

Bayesian Sparse Partial Least Squares

Diego Vidaurre

diego.vidaurre@ohba.ox.ac.uk

*Oxford Centre for Human Brain Activity, University of Oxford,
Oxford OX3 7JX, U.K.*

Marcel A. J. van Gerven

m.vangerven@donders.ru.nl

*Radboud University Nijmegen, Donders Institute for Brain,
Cognition and Behavior, Nijmegen 6525H, Netherlands*

Concha Bielza

mcbielza@fi.upm.es

Pedro Larrañaga

pedro.larranaga@fi.upm.es

*Computational Intelligence Group, Universidad Politécnica Madrid,
Madrid 28660, Spain*

Tom Heskes

t.heskes@science.ru.nl

*Radboud University Nijmegen, Institute for Computing and Information
Science, Intelligent Systems, Nijmegen 6525H, Netherlands*

Partial least squares (PLS) is a class of methods that makes use of a set of latent or unobserved variables to model the relation between (typically) two sets of input and output variables, respectively. Several flavors, depending on how the latent variables or components are computed, have been developed over the last years. In this letter, we propose a Bayesian formulation of PLS along with some extensions. In a nutshell, we provide sparsity at the input space level and an automatic estimation of the optimal number of latent components. We follow the variational approach to infer the parameter distributions. We have successfully tested the proposed methods on a synthetic data benchmark and on electrocorticogram data associated with several motor outputs in monkeys.

1 Introduction ---

Partial least squares (PLS) (Wold, Sjöström, & Eriksson, 2001) is a family of techniques originally devised for modeling two sets of observed variables, which we shall call input and output components, by means of some

(typically low-dimensional) set of latent or unobserved components. This model can also be extended to deal with more than two sets of components (Wangen & Kowalsky, 1989). It is commonly used for regression but is also applicable to classification (Barker & Rayens, 2003). In this letter, we focus on the regression paradigm. Latent components are generated by some linear transformation of the input components, while the output components are assumed to be generated by some linear transformation of the latent components.

The difference between PLS and related techniques lies in how the latent components are estimated (Hastie, Tibshirani, & Friedman, 2008; Rosipal & Krämer, 2006). Unlike PLS, principal components regression (PCR), for example, does not consider the output when constructing the latent components. Also, PLS differs from canonical correlation analysis (CCA) in that CCA treats the input and output spaces symmetrically (Hardoon, Szedmak, & Shawe-Taylor, 2004). A complete comparison between PLS, PCR and classical shrinkage regression from a statistical perspective is given by Frank and Friedman (1993), and further insight into the shrinkage properties of PLS can be found, for instance, by Goutis (1996). There exist Bayesian formulations of some latent component models in the literature, such as PCA (Bishop, 1998; Nakajima, Sugiyama, & Babacan, 2011), CCA (Fujiwara, Miyawaki, & Kamitani, 2009; Virtanen, Klami, & Kaski, 2011; Wang, 2007) and mixtures of factor analyzers (Beal, 2003; Ghahramani & Beal, 2000). To our knowledge, however, a Bayesian version of PLS has not yet been proposed.

Different varieties of PLS regression arise by the way they extract latent components (Rosipal & Krämer, 2006). In its classic form, PLS aims to maximize the covariance among the latent components, which are constrained to be orthogonal, using the nonlinear iterative partial least squares (NIPALS) algorithm (Wold, 1975). This is more an algorithmic than a traditional statistical approach, and, hence, the analysis of its properties is less obvious. A more rigorous approach (from a statistical perspective) is taken by de Jong (1993), who directly formulates the latent space as a linear projection of the input space and solves the resulting optimization problem by the so-called SIMPLS algorithm. The SIMPLS algorithm is equivalent to NIPALS only when the output space is unidimensional. Sparsifying accounts of PLS are proposed by van Gerwen, Chao, and Heskes (2012) and Chun and Keleş (2010). A kernelized approach has been introduced by Lindgren, Geladi, and Wold (1993) and Rosipal and Trejo (2001).

The main goal of this letter is to develop a Bayesian approach for PLS regression. We use variational inference (Jaakkola, 2001) for estimating the parameters. Let \mathbf{X} be an $N \times p$ input matrix and \mathbf{Y} be an $N \times q$ output matrix, with elements x_{ni} and y_{nj} and rows \mathbf{x}_n and \mathbf{y}_n . Assuming centered data, we follow the definition of PLS given by

$$\mathbf{Z} = \mathbf{X}\mathbf{P} + \boldsymbol{\epsilon}_Z, \quad \mathbf{Y} = \mathbf{Z}\mathbf{Q} + \boldsymbol{\epsilon}_Y,$$

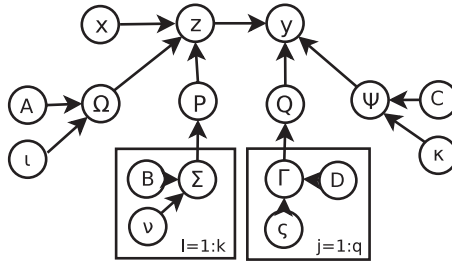


Figure 1: Bayesian hierarchy of the proposed Bayesian PLS model, where Z lies in the latent space.

where P and Q are, respectively, $p \times k$ and $k \times q$ loading matrices, Z is the $N \times k$ latent score matrix, with elements z_{il} and rows $z_{i\cdot}$, and ϵ_Z and ϵ_Y are the matrices of residuals. We use an intermediate k -dimensional latent space, k typically being lower than p and q . We consider a Bayesian hierarchy defined through several normal Wishart distributions for the latent and output variables, as well as for the loading matrices:

$$\begin{aligned}
 z' &\sim \mathcal{N}(x'P, \Omega), \quad \Omega^{-1} \sim \mathcal{W}(A, \iota), \quad p_l \sim \mathcal{N}(0, \Sigma_l), \quad \Sigma_l^{-1} \sim \mathcal{W}(B_l, \nu_l), \\
 y' &\sim \mathcal{N}(z'Q, \Psi), \quad \Psi^{-1} \sim \mathcal{W}(C, \kappa), \quad q_j \sim \mathcal{N}(0, \Gamma_j), \quad \Gamma_j^{-1} \sim \mathcal{W}(D_j, \varsigma_j),
 \end{aligned}
 \tag{1.1}$$

with $l = 1, \dots, k$ and $j = 1, \dots, q$, and where p_l is the l th column of P and q_j is the j th column of Q . A, B_l, C, D_j are the scale matrix hyperparameters of the Wishart prior distributions, and $\iota, \nu_l, \kappa, \varsigma_j$ are the corresponding degrees of freedom. (See Figure 1 for a graphical representation using plate notation.) In the remainder, we suppress the hyperparameters in our notation when it is clear from the context.

By imposing separate gaussian priors for each column of P (and Q), we are allowing different input (and latent) variable couplings for each component $l = 1, \dots, k$ (and $j = 1, \dots, q$). An obvious simplification, which we take in the rest of the letter, is to let $\Sigma_l = \Sigma$ (and $\Gamma_j = \Gamma$), so that information is borrowed among the latent components and responses and the number of parameters is reduced. The derivation of the general case is straightforward.

It might appear that for large p scenarios, a large number of parameters is associated with the full precision matrix Σ^{-1} . However, as we will show, the estimation of these matrices is low rank, so that the effective number of parameters is kept reasonably low.

Note that $E[y'|x'] = x'PQ$. Since PQ has at most rank k , this formulation is related to reduced rank methods (Izenman, 1975) and approaches that penalize the nuclear norm of the coefficient matrix (Yuan, Ekici, Lu,

& Monteiro, 2007). There is also a connection with the multivariate group Lasso (Obozinski, Wainwright, & Jordan, 2011), which imposes an L_1/L_2 -penalty on the coefficient matrix so that a common sparsity pattern is shared by all responses. However, the multivariate group Lasso does not account for correlated errors. The sparse multivariate regression with covariance estimation approach (Rothman, Levina, & Zhu, 2010), on the other hand, does consider correlation between the responses by simultaneously estimating the coefficient matrix and the (full) inverse covariance matrix of the response variables. The coefficient matrix is L_1 -regularized, and then the sparsity pattern can vary for each response. Our approach is also somewhat related to the multitask feature learning problem, where each task has a different set of inputs and the goal is to find some shared structural parameterization that is beneficial for the individual tasks. For example, the method proposed by Argyriou, Evgeniou, and Pontil (2006) seeks a low-rank linear transformation such that the outputs are encouraged to share a common input sparsity pattern. The method introduced by Ando and Zhang (2005) is formulated so that unlabeled data can be used for learning common underlying predictive functional structures. However, these approaches do not build on a generative model and it is not possible to express a Bayesian formulation that leads to the same estimator.

The rest of the letter is organized as follows. Section 2 introduces the variational approximation in the basic setting. Section 3 describes how to achieve a sparse solution. Section 4 proposes an improved model that aims to estimate the optimal number of latent components and increase the accuracy. Section 5 presents a simulation study with comparisons to other methods. Section 6 provides some results for real neural signal decoding. Finally, section 7 provides conclusions and directions for future work.

2 Variational Parameter Inference

We are interested in the posterior distribution $Pr(P, Q | X, Y)$, given by

$$\int P(P, Q, Z, \Omega^{-1}, \Sigma^{-1}, \Psi^{-1}, \Gamma^{-1} | X, Y) dZ d\Omega^{-1} d\Sigma^{-1} d\Psi^{-1} d\Gamma^{-1}.$$

For computational reasons, we approximate the posterior distribution of the parameters given Y by a variational distribution with the following factorization:

$$\begin{aligned} P(P, Q, Z, \Omega^{-1}, \Sigma^{-1}, \Psi^{-1}, \Gamma^{-1} | X, Y) \\ \approx F(P, Q, Z, \Omega^{-1}, \Sigma^{-1}, \Psi^{-1}, \Gamma^{-1}) = F(Z)F(P, \Omega^{-1}, \Sigma^{-1}, Q, \Psi^{-1}, \Gamma^{-1}). \end{aligned}$$

The variational approximation automatically (i.e., without the need for further assumptions) factorizes into $F(Z)F(P, \Omega^{-1}, \Sigma^{-1})F(Q, \Psi^{-1}, \Gamma^{-1})$.

This can be easily verified by inspecting the functional $F(P, \Omega^{-1}, \Sigma^{-1}, Q, \Psi^{-1}, \Gamma^{-1})$, defined as the log of the joint distribution when we take the expectation with respect to Z (Beal, 2003). Since Z separates $P, \Omega^{-1}, \Sigma^{-1}$ from $Q, \Psi^{-1}, \Gamma^{-1}$ in the Bayesian hierarchy, the resulting expression for $F(P, \Omega^{-1}, \Sigma^{-1}, Q, \Psi^{-1}, \Gamma^{-1})$ does not have interaction terms between the two groups of variables.

Also on computational grounds, we assume $F(P, \Omega^{-1}, \Sigma^{-1}) = F(P)F(\Omega^{-1}, \Sigma^{-1})$ and, analogously, $F(Q, \Psi^{-1}, \Gamma^{-1}) = F(Q)F(\Psi^{-1}, \Gamma^{-1})$. From this, we have an automatic factorization between Ω^{-1} and Σ^{-1} (Ψ^{-1} and Γ^{-1}). Finally, $F(Z)$ automatically factorizes into $\prod_{n=1}^N F(z_n)$, so that up to a constant, we have

$$\begin{aligned} F(z_n) &= E_{P,Q,\Omega^{-1},\Psi^{-1}}[\log P(z_n|X, P, \Omega^{-1}) + \log P(y|z_n, Q, \Psi^{-1})] \\ &= -\frac{1}{2}z'_n(E[\Omega^{-1}] + E[Q\Psi^{-1}Q'])z_n + (x'_n\mu_P E[\Omega^{-1}] + y'_n E[\Psi^{-1}]\mu_Q)z_n, \end{aligned}$$

where μ_P and μ_Q are the expectations of, respectively, P and Q . Expectations are with regard to the variational distribution. Completing the square, we have

$$F(z_n) = \mathcal{N}(z_n; \mu_{z_n}, S_{z_n}) \tag{2.1}$$

with $S_{z_n} = (E[\Omega^{-1}] + E[Q\Psi^{-1}Q'])^{-1}$ and $\mu_{z_n} = S_{z_n} (E[\Omega^{-1}]\mu'_P x_n + \mu_Q E[\Psi^{-1}]y_n)$. We can compute

$$E[Q\Psi^{-1}Q'] = \mu_Q E[\Psi^{-1}]\mu'_Q + \sum_{j_1=1}^q \sum_{j_2=1}^q E[\Psi^{-1}_{j_1 j_2}] S_{Q_{j_1 j_2}},$$

where $S_{Q_{j_1 j_2}}$ denotes the $k \times k$ cross-covariance matrix relative to the j_1 th and j_2 th columns of the loading matrix Q .

For Ω^{-1} , we have, up to a constant,

$$\begin{aligned} \log F(\Omega^{-1}) &= E_{Z,P}[\log P(Z|X, P, \Omega^{-1}) + \log P(\Omega^{-1})] \\ &= \frac{\iota - k - N - 1}{2} \log |\Omega| - \frac{1}{2} \text{Tr}((E[Z'Z] + E[P'X'XP] \\ &\quad - \mu'_Z X \mu_P - \mu'_P X' \mu_Z + A^{-1}) \Omega^{-1}), \end{aligned}$$

where we have used standard properties of the trace operator. Here, we can identify a Wishart distribution,

$$F(\Omega^{-1}) = \mathcal{W}(\Omega^{-1}; \tilde{A}^{-1}, \tilde{\iota}), \tag{2.2}$$

with $\tilde{\iota} = \iota + N$ and $\tilde{\mathbf{A}}^{-1} = (E[\mathbf{Z}'\mathbf{Z}] + E[\mathbf{P}'\mathbf{X}'\mathbf{X}\mathbf{P}] - \boldsymbol{\mu}'_Z\mathbf{X}\boldsymbol{\mu}_P - \boldsymbol{\mu}'_P\mathbf{X}\boldsymbol{\mu}_Z + \mathbf{A}^{-1})^{-1}$, where $E[\mathbf{Z}'\mathbf{Z}] = \sum_{n=1}^N E[z_n z'_n] = \sum_{n=1}^N (\mathbf{S}_{z_n} + \boldsymbol{\mu}_{z_n} \boldsymbol{\mu}'_{z_n})$.

If we set $\boldsymbol{\Omega}^{-1}$ to be diagonal, then the variational distribution $F(\mathbf{P})$ factorizes over columns and we get a gamma distribution for each element Ω_{ll}^{-1} :

$$F(\Omega_{ll}^{-1}) = \mathcal{G}(\Omega_{ll}^{-1}; \tilde{\iota}, \tilde{\mathbf{A}}_{ll}^{-1}) \tag{2.3}$$

with $\tilde{\iota} = \iota + \frac{N}{2}$ and $\tilde{\mathbf{A}}_{ll}^{-1} = \frac{1}{2}(E[\mathbf{Z}_l \mathbf{Z}'_l] + E[\mathbf{p}'_l \mathbf{X}' \mathbf{X} \mathbf{p}_l] - 2\boldsymbol{\mu}'_{z_l} \mathbf{X} \boldsymbol{\mu}_{p_l}) + \mathbf{A}_{ll}^{-1}$, where \mathbf{Z}_l denotes the l th column of \mathbf{Z} .

If we do not factorize $F(\mathbf{P})$, that is, if $\boldsymbol{\Omega}^{-1}$ is not chosen to be diagonal, then we have, up to a constant,

$$\begin{aligned} \log F(\mathbf{P}) &= E_{\mathbf{Z}, \boldsymbol{\Sigma}^{-1}, \boldsymbol{\Omega}^{-1}} \left[\log P(\mathbf{Z}|\mathbf{X}, \mathbf{P}, \boldsymbol{\Omega}^{-1}) + \sum_{l=1}^k \log P(\mathbf{p}_l|\boldsymbol{\Sigma}) \right] \\ &= - \sum_{n=1}^N \left(\frac{1}{2} \mathbf{x}'_n \mathbf{P} E[\boldsymbol{\Omega}^{-1}] \mathbf{P}' \mathbf{x}_n - \mathbf{x}'_n \mathbf{P} E[\boldsymbol{\Omega}^{-1}] \boldsymbol{\mu}_{z_n} \right) - \frac{1}{2} \sum_{l=1}^k \mathbf{p}'_l E[\boldsymbol{\Sigma}^{-1}] \mathbf{p}_l. \end{aligned}$$

We define $\tilde{\mathbf{p}}$ as the concatenation of the rows of \mathbf{P} and $\tilde{\mathbf{p}}^*$ as the concatenation of the rows of the $p \times k$ least-squares solution $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\mu}_Z$, so that after some algebra, we can identify a pk -dimensional gaussian distribution,

$$F(\mathbf{P}) = \mathcal{N}(\mathbf{P}; \boldsymbol{\mu}_P, \mathbf{S}_P) \tag{2.4}$$

with $\mathbf{S}_{\tilde{\mathbf{p}}} = (E[\boldsymbol{\Sigma}^{-1}] \otimes \mathbf{I}_k + \mathbf{X}'\mathbf{X} \otimes E[\boldsymbol{\Omega}^{-1}])^{-1}$ and $\boldsymbol{\mu}_{\tilde{\mathbf{p}}} = \mathbf{S}_{\tilde{\mathbf{p}}}(\mathbf{X}'\mathbf{X} \otimes E[\boldsymbol{\Omega}^{-1}])\tilde{\mathbf{p}}^*$, where \mathbf{I}_k is the $k \times k$ identity matrix and \otimes denotes the Kronecker product. From this expression, we can reconstruct $\boldsymbol{\mu}_P$ and \mathbf{S}_P .

When $\boldsymbol{\Omega}^{-1}$ is diagonal, we can simplify $F(\mathbf{P}) = \prod_{l=1}^k F(\mathbf{p}_l)$. For each factor, we have

$$F(\mathbf{p}_l) = \mathcal{N}(\mathbf{p}_l; \boldsymbol{\mu}_{p_l}, \mathbf{S}_{p_l}) \tag{2.5}$$

with $\mathbf{S}_{p_l} = (E[\boldsymbol{\Sigma}^{-1}] + E[\Omega_{ll}^{-1}]\mathbf{X}'\mathbf{X})^{-1}$ and $\boldsymbol{\mu}_{p_l} = E[\Omega_{ll}^{-1}]\mathbf{S}_{p_l}\mathbf{X}'\boldsymbol{\mu}_{z_l}$.

For $\boldsymbol{\Sigma}^{-1}$, we have, up to a constant,

$$\begin{aligned} \log F(\boldsymbol{\Sigma}^{-1}) &= E_{p_l} \left[\log P(\boldsymbol{\Sigma}^{-1}) + \sum_{l=1}^k \log P(\mathbf{p}_l|\boldsymbol{\Sigma}^{-1}) \right] \\ &= \frac{v_l - p - k - 1}{2} \log |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \text{Tr}((E[\mathbf{P}\mathbf{P}'] + \mathbf{B}_l^{-1})\boldsymbol{\Sigma}^{-1}), \end{aligned}$$

where we can identify a Wishart distribution:

$$F(\Sigma^{-1}) = \mathcal{W}(\Sigma^{-1}; \tilde{\mathbf{B}}^{-1}, \tilde{\mathbf{v}}), \quad (2.6)$$

with $\tilde{\mathbf{B}}^{-1} = (E[\mathbf{P}\mathbf{P}'] + \mathbf{B}^{-1})^{-1}$ and $\tilde{\mathbf{v}} = \nu + k$.

Note that the matrix $E[\mathbf{P}\mathbf{P}']$ is not full rank as far as $p > k$, which is typically the case. It has, in fact, rank k . Then the effective number of parameters of this matrix is not $p(p-1)/2$, but, at most, $pk + 1 - k(k-1)/2$. When p is high relative to N , it becomes necessary to borrow information between the components $l = 1, \dots, k$, suggesting the choice $\Sigma_l^{-1} = \Sigma^{-1}$.

Calculations for \mathbf{Q} and dependent distributions are similar to those of \mathbf{P} and are given in appendix A. Brown & Zidek (1980) theoretically showed that an adaptive joint estimation dominates an independent estimation for each of the outputs separately when the number of inputs is considerably larger than the number of outputs, but this domination breaks down when the number of outputs approaches the number of inputs. In our situation, when estimating \mathbf{Q} , the number of outputs q typically even exceeds the number of hidden units k . This suggests that the factorization $F(\mathbf{Q}) = \prod_{j=1}^q F(q_j)$, mimicking independent estimation, is the sensible choice for $q > k$, which is often the case.

In short, the proposed approach proceeds as follows:

1. Initialize \mathbf{Z} to the k first principal components of \mathbf{Y} .
2. Compute the distributions of Ω^{-1} , \mathbf{P} and Σ^{-1} using equations 2.2, 2.4, and 2.6.
3. Compute the distributions of Ψ^{-1} , \mathbf{Q} and Γ^{-1} using equations A.1, A.3, and A.5.
4. Compute the distribution of \mathbf{Z} using equation 2.1.
5. Repeat steps 2 to 4 until convergence.

This grouping of the updates is motivated by the structure of the Bayesian hierarchy and the variational factorization in equation 2.1. A variant of the basic scheme, by assuming diagonality of Ω^{-1} and Ψ^{-1} , arises by substituting equation 2.2 by 2.3, 2.4 by 2.5, A.1 by A.2, and A.3 by A.4. A variational lower bound of the evidence is given in appendix B.

3 Sparsity in \mathbf{P} and \mathbf{Q}

For achieving sparsity on the input variables, we may impose a group-sparsifying prior on \mathbf{P} , so that the groups are the k -dimensional rows of \mathbf{P} . By setting an individual regularization parameter on each group and integrating out \mathbf{P} , the maximum likelihood value of such regularization parameters will effectively discard some groups. This is an example of groupwise automatic relevance determination (see, for example, Virtanen et al., 2011).

To achieve this objective, we set priors

$$P_i \sim \mathcal{N}(\mathbf{0}, \sigma_i^2 \mathbf{I}_k),$$

where P_i is the i th row of P . This way, we will effectively drop the useless inputs (rows of P). We define the precisions σ_i^{-2} to be gamma distributed. The variational approximation of the posterior of p_i becomes a gaussian distribution with parameters $S_{p_i} = (\text{diag}(E[\sigma^{-2}]) + E[\Omega_{ll}^{-1}]X'X)^{-1}$ and $\mu_{p_i} = E[\Omega_{ll}^{-1}]S_{p_i}X'\mu_{Z_i}$.

The derivation for nonfactorized P is straightforward. For σ_i^{-2} , we have

$$E[\sigma_i^{-2}] = \frac{2\nu + k}{E[P_i'P_i] + 2B_{ii}^{-1}} = \frac{2\nu + k}{\sum_{l=1}^k (S_{p_{i,l}} + \mu_{p_{i,l}}^2) + 2B_{ii}^{-1}}.$$

Also, we impose a similar groupwise prior on Q :

$$Q_l \sim \mathcal{N}(\mathbf{0}, \gamma_l^2 \mathbf{I}_q),$$

with the precision γ_l^{-2} being gamma distributed. The idea is to obtain a data-driven estimation of the importance of each latent component when estimating Q . The variational approximation of q_j is a gaussian distribution with parameters $S_{q_j} = (\text{diag}(E[\gamma^{-2}]) + E[\Psi_{jj}^{-1}]E[Z'Z])^{-1}$ and $\mu_{q_j} = E[\Psi_{jj}^{-1}]S_{q_j}\mu_Z'Y_{.j}$.

Again, we follow a variational approximation to obtain

$$E[\gamma_l^{-2}] = \frac{2\zeta + q}{E[Q_l'Q_l] + 2D_{ll}^{-1}} = \frac{2\zeta + q}{\sum_{j=1}^q (S_{Q_{l,j}} + \mu_{Q_{l,j}}^2) + 2D_{ll}^{-1}}.$$

Then a value σ_i^{-2} (γ_l^{-2}) close to zero means that the corresponding input (latent) variable can be considered irrelevant.

4 Adaptive P and Automatic Selection of k ---

The proposed estimation of P disregards Y and uses only the current state of Z . Put differently, equal attention is paid to all latent components when estimating P , no matter the contribution of each latent component to the prediction of Y . A possible improvement would be to focus on those latent components that are more useful for modeling Y , regularizing more the latent components that are relatively useless—those whose coefficients in Q are lower.

Here, we impose groupwise priors on P for both rows and columns so as to achieve two goals. First, we provide a more adaptive estimation of

P , which we hope will reduce the bias. Second, starting with a reasonably high value of k , we will be better able to obtain an estimate of the adequate dimension for the latent space. Hence, the purpose of this section is to improve the model performance for a given number of latent components. A proper model selection scheme, which would compare the model evidence (see appendix B) for different values of k , would be an alternative (or complementary) way to go.

We consider the variational approximation with P and Q factorized over columns. We let γ^{-2} also influence P so that it indeed controls the latent space for both projection matrices. We balance the regularization effect of γ^{-2} with an additional variable ϕ^{-2} , so that the relative contributions of σ^{-2} and γ^{-2} are adaptively estimated from data. We choose ϕ^{-2} to be gamma distributed for keeping conjugacy. Specifically, we propose

$$P_{ii} \sim \mathcal{N}(\mathbf{0}, (\sigma_i^{-2} + \phi^{-2}\gamma_l^{-2})^{-1}), \quad \phi^{-2} \sim \mathcal{G}(e, \varphi).$$

The variational distribution of p_l is a gaussian distribution whose parameters are $S_{p_l} = (\text{diag}(E[\sigma^{-2}]) + E[\phi^{-2}]E[\gamma_l^{-2}]\mathbf{I}_p + E[\Omega_{ll}^{-1}]\mathbf{X}'\mathbf{X})^{-1}$ and $\mu_{p_l} = E[\Omega_{ll}^{-1}]S_{p_l}\mathbf{X}'\mu_{z_l}$.

For γ^{-2} , we have

$$E[\gamma_l^{-2}] = \frac{2\zeta + p + q}{\phi^{-2} \sum_{i=1}^p (S_{P_{i,i}} + \mu_{P_{i,i}}^2) + \sum_{j=1}^q (S_{Q_{j,j}} + \mu_{Q_{j,j}}^2) + 2d_l^{-1}},$$

and, for ϕ^{-2} ,

$$E[\phi^{-2}] = \frac{2\varphi + pk}{\sum_{i=1}^p \sum_{l=1}^k \gamma_l^{-2} (S_{P_{i,i}} + \mu_{P_{i,i}}^2) + 2e^{-1}}.$$

Then, by automatic relevance determination, the number of coordinates of γ^{-2} significantly greater than zero at convergence is an estimation of the optimal number of latent variables k .

5 Synthetic Experiments

In order to test in practice the performance of the algorithms, we now present a synthetic simulation study where we compare the sparse approach proposed in section 3 (SPLS for short) with other methods. For simplicity, we do not present results for the nonfactorized estimation of P and Q or for full matrices Σ^{-1} and Γ^{-1} .

We have tested some alternative PLS methods: the NIPALS and SIMPLS algorithms and the frequentist sparse PLS from Chun and Keleş (2010) (sSIMPLS). All the PLS methods have been provided with four latent

components. Also, we have tested a number of non-PLS univariate and multivariate regression techniques. Univariate methods are separately applied for each response; these include ordinary least squares regression (OLS), ridge regression, and the Lasso (Hastie et al., 2008). From the multivariate side, we have tested the multivariate group lasso (MGL) (Obozinski et al., 2011) and the sparse multivariate regression with covariance estimation approach (MRCE) (Rothman et al., 2010). The regularization parameters for the non-Bayesian methods have been chosen by cross-validation.

The design of the simulation study is as follows. In all experiments, the number of input variables is $p = 50$, and the number of responses is $q = 8$. We have covered several different situations, varying the number of hidden components, which we denote as k_0 , and the amount of training data ($N = 100, 500$). We have set $k_0 = 1, 2, 4, 8$. Within each situation, we generated 100 random replications. For each of them, we have sampled 1000 testing data points. Then, for each situation and each model, reports are shown over $100 \times 1000 \times 8 = 800,000$ residuals.

For each random replication, the input matrix was generated according to a multivariate gaussian distribution with zero mean and covariance matrix M , whose elements were set as $M_{i_1 i_2} = r_1^{|i_1 - i_2|}$, and r_1 was sampled from the uniform distribution with support $[0, 1]$. Hence, depending on r_1 , the degree of collinearity in the input matrix is different among the replications. Of course, the testing input matrix was sampled using the same covariance matrix M .

When $k_0 < 8$, the latent components were generated as $Z = XP + \epsilon_Z$. Sparsity was enforced by setting each row of P to zero with probability 0.8 (ensuring at least two relevant input variables). The rest of the rows were sampled from a standard normal distribution. The noise ϵ_Z was generated from a gaussian distribution with zero mean and diagonal covariance matrix Ω , with $\Omega_{ll} = r_2 \text{sd}(Xp_l)$, where $\text{sd}(\cdot)$ denotes the sample standard deviation and the value r_2 was sampled from the uniform distribution with support $[0.01, 0.1]$ (separately for each diagonal element).

We generated the responses as $Y = ZQ + \epsilon_Y$. We did not consider sparsity in Q , whose elements were sampled from a standard normal distribution. The noise ϵ_Y was sampled from a gaussian distribution with zero mean and diagonal covariance matrix Ψ , with diagonal elements $\Psi_{jj} = r_3 \text{sd}(Zq_j)$, where r_3 was sampled from the uniform distribution with support $[0.25, 0.5]$.

When $k_0 = 8$, the response was computed as $Y = XF + \epsilon_Y$, where F has rank q . Since $q > k$, F has a higher number of effective parameters than PQ in the factorized ($k_0 < 8$) case. F was sampled from a normal standard distribution, considering sparsity as before.

Figures 2, 3, 4, and 5 show box plots with the performance of the different methods for all situations. Results are in terms of the explained variance R^2 . In short, these graphs illustrate how the methods behave for different

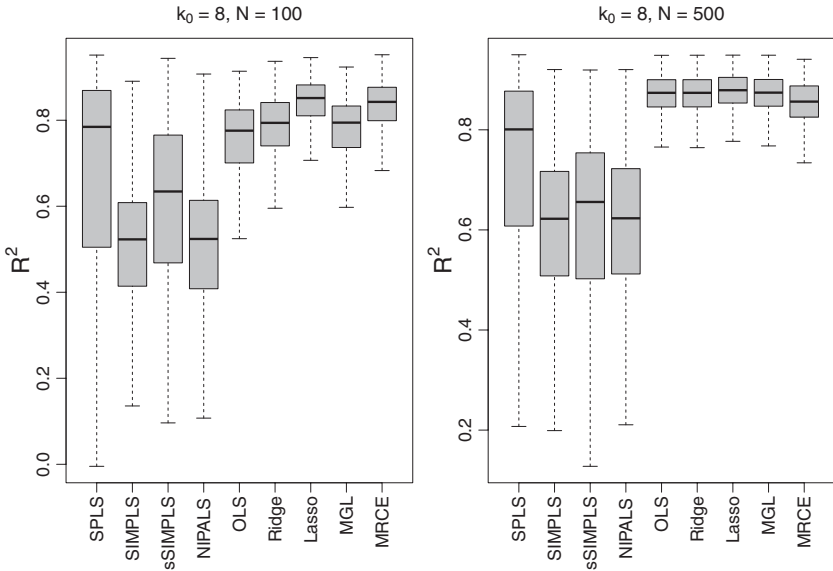


Figure 2: Box plot of the R^2 values for $k_0 = 8$ and $N = 100, 500$.

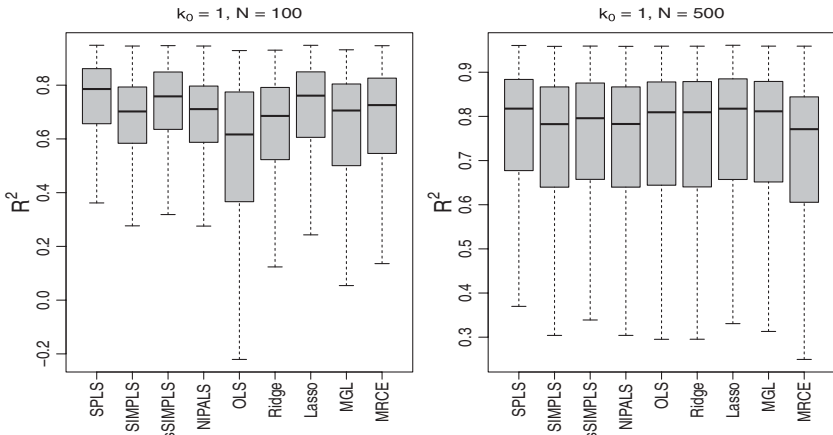


Figure 3: Box plot of the R^2 values for $k_0 = 1$ and $N = 100, 500$.

complexities of the true response function and amounts of training input data.

When $k_0 = 8$, the response function does not factorize and has the highest number of parameters. In this case, the PLS methods have fewer parameters than the actual true response function and tend to underfit. This is in

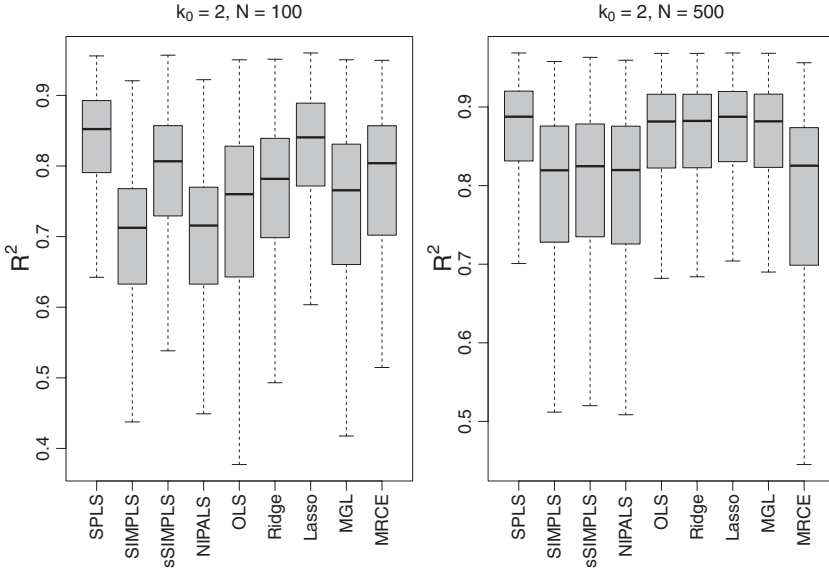


Figure 4: Box plot of the R^2 values for $k_0 = 2$ and $N = 100, 500$.

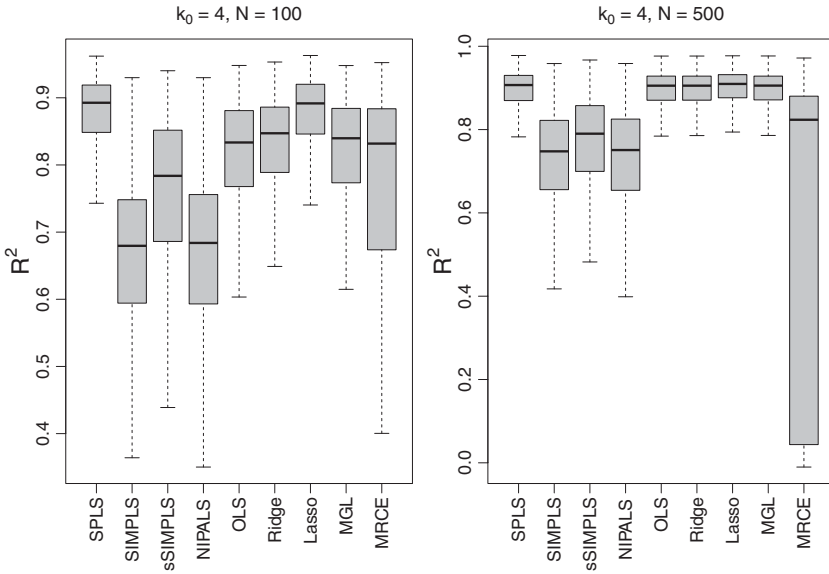


Figure 5: Box plot of the R^2 values for $k_0 = 4$ and $N = 100, 500$.

particular the case for $N = 500$, when there are sufficient data for a reliable estimation of all the parameters. Interestingly, SPLS is the only PLS method that does not perform much worse than ridge regression, the Lasso, MGL, and MRCE for $N = 100$ and also for $N = 500$.

The $k_0 = 1$ case is the opposite extreme. Whereas the full-rank methods aim to estimate $pq = 400$ parameters, the PLS methods, which are set to use four latent components, estimate $4p + 4q = 232$ parameters. Still, this is much more than the true number of parameters ($p + q = 58$). For this reason, the PLS methods are not much better than the others. Indeed, SPLS and the Lasso appear to be the more effective in controlling the complexity of the model. Note that the differences between the methods are very subtle for $N = 500$.

When $k_0 = 2$, the true response function has 116 parameters. In this case, for $N = 100$, SPLS and the Lasso clearly outperform the other methods. SPLS is slightly better than the Lasso. Surprisingly, SIMPLS and NIPALS do not perform better than OLS. For $N = 500$, SPLS, OLS, ridge regression, the Lasso, and MGL are almost indistinguishable and better than the others.

Finally, when $k_0 = 4$, the true number of parameters (232) matches that of the PLS methods. Surprisingly, for $N = 100$, SIMPLS, sSIMPLS, and NIPALS are again less accurate than any of the univariate regression methods, MGL and MRCE. On the contrary, SPLS, followed by the Lasso, are the better methods. For $N = 500$, SPLS, the univariate regression methods (including OLS) and MGL have the same performance, which is not far from the optimal Bayes error.

It is noticeable that SPLS is the only PLS method that works as well as the Lasso, and even beats it in some situations ($N = 100$ and $k_0 = 1, 2$). The good performance of the univariate regression techniques (in particular, ridge regression and the Lasso) is very likely because Ψ is diagonal, even when there is a coupling due to Ω . The poor performance of MRCE probably comes from the estimation of a full inverse covariance matrix of the response variables.

In these experiments, the performance of APLS (not shown) is not very different from that of SPLS. This is not surprising, because the synthetic data sets were generated according to equations 1.1, which correspond to the SPLS model. In the next section, we shall see an example where the adaptive version is clearly better.

6 Electrocardiogram Data Decoding

In this section we describe some experimental results on a neuroscientific data set. In particular, we aim to decode the motor output from electrocardiogram (ECoG) signals collected in monkeys. ECoG signals were recorded at a sampling rate of 1 kHz per channel, for 32 electrodes implanted in the right hemisphere, and filtered. Afterward, the signals were bandpass-filtered from 0.3 to 500 Hz. A time-frequency representation of

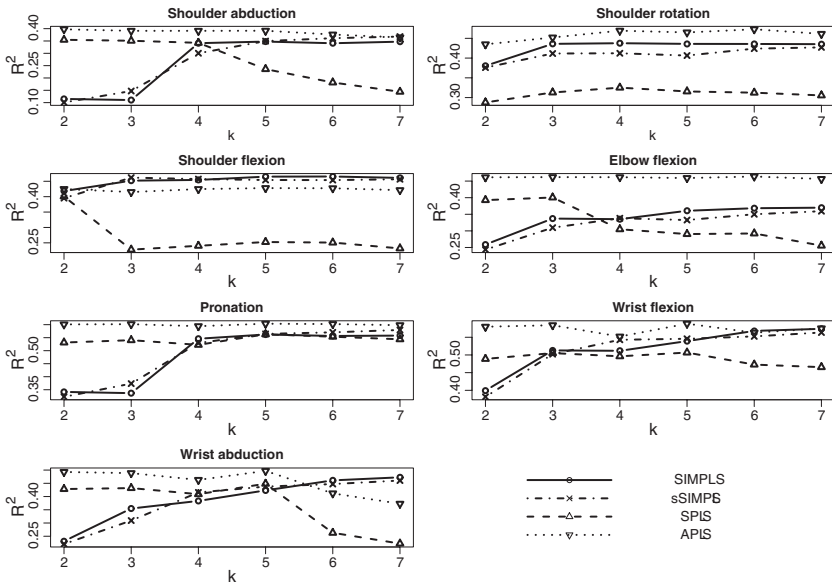


Figure 6: R^2 value for each response for a different number of latent components $k = 2, \dots, 7$.

the ECoG signals containing $p = 1600$ variables (32 electrodes, 10 frequency bins, and 5 time lags) was used as the input for the algorithms. The monkey’s movements were captured at a sampling rate of 120 Hz. Seven degrees of freedom of the monkey’s movements ($q = 7$) are to be modeled and predicted. From the available data, we have used $N = 5982$ time points for training and 5982 more for testing. More details about the data can be found in van Gerven et al. (2012).

Given the high dimensionality of the data, we have run the simplest version of the algorithms; we have factorized P and Q , and Σ^{-1} and Γ^{-1} are taken to be diagonal.

Figure 6 illustrates the performance (explained variance R^2 over testing data) of the proposed approaches compared to SIMPLS and sSIMPLS. Although not included in the plot, the results of NIPALS are almost indistinguishable from SIMPLS. It is remarkable that APLS performs better than the other methods for most of the responses and is always better than the nonadaptive algorithm. Note that APLS appears to need fewer latent components to give a reasonable estimation. Moreover, APLS is more robust to the choice of k than the other algorithms (including SPLS). Only for the wrist abduction response is the APLS accuracy clearly decreased when $k > 5$. On the other hand, SPLS performs worse in general for high values of k .

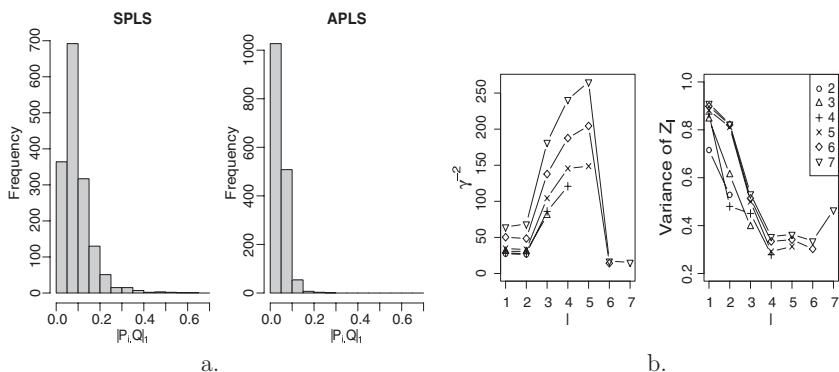


Figure 7: (a) Histogram of values $\sum_{j=1}^q |P'_{i,j} q_j|$ for SPLS and APLS. (b) Vector $\boldsymbol{\gamma}^{-2}$ and vector of individual latent component variances for APLS. Each line, labeled with a different type of symbol, represents a different run with a different maximum number of latent components.

Figure 7a shows, for SPLS and APLS, a histogram with the values of $\sum_{j=1}^q |P'_{i,j} q_j|$ as a measure of the importance of each input variable for the output prediction. (Note that σ^{-2} reflects the importance of each input for predicting the latent variables and hence is not a good measure of importance of each input variable for the output prediction.) It is worth noting that APLS yields sparser solutions than SPLS in this sense.

Figure 7b shows the values of $\boldsymbol{\gamma}^{-2}$ and the variance of the latent components for six executions of APLS, each with a different number of latent components $k = 2, \dots, 7$. Each line, labeled with a different type of symbol, corresponds to a different value k . Note that latent variables that have a high value γ_l^{-2} (left graph) or exhibit a low variance (right graph) are given less importance. From the graph, we can conclude that the first two latent components are the most important for the prediction. For example, the line with the \times symbol corresponds to $k = 5$ and has five components (symbols). Each component of this line corresponds to a latent component in the model. The lowest value of γ_l^{-2} (or the highest variance of Z_l) for this model pertains to the first two components. Hence, for this model, the two relevant components are $l = 1, 2$. For the models with $k = 6, 7$ components (whose lines have the \diamond and ∇ symbols), however, the last components have the lowest values for $\boldsymbol{\gamma}^{-2}$. The accuracy in these cases is worse than the accuracy that can be obtained with lower k values (see Figure 6). In summary, for all runs, there are two predominant latent components (either the first two or the last two), which indicates that $k = 2$ is a reasonable estimate of the optimal value of latent components.

Finally, Figure 8 demonstrates the trajectories decoded by SIMPLS and APLS. We have used $k = 7$ for SIMPLS and $k = 2$ for APLS. The two

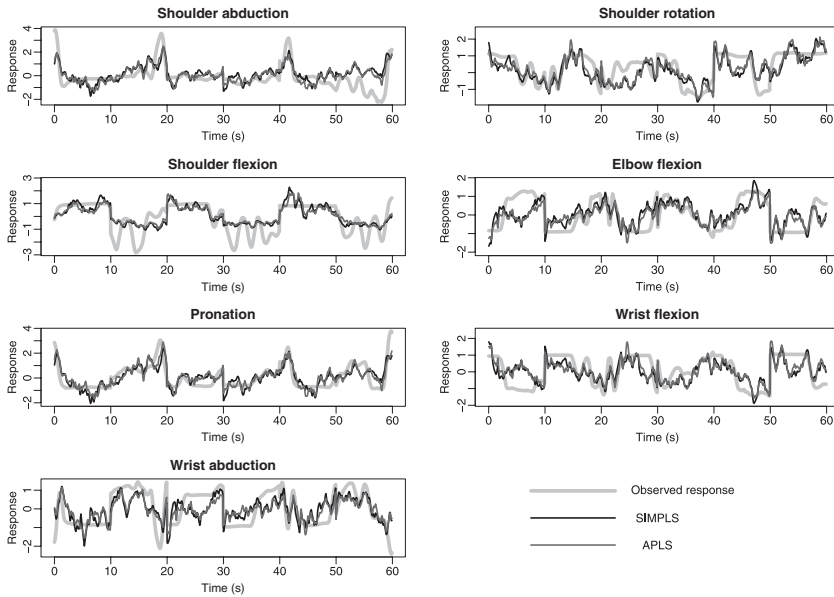


Figure 8: Observed and estimated responses in a 60 second time window, as predicted by SIMPLS and APLS. Two latent variables for APLS and seven latent variables for SIMPLS have been used.

algorithms turn out to do well. APLS is a bit smoother (less noisy) than SIMPLS. Again, the outcome for NIPALS is very similar to SIMPLS.

Table 1 illustrates the results of APLS versus OLS, ridge regression, the Lasso, MGL, and MRCE. Interestingly, the performance of APLS is higher than the other methods for all responses. Most differences are statistically significant according to a *t*-test. MGL does worse than ridge regression and the Lasso. MRCE, however, is also quite competitive.

7 Discussion

We have proposed a Bayesian formulation of PLS, with extensions for sparsity, adaptive modeling of P , and automatic determination of k , and we empirically showed that they perform well on ECoG data decoding.

The proposed approximation relies on the Bayesian paradigm, and, hence, regularization is performed in a data-driven fashion with low risk of overfitting. An advantage is interpretability of the model: using diagonal matrices $\Sigma^{-1} = \text{diag}(\sigma^{-2})$ and $\Gamma^{-1} = \text{diag}(\gamma^{-2})$, automatic relevance determination provides a measure of the relevance of each input and latent component. If, in addition, we use the adaptive extension proposed in section 4, we can obtain a reasonable estimate of the optimal number k of latent

Table 1: Mean Absolute Error (and Standard Deviations) for APLS ($k = 2$), OLS, Univariate Ridge, Univariate Lasso, Multivariate Group Lasso (MGL), and Multivariate Regression with Covariance Estimation (MRCE).

	APLS	OLS	Ridge	Lasso	MGL	MRCE
SA	0.59 (± 0.003)	0.92(± 0.004)	0.70(± 0.003)	0.77(± 0.003)	0.90(± 0.004)	0.60(± 0.003)
SF	0.64 (± 0.003)	0.92(± 0.004)	0.72(± 0.003)	0.79(± 0.003)	0.91(± 0.004)	0.72(± 0.002)
P	0.51 (± 0.002)	0.70(± 0.003)	0.56(± 0.002)	0.60(± 0.002)	0.70(± 0.003)	0.51 (± 0.002)
WA	0.57 (± 0.002)	0.84(± 0.004)	0.62(± 0.003)	0.74(± 0.003)	0.84(± 0.003)	0.68(± 0.002)
SR	0.55 (± 0.002)	0.81(± 0.003)	0.61(± 0.002)	0.70(± 0.003)	0.80(± 0.003)	0.57(± 0.002)
EF	0.53 (± 0.002)	0.84(± 0.004)	0.70(± 0.003)	0.68(± 0.003)	0.81(± 0.003)	0.66(± 0.002)
WF	0.54 (± 0.002)	0.75(± 0.003)	0.58(± 0.002)	0.65(± 0.002)	0.74(± 0.003)	0.66(± 0.002)

Notes: Each row corresponds to a motor output: shoulder abduction (SA), shoulder flexion (SF), pronation (P), wrist abduction (WA), shoulder rotation (SR), elbow flexion (EF), and wrist flexion (WF). The best method is highlighted in bold.

components from $\boldsymbol{\gamma}^{-2}$. Unlike other PLS formulations, the adaptive model appears to be robust to the choice of k .

For automatically selecting k , we can run the algorithm with several values of k and then select the one that reaches the highest model evidence (see appendix B). A more sophisticated solution is to move toward a nonparametric method, where a proper prior on \mathbf{Z} would automatically select the optimal number of latent components. However, the model would probably lose its conjugacy, so that variational inference would no longer be practicable, and we would have to resort to sampling methods of inference.

Note that up to permutation of the latent components and sign flips, the model is identifiable thanks to the priors over \mathbf{P} and \mathbf{Q} , even when neither \mathbf{Q} nor the latent components are forced to be orthonormal. Furthermore, although the model is unidentifiable with respect to permutations of the latent components, due to the initialization of \mathbf{Z} (step 1 of the algorithm in section 2), the method will always produce the same order in the latent components across different runs.

Future developments can involve a Markovian consideration of the time dynamics. By means of this extension, a connection between PLS and an input-output linear dynamical system (Beal, 2003) can be established.

Appendix A

We now formulate the variational update equations for $\boldsymbol{\Psi}^{-1}$, \mathbf{Q} , $\boldsymbol{\Gamma}^{-1}$. For a full matrix $\boldsymbol{\Psi}^{-1}$, we have a Wishart distribution,

$$F(\boldsymbol{\Psi}^{-1}) = \mathcal{W}(\boldsymbol{\Psi}^{-1}; \tilde{\mathbf{C}}^{-1}, \tilde{\kappa}), \quad (\text{A.1})$$

with $\tilde{\mathbf{C}}^{-1} = (\mathbf{Y}'\mathbf{Y} + E[\mathbf{Q}'\mathbf{Z}'\mathbf{Z}\mathbf{Q}] - \boldsymbol{\mu}'_{\mathbf{Z}}\boldsymbol{\mu}_{\mathbf{Q}} - \boldsymbol{\mu}'_{\mathbf{Q}}\boldsymbol{\mu}'_{\mathbf{Z}}\mathbf{Y} + \mathbf{C}^{-1})^{-1}$ and $\tilde{\kappa} = \kappa + N$.

For diagonal Ψ^{-1} , we have the diagonal components to be gamma distributed:

$$F(\Psi_{jj}^{-1}) = \mathcal{G}(\Psi_{jj}^{-1}; \tilde{\kappa}, \tilde{C}_{jj}^{-1}), \tag{A.2}$$

with $\tilde{\kappa} = \kappa + \frac{N}{2}$ and $\tilde{C}_{jj}^{-1} = \frac{1}{2}(\mathbf{Y}_{\cdot j}\mathbf{Y}'_{\cdot j} + E[\mathbf{q}'_j\mathbf{Z}'\mathbf{Z}\mathbf{q}_j] - 2\mathbf{Y}'_{\cdot j}\boldsymbol{\mu}_Z\boldsymbol{\mu}_{q_j}) + C_{jj}^{-1}$.

For \mathbf{Q} , we have, in the general case, a kq -dimensional gaussian distribution,

$$F(\mathbf{Q}) = \mathcal{N}(\mathbf{Q}; \boldsymbol{\mu}_Q, \mathbf{S}_Q), \tag{A.3}$$

whose parameters can be reconstructed from $\mathbf{S}_{\tilde{q}} = (E[\boldsymbol{\Gamma}^{-1}] \otimes \mathbf{I}_q + E[\mathbf{Z}'\mathbf{Z}] \otimes E[\Psi^{-1}])^{-1}$ and $\boldsymbol{\mu}_{\tilde{q}} = \mathbf{S}_{\tilde{q}}(E[\mathbf{Z}'\mathbf{Z}] \otimes E[\Psi^{-1}])\tilde{\mathbf{q}}^*$, where $\tilde{\mathbf{q}}$ is defined like $\tilde{\mathbf{p}}$ and $\tilde{\mathbf{q}}^*$ is the concatenation of the rows of $(\boldsymbol{\mu}'_Z\mathbf{X})^{-1}\boldsymbol{\mu}'_Z\mathbf{Y}$. If we choose to factorize $F(\mathbf{Q})$, which follows from taking Ψ^{-1} to be diagonal, we have

$$F(\mathbf{q}_j) = \mathcal{N}(\mathbf{q}_j; \boldsymbol{\mu}_{q_j}, \mathbf{S}_{q_j}), \tag{A.4}$$

with $\mathbf{S}_{q_j} = (E[\boldsymbol{\Gamma}^{-1}] + E[\Psi_{jj}^{-1}]E[\mathbf{Z}'\mathbf{Z}])^{-1}$ and $\boldsymbol{\mu}_{q_j} = E[\Psi_{jj}^{-1}]\mathbf{S}_{q_j}\boldsymbol{\mu}'_Z\mathbf{Y}_{\cdot j}$

With regard to $\boldsymbol{\Gamma}_j^{-1}$, we have

$$F(\boldsymbol{\Gamma}^{-1}) = \mathcal{W}(\boldsymbol{\Gamma}^{-1}; \tilde{\mathbf{D}}^{-1}, \tilde{\zeta}), \tag{A.5}$$

with $\tilde{\mathbf{D}}^{-1} = (E[\mathbf{Q}\mathbf{Q}'] + \mathbf{D}^{-1})^{-1}$ and $\tilde{\zeta} = \zeta + q$.

Appendix B

In this appendix, we derive a variational lower bound of the evidence from model 1. This can be used to monitor the inference process and check convergence. The lower bound is defined as

$$\begin{aligned} \mathcal{L} = & E[\ln P(\boldsymbol{\Sigma}^{-1})] + \sum_{l=1}^k E[\ln P(\mathbf{p}_l|\boldsymbol{\Sigma}^{-1})] + E[\ln P(\boldsymbol{\Omega}^{-1})] \\ & + \sum_{n=1}^N E[\ln P(z_n|\mathbf{x}_n, \mathbf{P}, \boldsymbol{\Omega}^{-1})] + E[\ln P(\boldsymbol{\Gamma}^{-1})] + \sum_{j=1}^q E[\ln P(\mathbf{q}_j|\boldsymbol{\Gamma}^{-1})] \\ & + E[\ln P(\Psi^{-1})] + \sum_{n=1}^N E[\ln P(\mathbf{y}_n|z_n, \mathbf{Q}, \Psi^{-1})] + \sum_{l=1}^k E[\ln F(\mathbf{p}_l)] \end{aligned}$$

$$\begin{aligned}
 &+ E[\ln F(\boldsymbol{\Omega}^{-1}, \boldsymbol{\Sigma}^{-1})] + \sum_{n=1}^N E[\ln F(z_n)] + \sum_{j=1}^q E[\ln F(q_j)] \\
 &+ E[\ln F(\boldsymbol{\Psi}^{-1}, \boldsymbol{\Gamma}^{-1})]. \tag{B.1}
 \end{aligned}$$

Particularizing for full matrices $\boldsymbol{\Omega}^{-1}$ and $\boldsymbol{\Psi}^{-1}$, we have

$$\begin{aligned}
 &E[\ln P(\boldsymbol{\Sigma}^{-1})] \\
 &= -\frac{\nu}{2} \ln |\mathbf{B}^{-1}| - \ln \left(2^{\nu p/2} \pi^{p(p-1)/4} \prod_{i=1}^p G\left(\frac{\nu+1-i}{2}\right) \right) \\
 &\quad + \frac{\nu-p-1}{2} \left(\sum_{i=1}^p \psi\left(\frac{\tilde{\nu}+1-i}{2}\right) + p \ln 2 + \ln |\tilde{\mathbf{B}}^{-1}| \right) - \frac{\tilde{\nu}}{2} \text{Tr}(\mathbf{B}\tilde{\mathbf{B}}^{-1}),
 \end{aligned}$$

$$\begin{aligned}
 &E[\ln P(\mathbf{p}_l | \boldsymbol{\Sigma}^{-1})] \\
 &= -\frac{p}{2} \ln(2\pi) + \frac{1}{2} \left(\sum_{i=1}^p \psi\left(\frac{\tilde{\nu}+1-i}{2}\right) + p \ln 2 + \ln |\tilde{\mathbf{B}}^{-1}| \right) \\
 &\quad - \frac{\tilde{\nu}}{2} \boldsymbol{\mu}'_{p_l} \tilde{\mathbf{B}}^{-1} \boldsymbol{\mu}_{p_l} - \frac{\nu+k}{2} \text{Tr}(\mathbf{S}_{p_l} \tilde{\mathbf{B}}^{-1}),
 \end{aligned}$$

$$\begin{aligned}
 &E[\ln P(\boldsymbol{\Omega}^{-1})] \\
 &= -\frac{\iota}{2} \ln |\tilde{\mathbf{A}}^{-1}| - \ln \left(2^{ik/2} \pi^{k(k-1)/4} \prod_{l=1}^k G\left(\frac{\iota+1-l}{2}\right) \right) \\
 &\quad + \frac{\iota-k-1}{2} \left(\sum_{l=1}^k \psi\left(\frac{\tilde{\iota}+1-l}{2}\right) + k \ln 2 + \ln |\tilde{\mathbf{A}}^{-1}| \right) - \frac{\tilde{\iota}}{2} \text{Tr}(\mathbf{A}\tilde{\mathbf{A}}^{-1}),
 \end{aligned}$$

$$\begin{aligned}
 &E[\ln P(z_n | \mathbf{x}_n, \mathbf{P}, \boldsymbol{\Omega}^{-1})] \\
 &= -\frac{p}{2} \ln(2\pi) + \frac{1}{2} \left(\sum_{l=1}^k \psi\left(\frac{\tilde{\iota}+1-l}{2}\right) + k \ln 2 + \ln |\tilde{\mathbf{A}}^{-1}| \right) \\
 &\quad - \frac{\tilde{\iota}}{2} \mathbf{x}'_n \left(\boldsymbol{\mu}_p \tilde{\mathbf{A}}^{-1} \boldsymbol{\mu}'_p + \sum_{l=1}^k \tilde{\mathbf{A}}_{ll}^{-1} \mathbf{S}_{p_l} \right) \mathbf{x}_n,
 \end{aligned}$$

$$\begin{aligned}
 &E[\ln P(\mathbf{y}_n | z_n, \mathbf{Q}, \boldsymbol{\Psi}^{-1})] \\
 &= -\frac{q}{2} \ln(2\pi) + \frac{1}{2} \left(\sum_{j=1}^q \psi\left(\frac{\tilde{\kappa}+1-j}{2}\right) + q \ln 2 + \ln |\tilde{\mathbf{C}}^{-1}| \right)
 \end{aligned}$$

$$\begin{aligned}
 & -\frac{\tilde{k}}{2} \boldsymbol{\mu}'_{z_n} \left(\boldsymbol{\mu}_Q \tilde{\mathbf{C}}^{-1} \boldsymbol{\mu}'_Q + \sum_{j=1}^q \tilde{\mathbf{C}}_{jj}^{-1} \mathbf{S}_{q_j} \right) \boldsymbol{\mu}_{z_n} \\
 & -\frac{\tilde{k}}{2} \text{Tr} \left(\mathbf{S}_{z_n} \left(\boldsymbol{\mu}_Q \tilde{\mathbf{C}}^{-1} \boldsymbol{\mu}'_Q + \sum_{j=1}^q \tilde{\mathbf{C}}_{jj}^{-1} \mathbf{S}_{q_j} \right) \right) - \frac{\tilde{k}}{2} \mathbf{y}'_n \tilde{\mathbf{C}}^{-1} \mathbf{y}_n + \tilde{k} \mathbf{y}'_n \tilde{\mathbf{C}}^{-1} \boldsymbol{\mu}'_Q \boldsymbol{\mu}_{z_n},
 \end{aligned}$$

where $G(\cdot)$ and $\psi(\cdot)$ are the gamma and digamma functions. The expressions for $E[\ln P(\boldsymbol{\Gamma}^{-1})]$, $E[\ln P(\mathbf{q}_j | \boldsymbol{\Gamma}^{-1})]$ and $E[\ln P(\boldsymbol{\Psi}^{-1})]$ are analogous to $E[\ln P(\boldsymbol{\Sigma}^{-1})]$, $E[\ln P(\mathbf{p}_l | \boldsymbol{\Sigma}^{-1})]$, and $E[\ln P(\boldsymbol{\Omega}^{-1})]$, respectively, and are not shown.

The rest of the terms in equation B.1 correspond to the negative entropies of the $F(\cdot)$ distributions—for example:

$$\begin{aligned}
 E[\ln F(\mathbf{p}_l)] &= \frac{1}{2} \ln |\mathbf{S}_{p_l}| + \frac{p}{2} (1 + \ln(2\pi)), \\
 E[\ln F(\boldsymbol{\Omega}^{-1}, \boldsymbol{\Sigma}^{-1})] &= -\frac{l}{2} \ln |\tilde{\mathbf{A}}^{-1}| - \ln \left(2^{0.5lk} \pi^{k(k-1)/4} \prod_{l=1}^k G\left(\frac{\tilde{l} + 1 - l}{2}\right) \right) \\
 & - \frac{\tilde{l} - k - 1}{2} \left(\sum_{l=1}^k \psi\left(\frac{\tilde{l} + 1 - l}{2}\right) + k \ln 2 + \ln |\tilde{\mathbf{A}}^{-1}| \right) + \frac{\tilde{l}k}{2} \\
 & - \frac{\nu}{2} \ln |\tilde{\mathbf{B}}^{-1}| - \ln \left(2^{0.5\nu p} \pi^{p(p-1)/4} \prod_{i=1}^p G\left(\frac{\tilde{\nu} + 1 - i}{2}\right) \right) \\
 & - \frac{\tilde{\nu} - p - 1}{2} \left(\sum_{i=1}^p \psi\left(\frac{\tilde{\nu} + 1 - i}{2}\right) + p \ln 2 + \ln |\tilde{\mathbf{B}}^{-1}| \right) + \frac{\tilde{\nu}p}{2}, \\
 E[\ln F(\mathbf{z}_n)] &= \frac{1}{2} \ln |\mathbf{S}_{z_n}| + \frac{k}{2} (1 + \ln(2\pi)).
 \end{aligned}$$

The variational lower bound for other variations of the method can be easily computed following the same line of argument.

Acknowledgments

This work has been partially supported by projects TIN2010-20900-C04-04 and Cajal Blue Brain of the Spanish Ministry of Science and Innovation (MINECO).

References

- Ando, R., & Zhang, T. (2005). A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6, 1817–1853.
- Argyriou, A., Evgeniou, T., & Pontil, M. (2006). Multi-task feature learning. In B. Schölkopf, B. Platt, & T. Hoffmann (Eds.), *Advances in neural information processing systems*, 19 (pp. 41–48). Cambridge, MA: MIT Press.
- Barker, M., & Rayens, W. (2003). Partial least squares for discrimination. *Journal of Chemometrics*, 17, 166–173.
- Beal, M. J. (2003). *Variational algorithms for approximate Bayesian inference*. Unpublished doctoral dissertation, University College London.
- Bishop, C. (1998). Bayesian principal components. In M. S. Keams, S. Sola, & D. A. Cohn (Eds.), *Advances in neural information processing systems*, 11. Cambridge, MA: MIT Press.
- Brown, P. J., & Zidek, J. V. (1980). Adaptive multivariate ridge regression. *Annals of Statistics*, 8, 64–74.
- Chun, H., & Keleş, S. (2010). Sparse partial least squares regression for simultaneous dimension reduction and variable selection. *Journal of the Royal Statistical Society: Series B*, 72, 3–25.
- de Jong, S. (1993). SIMPLS: An alternative approach to partial least squares regression. *Chemometrics and Intelligent Laboratory Systems*, 18, 251–263.
- Frank, I. E., & Friedman, J. H. (1993). A statistical view of some chemometrics regression tools. *Technometrics*, 35, 109–135.
- Fujiwara, Y., Miyawaki, Y., & Kamitani, Y. (2009). Estimating image bases for visual image reconstruction from human brain activity. In Y. Bengio, P. Schuurmans, J. Lafferty, C. Williams, & A. Culotta (Eds.), *Advances in neural information processing systems*, 22. Cambridge, MA: MIT Press.
- Ghahramani, Z., & Beal, M. J. (2000). Variational inference for Bayesian mixtures of factor analysers. In S. A. Solla, T. K. Leen, & K.-R. Müller (Eds.), *Advances in neural information processing systems*, 12. Cambridge, MA: MIT Press.
- Goutis, C. (1996). Partial least squares yields shrinkage estimators. *Annals of Statistics*, 24, 816–824.
- Hardoon, D. R., Szedmak, S., & Shawe-Taylor, J. (2004). Canonical correlation analysis: An overview with application to learning methods. *Neural Computation*, 16, 2639–2664.
- Hastie, T., Tibshirani, R., & Friedman, J. (2008). *The elements of statistical learning: Data mining, inference, and predictions* (2nd ed.). New York: Springer.
- Izenman, A. (1975). Reduced-rank regression for the multivariate linear model. *Journal of Multivariate Analysis*, 5, 248–264.
- Jaakkola, T. S. (2001). Tutorial on variational approximation methods. In M. Opper & D. Saad (Eds.), *Advanced mean field methods: Theory and practice*. Cambridge, MA: MIT Press.
- Lindgren, F., Geladi, P., & Wold, S. (1993). The kernel algorithm for PLS. *Journal of Chemometrics*, 7, 45–59.
- Nakajima, S., Sugiyama, M., & Babacan, D. (2011). On Bayesian PCA: Automatic dimensionality selection and analytic solution. In *37th International Conference on Machine Learning*. International Machine Learning Society.

- Obozinski, G., Wainwright, M. J., & Jordan, M. I. (2011). Support union recovery in high-dimensional multivariate regression. *Annals of Statistics*, 39, 1–47.
- Rosipal, R., & Krämer, N. (2006). Overview and recent advances in partial least squares. *Lecture Notes in Computer Science*, 3940, 34–51.
- Rosipal, R., & Trejo, L. J. (2001). Kernel partial least squares regression in reproducing kernel Hilbert space. *Journal of Machine Learning Research*, 2, 97–123.
- Rothman, A. J., Levina, E., & Zhu, J. (2010). Sparse multivariate regression with covariance estimation. *Journal of Computational and Graphical Statistics*, 19, 947–962.
- van Gerven, M. A. J., Chao, Z. C., & Heskes, T. (2012). On the decoding of intracranial data using sparse orthonormalized partial least squares. *Journal of Neural Engineering*, 9, 026017.
- Virtanen, S., Klami, A., & Kaski, S. (2011). Bayesian CCA via group sparsity. In *37th International Conference on Machine Learning*. International Machine Learning Society.
- Wang, C. (2007). Variational Bayesian approach to canonical correlation analysis. *IEEE Transactions on Neural Networks*, 18, 905–910.
- Wangen, L. E., & Kowalsky, B. R. (1989). A multiblock partial least squares algorithm for investigating complex chemical systems. *Journal of Chemometrics*, 3, 3–20.
- Wold, H. (1975). Path models with latent variables: The NIPALS approach. In H. M. Blalock (Ed.), *Quantitative sociology: International perspectives on mathematical and statistical model building* (pp. 307–357). New York: Academic Press.
- Wold, S., Sjöström, M., & Eriksson, L. (2001). PLS-regression: A basic tool of chemometrics. *Chemometrics and Intelligent Laboratory Systems*, 58, 109–130.
- Yuan, M., Ekici, A., Lu, Z., & Monteiro, R. (2007). Dimension reduction and coefficient estimation in multivariate linear regression. *Journal of the Royal Statistical Society: Series B*, 69, 329–346.