for $\gamma > \gamma^*$ or for $\rho < \rho^*$. The solid line for $0 < \rho < \rho^*$ in Fig. 1 indicates a smooth (second order) phase transition from $m = 0$ to $m = 1$. For $\rho > \rho^*$, the transition from $m = 0$ to $m = 1$ is discontinuous: for each ρ there is a range of values of γ where two variational solutions $m \approx 0$ and $m \approx 1$ co-exist. For comparison, we also show the line $\gamma = -p\rho/2$ that separates the solution $s = 0$ and $s = 1$ according the exact (non-variational) solution equation (14).

The multi-valued variational solution results in a hysteresis effect. When the solution is computed for increasing γ , the $m \approx 0$ solution is obtained until it no longer exists. If the

Fig. 2 Uni-variate solution for different regression methods. All methods yield a shrunk solution (deviation from *diagonal line*). Variational Garrote (VG) with $\gamma = -10$, $p = 100$ and $\sigma_y^2 = 1$. Ridge regression with $\lambda = 0.5$. Garrote with $\gamma = 1/4$. Lasso with $\gamma = 1/2$



sequence of solutions is computed for decreasing γ the $m \approx 1$ solution is obtained for values of γ where previously the $m \approx 0$ solution was obtained.

From this simple one-dimensional case we may infer that the variational approximation is relatively easy to compute in the uni-modal region (small ρ or γ not too negative) and becomes more inaccurate in the region where multiple optima exist (region between the dot-dashed and dashed lines in Fig. 1).

It is interesting to compare the uni-variate solution of the Variational Garrote with ridge regression, lasso or Breiman’s Garrote, which was previously done for the latter three methods in Tibshirani (1996). Suppose that data are generated from the model $y = wx + \xi$ with $\langle \xi^2 \rangle = \langle x^2 \rangle = 1$. We compare the solutions as a function of w . The OLS solution is approximately given by $w_{ols} \approx \langle xy \rangle = w$, where we ignore the statistical deviations of order $1/p$ due to the finite data set size. Similarly, the ridge regression solution is given by $w_{ridge} \approx \lambda w$, with $0 < \lambda < 1$ depending on the ridge prior. The lasso solution (for non-negative w) is given by $w_{lasso} = (w - \gamma)^+$ (Tibshirani 1996), with γ depending on the ℓ_1 constraint. Breiman’s Garrote solution is given by $w_{garrote} = (1 - \frac{\gamma}{w^2})^+ w$ (Tibshirani 1996), with γ depending on the ℓ_1 constraint. The VG solution is given by $w_{vg} = mw$, with m the solution of Eq. (15). Note, that the VG solution depends, in addition to w , γ , on the unexplained variance σ_y^2 and the number of samples p , whereas the other methods do not.

The qualitative difference of the solutions is shown in Fig. 2. The ridge regression solution is off by a constant multiplicative factor. The lasso solution is zero for small w and for larger w gives a solution that is shifted downwards by a constant factor. Breiman’s Garrote is identical to the lasso for small w and shrinks less for larger w . The VG gives an almost ideal behavior and can be interpreted as a soft version of variable selection: For small w the solution is close to zero and the variable is ignored, and above a threshold it is identical to the OLS solution.

The qualitative nature of the phase plot Fig. 1 and the input-output behavior Fig. 2 extends to the multi-variate orthogonal case. The symmetry breaking of feature i is independent of all other features, except for the term $\delta = \sum_{j \neq i} b_j^2 m_j$ that enters through β . If we increase γ , δ increases in steps each time that one of the features j switches from $m_j \approx 0$ to $m_j \approx 1$. Thus δ is constant almost always, except at the step points. Since the critical values

of ρ and γ depend in a simple way on δ , the phase plot for the multivariate orthogonal case is qualitatively the same as for the uni-variate case.

5 Numerical examples

In the following examples, we compare the VG with lasso, ridge regression and in some cases, with the paired mean field approach (PMF) (Titsias and Lázaro-Gredilla 2011). We show that the VG and PMF significantly outperform the lasso and ridge regression on a large number of different examples both in terms of the accuracy of the solution, as well as in prediction error. In addition, we show that the VG does not suffer from the inconsistency of the lasso method when the input correlations are large. We finally show how all methods compare as a function of the level of noise, the sparsity of the target solution, the number of samples and the number of irrelevant predictors.

For most of the examples, we generate training, validation and testing sets. Inputs are generated from a zero mean multi-variate Gaussian distribution with specified covariance structure. We generate outputs $y^\mu = \sum_i \hat{w}_i x_i^\mu + d\xi^\mu$ with $d\xi^\mu \in \mathcal{N}(0, \hat{\sigma})$ and \hat{w}_i depending on the problem.

For VG, ridge regression and lasso, we optimize the model parameters on the training set and, when necessary, optimize the hyper-parameters (γ for VG, λ for ridge regression and lasso) that minimize the quadratic error on the validation set. For the lasso, we used the method described in Friedman et al. (2010).⁴

Comparison with PMF is performed using the software available online for the regression case with one-dimensional output.⁵ For PMF, we merge both training and validation sets and the resulting dataset is used as input for the PMF method. This ensures that all methods use the same data for parameter estimation.

We also consider a modified version of PMF which replaces the update of π in the M-Step with a sequential annealing-reheating procedure such as the one proposed for γ in the VG. We observed empirically that the best strategy is to perform a sweep from sparse to dense $\pi_{0 \rightarrow 1}$ solutions (forward pass) followed by a sweep from dense to sparse $\pi_{1 \rightarrow 0}$ solutions (backward pass) and select the solution with maximum bound value (or minimum negative bound as we report here) in the backward pass. PMF does not over-fit as a function of π and thus does not require the use of a validation set. We refer to such variant of PMF as PMF-ANNEAL.

We define the solution vector for a given method as \mathbf{v} . For VG, the components are $v_i \equiv m_i w_i$. In the case of PMF and PMF-ANNEAL, m_i corresponds to the spike-and-slab variational posterior and w_i to the variational mean for the weights.⁶ For ridge and lasso $v_i \equiv w_i$.

5.1 Small Example 1

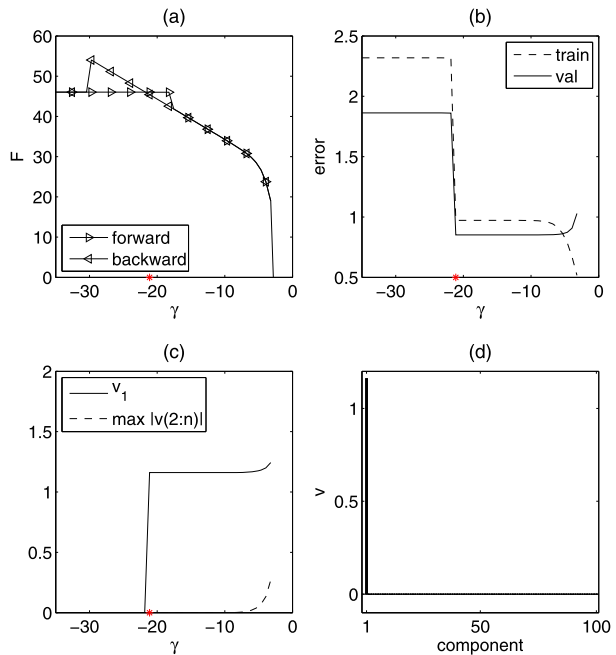
In the first example, we take independent inputs $x_i^\mu \in \mathcal{N}(0, 1)$ and a teacher weight vector with only one non-zero entry: $\hat{w} = (1, 0, \dots, 0)$, $n = 100$ and $\hat{\sigma} = 1$. The training set size $p = 50$, validation set size $p_v = 50$ and test set size $p_t = 400$. We choose $\epsilon = 0.001$ in Eq. (13), $\gamma_{\max} = 0.02\gamma_{\min}$, $\Delta\gamma = -0.02\gamma_{\min}$ (see Algorithm 1 for details).

⁴<http://www-stat.stanford.edu/~tibs/glmnet-matlab/>.

⁵<http://www.well.ox.ac.uk/~mtitsias/software.html>.

⁶The notation in Titsias and Lázaro-Gredilla (2011) uses \tilde{w}_i for w_i and γ_i for m_i .

Fig. 3 *Top left (a)*: minimal variational free energy versus γ . The *two curves* correspond to warm start solution from small to large γ ('forward') and from large to small γ ('backward') (see also Algorithm 1). *Top right (b)*: training and validation error versus γ . The optimal γ minimizes the validation error. *Bottom left (c)*: solution $v_1 = m_1 w_1$ and $\max_{i=2:n} |m_i w_i|$. The correct solution is found in the range $\gamma \approx -20$ to $\gamma \approx -5$. *Bottom right (d)*: optimal solution $v_i = m_i w_i$ versus i



Results for a single run of the VG are shown in Fig. 3. In Fig. 3a, we plot the minimal variational free energy F versus γ for both the forward and backward run. Note, the hysteresis effect due to the local minima. For each γ , we use the solution with the lowest F . In Fig. 3b, we plot the training error and validation error versus γ . The optimal $\gamma \approx -21$ is denoted by a star and the corresponding $\sigma = 1/\sqrt{\beta} = 1.05$. In Fig. 3c, we plot the non-zero component $v_1 = m_1 w_1$ and the maximum absolute value of the remaining components versus γ . Note the robustness of the VG solution in the sense of the large range of γ values for which the correct solution is found. In Fig. 3d, we plot the optimal solution $v_i = m_i w_i$ versus i .

In Fig. 4 we show the lasso (top row) and ridge regression (bottom row) results for the same data set. The optimal value for λ minimizes the validation error (star). In Fig. 4b, c we see that the lasso selects a number of incorrect features as well. Figure 4b also shows that the lasso solution with a larger λ in the range $0.45 < \lambda < 0.95$ could select the single correct feature, but would then estimate \hat{w}_1 too small due to the large shrinkage effect. Ridge regression gives very bad results. The non-zero feature is too small and the remaining features have large values. Note from Fig. 4e, that ridge regression yields a non-sparse solution for all values of λ .

Table 1 shows that the VG significantly outperforms the lasso method and ridge regression both in terms of prediction error, the accuracy of the estimation of the parameters and the number of non-zero parameters. In this simple example, there is no significant difference in the prediction error of lasso, PMF and VG, but the lasso solution is significantly less sparse. There is no significant difference between the solutions found by PMF and VG.

5.2 Small Example 2

In the second example, we consider the effect of correlations in the input distribution. Following Tibshirani (1996) we generate input data from a multi-variate Gaussian distribution

Table 1 Results for Example 1 averaged over 20 instances. Train is mean squared error (MSE) on the training set. Val is MSE on the validation set. Test is MSE on the test set. # non-zero is the number of non-zero elements in the lasso solution and $\sum_{i=1}^n (m_i > 0.5)$ for VG and PMF. $\|\delta v\|_1 = \sum_{i=1}^n |v_i - \hat{w}_i|$

	Train	Val	Test	# non-zero	$\ \delta v\ _1$
Ridge	0.60 ± 0.43	1.72 ± 0.39	1.80 ± 0.12	–	3.97 ± 1.23
Lasso	0.78 ± 0.26	1.07 ± 0.20	1.17 ± 0.20	8.65 ± 6.75	0.80 ± 0.57
PMF	–	–	1.02 ± 0.10	1.5 ± 1.19	0.33 ± 0.37
VG	0.85 ± 0.22	0.96 ± 0.17	1.01 ± 0.10	1.20 ± 0.52	0.31 ± 0.30
True	0.93 ± 0.14	0.87 ± 0.20	0.98 ± 0.04	1	0

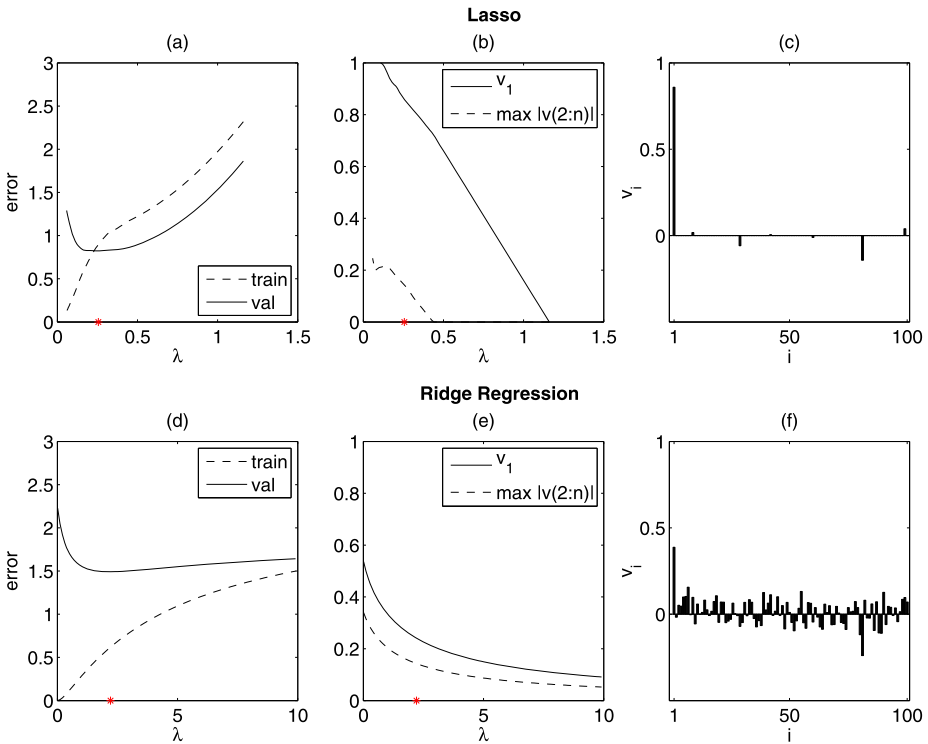


Fig. 4 Regression solution for lasso and ridge regression for same data set as in Fig. 3. *Top row (a)–(c): lasso. Bottom row (d)–(f): ridge regression. Left column (a), (d): training and validation errors versus λ . Middle column (b), (e): solution for the non-zero feature v_1 and the zero-features $\max_{i=2:n} |v_i|$. Right column (c), (f): optimal lasso and ridge regression solution v_i versus i*

with covariance matrix $\chi_{ij} = \zeta^{|i-j|}$, with $\zeta = 0.5$. In addition, we choose multiple features non-zero: $\hat{w}_i = 1, i = 1, 2, 5, 10, 50$ and all other $\hat{w}_i = 0$. We use $n = 100, \hat{\sigma} = 1$ and $p/p_v/p_r = 50/50/400$. In Table 2 we compare the performance of the VG, lasso, ridge regression and PMF on 20 random instances. We see that the VG and PMF significantly outperform the lasso method and ridge regression both in terms of prediction error and accuracy of the estimation of the parameters. Again, there is no significant difference between PMF and VG.

Table 2 Results for Example 2. For definitions see caption of Table 1

	Train	Val	Test	# non-zero	$\ \delta \mathbf{v}\ _1$
Ridge	0.32 ± 0.27	3.30 ± 0.67	3.46 ± 0.31	—	11.09 ± 0.93
Lasso	0.75 ± 0.37	1.39 ± 0.37	1.48 ± 0.29	16.30 ± 6.60	2.08 ± 0.87
PMF	—	—	1.06 ± 0.11	5.15 ± 0.49	0.67 ± 0.35
VG	0.80 ± 0.25	1.13 ± 0.31	1.15 ± 0.21	5.05 ± 0.51	0.83 ± 0.54
True	0.93 ± 0.14	0.87 ± 0.20	0.98 ± 0.04	5	0

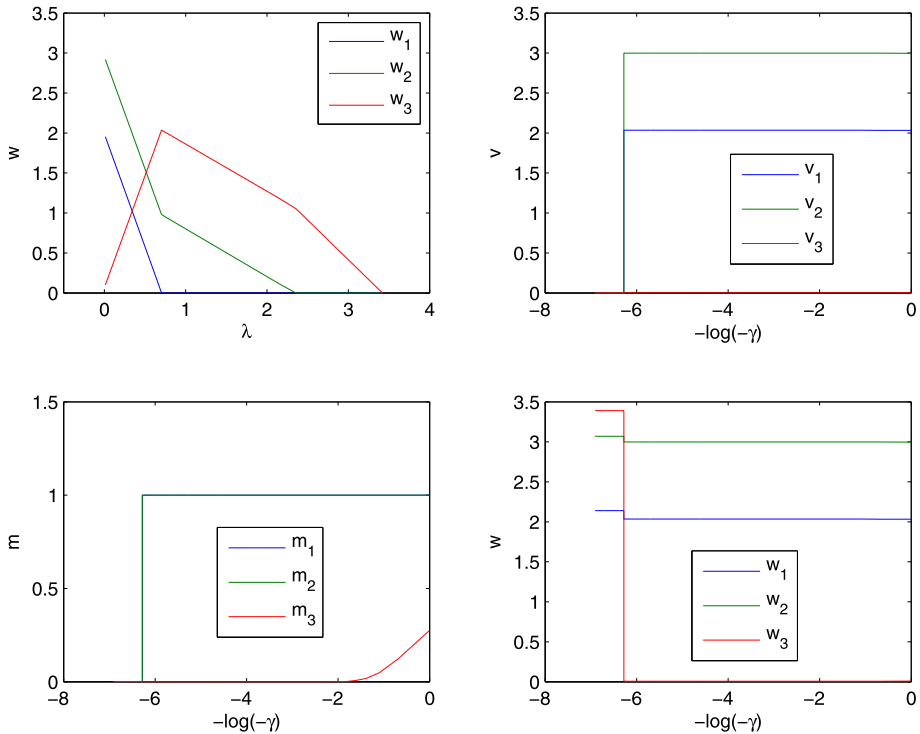


Fig. 5 Lasso and VG solution for the inconsistent Example a of Zhao and Yu (2006). *Top left*: lasso solution versus λ is called inconsistent because it does not contain a λ for which the correct sparsity ($w_{1,2} \neq 0, w_3 = 0$) is obtained. *Top right*: the VG solution for \mathbf{v} versus γ contains large range of γ for which the correct solution is obtained. *Bottom left*: VG solution for \mathbf{m} (curves for $m_{1,2}$ are identical). *Bottom right*: VG solution for \mathbf{w} (Color figure online)

5.3 Analysis of consistency: VG vs lasso

It is well-known that the lasso method may yield inconsistent results when input variables are correlated. In Zhao and Yu (2006), necessary and sufficient conditions for consistency are derived. In addition, they give a number of examples where lasso gives inconsistent results. Their simplest example has three input variables, x_1, x_2, x_3 . x_1, x_2, ξ, e are independent and Normal distributed random variables, $x_3 = 2/3x_1 + 2/3x_2 + \xi$ and $y = \sum_{i=1}^3 \hat{w}_i x_i + e$, $p = 1000$. When $\hat{w} = (-2, 3, 0)$ (Example b) this example is consis-

Table 3 Accuracy of ridge, lasso and VG for Example 1a, b from Zhao and Yu (2006). $p = p_v = 1000$. Parameters λ (ridge and lasso) and γ (VG) optimized through cross validation. $\|\delta\mathbf{v}\|_1$ as before, $\max(|v_3|)$ is maximum over 100 trials of the absolute value of v_3 . Example a is inconsistent for lasso and yields much larger errors than the VG. Example b is consistent and the quality of the lasso and VG are similar. Ridge regression is bad for both examples

	Example a		Example b	
	$\ \delta\mathbf{v}\ _1$	$\max(v_3)$	$\ \delta\mathbf{v}\ _1$	$\max(v_3)$
Ridge	0.64 ± 0.18	0.48	0.02 ± 0.02	0.27
Lasso	0.19 ± 0.14	0.30	0.00 ± 0.00	0.00
VG	0.05 ± 0.03	0.00	0.00 ± 0.00	0.00

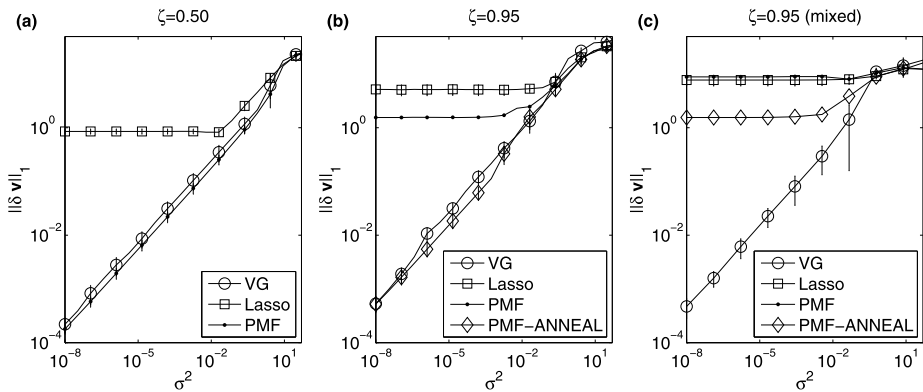


Fig. 6 Accuracy as a function of the noise. $n = 100, p = 100, p_v = 20$ and $\hat{w}_i = 1$ for 10 randomly chosen components i . (a) For weakly correlated inputs $\zeta = 0.5$ VG and PMF show comparable performance superior to lasso. (b) For strongly correlated inputs $\zeta = 0.95$ VG performs better than PMF (errorbars for PMF are not shown for clarity) but similarly to PMF-ANNEAL. (c) For $\zeta = 0.95, p = 60, p_v = 5$, and $\hat{w}_i \pm 1$ and mixed input correlations VG outperforms all methods on average

tent, but when $\hat{w} = (2, 3, 0)$ (Example a) this example violates the consistency condition. The lasso and VG solution for Example a for different values of λ and γ are shown in Fig. 5a, b, respectively. The VG solution $v_i = m_i w_i$ in terms of m_i and w_i is shown in Fig. 5c, d. The average results over 100 instances for Example a and Example b are shown in Table 3. We see that the VG does not suffer from inconsistency and always finds the correct solution, avoiding sub-optimal local minima.

5.4 Effect of the noise

In this subsection we show the accuracy VG, lasso and PMF as a function of the noise $\hat{\sigma}^2$. We generate data with $n = 100, p = 100, p_v = 20$ and $\hat{w}_i = 1$ for 10 randomly chosen components i . We vary $\hat{\sigma}^2$ in the range 10^{-8} to 10 for two values of the correlation strength in the inputs $\zeta = 0.5, 0.95$.

For weakly correlated inputs, Fig. 6a, we distinguish three noise domains: for large noise all methods produce errors of $\mathcal{O}(1)$ and fail to find the predictive features. For intermediate and low noise levels, $10^0 > \hat{\sigma}^2 > 10^{-2}$, VG and PMF perform significantly better than lasso. In the limit of zero noise, the error of VG and PMF keeps on decreasing whereas the lasso error saturates to a constant value.

Table 4 Comparison of VG and PMF in the Boston-housing dataset in terms of approximating the ground-truth \hat{w} . Average errors $\|\delta\mathbf{v}\|_1 = \sum_{i=1}^n |v_i - \hat{w}_i|$, with v_i the approximation of VG or PMF, together with 95 % confidence intervals (given by percentiles) obtained after 300 random initializations for both soft and extreme initializations

	<i>Soft-error</i>	<i>Extreme-error</i>
PMF (Titsias and Lázaro-Gredilla 2011)	0.208 [0.002, 0.454]	0.204 [0.002, 0.454]
PMF	0.237 [0.001, 0.454]	0.209 [0.001, 0.454]
VG	0.006 [0.006, 0.006]	0.006 [0.006, 0.006]

For strongly correlated inputs, Fig. 6b, we observe that whereas the error of VG scales approximately as before, PMF gets stuck in local minima in some instances, yielding worse average performance than VG. In contrast, the annealed version of PMF is able to avoid these sub-optimal solutions, resulting in average performance comparable to VG.

Finally, we consider a more challenging problem in which the weights have mixed signs $\hat{w}_i = \pm 1$, inputs are positively and negatively correlated, and a small number of samples is available ($p = 60$, $p_v = 5$). To generate negatively correlated inputs, we select a subset of the predictors and for each predictor we first obtain the indices (sample numbers) of their values sorted in ascending order. Then, we replace the predictor values with the values sorted in descending order using the previous indices. Average results for this setup are shown Fig. 6c. In this case, VG error scales as before, whereas PMF-ANNEAL gets stuck in sub-optimal solutions in some instances.

We can thus conclude that the use of annealing-reheating in the hyper-parameter optimization only explains partially the better performance of the VG compared to PMF. The results on the mixed problem suggest that the combination of a naive mean field variational approximation with a MAP step also helps to avoid local minima.

5.5 Boston-housing dataset: VG vs PMF

We now focus on comparing in more detail the performance of VG with PMF. In Titsias and Lázaro-Gredilla (2011), the Boston-housing dataset⁷ is used to test the accuracy of the PMF approximation compared to a naive mean field approximation.

This is a linear regression problem that consists of 456 training examples with one-dimensional response variable y and 13 predictors that include housing values. We use here the same setup as in Titsias and Lázaro-Gredilla (2011) to compare VG with PMF. For PMF, hyper-parameters were fixed to values $\sigma^2 = 0.1 \times \text{var}(y)$, $\pi = 0.25$, $\sigma_w^2 = 1$ where $\text{var}(y)$ denotes the output variance. For the VG, we use $\beta = 1/\sigma^2$, $\gamma = \log(\pi/(1 - \pi))$. Since γ and β are given, the VG algorithm reduces to iterate Eqs. (8) and (9) starting from a random \mathbf{m} . Similarly, the PMF reduces to perform an E-step given the fixed hyperparameter values.

As in Titsias and Lázaro-Gredilla (2011), we use random initial values for the variational parameters *between* 0 and 1 (*soft* initialization) and random values *equal to* 0 or 1 (*hard* initialization). We considered as ground truth $\hat{w} \equiv \mathbf{w}^{\text{tr}}$ the result of the efficient paired Gibbs sampler developed in Titsias and Lázaro-Gredilla (2011).

Table 4 shows the results. The first and second rows show the errors reported in Titsias and Lázaro-Gredilla (2011) and the errors that we obtain using their software, respectively.

⁷<http://archive.ics.uci.edu/ml/datasets/Housing>.

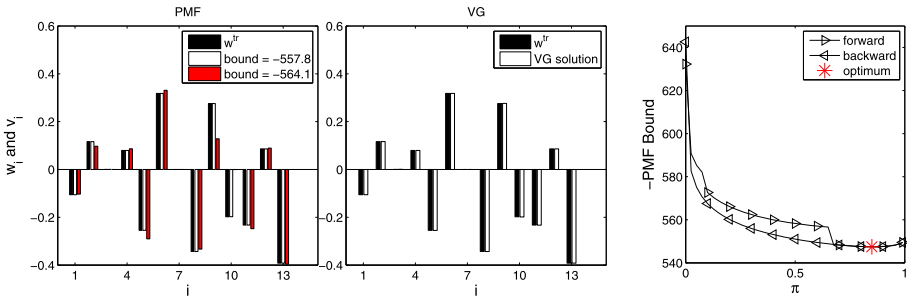


Fig. 7 Boston-housing results. w^{tr} are the true weights. *Left*: PMF finds two solutions (in white and red). The red one is suboptimal (predictor 10 is a false negative and predictor 9 is underestimated). *Middle*: VG always finds the same optimum. *Right*: hysteresis effect for π in PMF. PMF is initially trapped in the local optimum (beginning of forward pass). The global optimum is found for $\pi > \pi^*$ ($\pi^* \approx 0.7$) and continued to be the solution in the backward pass (Color figure online)

We observe a small discrepancy in the average errors. However, if we consider the percentiles, the results are consistent.

PMF finds two local optima depending on the initialization: one is the correct solution (error $\approx 10^{-3}$) whereas the other has error 0.454. These two solutions are found equally often for both soft or hard initializations, showing no dependence on the type of initialization, in agreement with Titsias and Lázaro-Gredilla (2011), and they are illustrated in Fig. 7(left).

The results of VG are shown on the third row of Table 4 and in Fig. 7(middle). Contrary to PMF, the VG shows no dependence on the initialization and always finds a solution with an error of order 10^{-3} .

The result of the annealed version of PMF for a case in which PMF converged to the suboptimal solution is illustrated in Fig. 7(right). The global optimum is found for $\pi > \pi^*$, ($\pi^* \approx 0.7$) during the forward pass and continued to be the solution in the backward pass, showing the hysteresis effect mentioned for γ in the VG. This means that for $\pi > \pi^*$, conventional PMF always converges to the global optimum, but that may not be case for $\pi < \pi^*$, depending on the initialization of the weights.

We also perform a similar experiment with VG for fixed values of γ in the corresponding range $\gamma = \log(\frac{\pi}{1-\pi})$ and for each value of γ we run VG using 100 random initial values. VG never finds a suboptimal solution and always converges to the same solution regardless of the fixed value of γ and the initialization. We thus conclude that the naive mean field variational approximation in combination with the MAP procedure do not suffer from local optima effect in this dataset.

5.6 Dependence on the number of samples

We now analyze the performance of all considered methods as a function of the number of samples available. We first analyze the case when inputs are not correlated and then consider correlations of practical relevance that appear in genetic datasets.

For these experiments, we generate the data for dimension $n = 500$ and noise level $\beta = 1$. We explore two scenarios: very sparse problems with only 10 % of active predictors and denser problems with 25 % of active predictors. The weights of the nonzero elements take integer values in increasing order starting from 1, i.e. in the sparse case, they take values from 1 to 50. We choose the validation set sizes very small ($p_v = p/10$). Choosing larger validation set sizes worsens the performance of VG compared to the PMF variants. This is

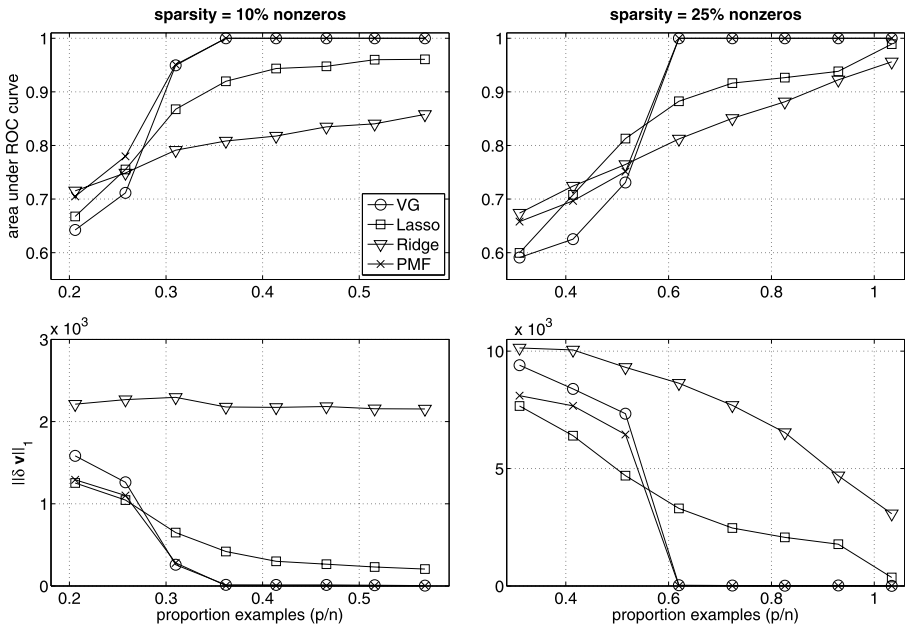


Fig. 8 *Uncorrelated case*: performance as a function of number of training samples p for two levels of sparsity (10 % and 25 % of non-zero entries). For each value averages over 20 runs are plotted. *Top*: area under the ROC curves (see text for definition). *Bottom*: reconstruction error, defined as $\|\delta v\|_1 = \sum_{i=1}^n |v_i - \hat{w}_i|$

due to the difference between using a cross validation or a Bayesian approach (PMF variants use both training and validation sets for learning).

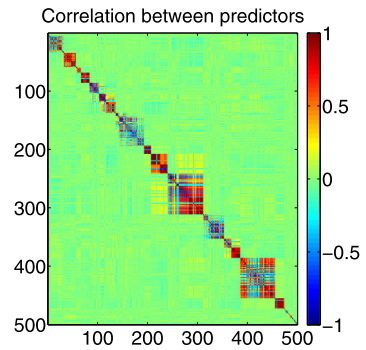
5.6.1 Uncorrelated case

Figure 8 shows results of performance for uncorrelated inputs. Top panels show the area under the Receiver Operating Characteristic (ROC) curve. The ROC curve is calculated by thresholding the weight estimates. Those weights that lie above (below) the threshold are considered as active (inactive) predictors. The ROC curve plots the fraction of true positives versus the fraction of false positives for all threshold values. The area under the curve measures the ability of the method to correctly classify those predictors that are and are not active. A value of 1 for the area represents a perfect classification whereas 0.5 represents random classification. The ROC is plotted as a function of the fraction of samples relative to the number of inputs: p/n .

For both VG and PMF, we observe in all performance measures a transition from a regime where solutions are poor to a regime with almost perfect recovery. This transition, not noticeable in the other (convex) methods, occurs at around 35 % of examples for 10 % of sparsity (left column) and shifts to higher values for denser problems (≈ 60 % for 25 % of sparsity, right column).

If we compare VG with PMF we see that PMF performs slightly better than VG in terms of area under the ROC curve and reconstruction error in the small sample size limit. Above the threshold, VG and PMF show equivalent performance. We observe no difference between PMF and PMF-ANNEAL (results not shown).

Fig. 9 Example of input correlation matrix in the genetic dataset (Color figure online)



We also see that lasso performs better than ridge regression, but the difference between both methods tends to be smaller for denser problems. Both lasso and ridge regression are significantly worse than VG and PMF.

5.6.2 Correlated case: genetic dataset

We now consider input data obtained from a genetic domain, where inputs x_i denote single nucleotide polymorphisms (SNPs) that have values $x_i \in \{0, 1, 2\}$. SNPs typically show correlations structured in blocks, where nearby SNPs are highly correlated, but show no dependence on distant SNPs. An example of such correlation matrix can be seen in Fig. 9. The raw genetic dataset for that experiment included 928 samples of 2399 three-valued SNP predictors $\{0, 1, 2\}$. To generate the dataset used in the analysis, we keep the original correlation structure of the input data but generate the outputs artificially using a randomly chosen set of active/inactive predictors. This allows to quantify the error of the different methods.

First, we filter out the less informative predictors (with entropy smaller than $\epsilon_e = 0.9$). This step removes 877 predictors. From the remaining set of 1522 predictors, we select incrementally the active ones checking that at each step the correlation between a new active predictor and the rest of active predictors is at most $\epsilon_\zeta = 0.9$. Once the active predictors have been selected, we select randomly the remaining (inactive) predictors to form a set of $n = 500$ total predictors. The values for n , ϵ_e and ϵ_ζ are chosen in a way that permits the analysis in terms of size of the training and validation sets.

Figure 10 shows the results. Contrary to the uncorrelated case, the existence of strong correlations between some of the predictors prevents a clear distinction between solution regimes as a function of sample size.

We observe, as before, that both VG and PMF are the preferable methods for sufficiently large training set size. The difference between ridge regression and lasso is more remarkable and ridge regression can even be a preferable choice than lasso for denser problems when a large number samples is available.

In all performance measures considered, VG performs better or comparable to PMF. In particular, VG significantly outperforms PMF for denser problems, which are harder due to the presence of more local minima. PMF-ANNEAL significantly improves the results of conventional PMF for both sparsity levels. From these results we can conclude that VG shows better or comparable performance than any other method considered.

5.7 Scaling with dimension n

We conclude our empirical study by analyzing how the methods scale, both in terms of the quality of the solution as in terms of CPU times, as a function of the number of features n

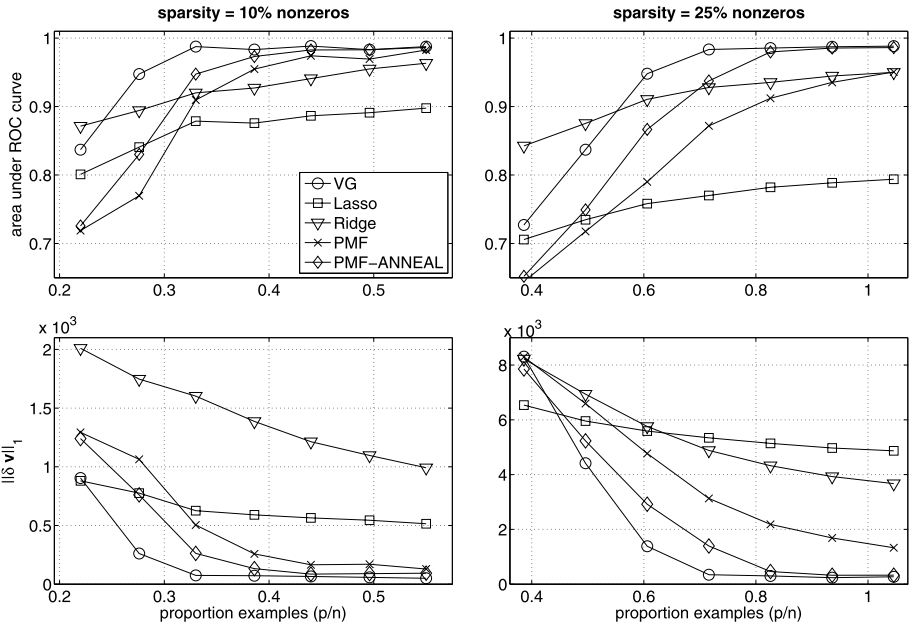


Fig. 10 *Correlated case*: performance as a function of number of training samples p for two levels of sparsity (10 % and 25 % of non-zero entries). For each value averages over 20 runs are plotted. *Top*: area under the ROC curves (see text for definition). *Bottom*: reconstruction error, defined as $\|\delta v\|_1 = \sum_{i=1}^n |v_i - \hat{w}_i|$

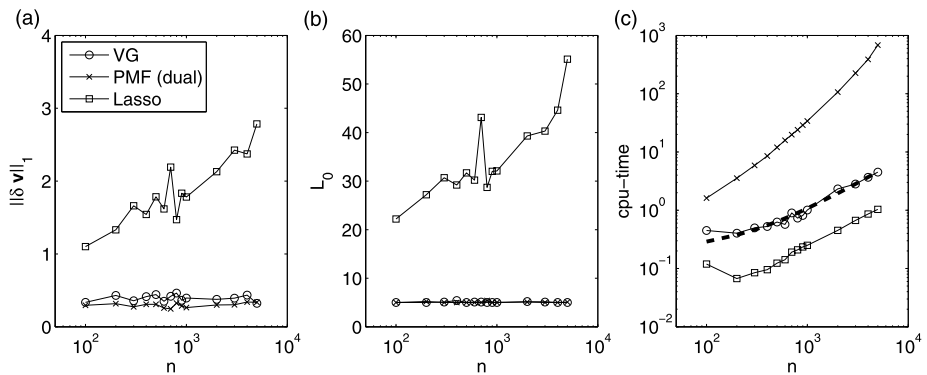


Fig. 11 *Scaling with n* : performance of VG (dual version), PMF and lasso as a function of the number of features n . (a) Error of the solution vector. (b) ℓ_0 of the solution vector. (c) cpu-time in seconds (*dashed line* corresponds to a linear fit). Data are generated as in Example 2. $p = 100$, $p_v = 100$, $\beta = 2$, $\zeta = 0$

for a constant number of samples. We use the data as in Example 2 above, with uncorrelated inputs.

Figure 11 shows the results for VG, PMF and lasso. For the VG, we use the dual method described in Appendix B. Figure 11a shows that the VG and PMF have constant quality in terms of the error $\|\delta v\|_1$, whereas the quality of the lasso deteriorates with n . Figure 11b shows that the VG and PMF have close to optimal norms $\ell_0 = 5$ and that the ℓ_0 norm of

the lasso deteriorates with n . Figure 11c shows that the computation time of all methods scales approximately linear with n . Lasso is significantly faster than VG and PMF, and VG is significantly faster than PMF. Note, however that the VG and the PMF methods are implemented in Matlab whereas the lasso method uses an optimized Fortran implementation.

6 Discussion

In this paper, we have analyzed the variational method for sparse regression using ℓ_0 penalty. We have presented a minimal version of the model with no (hierarchical) prior distributions to highlight some important features: the variational ridge term that dynamically regularizes the regression; the input-output behavior as a smoothed version of hard feature selection; a phase plot that shows when the variational solution is unique in the orthogonal design case for different p, ρ, γ .

The VG suffers from local minima as can be expected for any method that needs to solve a non-convex problem. We have shown evidence that the combined variational/MAP approach together with the annealing procedure that results from increasing γ , followed by a “heating” phase to detect hysteresis works well in practice, helping to avoid local minima. In particular, we have shown that VG can outperform a more complex model such as PMF precisely because of that reason. Further, we also have observed that VG can be still preferable to an improved version of PMF (PMF-ANNEAL) in a practical scenario with strongly correlated inputs and/or moderately sparse problems.

As mentioned in Sect. 3, the approach of Carbonetto and Stephens (2012) shares many similarities with the VG. It would be of interest to compare both approaches. We leave this comparison and other more powerful approximations, such as structured mean field approximation or belief propagation for future work.

We have seen that the performance of the VG is excellent in the zero noise limit. In this limit, the regression problem reduces to a compressed sensing problem (Candes and Tao 2005; Donoho 2006). The performance of compressed sensing with ℓ_q sparseness penalty was analyzed theoretically in Kabashima et al. (2009), showing the superiority of the ℓ_1 penalty in comparison to the ℓ_2 penalty and suggesting the optimality of the ℓ_0 penalty. Our numerical results are in agreement with this finding.

Our implementation uses parallel updating of Eqs. (8)–(10) or for the dual formulation equations (8), (22), (25), (28). One may consider also a sequential updating. This was done successfully for the lasso based on the idea of the Gauss-Seidel algorithm (Friedman et al. 2010). The advantage of such an approach is that each update is linear in both n and p , since only the non-zero components need to be updated. However, the number of updates to converge will be larger. The proof of convergence for such a coordinate descend method for the VG is likely to be more complex than for the lasso due to non-convexity. As a result, a smoothing parameter $\eta \neq 1$ (see Algorithm 1) may still be required.

Acknowledgements We would like to thank M. Titsias for providing the code of PMF and specially the Boston Housing files. We also thank Kevin Sharp and Wim Wiegenrick for useful discussions and anonymous reviewers for helping on improving the manuscript.

Open Access This article is distributed under the terms of the Creative Commons Attribution License which permits any use, distribution, and reproduction in any medium, provided the original author(s) and the source are credited.

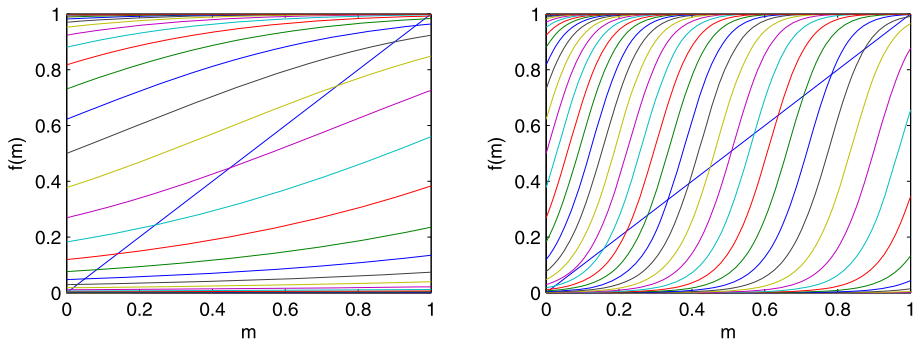


Fig. 12 $f(m)$ vs m . *Left (a)*: $p = 100, \gamma = -10$, different lines correspond to different values of $0 < \rho < 1$ (higher lines are higher ρ). The solution for m is given by the intersection f with the diagonal line. The solution for m is unique and increases with increasing ρ . *Right (b)*: same as left, but with $p = 100, \gamma = -30$. Depending on ρ , there are one or three solutions for m . The solutions close to $m \approx 0, 1$ correspond to local minima of F . The intermediate solution corresponds to a local maximum of F (Color figure online)

Appendix A: Phase plot computation for the orthogonal case

In the uni-variate case, $f(m)$ in Eq. (15) is an increasing function of m and crosses the line m either 1 or three times, depending on the values of p and γ (see Fig. 12). In the multivariate orthogonal case, this is still true, since the influence of other features is only through β . We can thus write $\beta^{-1} = \sigma_y^2(1 - \rho m - \delta)$, where $0 \leq \delta < 1$ is a function of the variational parameters of the other features. Thus, there are regions of parameter space γ, p, ρ where the uni-variate solution is unique and others for which there are two stable solutions.

The transition between these two regions is when $f'(m) = 1$ and $f(m) = m$. These two equations imply

$$\left(1 + \frac{p}{2}\right)\rho^2 m^2 - \left(2\rho(1 - \delta) + \frac{p}{2}\rho^2\right)m + (1 - \delta)^2 = 0 \tag{17}$$

This quadratic equation in m has either zero, one or two solutions, corresponding to no touching, touching once and touching twice, respectively. Denote $a = \left(1 + \frac{p}{2}\right)\rho^2, b = 2\rho(1 - \delta) + \frac{p}{2}\rho^2$. The critical value for ρ, p is when Eq. (17) has one solution for m , which occurs when

$$D = b^2 - 4a(1 - \delta)^2 = \frac{p}{2}\rho^2(\rho - \rho^*)\left(\frac{p}{2}\rho + 2(1 - \delta) + 2(1 - \delta)\sqrt{1 + \frac{p}{2}}\right) = 0 \tag{18}$$

$$\rho^* = \frac{4}{p}(1 - \delta)\left(\sqrt{1 + \frac{p}{2}} - 1\right)$$

Thus, D is positive when $\rho > \rho^*$ and Eq. (17) has two solutions for m . We denote these solutions by $m_{1,2} = \frac{b \pm \sqrt{D}}{2a}$. Note, that the solutions in these critical points only depend on ρ, p . For each of these solutions we must find a γ such that $f(m) = m$, which is given by

$$\gamma_i = \log \frac{m_i}{1 - m_i} - \frac{p}{2} \frac{\rho}{1 - \rho m_i} \quad i = 1, 2 \tag{19}$$

It is easy to see that the smallest of these solutions $m_1 < m_2$ corresponds to a local maximum of the free energy and can be discarded. Thus, when $\rho > \rho^*$ and $\gamma_2 < \gamma < \gamma_1$ two stable variational solutions $m \approx 0, 1$ co-exist.

When $\rho < \rho^*$, Eq. (17) has no solutions for m . In this case the conditions $f'(m) = 1$ and $f(m) = m$ cannot be jointly satisfied and the variational solution is unique.

From Eq. (18) we see that ρ^* is a decreasing function of p and when $p \gg 1$, $\rho^* \approx 2\sqrt{\frac{2}{p}}$. In the critical point, where $\rho = \rho^*(p)$, $m = b/2a \approx \frac{1}{2}(1 + \sqrt{\frac{2}{p}})$ and

$$\gamma^* \approx -\sqrt{2p}(1 - \delta) \tag{20}$$

When $\rho < \rho^*$ or $\gamma > \gamma^*$ the variational solution is unique. We illustrate the phase plot ρ, γ for $p = 100$ in Fig. 1a.

Appendix B: Dual formulation

The solution of the system of Eqs. (8)–(10) by fixed point iteration requires the repeated solution of the n dimensional linear system $\chi'w = b$. When $n > p$, we can obtain a more efficient method using a dual formulation.

We define new variables $z^\mu = \sum_i m_i w_i x_i^\mu$ and add Lagrange multipliers λ^μ :

$$\begin{aligned} F = & -\frac{p}{2} \log \frac{\beta}{2\pi} + \frac{\beta}{2} \sum_{\mu}^p (z^\mu - y^\mu)^2 + \frac{\beta p}{2} \sum_i m_i (1 - m_i) w_i^2 \chi_{ii} \\ & - \gamma \sum_{i=1}^n m_i + \sum_{i=1}^n (m_i \log m_i + (1 - m_i) \log(1 - m_i)) \\ & + \sum_{\mu} \lambda^\mu \left(z^\mu - \sum_i m_i w_i x_i^\mu \right) \end{aligned} \tag{21}$$

We compute the derivatives of Eq. (21):

$$\begin{aligned} \frac{\partial F}{\partial w_i} &= m_i \left(\beta p (1 - m_i) \chi_{ii} w_i - \sum_{\mu} \lambda^\mu x_i^\mu \right) \\ \frac{\partial F}{\partial z^\mu} &= \beta (z^\mu - y^\mu) + \lambda^\mu \\ \frac{\partial F}{\partial \beta} &= -\frac{p}{2\beta} + \frac{1}{2} \sum_{\mu}^p (z^\mu - y^\mu)^2 + \frac{p}{2} \sum_i m_i (1 - m_i) w_i^2 \chi_{ii} \\ \frac{\partial F}{\partial m_i} &= \frac{\beta p}{2} (1 - 2m_i) w_i^2 \chi_{ii} - \gamma + \sigma^{-1}(m_i) - \sum_{\mu} \lambda_{\mu} w_i x_i^\mu \\ \frac{\partial F}{\partial \lambda^\mu} &= z^\mu - \sum_i m_i w_i x_i^\mu \end{aligned}$$

By setting $\frac{\partial F}{\partial w_i} = \frac{\partial F}{\partial z^\mu} = 0$ we obtain

$$w_i = \frac{1}{\beta p \chi_{ii}} \frac{1}{1 - m_i} \sum_{\mu} \lambda^\mu x_i^\mu \tag{22}$$

and $z^\mu = y^\mu - \frac{1}{\beta} \lambda^\mu$. Setting the remaining derivatives to zero, and eliminating w_i and z^μ we obtain Eq. (8) and

$$\beta = \frac{1}{p} \sum_{\mu\nu} \lambda_\mu \lambda_\nu A_{\mu\nu} \tag{23}$$

$$\beta y^\mu = \sum_{\nu} A_{\mu\nu} \lambda^\nu \tag{24}$$

with $A_{\mu\nu}$ given by

$$A_{\mu\nu} = \delta_{\mu\nu} + \frac{1}{p} \sum_i \frac{m_i}{1 - m_i} \frac{x_i^\mu x_i^\nu}{\chi_{ii}} \tag{25}$$

For given $A_{\mu\nu}$, let \hat{y} denote the solution of

$$\sum_{\nu=1}^p A_{\mu\nu} \hat{y}^\nu = y^\mu \tag{26}$$

Then it is easy to verify that

$$\frac{1}{\beta} = \frac{1}{p} \sum_{\mu} \hat{y}^\mu y^\mu \tag{27}$$

$$\lambda^\mu = \beta \hat{y}^\mu \tag{28}$$

solve the system of Eqs. (23)–(24).

References

Titsias, M., & Lázaro-Gredilla, M. (2011). Spike and slab variational inference for multi-task and multiple kernel learning. In *Advances in neural information processing systems* (Vol. 24, pp. 2339–2347).

Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B, Methodological*, 58(1), 267–288.

Friedman, J. H., Hastie, T., & Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1), 1–22.

Frank, I. E., & Friedman, J. H. (1993). A statistical view of some chemometrics regression tools. *Technometrics*, 35(2), 109–135.

George, E. I., & McCulloch, R. E. (1993). Variable selection via Gibbs sampling. *Journal of the American Statistical Association*, 88(423), 881–889.

Mitchell, T. J., & Beauchamp, J. J. (1988). Bayesian variable selection in linear regression. *Journal of the American Statistical Association*, 83(404), 1023–1032.

Brown, P. J., Vannucci, M., & Fearn, T. (1998). Multivariate Bayesian variable selection and prediction. *Journal of the Royal Statistical Society, Series B, Statistical Methodology*, 60(3), 627–641.

Clyde, M., & George, E. I. (2004). Model uncertainty. *Statistical Science*, 19(1), 81–94.

Ishwaran, H., & Rao, J. S. (2005). Spike and slab variable selection: frequentist and Bayesian strategies. *The Annals of Statistics*, 33(2), 730–773.

- Carbonetto, P., & Stephens, M. (2012). Scalable variational inference for Bayesian variable selection in regression, and its accuracy in genetic association studies. *Bayesian Analysis*, 7(1), 73–108.
- Logsdon, B. A., Hoffman, G. E., & Mezey, J. G. (2010). A variational Bayes algorithm for fast and accurate multiple locus genome-wide association analysis. *BMC Bioinformatics*, 11, 58.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., & Saul, L. K. (1999). An introduction to variational methods for graphical models. *Machine Learning*, 37(2), 183–233.
- Kappen, H. J., & Spanjers, J. J. (2000). Mean field theory for asymmetric neural networks. *Physical Review E*, 61, 5658–5663.
- Murphy, K. P., Weiss, Y., & Jordan, M. I. (1999). Loopy Belief Propagation for approximate inference: an empirical study. In *Proceedings of the 15th annual conference on uncertainty in artificial intelligence* (pp. 467–475). San Francisco: Morgan Kaufmann.
- Opper, M., & Saad, D. (2001). *Advanced mean field methods: theory and practice*. Cambridge: MIT Press.
- Wainwright, M. J., & Jordan, M. I. (2008). Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1–2), 1–305.
- Barber, D., & Wiering, W. (1999). Tractable variational structures for approximating graphical models. In *Advances in neural information processing systems II* (pp. 183–189). Cambridge: MIT Press.
- Breiman, L. (1995). Better subset regression using the nonnegative garrote. *Technometrics*, 37(4), 373–384.
- O’Hara, R. B., & Sillanpää, M. J. (2009). A review of Bayesian variable selection methods: what, how and which. *Bayesian Analysis*, 4(1), 85–118.
- Attias, H. (1999). Independent factor analysis. *Neural Computation*, 11(4), 803–851.
- Yoshida, R., & West, M. (2010). Bayesian learning in sparse graphical factor models via variational mean-field annealing. *Journal of Machine Learning Research*, 99, 1771–1798.
- Hernández-Lobato, D., Hernández-Lobato, J. M., Helleputte, T., & Dupont, P. (2010). Expectation propagation for Bayesian multi-task feature selection. In *Proceedings of the 2010 European conference on machine learning and knowledge discovery in databases: part I* (pp. 522–537). New York: Springer.
- Tibshirani, R. (2011). Regression shrinkage and selection via the lasso: a retrospective. *Journal of the Royal Statistical Society, Series B, Statistical Methodology*, 73(3), 273–282.
- Fan, J., & Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456), 1348–1360.
- Mazumder, R., Friedman, J. H., & Hastie, T. (2011). SparseNet: coordinate descent with nonconvex penalties. *Journal of the American Statistical Association*, 106(495), 1125–1138.
- Park, T., & Casella, G. (2008). The Bayesian lasso. *Journal of the American Statistical Association*, 103(482), 681–686.
- Griffin, J. E., & Brown, P. J. (2011). Bayesian hyper-lassos with non-convex penalization. *Australian & New Zealand Journal of Statistics*, 53(4), 423–442.
- Zhao, P., & Yu, B. (2006). On model selection consistency of lasso. *Journal of Machine Learning Research*, 7, 2541–2563.
- Candes, E. J., & Tao, T. (2005). Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12), 4203–4215.
- Donoho, D. L. (2006). Compressed sensing. *IEEE Transactions on Information Theory*, 52(4), 1289–1306.
- Kabashima, Y., Wadayama, T., & Tanaka, T. (2009). A typical reconstruction limit for compressed sensing based on L_p -norm minimization. *Journal of Statistical Mechanics: Theory and Experiment*, 2009(09), L09003.