



## Combining feature spaces for classification

Theodoros Damoulas\*, Mark A. Girolami

*Inference Research Group, Department of Computing Science, Faculty of Information and Mathematical Sciences, University of Glasgow, 18 Lilybank Gardens, Glasgow G12 8QQ, Scotland, UK*

### ARTICLE INFO

#### Article history:

Received 2 November 2007  
Received in revised form 24 January 2009  
Accepted 5 April 2009

#### Keywords:

Variational Bayes approximation  
Multiclass classification  
Kernel combination  
Hierarchical Bayes  
Bayesian inference  
Ensemble learning  
Multi-modal modelling  
Information integration

### ABSTRACT

In this paper we offer a variational Bayes approximation to the multinomial probit model for basis expansion and kernel combination. Our model is well-founded within a hierarchical Bayesian framework and is able to instructively combine available sources of information for multinomial classification. The proposed framework enables informative integration of possibly heterogeneous sources in a multitude of ways, from the simple summation of feature expansions to weighted product of kernels, and it is shown to match and in certain cases outperform the well-known ensemble learning approaches of combining individual classifiers. At the same time the approximation reduces considerably the CPU time and resources required with respect to both the ensemble learning methods and the full Markov chain Monte Carlo, Metropolis–Hastings within Gibbs solution of our model. We present our proposed framework together with extensive experimental studies on synthetic and benchmark datasets and also for the first time report a comparison between summation and product of individual kernels as possible different methods for constructing the composite kernel matrix.

© 2009 Elsevier Ltd. All rights reserved.

### 1. Introduction

Classification, or supervised discrimination, has been an active research field within the pattern recognition and machine learning communities for a number of decades. The interest comes from the nature of the problems encountered within the communities, a large number of which can be expressed as discrimination tasks, i.e. to be able to learn and/or predict in which category an object belongs. In the supervised and semi-supervised learning setting, [24,6], where labeled (or partially labeled) training sets are available in order to predict labels for unseen objects, we are in the classification domain and in the unsupervised learning case, where the machine is asked to learn categories from the available unlabeled data, we fall in the domain of clustering.

In both of these domains a recurring theme is how to utilize constructively all the information that is available for the specific problem in hand. This is specifically true when a large number of features from different sources are available for the same object and there is limited or no a priori knowledge of their significance and contribution to the classification or clustering task. In such cases, concatenating all the features into a single feature space does not

guarantee an optimum performance and it exacerbates the “curse of dimensionality” problem.

The problem of classification in the case of multiple feature spaces or sources has been mostly dealt with in the past by ensemble learning methods [12], namely combinations of individual classifiers. The idea behind that approach is to train one classifier in every feature space and then combine their class predictive distributions. Different ways of combining the output of the classifiers have been studied [18,34,19] and also meta-learning an overall classifier on these distributions [36] has been proposed.

The drawbacks of combining classifiers lie on their theoretical justification, on the processing loads incurred as multiple training has to be performed, and on the fact that the individual classifiers operate independently on the data. Their performance has been shown to significantly improve over the best individual classifier in many cases but the extra load of training multiple classifiers may possibly restrain their application when resources are limited. The typical combination rules for classifiers are ad hoc methods [5] that are based on the notions of the *linear* or *independent* opinion pools. However, as noted by Berger [5], “... it is usually better to approximate a “correct answer” than to adopt a completely ad hoc approach.” and this has a bearing, as we shall see, on the nature of classifier combination and the motivation for kernel combination.

In this paper we offer, for the classification problem, a multi-class kernel machine [30,31] that is able to combine kernel spaces in an informative manner while at the same time learning their significance. The proposed methodology is general and can also be

\* Corresponding author. Tel.: +44 141 3302421; fax: +44 141 3308627.

E-mail addresses: [theo@dcs.gla.ac.uk](mailto:theo@dcs.gla.ac.uk) (T. Damoulas), [girolami@dcs.gla.ac.uk](mailto:girolami@dcs.gla.ac.uk) (M.A. Girolami).

employed outside the “kernel-domain”, i.e. without the need to embed the features into Hilbert spaces, by allowing for combination of basis function expansions. That is useful in the case of multi-scale problems and wavelets [3]. The combination of kernels or dictionaries of basis functions as an alternative to classifier combination offers the advantages of reduced computational requirements, the ability to learn the significance of the individual sources and an improved solution based on the inferred significance.

Previous related work, includes the area of *kernel learning*, where various methodologies have been proposed in order to learn the possibly composite kernel matrix, by convex optimization methods [20,33,22,23], by the introduction of hyper-kernels [28] or by direct ad hoc construction of the composite kernel [21]. These methods operate within the context of support vector machine (SVM) learning and hence inherit the arguable drawbacks of non-probabilistic outputs and ad hoc extensions to the multiclass case. Furthermore, another related approach within the non-parametric Gaussian process (GP) methodology [27,29] has been proposed by Girolami and Zhong [16], where instead of kernel combination the integration of information is achieved via combination of the GP covariance functions.

Our work further develops the work of Girolami and Rogers [14], which we generalize to the multiclass case by employing a multinomial probit likelihood and introducing latent variables that give rise to efficient Gibbs sampling from the parameter posterior distribution. We bound the marginal likelihood or model evidence [25] and derive the variational Bayes (VB) approximation for the multinomial probit composite kernel model, providing a fast and efficient solution. We are able to combine kernels in a general way e.g. by summation, product or binary rules and learn the associated weights to infer the significance of the sources.

For the first time we offer, a VB approximation on an explicit multiclass model for kernel combination and a comparison between kernel combination methods, from summation to weighted product of kernels. Furthermore, we compare our methods against classifier combination strategies on a number of datasets and we show that the accuracy of our VB approximation is comparable to that of the full Markov chain Monte Carlo (MCMC) solution.

The paper is organized as follows. First we introduce the concepts of composite kernels, kernel combination rules and classifier combination strategies and give an insight on their theoretical underpinnings. Next we introduce the multinomial probit composite kernel model, the MCMC solution and the VB approximation. Finally, we present results on synthetic and benchmark datasets and offer discussion and concluding remarks.

## 2. Classifier versus kernel combination

The theoretical framework behind classifier combination strategies has been explored relatively recently<sup>1</sup> [12,18,19], despite the fact that ensembles of classifiers have been widely used experimentally before. We briefly review that framework in order to identify the common ground and the deviations from the classifier combination to the proposed kernel (or basis function) combination methodology.

Consider  $S$  feature spaces or sources of information.<sup>2</sup> An object  $n$  belonging to a dataset  $\{\mathbf{X}, \mathbf{t}\}$ , with  $\mathbf{X}$  the collection of all objects and  $\mathbf{t}$  the corresponding labels, is represented by  $S, D^s$ -dimensional feature vectors  $\mathbf{x}_n^s$  for  $s = 1, \dots, S$  and  $\mathbf{x}_n^s \in \mathfrak{R}^{D^s}$ . Let the number of classes be

<sup>1</sup> At least in the classification domain for machine learning and pattern recognition. The related statistical decision theory literature has studied opinion pools much earlier, see Berger [5] for a full treatment.

<sup>2</sup> Throughout this paper  $m$  denotes scalar,  $\mathbf{m}$  vector and  $\mathbf{M}$  a Matrix. If  $\mathbf{M}$  is a  $C \times N$  matrix then  $\mathbf{m}_c$  is the  $c$ th row vector,  $\mathbf{m}_n$  the  $n$ th column vector of that matrix and all other indices imply row vectors.

$C$  with the target variable  $t_n = c = 1, \dots, C$  and the number of objects  $N$  with  $n = 1, \dots, N$ . Then, the class posterior probability for object  $n$  will be  $P(t_n|\mathbf{x}_n^1, \dots, \mathbf{x}_n^S)$  and according to Bayesian decision theory we would assign the object  $n$  with the class that has the maximum a posteriori probability.

The classifier combination methodologies are based on the following assumption, which as we shall see is not necessary for the proposed kernel combination approach. The assumption is that although “it is essential to compute the probabilities of various hypotheses by considering all the measurements simultaneously ... it may not be a practical proposition.” [18] and it leads to a further approximation on the (noisy) class posterior probabilities which are now broken down to individual contributions from classifiers trained on each feature space  $s$ .

The *product combination* rule, assuming equal a priori class probabilities:

$$P_{\text{prod}}(t_n|\mathbf{x}_n^1, \dots, \mathbf{x}_n^S) = \frac{\prod_{s=1}^S (P(t_n|\mathbf{x}_n^s) + \varepsilon^s)}{\sum_{c=1}^C \prod_{s=1}^S (P(t_n|\mathbf{x}_n^s) + \varepsilon^s)} \quad (1)$$

The *mean combination* rule:

$$P_{\text{mean}}(t_n|\mathbf{x}_n^1, \dots, \mathbf{x}_n^S) = \frac{1}{S} \sum_{s=1}^S (P(t_n|\mathbf{x}_n^s) + \varepsilon^s) \quad (2)$$

with  $\varepsilon^s$  the prediction error made by the individual classifier trained on feature space  $s$ .

The theoretical justification for the product rule comes from the *independence* assumption of the feature spaces, where  $\mathbf{x}_n^1, \dots, \mathbf{x}_n^S$  are assumed to be uncorrelated; and the mean combination rule is derived on the opposite assumption of extreme correlation. As it can be seen from Eqs. (1) and (2), the individual errors  $\varepsilon^s$  of the classifiers are either added or multiplied together and the hope is that there will be a synergetic effect from the combination that will cancel them out and hence reduce the overall classification error.

Instead of combining classifiers we propose to combine the feature spaces and hence obtain the class posterior probabilities from the model

$$P_{\text{spaces}}(t_n|\mathbf{x}_n^1, \dots, \mathbf{x}_n^S) = P(t_n|\mathbf{x}_n^1, \dots, \mathbf{x}_n^S) + \varepsilon$$

where now we only have one error term  $\varepsilon$  from the single classifier operating on the composite feature space. The different methods to construct the composite feature space reflect different approximations to the  $P(t_n|\mathbf{x}_n^1, \dots, \mathbf{x}_n^S)$  term without incorporating individual classifier errors into the approximation. Still though, underlying assumptions of independence for the construction of the composite feature space are implied through the kernel combination rules and as we shall see this will lead to the need for *diversity* as it is commonly known for other ensemble learning approaches.

Although our approach is general and can be applied to combine basis expansions or kernels, we present here the composite kernel construction as our main example. Embedding the features into Hilbert spaces via the kernel trick [30,31] we can define<sup>3</sup> :

The  $N \times N$  *mean composite* kernel as

$$\mathbf{K}^{\beta\Theta}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{s=1}^S \beta_s K^{s\Theta_s}(\mathbf{x}_i^s, \mathbf{x}_j^s) \quad \text{with} \quad \sum_{s=1}^S \beta_s = 1 \quad \text{and} \quad \beta_s \geq 0 \quad \forall s$$

<sup>3</sup> Superscripts denote “function of”, i.e.  $\mathbf{K}^{\beta\Theta}$  denotes that the kernel  $\mathbf{K}$  is a function of  $\beta$  and  $\Theta$ , unless otherwise specified.

The  $N \times N$  product composite kernel<sup>4</sup> as

$$\mathbf{K}^{\beta\Theta}(\mathbf{x}_i, \mathbf{x}_j) = \prod_{s=1}^S K^{s\theta_s}(\mathbf{x}_i^s, \mathbf{x}_j^s)^{\beta_s} \quad \text{with } \beta_s \geq 0 \quad \forall s$$

And the  $N \times N$  binary composite kernel as

$$\mathbf{K}^{\beta\Theta}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{s=1}^S \beta_s K^{s\theta_s}(\mathbf{x}_i^s, \mathbf{x}_j^s) \quad \text{with } \beta_s \in \{0, 1\} \quad \forall s$$

where  $\Theta$  are the kernel parameters,  $\beta$  the combinatorial weights and  $K$  is the kernel function employed, typically in this study, a Gaussian or polynomial function. We can now weigh accordingly the feature spaces by the parameter  $\beta$  which, as we shall see further on, can be inferred from the observed evidence. Learning these combinatorial parameters is the main objective of multiple kernel learning methods such as the proposed ones in this article.

### 3. The multinomial probit model

In accordance with Albert and Chib [1], we introduce auxiliary variables  $\mathbf{Y} \in \mathfrak{R}^{C \times N}$  that we regress onto with our composite kernel and the parameters (regressors)  $\mathbf{W} \in \mathfrak{R}^{C \times N}$ . The intuition is that the regressors express the weight with which a data point “votes” for a specific class  $c$  and the auxiliary variables are continuous “target” values that are related to class membership ranking for a specific point. Following the standardized noise model  $\varepsilon \sim \mathcal{N}(0, 1)$  which results in  $y_{cn} = \mathbf{w}_c \mathbf{k}_n^{\beta\Theta} + \varepsilon$ , with  $\mathbf{w}_c$  the  $1 \times N$  row vector of class  $c$  regressors and  $\mathbf{k}_n^{\beta\Theta}$  the  $N \times 1$  column vector of inner products for the  $n$ th element, leads to the following Gaussian probability distribution:

$$p(y_{cn} | \mathbf{w}_c, \mathbf{k}_n^{\beta\Theta}) = \mathcal{N}_{y_{cn}}(\mathbf{w}_c \mathbf{k}_n^{\beta\Theta}, 1) \quad (3)$$

The link from the auxiliary variable  $y_{cn}$  to the discrete target variable of interest  $t_n \in \{1, \dots, C\}$  is given by  $t_n = i \iff y_{in} > y_{jn} \quad \forall j \neq i$  and by the following marginalization  $P(t_n = i | \mathbf{W}, \mathbf{k}_n^{\beta\Theta}) = \int P(t_n = i | \mathbf{y}_n) P(\mathbf{y}_n | \mathbf{W}, \mathbf{k}_n^{\beta\Theta}) d\mathbf{y}_n$ , where  $P(t_n = i | \mathbf{y}_n)$  is a delta function, results in the multinomial probit likelihood as

$$P(t_n = i | \mathbf{W}, \mathbf{k}_n^{\beta\Theta}) = \mathcal{E}_{p(u)} \left\{ \prod_{j \neq i} \Phi(u + (\mathbf{w}_i - \mathbf{w}_j) \mathbf{k}_n^{\beta\Theta}) \right\} \quad (4)$$

where  $\mathcal{E}$  is the expectation taken with respect to the standardized normal distribution  $p(u) = \mathcal{N}(0, 1)$  and  $\Phi$  is the cumulative density function. We can now consider the prior and hyper-prior distributions to be placed on our model.

#### 3.1. Prior probabilities

Having derived the multinomial probit likelihood we complete our Bayesian model by introducing prior distributions on the model parameters. The choice of prior distributions is justified as follows: for the regressors  $\mathbf{W}$  we place a product of zero mean Gaussian distributions with scale  $\zeta_{cn}$  that reflects independence (product) and lack of prior knowledge (zero mean normal). The only free parameter on this distribution is the scale, on which we place a gamma prior distribution that ensures a positive value, is a conjugate pair with the normal and is controlled via the hyper-parameters  $\tau, \nu$ . This prior setting propagates the uncertainty to a higher level and it is commonly employed in the statistics literature [11] for the above-mentioned reasons of conjugacy and (lack of) prior knowledge.

<sup>4</sup> In this case only the superscript  $\beta_s$  in  $K^{s\theta_s}(\mathbf{x}_i^s, \mathbf{x}_j^s)^{\beta_s}$  denotes power. The superscript  $s$  in  $K^s(\dots)$  or generally in  $\mathbf{K}^s$  indicates that this is the  $s$ th base kernel.

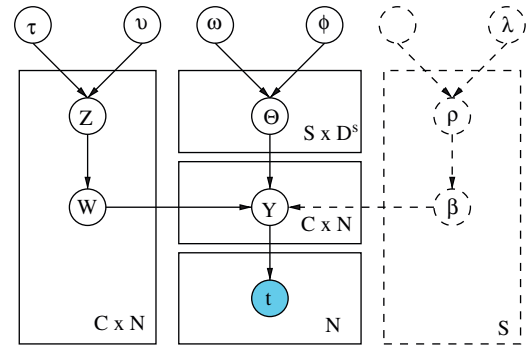


Fig. 1. Plates diagram of the model for mean composite kernel.

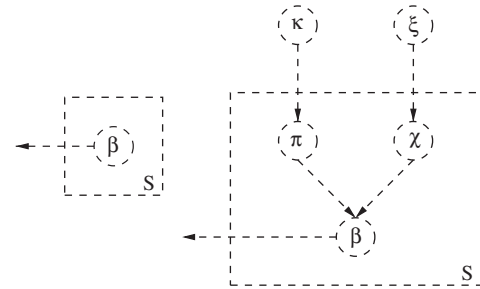


Fig. 2. Modification for binary composite (left) and product composite kernel (right).

The hyper-parameter prior setting of  $\tau, \nu$  dictates the prior form of the gamma distribution and can be set to ensure an uninformative distribution (as it will be done in this work) or we can induce sparsity via setting a flat prior ( $\tau, \nu \rightarrow 0$ ) as in Damoulas et al. [9]. This leads to the empirical Bayes method of type-II maximum likelihood as employed in the construction of relevant vector machines [35].

Furthermore, we place a gamma distribution on each kernel parameter since  $\theta_{sd} \in \mathfrak{R}^+$ . In the case of the mean composite kernel, a Dirichlet distribution with parameters  $\rho$  is placed on the combinatorial weights in order to satisfy the constraints imposed on the possible values which are defined on a simplex and assists in statistical identifiability. A further gamma distribution, again to ensure positive values, is placed on each  $\rho_s$  with associated hyper-parameters  $\mu, \lambda$ .

In the product composite kernel case we employ the right dashed plate in Fig. 2 which places a gamma distribution on the combinatorial weights  $\beta$ , that do not need to be defined on a simplex anymore, with an exponential hyper-prior distribution on each of the parameters  $\pi_s, \chi_s$ .

Finally, in the binary composite kernel case we employ the left dashed plate in Fig. 2 which places a binomial distribution on each  $\beta_s$  with equal probability of being 1 or zero (unless prior knowledge says otherwise). The small size of the possible  $2^S$  states of the  $\beta$  vector allows for their explicit consideration in the inference procedure and hence there is no need to place any hyper-prior distributions.

The model, for case of the mean composite kernel, is depicted graphically in Fig. 1 where the conditional relation of the model parameters and associated hyper-parameters can be seen. The accompanied variations for the binary and product composite kernel are given in Fig. 2.

The intuition behind the hierarchical construction of the proposed model is that uncertainty is propagated into a higher level of prior distribution setting. In other words, we place further hyper-prior distributions on parameters of the prior distributions and we let the evidence guide us to the required posteriors. In that way the extra levels allow less sensitivity for initial settings (in the case where

fixed values are used for the hyper-parameters instead of empirical Bayes) since a specific hyper-parameter value will still lead to a (prior) *distribution* of possible parameter values which can still be “corrected” by the evidence of the data towards the appropriate final posterior distribution.

### 3.2. Exact Bayesian inference: Gibbs sampling

The standard statistical solution for performing Bayesian inference is via Markov chain Monte Carlo sampling methods that are powerful but computationally expensive. These approaches still approximate the desired posterior distribution but only due to the limited number of drawn samples, hence the name exact inference. From these methods, Gibbs sampling is the preferred alternative as it avoids the need for tuning acceptance ratios as is typical in the original Metropolis–Hastings (MH) algorithms [17]. Such a sampling scheme is only possible when the hard-to-sample-from joint posterior distribution can be broken down to conditional distributions from which we can easily draw samples.

In our model such a sampler naturally arises by exploiting the conditional distributions of  $\mathbf{Y}|\mathbf{W}, \dots$  and  $\mathbf{W}|\mathbf{Y}, \dots$ , so that we can draw samples from the parameter posterior distribution  $P(\Psi|\mathbf{t}, \mathbf{X}, \Xi)$  where  $\Psi = \{\mathbf{Y}, \mathbf{W}, \beta, \Theta, \mathbf{Z}, \rho\}$  and  $\Xi$  the aforementioned hyper-parameters of the model. In the case of non-conjugate pair of distributions, Metropolis–Hastings sub-samplers can be employed for posterior inference. More details of the Gibbs sampler are provided analytically in [8] and summarized in Appendix A.

### 3.3. Approximate inference: variational approximation

Exact inference via the Gibbs sampling scheme is computationally expensive and in this article we derive a deterministic variational approximation which is an efficient method with comparable accuracy to the full sampling scheme as it will be further demonstrated. The variational methodology, see Beal [4] for a recent treatment, offers a lower bound on the model evidence using an ensemble of factored posteriors to approximate the joint parameter posterior distribution. Although the factored ensemble implies independence of the approximate posteriors, which is a typical strong assumption of variational methods, there is weak coupling through the current estimates of parameters as it can be seen below.

Considering the joint likelihood of the model<sup>5</sup> defined as  $p(\mathbf{t}, \Psi|\mathbf{X}, \Xi) = p(\mathbf{t}|\mathbf{Y})p(\mathbf{Y}|\mathbf{W}, \beta, \Theta)p(\mathbf{W}|\mathbf{Z})p(\mathbf{Z}|\tau, v)p(\beta|\rho)p(\Theta|\omega, \phi)p(\rho|\mu, \lambda)$  and the factorable ensemble approximation of the required posterior  $p(\Psi|\Xi, \mathbf{X}, \mathbf{t}) \approx Q(\Psi) = Q(\mathbf{Y})Q(\mathbf{W})Q(\beta)Q(\Theta)Q(\mathbf{Z})Q(\rho)$  we can bound the model evidence using Jensen’s inequality:

$$\log p(\mathbf{t}) \geq \mathcal{E}_{Q(\Psi)}\{\log p(\mathbf{t}, \Psi|\Xi)\} - \mathcal{E}_{Q(\Psi)}\{\log Q(\Psi)\} \quad (5)$$

and minimize it with distributions  $Q(\Psi_i) \propto \exp(\mathcal{E}_{Q(\Psi_{-i})}\{\log p(\mathbf{t}, \Psi|\Xi)\})$  where  $Q(\Psi_{-i})$  is the factorable ensemble with the  $i$ th component removed. The above standard steps describe the adoption of the *mean-field* or variational approximation theory for our specific model.

The resulting approximate posterior distributions are given below with full details of the derivations in Appendix B. First, the approximate posterior over the auxiliary variables is given by

$$Q(\mathbf{Y}) \propto \prod_{n=1}^N \delta(y_{i,n} > y_{k,n} \forall k \neq i) \delta(t_n = i) \mathcal{N}_{\mathbf{y}_n}(\tilde{\mathbf{W}}\mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}}, \mathbf{1}) \quad (6)$$

<sup>5</sup> The *mean composite* kernel is considered as an example. Modifications for the *binary* and *product composite* kernel are trivial and respective details are given in the appendices.

which is a product of  $N$   $C$ -dimensional conically truncated Gaussians demonstrating independence across samples as expected from our initial i.i.d assumption. The shorthand tilde notation denotes posterior expectations in the usual manner, i.e.  $f(\tilde{\beta}) = \mathcal{E}_{Q(\beta)}\{f(\beta)\}$ , and the posterior expectations (details in Appendix C) for the auxiliary variable follow as

$$\tilde{y}_{cn} = \tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} - \frac{\mathcal{E}_{p(u)}\{\mathcal{N}_u(\tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} - \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}}, 1)\Phi_u^{n,i,c}\}}{\mathcal{E}_{p(u)}\{\Phi(u + \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} - \tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}})\Phi_u^{n,i,c}\}} \quad (7)$$

$$\tilde{y}_{in} = \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} - \left( \sum_{c \neq i} \tilde{y}_{cn} - \tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} \right) \quad (8)$$

where  $\Phi$  is the standardized cumulative distribution function (CDF) and  $\Phi_u^{n,i,c} = \prod_{j \neq i,c} \Phi(u + \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} - \tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}})$ . Next, the approximate posterior for the regressors can be expressed as

$$Q(\mathbf{W}) \propto \prod_{c=1}^C \mathcal{N}_{\mathbf{w}_c}(\tilde{\mathbf{y}}_c \mathbf{K}^{\tilde{\beta}\tilde{\Theta}} \mathbf{V}_c, \mathbf{V}_c) \quad (9)$$

where the covariance is defined as

$$\mathbf{V}_c = \left( \sum_{i=1}^S \sum_{j=1}^S \tilde{\beta}_i \tilde{\beta}_j \mathbf{K}^{i\tilde{\theta}_i} \mathbf{K}^{j\tilde{\theta}_j} + (\tilde{\mathbf{Z}}_c)^{-1} \right)^{-1} \quad (10)$$

and  $\tilde{\mathbf{Z}}_c$  is a diagonal matrix of the expected variances  $\tilde{\zeta}_1, \dots, \tilde{\zeta}_N$  for each class. The associated posterior mean for the regressors is therefore  $\tilde{\mathbf{w}}_c = \tilde{\mathbf{y}}_c \mathbf{K}^{\tilde{\beta}\tilde{\Theta}} \mathbf{V}_c$  and we can see the coupling between the auxiliary variable and regressor posterior expectation.

The approximate posterior for the variances  $\mathbf{Z}$  is an updated product of inverse-gamma distributions (gamma on the scales) and the posterior mean for the scale is given by

$$\frac{\tau + \frac{1}{2}}{v + \frac{1}{2} \mathbf{w}_{cn}^2} \quad (11)$$

for details see Appendix B.3 or Denison et al. [11]. Finally, the approximate posteriors for the kernel parameters  $Q(\Theta)$ , the combinatorial weights  $Q(\beta)$  and the associated hyper-prior parameters  $Q(\rho)$ , or  $Q(\pi)$ ,  $Q(\chi)$  in the product composite kernel case, can be obtained by importance sampling [2] in a similar manner to Girolami and Rogers [15] since no tractable analytical solution can be offered. Details are provided in Appendix B.

Having described the approximate posterior distributions of the parameters and hence obtained the posterior expectations we turn back to our original task of making class predictions  $\mathbf{t}^*$  for  $N_{test}$  new objects  $\mathbf{X}^*$  that are represented by  $S$  different information sources  $\mathbf{X}^{s*}$  embedded into Hilbert spaces as base kernels  $\mathbf{K}^{s\theta_s, \beta_s}$  and combined into a composite *test* kernel  $\mathbf{K}^{*\Theta, \beta}$ . The predictive distribution for a single new object  $\mathbf{x}^*$  is given by  $p(t^* = c|\mathbf{x}^*, \mathbf{X}, \mathbf{t}) = \int p(t^* = c|\mathbf{y}^*)p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{t})d\mathbf{y}^* = \int \delta_c^* p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{t})d\mathbf{y}^*$  which ends up, see Appendix D for complete derivation, as

$$p(t^* = c|\mathbf{x}^*, \mathbf{X}, \mathbf{t}) = E_{p(u)} \left\{ \prod_{j \neq c} \Phi \left[ \frac{1}{v_j^*} (u \tilde{v}_c^* + \tilde{m}_c^* - \tilde{m}_j^*) \right] \right\} \quad (12)$$

where, for the general case of  $N_{test}$  objects,  $\tilde{\mathbf{m}}_c^* = \tilde{\mathbf{y}}_c \mathbf{K}(\mathbf{K}^* \mathbf{K}^{*T} + \mathbf{V}_c^{-1})^{-1} \mathbf{K}^* \mathcal{V}_c^*$  and  $\mathcal{V}_c^* = (\mathbf{I} + \mathbf{K}^{*T} \mathbf{V}_c \mathbf{K}^*)$  while we have dropped the notation for the dependance of the train  $\mathbf{K}(N \times N)$  and test  $\mathbf{K}^*(N \times N_{test})$  kernels on  $\Theta, \beta$  for clarity. In Algorithm 1 we summarize the VB approximation in a pseudo-algorithmic fashion.

**Algorithm 1.** VB multinomial probit composite kernel regression

- 1: Initialize  $\Xi$ , sample  $\Psi$ , create  $\mathbf{K}_s|\beta_s, \theta_s$  and hence  $\mathbf{K}|\beta, \Theta$
- 2: **while** Lower Bound changing **do**
- 3:  $\tilde{\mathbf{w}}_c \leftarrow \tilde{\mathbf{y}}_c \mathbf{K} \mathbf{V}_c$
- 4:  $\tilde{\mathbf{y}}_{cn} \leftarrow \tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} - \frac{\mathcal{E}_{p(u)}\{\mathcal{N}_u(\tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} - \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}}, 1)\Phi_u^{n,i,c}\}}{\mathcal{E}_{p(u)}\{\Phi(u + \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} - \tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}})\Phi_u^{n,i,c}\}}$
- 5:  $\tilde{\mathbf{y}}_{in} \leftarrow \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} - \left(\sum_{j \neq i} \tilde{\mathbf{y}}_{jn} - \tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}}\right)$
- 6:  $\tilde{\zeta}_{cn}^{-1} \leftarrow \frac{\tau + \frac{1}{2}}{\nu + \frac{1}{2} \tilde{\mathbf{w}}_{cn}^2}$
- 7:  $\tilde{\rho}, \tilde{\beta}, \tilde{\Theta} \leftarrow \tilde{\rho}, \tilde{\beta}, \tilde{\Theta}|\tilde{\mathbf{w}}_c, \tilde{\mathbf{y}}_n$  by importance sampling
- 8: Update  $\mathbf{K}|\beta, \Theta$  and  $\mathbf{V}_c$
- 9: **end while**
- 10: Create composite test kernel  $\mathbf{K}^*|\tilde{\beta}, \tilde{\Theta}$
- 11:  $\tilde{\mathcal{V}}_c^* \leftarrow (\mathbf{I} + \mathbf{K}^* \mathbf{V}_c \mathbf{K}^*)^{-1}$
- 12:  $\tilde{\mathbf{m}}_c^* \leftarrow \tilde{\mathbf{y}}_c \mathbf{K} (\mathbf{K}^* \mathbf{K}^{*T} + \mathbf{V}_c^{-1})^{-1} \mathbf{K}^* \tilde{\mathcal{V}}_c^*$
- 13: **for**  $n = 1$  to  $N_{test}$  **do**
- 14: **for**  $c = 1$  to  $C$  **do**
- 15: **for**  $i = 1$  to  $K$  Samples **do**
- 16:  $u_i \leftarrow \mathcal{N}(0, 1), p_{cn}^i \leftarrow \prod_{j \neq c} \Phi\left[\frac{1}{v_j^*} (u_i v_c^* + \tilde{m}_c^* - \tilde{m}_j^*)\right]$
- 17: **end for**
- 18: **end for**
- 19:  $P(t_n^* = c | \mathbf{x}_n^*, \mathbf{X}, \mathbf{t}) = \frac{1}{K} \sum_{i=1}^K p_{cn}^i$
- 20: **end for**

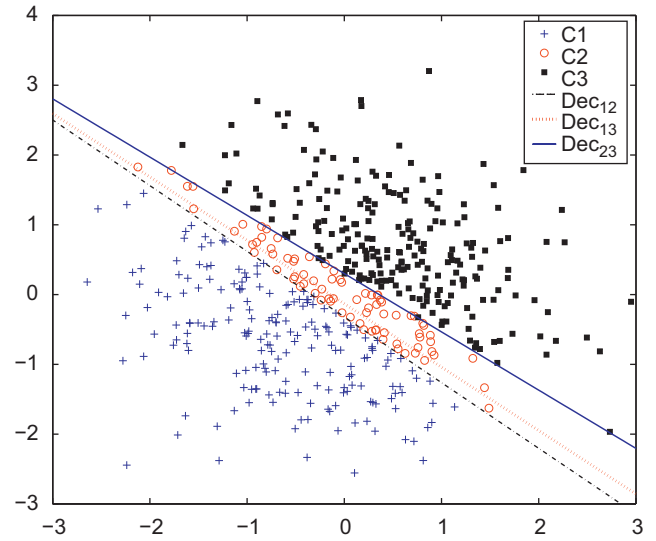
**4. Experimental studies**

This section presents the experimental findings of our work with emphasis first on the performance of the VB approximation with respect to the full MCMC Gibbs sampling solution and next on the efficiency of the proposed feature space combination method versus the well known classifier combination strategies. We employ two artificial low-dimensional datasets, a linearly and a non-linearly separable one introduced by Neal [27], for the comparison between the approximation and the full solution; standard UCI<sup>6</sup> multinomial datasets for an assessment of the VB performance; and finally, two well known large benchmark datasets to demonstrate the utility of the proposed feature space combination against ensemble learning and assess the performance of the different kernel combination methods.

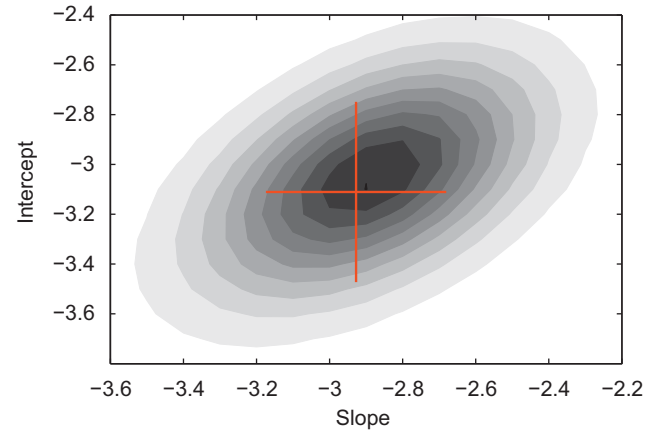
Unless otherwise specified, all the hyper-parameters were set to uninformative values and the Gaussian kernel parameters were fixed to  $1/D$  where  $D$  the dimensionality of the vectors. Furthermore, the convergence of the VB approximation was determined by monitoring the lower bound and the convergence occurred when there was less than 0.1% increase in the bound or when the maximum number of VB iterations was reached. The burn-in period for the Gibbs sampler was set to 10% of the total 100,000 of samples. Finally, all the CPU times reported in this study are for a 1.6GHz Intel based PC with 2Gb RAM running unoptimized Matlab codes and all the  $p$ -values reported are for a two sample  $t$ -test with the null hypothesis that the distributions have a common mean.

**4.1. Synthetic datasets**

In order to illustrate the performance of the VB approximation against the full Gibbs sampling solution, we employ two low-



**Fig. 3.** Linearly separable dataset with known regressors defining the decision boundaries.  $C_n$  denotes the members of class  $n$  and  $\text{Dec}_{ij}$  is the decision boundary between classes  $i$  and  $j$ .



**Fig. 4.** Gibbs posterior distribution of a decision boundary's ( $\text{Dec}_{12}$ ) slope and intercept for a chain of 100,000 samples. The cross shows the original decision boundary employed to sample the dataset.

dimensional datasets which enable us to visualize the decision boundaries and posterior distributions produced by either method. First we consider a linearly separable case in which we construct the dataset by fixing our regressors  $\mathbf{W} \in \mathfrak{R}^{C \times D}$ , with  $C=3$  and  $D=3$ , to known values and sample two-dimensional covariates  $\mathbf{X}$  plus a constant term. In that way, by knowing the true values of our regressors, we can examine the accuracy of both the Gibbs posterior distribution and the approximate posterior estimate of the VB. In Fig. 3 the dataset together with the optimal decision boundaries constructed by the known regressor values can be seen.

In Figs. 4 and 5 we present the posterior distributions of one decision boundary's ( $\text{Dec}_{12}$ ) slope and intercept based on both our obtained Gibbs samples and the approximate posterior of the regressors  $\mathbf{W}$ . As we can see, the variational approximation is in agreement with the mass of the Gibbs posterior and it successfully captures the pre-determined regressors values.

However, as it can be observed the approximation is overconfident in the prediction and produces a smaller covariance for the posterior distribution as expected [10]. Furthermore, the

<sup>6</sup> A. Asuncion, D.J. Newman, 2007. UCI Machine Learning Repository, Department of Information and Computer Science, University of California, Irvine, CA [http://www.ics.uci.edu/~mllearn/MLRepository.html]

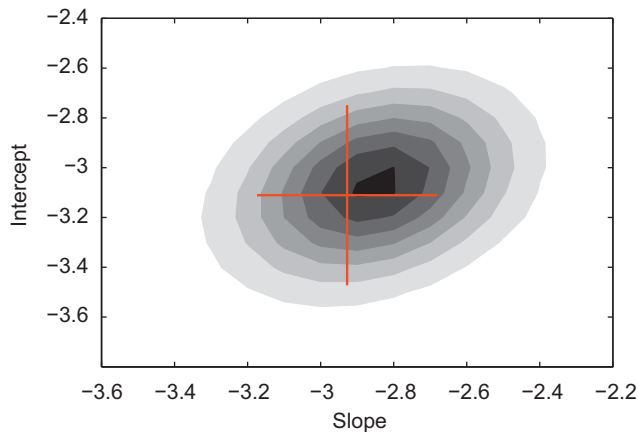


Fig. 5. The variational approximate posterior distribution for the same case as above. Employing 100,000 samples from the approximate posterior of the regressors  $\mathbf{W}$  in order to estimate the approximate posterior of the slope and intercept.

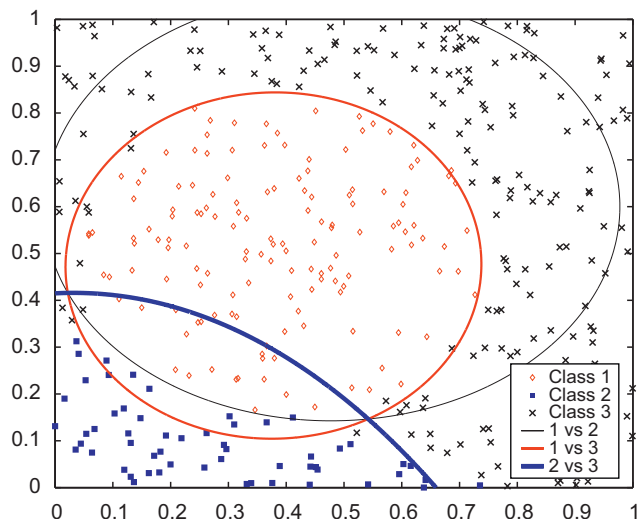


Fig. 7. Decision boundaries from the VB approximation.

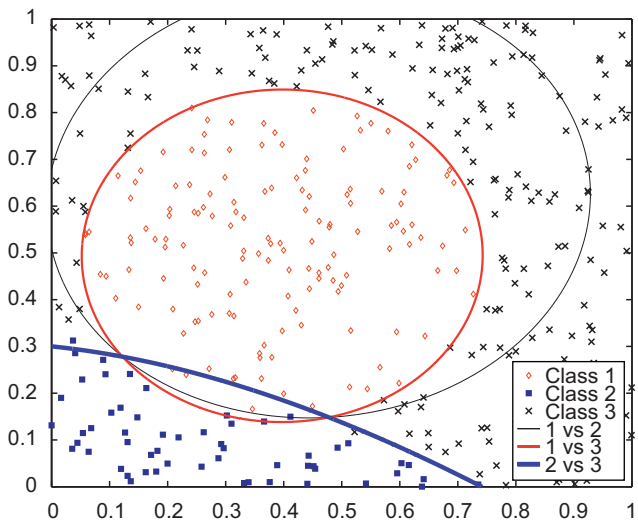


Fig. 6. Decision boundaries from the Gibbs sampling solution.

probability mass is concentrated in a very small area due to the very nature of VB approximation and similar mean field methods that make extreme “judgments” as they do not explore the posterior space by Markov chains:

$$C_{\text{Gibbs}} = \begin{bmatrix} 0.16 & 0.18 \\ 0.18 & 0.22 \end{bmatrix}, \quad C_{\text{VB}} = \begin{bmatrix} 0.015 & 0.015 \\ 0.015 & 0.018 \end{bmatrix} \quad (13)$$

The second synthetic dataset we employ is a four-dimensional three-class dataset  $\{\mathbf{X}, \mathbf{t}\}$  with  $N = 400$ , first described by Neal [27], which defines the first class as points in an ellipse  $\alpha > x_1^2 + x_2^2 > \beta$ , the second class as points below a line  $\alpha x_1 + \beta x_2 < \gamma$  and the third class as points surrounding these areas, see Fig. 6.

We approach the problem by: (1) introducing a second order polynomial expansion on the original dataset  $F(\mathbf{x}_n) = [1 \ x_{n1} \ x_{n2} \ x_{n1}^2 \ x_{n1}x_{n2} \ x_{n2}^2]$  while disregarding the uninformative dimensions  $x_3, x_4$ , (2) modifying the dimensionality of our regressors  $\mathbf{W}$  to  $C \times D$  and analogously the corresponding covariance and (3) substituting  $F(\mathbf{X})$  for  $\mathbf{K}$  in our derived methodology. Due to our expansion we now have a  $2 - D$  decision plane that we can plot and a six-dimensional regressor  $\mathbf{w}$  per class. In Fig. 6 we plot the decision boundaries produced from the full Gibbs solution by averaging over the posterior parameters after 100,000 samples and in Fig. 7

Table 1  
CPU time (s) comparison for 100,000 Gibbs samples versus 100 VB iterations.

Gibbs	VB
41,720 (s)	120.3 (s)

Notice that the number of VB iterations needed for the lower bound to converge is typically less than 100.

the corresponding decision boundaries from the VB approximation after 100 iterations.

As it can be seen, both the VB approximation and the MCMC solution produce similar decision boundaries leading to good classification performances—2% error for both Gibbs and VB for the above decision boundaries—depending on training size. However, the Gibbs sampler produces tighter boundaries due to the Markov chain exploring the parameter posterior space more efficiently than the VB approximation.

The corresponding CPU times are given in Table 1.

#### 4.2. Multinomial UCI datasets

To further assess the VB approximation of the proposed multinomial probit classifier, we explore a selection of UCI multinomial datasets. The performances are compared against reported results [26] from well-known methods in the pattern recognition and machine learning literature. We employ an RBF (VB RBF), a second order polynomial (VB P) and a linear kernel (VB L) with the VB approximation and report 10-fold cross-validated (CV) error percentages, in Table 2, and CPU times, in Table 3, for a maximum of 50 VB iterations unless convergence has already occurred. The comparison with the  $K$ -nn and PK-nn is for standard implementations of these methods, see Manocha and Girolami [26] for details.

As we can see the VB approximation to the multinomial probit outperforms in most cases both the  $K$ -nn and PK-nn although not offering statistical significant improvements as the variance of the 10-fold CV errors is quite large in most cases.

#### 4.3. Benchmark datasets

##### 4.3.1. Handwritten numerals classification

In this section we report results demonstrating the efficiency of our kernel combination approach when multiple sources of

**Table 2**  
Ten-fold cross-validated error percentages (mean ± std.) on standard UCI multinomial datasets.

Dataset	VB RBF	VB L	VB P	K-nn	PK-nn
Balance	8.8 ± 3.6	12.2 ± 4.2	<b>7.0 ± 3.3</b>	11.5 ± 3.0	10.2 ± 3.0
Crabs	23.5 ± 11.3	<b>13.5 ± 8.2</b>	21.5 ± 9.1	15.0 ± 8.8	19.5 ± 6.8
Glass	27.9 ± 10.1	35.8 ± 11.8	28.4 ± 8.9	29.9 ± 9.2	<b>26.7 ± 8.8</b>
Iris	<b>2.7 ± 5.6</b>	11.3 ± 9.9	4.7 ± 6.3	5.3 ± 5.2	4.0 ± 5.6
Soybean	6.5 ± 10.5	6 ± 9.7	<b>4 ± 8.4</b>	14.5 ± 16.7	4.5 ± 9.6
Vehicle	<b>25.6 ± 4.0</b>	29.6 ± 3.3	26 ± 6.1	36.3 ± 5.2	37.2 ± 4.5
Wine	4.5 ± 5.1	2.8 ± 4.7	<b>1.1 ± 2.3</b>	3.9 ± 3.8	3.4 ± 2.9

The bold letters indicate top mean accuracy.

**Table 3**  
Running times (seconds) for computing 10-fold cross-validation results.

Dataset	Balance	Crabs	Glass	Iris	Soybean	Vehicle	Wine
VB CPU time (s)	2285	270	380	89	19	3420	105

**Table 4**  
Classification percentage error (mean ± std.) of individual classifiers trained on each feature space.

FR	KL	PX	ZM
Full MCMC Gibbs sampling—single FS			
27.3 ± 3.3	11.0 ± 2.3	7.3 ± 2	25.2 ± 3

**Table 5**  
Combinations of the individual classifiers based on four widely used rules.

Prod	Sum	Max	Maj
Full MCMC Gibbs sampling—Comb. classifiers			
5.1 ± 1.7	5.3 ± 2	8.4 ± 2.3	8.45 ± 2.2

Product (Prod), Mean (Sum), Majority voting (Maj) and Maximum (Max).

**Table 6**  
Combination of feature spaces.

Bin	FixSum	WSum	FixProd	WProd
Full MCMC Gibbs sampling—Comb. FS				
5.7 ± 2	5.5 ± 2	5.8 ± 2.1	5.2 ± 1.8	5.9 ± 1.2

Gibbs sampling results for four feature space combination methods: Binary (Bin), Mean with fixed weights (FixSum), Mean with inferred weights (WSum), product with fixed weights (FixProd) and product with inferred weights (WProd).

information are available. Comparisons are made against combination of classifiers and also between different ways to combine the kernels as we have described previously. In order to assess the above, we make use of two datasets that have multiple feature spaces describing the objects to be classified. The first one is a large  $N=2000$ , multinomial  $C=10$  and “multi-featured”  $S=4$  UCI dataset, named “Multiple Features”. The objects are handwritten numerals, from 0 to 9, and the available features are the Fourier descriptors (FR), the Karhunen–Loève features (KL), the pixel averages (PX) and the Zernike moments (ZM). Tax et al. [34] have previously reported results on this problem by combining classifiers but have employed a different test set which is not publicly available. Furthermore, we allow the rotation invariance property of the ZM features to cause problems in the distinction between digits 6 and 9. The hope is that the remaining feature spaces can compensate on the discrimination.

In Tables 4–7, we report experimental results over 50 repeated trials where we have randomly selected 20 training and 20 testing objects from each class. For each trial we employ (1) a single classifier on each feature space, (2) the proposed classifier on the composite feature space. This allows us to examine the performance of

**Table 7**  
Combination of feature spaces.

Bin	FixSum	WSum	FixProd	WProd
VB approx.—Comb. FS				
5.53 ± 1.7	4.85 ± 1.5	6.1 ± 1.6	5.35 ± 1.4	6.43 ± 1.8

VB approximation.

combinations of classifiers versus combination of feature spaces. It is worth noting that a concatenation of feature spaces for this specific problem has been found to perform as good as ensemble learning methods [8]. We report results from both the full MCMC solution and the VB approximation in order to analyse the possible degradation of the performance for the approximation. In all cases we employ the multinomial probit kernel machine using a Gaussian kernel with fixed parameters.

As we can see from Tables 5 and 4 the two best performing ensemble learning methods outperform all of the individual classifiers trained on separate feature spaces. With a  $p$ -value of  $1.5e^{-07}$  between the Pixel classifier and the Product rule, it is a statistical significant difference. At the same time, all the kernel combination methods in Table 6 match the best performing classifier combination approaches, with a  $p$ -value of 0.91 between the Product classifier combination rule and the product kernel combination method. Finally, from Table 7 it is obvious that the VB approximation performs very well compared with the full Gibbs solution and even when combinatorial weights are inferred, which expands the parameter space, there is no statistical difference, with a  $p$ -value of 0.47, between the MCMC and the VB solution.

Furthermore, the variants of our method that employ combinatorial weights offer the advantage of inferring the significance of the contributing sources of information. In Fig. 8, we can see that the pixel (PX) and Zernike moments (ZM) feature spaces receive large weights and hence contribute significantly in the composite feature space. This is in accordance with our expectations, as the pixel feature space seems to be the best performing individual classifier and complementary predictive power mainly from the ZM channel is improving on that.

The results indicate that there is no benefit in the classification error performance when weighted combinations of feature spaces are employed. This is in agreement with past work by Lewis et al. [23] and Girolami and Zhong [16]. The clear benefit, however, remains the ability to infer the relative significance of the sources and hence gain a better understanding of the problem.

#### 4.3.2. Protein fold classification

The second multi-feature dataset we examine was first described by Ding and Dubchak [13] and is a protein fold classification problem with  $C=27$  classes and  $S=6$  available feature spaces (approx.

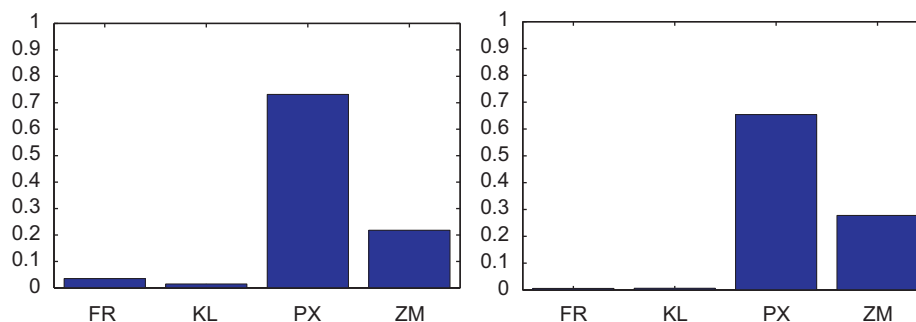


Fig. 8. Typical combinatorial weights from Mean (left) and Product (right) composite kernel.

Table 8

Classification error percentage on the individual feature spaces.

Comp	Hydro	Pol	Polz	Str	Vol
VB approximation—Single FS					
<b>49.8 ± 0.6</b>	63.4 ± 0.6	65.2 ± 0.4	66.3 ± 0.4	60.0 ± 0.25	65.6 ± 0.3

(1) Amino Acid composition (Comp); (2) Hydrophobicity profile (HP); (3) Polarity (Pol); (4) Polarizability (Polz); (5) Secondary structure (Str); and (6) van der Waals volume profile (Vol).

The bold letters indicate statistical significant improvements.

Table 9

Classification error percentages by combining the predictive distributions of the individual classifiers.

Prod	Sum	Max	Maj
VB approximation—Comb. classifiers			
49.8 ± 0.6	<b>47.6 ± 0.4</b>	53.7 ± 1.4	54.1 ± 0.5

The bold letters indicate statistical significant improvements.

Table 10

Classification error percentages by feature space combination.

Bin	FixSum	WSum	FixProd	WProd
VB approximation—Comb. FS				
40.7 ± 1.2	<b>40.1 ± 0.3</b>	44.4 ± 0.6	43.8 ± 0.14	43.2 ± 1.1

The bold letters indicate statistical significant improvements.

Table 11

Typical running times (seconds).

Method	Classifier combination	FS combination
CPU time (s)	11,519	1256

Ten-fold reduction in CPU time by combining feature spaces.

20-D each). The dataset is available online<sup>7</sup> and it is divided into a train set of  $N = 313$  size and an independent test set with  $N_{test} = 385$ . We use the original dataset and do not consider the feature modifications (extra features, modification of existing ones and omission of four objects) suggested recently in the work by Shen and Chou [32] in order to increase prediction accuracy. In Table 8 we report the performances of the individual classifiers trained on each feature space, in Table 9 the classifier combination performances and in Table 10 the kernel combination approaches. Furthermore, in Table 11 we give the corresponding CPU time requirements between

classifier and kernel combination methods for comparison. It is worth noting that a possible concatenation of all features into a single representation has been examined by Ding and Dubchak [13] and was found not to be performing competitively.

The experiments are repeated over five randomly initialized trials for a maximum of 100 VB iterations by monitoring convergence with a 0.1% increase threshold on the lower bound progression. We employ second order polynomial kernels as they were found to give the best performances across all methods.

As it can be observed the best of the classifier combination methods outperform any individual classifier trained on a single feature space with a  $p$ -value between the Sum combination rule and the composition classifier of  $2.3e - 02$ . However, all the kernel combination methods perform better than the best classifier combination with a  $p$ -value of  $1.5e - 08$  between the latter (FixSum) and the best performing classifier combination method (Sum). The corresponding CPU times show a 10-fold reduction by the proposed feature space combination method from the classifier combination approaches.

Furthermore, with a top performance of 38.3 error % for a single binary combination run (Bin) we match the state-of-the-art reported on the original dataset by Girolami and Zhong [16] employing Gaussian process priors. The best reported performance by Ding and Dubchak [13] was a 56% accuracy (44%) error employing an all-vs-all method with  $S \times C \times (C - 1)/2 = 8240$  binary SVM classifiers. We offer instead a single multiclass kernel machine able to combine the feature spaces and achieve a 60–62% accuracy in reduced computational time. It is worth noting that recent work by the authors [7] has achieved a 70% accuracy by considering additional state-of-the-art string kernels.

Finally, in Fig. 9 the combinatorial weights can be seen for the case of the mean and the product<sup>8</sup> composite kernel. The results are in agreement with the past work identifying the significance of the amino acid composition, hydrophobicity profile and secondary structure feature spaces for the classification task.

## 5. Discussion

A multinomial probit composite kernel classifier, based on a well-founded hierarchical Bayesian framework and able to instructively combine feature spaces has been presented in this work. The full MCMC MH within Gibbs sampling approach allows for exact inference to be performed and the proposed variational Bayes approximation reduces significantly computational requirements while retaining the same levels of performance. The proposed methodology enables inference to be performed in three significant levels by

<sup>7</sup> <http://crd.lbl.gov/~cding/protein/>

<sup>8</sup> Notice that for the product composite kernel the combinatorial weights are not restricted in the simplex.



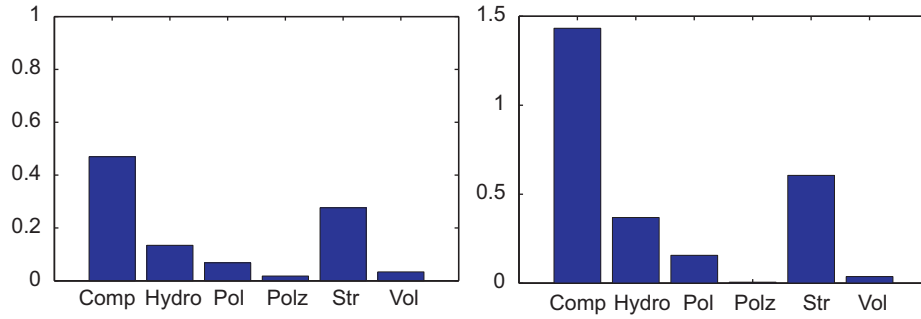


Fig. 9. Typical combinatorial weights from Mean (left) and Product (right) composite kernel for the protein fold dataset.

learning the combinatorial weights associated to each feature space, the kernel parameters associated to each dimension of the feature spaces and finally the parameters/regressors associated with the elements of the composite feature space.

An explicit comparison between the Gibbs sampling and the VB approximation has been presented. The approximate solution was found not to worsen significantly the classification performance and the resulting decision boundaries were similar to the MCMC approach, though the latter is producing tighter descriptions of the classes. The VB posteriors over the parameters were observed to be narrow as expected and all the probability mass was concentrated on a small area, however, within the posterior distributions produced by the Gibbs solution.

The proposed feature space combination approach was found to be as good as, or outperforming the best classifier combination rule examined and achieving the state-of-the-art while at the same time offering a 10-fold reduction in computational time and a better scaling when multiple spaces are available. At the same time, contrary to the classifier combination approach, inference can be performed on the combinatorial weights of the composite feature space, enabling a better understanding of the significance of the sources. In comparison with previously employed SVM combinations for the multiclass problem of protein fold prediction, our method offers a significant improvement of performance while employing a single classifier instead of thousands, reducing the computational time from approximately 12 to 0.3 CPU hours.

Finally, with respect to the two main problems considered, we can conclude that there is greater improvement on zero-one loss when diverse sources of information are used. This is a well known observation in ensemble learning methodology and kernel combination methods follow the same phenomenon. Theoretical justification for this can be easily provided via a simple bias-variance-covariance decomposition of the loss which follows by the same analysis as the one used for an ensemble of regressors. We are currently investigating a theoretical justification for another well known phenomenon in multiple kernel learning, that of zero-one loss equivalence between the average and the weighted kernel combination rules on specific problems. Recently we have offered in Damoulas et al. [9] a further deterministic approximation based on an EM update scheme which induces sparsity and in fact generalizes the relevance vector machine to the multi-class and multi-kernel setting.

### Acknowledgments

This work is sponsored by NCR Financial Solutions Group Ltd. The authors would particularly like to acknowledge the help and support from Dr. Gary Ross and Dr. Chao He of NCR Labs. The second author is supported by an EPSRC Advanced Research Fellowship (EP/E052029/1).

### Appendix A. Gibbs sampler

The conditional distribution for the auxiliary variables  $\mathbf{Y}$  is a product of  $N$   $C$ -dimensional conically truncated Gaussians given by

$$\prod_{n=1}^N \delta(y_{in} > y_{jn} \forall j \neq i) \delta(t_n = i) \mathcal{N}_{\mathbf{y}_n}(\mathbf{W} \mathbf{k}_n^{\beta \Theta}, \mathbf{I})$$

and for the regressors  $\mathbf{W}$  is a product of Gaussian distributions  $\prod_{c=1}^C \mathcal{N}_{\mathbf{w}_c}(\mathbf{m}_c, \mathbf{V}_c)$  where

$$\mathbf{m}_c = \mathbf{y}_c \mathbf{K}^{\beta \Theta} \mathbf{V}_c \text{ (row vector) and } \mathbf{V}_c = (\mathbf{K}^{\beta \Theta} \mathbf{K}^{\beta \Theta} + \mathbf{Z}_c^{-1})^{-1}$$

hence by iteratively sampling from these distributions will lead the Markov chain to sample from the desired posterior distribution. The typical MH subsamplers employed, in the case of the *mean composite* kernel model have acceptance ratios

$$A(\beta^i, \beta^*) = \min \left( 1, \frac{\prod_{c=1}^C \prod_{n=1}^N \mathcal{N}_{\mathbf{y}_{cn}}(\mathbf{w}_c \mathbf{k}_n^{\beta^* \Theta}, 1) \prod_{s=1}^S \beta_s^{*\rho_s-1}}{\prod_{c=1}^C \prod_{n=1}^N \mathcal{N}_{\mathbf{y}_{cn}}(\mathbf{w}_c \mathbf{k}_n^{\beta^i \Theta}, 1) \prod_{s=1}^S \beta_s^{i\rho_s-1}} \right)$$

$$A(\rho^i, \rho^*) = \min \left( 1, \frac{\Phi(\rho^*) \prod_{s=1}^S \beta_s^{\rho_s^*-1} \prod_{s=1}^S \rho_s^{*\mu-1} e^{-\lambda \rho_s^*}}{\Phi(\rho^i) \prod_{s=1}^S \beta_s^{i\rho_s-1} \prod_{s=1}^S \rho_s^{i\mu-1} e^{-\lambda \rho_s^i}} \right)$$

$$\text{where } \Phi(\rho) = \frac{\Gamma(\sum_{s=1}^S \rho_s)}{\prod_{s=1}^S \Gamma(\rho_s)}$$

with the proposed move symbolized by  $*$  and the current state with  $i$ .

For the *binary composite* kernel, an extra Gibbs step is introduced in our model as  $p(\beta_i = 0 | \beta_{-i}, \mathbf{Y}, \mathbf{K}^{\beta \Theta_s} \forall s \in \{1, \dots, S\})$  which depends on the marginal likelihood given by  $p(\mathbf{Y} | \beta, \mathbf{K}^{\beta \Theta_s} \forall s \in \{1, \dots, S\}) = \prod_{c=1}^C (2\pi)^{-N/2} |\Omega_c|^{-1/2} \exp\{-\frac{1}{2} \mathbf{y}_c \Omega_c \mathbf{y}_c^T\}$  with  $\Omega_c = \mathbf{I} + \mathbf{K}^{\beta \Theta} \mathbf{Z}_c^{-1} \mathbf{K}^{\beta \Theta}$ .

The *product composite* kernel employs two MH subsamplers to sample  $\beta$  (from a gamma distribution this time) and the hyperparameters  $\pi, \chi$ . The kernel parameters  $\Theta$  are inferred in all cases via an extra MH subsampler with acceptance ratio

$$A(\Theta^i, \Theta^*) = \min \left( 1, \frac{\prod_{s=1}^S \prod_{d=1}^D \theta_{sd}^* \prod_{c=1}^C \prod_{n=1}^N \mathcal{N}_{\mathbf{y}_{cn}}(\mathbf{w}_c \mathbf{k}_n^{\beta \Theta^*}, 1) \theta_{sd}^{*\omega-1} e^{-\phi \theta_{sd}^*}}{\prod_{s=1}^S \prod_{d=1}^D \theta_{sd}^i \prod_{c=1}^C \prod_{n=1}^N \mathcal{N}_{\mathbf{y}_{cn}}(\mathbf{w}_c \mathbf{k}_n^{\beta \Theta^i}, 1) \theta_{sd}^{i\omega-1} e^{-\phi \theta_{sd}^i}} \right)$$

Finally, the Monte Carlo estimate of the predictive distribution is used to assign the class probabilities according to a number of samples  $L$  drawn from the predictive distribution which, considering

Eq. (4) and substituting for the test composite kernel  $\mathbf{k}_n^{*\beta\Theta}$  defines the estimate of the predictive distribution as

$$P(t_* = c | \mathbf{x}_*) = \frac{1}{L} \sum_{l=1}^L \mathcal{E}_{p(u)} \left\{ \prod_{j \neq c} \Phi(u + (\mathbf{w}_c^l - \mathbf{w}_j^l) \mathbf{k}_n^{*\beta\Theta}) \right\}$$

## Appendix B. Approximate posterior distributions

We give here the full derivations of the approximate posteriors for the model under the VB approximation. We employ a first order approximation for the kernel parameters  $\Theta$ , i.e.  $\mathcal{E}_{Q(\Theta)}\{\mathbf{K}^{i\theta_i} \mathbf{K}^{j\theta_j}\} \approx \mathbf{K}^{i\tilde{\theta}_i} \mathbf{K}^{j\tilde{\theta}_j}$ , to avoid nonlinear contributions to the expectation. The same approximation is applied to the *product composite* kernel case for the combinatorial parameters  $\beta$  where  $\mathcal{E}_{Q(\beta)}\{\mathbf{K}^{i\beta_i} \mathbf{K}^{j\beta_j}\} \approx \mathbf{K}^{i\tilde{\beta}_i} \mathbf{K}^{j\tilde{\beta}_j}$

### B.1. Q(Y)

$$\begin{aligned} Q(\mathbf{Y}) &\propto \exp\{E_{Q(\mathbf{W})Q(\beta)Q(\Theta)}\{\log J \cdot L\}\} \\ &\propto \exp\{E_{Q(\mathbf{W})Q(\beta)Q(\Theta)}\{\log p(\mathbf{t}|\mathbf{Y}) + \log p(\mathbf{Y}|\mathbf{W}, \beta, \Theta)\}\} \\ &\propto \prod_{n=1}^N \delta(y_{i,n} > y_{k,n} \forall k \neq i) \delta(t_n = i) \exp\{E_{Q(\mathbf{W})Q(\beta)Q(\Theta)} \log p(\mathbf{Y}|\mathbf{W}, \beta, \Theta)\} \end{aligned}$$

where the exponential term can be analysed as follows:

$$\begin{aligned} &\exp\{E_{Q(\mathbf{W})Q(\beta)Q(\Theta)} \log p(\mathbf{Y}|\mathbf{W}, \beta)\} \\ &= \exp \left\{ E_{Q(\mathbf{W})Q(\beta)Q(\Theta)} \log \prod_{n=1}^N \mathcal{N}_{\mathbf{y}_n}(\mathbf{W} \mathbf{k}_n^{\beta\Theta}, \mathbf{I}) \right\} \\ &= \exp \left\{ E_{Q(\mathbf{W})Q(\beta)Q(\Theta)} \left\{ \sum_{n=1}^N \log \frac{1}{(2\pi)^{C/2} |\mathbf{I}|^{1/2}} \right. \right. \\ &\quad \left. \left. \times \exp \left( -\frac{1}{2} (\mathbf{y}_n - \mathbf{W} \mathbf{k}_n^{\beta\Theta})^T (\mathbf{y}_n - \mathbf{W} \mathbf{k}_n^{\beta\Theta}) \right) \right\} \right\} \\ &= \exp \left\{ E_{Q(\mathbf{W})Q(\beta)Q(\Theta)} \left\{ \sum_{n=1}^N \left( -\frac{1}{2} (\mathbf{y}_n^T \mathbf{y}_n - 2 \mathbf{y}_n^T \mathbf{W} \mathbf{k}_n^{\beta\Theta} \right. \right. \right. \\ &\quad \left. \left. + (\mathbf{k}_n^{\beta\Theta})^T \mathbf{W}^T \mathbf{W} \mathbf{k}_n^{\beta\Theta}) \right) \right\} \right\} \\ &= \exp \left\{ \sum_{n=1}^N -\frac{1}{2} \left( \mathbf{y}_n^T \mathbf{y}_n - 2 \mathbf{y}_n^T \tilde{\mathbf{W}} \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}} \right. \right. \\ &\quad \left. \left. + \sum_{i=1}^S \sum_{j=1}^S \tilde{\beta}_i \tilde{\beta}_j (\mathbf{k}_n^{i\tilde{\theta}_i})^T \tilde{\mathbf{W}}^T \mathbf{W} \mathbf{k}_n^{j\tilde{\theta}_j} \right) \right\} \end{aligned}$$

where  $\mathbf{k}_n^{i\tilde{\theta}_i}$  is the  $n$ th  $N$ -dimensional column vector of the  $i$ th base kernel with kernel parameters  $\tilde{\theta}_i$ . Now from this exponential term we can form the posterior distribution as a Gaussian and reach to the final expression:

$$Q(\mathbf{Y}) \propto \prod_{n=1}^N \delta(y_{i,n} > y_{k,n} \forall k \neq i) \delta(t_n = i) \mathcal{N}_{\mathbf{y}_n}(\tilde{\mathbf{W}} \mathbf{k}_n^{\tilde{\beta}\tilde{\Theta}}, \mathbf{I}) \quad (\text{B.1})$$

which is a  $C$ -dimensional conically truncated Gaussian.

### B.2. Q(W)

$$\begin{aligned} Q(\mathbf{W}) &\propto \exp\{E_{Q(\mathbf{Y})Q(\beta)Q(\zeta)Q(\Theta)}\{\log p(\mathbf{Y}|\mathbf{W}, \beta) + \log p(\mathbf{W}|\zeta)\}\} \\ &= \exp \left\{ E_{Q(\mathbf{Y})Q(\beta)Q(\Theta)} \left\{ \log \prod_{c=1}^C \mathcal{N}_{\mathbf{y}_c}(\mathbf{w}_c \mathbf{K}^{\beta\Theta}, \mathbf{I}) \right\} \right. \\ &\quad \left. + E_{Q(\zeta)} \left\{ \log \prod_{c=1}^C \mathcal{N}_{\mathbf{w}_c}(\mathbf{0}, \mathbf{Z}_c) \right\} \right\} \\ &= \exp \left\{ E_{Q(\mathbf{Y})Q(\beta)Q(\Theta)} \left\{ \sum_{c=1}^C \log \frac{1}{(2\pi)^{N/2} |\mathbf{I}|^{1/2}} \right. \right. \\ &\quad \left. \left. \times \exp \left\{ -\frac{1}{2} (\mathbf{y}_c - \mathbf{w}_c \mathbf{K}^{\beta\Theta}) (\mathbf{y}_c - \mathbf{w}_c \mathbf{K}^{\beta\Theta})^T \right\} \right\} \right. \\ &\quad \left. + E_{Q(\zeta)} \left\{ \sum_{c=1}^C \log \frac{1}{(2\pi)^{N/2} |\mathbf{Z}_c|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{w}_c (\mathbf{Z}_c)^{-1} \mathbf{w}_c^T) \right\} \right\} \right\} \\ &= \exp \left\{ E_{Q(\mathbf{Y})Q(\beta)Q(\Theta)} \left\{ \sum_{c=1}^C -\frac{1}{2} (\mathbf{y}_c \mathbf{V}_c^T \right. \right. \\ &\quad \left. \left. - 2 \mathbf{y}_c \mathbf{K}^{\beta\Theta} \mathbf{w}_c^T + \mathbf{w}_c \mathbf{K}^{\beta\Theta} \mathbf{K}^{\beta\Theta} \mathbf{w}_c^T) \right\} \right. \\ &\quad \left. + E_{Q(\zeta)} \left\{ \sum_{c=1}^C -\frac{1}{2} \log \prod_{n=1}^N \zeta_{cn} - \frac{1}{2} \mathbf{w}_c (\mathbf{Z}_c)^{-1} \mathbf{w}_c^T \right\} \right\} \\ &= \exp \left\{ \sum_{c=1}^C -\frac{1}{2} \left\{ \widetilde{\mathbf{y}_c \mathbf{y}_c^T} - 2 \widetilde{\mathbf{y}_c} \mathbf{K}^{\tilde{\beta}\tilde{\Theta}} \widetilde{\mathbf{w}_c^T} \right. \right. \\ &\quad \left. \left. + \mathbf{w}_c \sum_{i=1}^S \sum_{j=1}^S \tilde{\beta}_i \tilde{\beta}_j \mathbf{K}^{i\tilde{\theta}_i} \mathbf{K}^{j\tilde{\theta}_j} \mathbf{w}_c^T + \sum_{n=1}^N \widetilde{\log \zeta_{cn} + \mathbf{w}_c (\tilde{\mathbf{Z}}_c)^{-1} \mathbf{w}_c^T} \right\} \right\} \end{aligned}$$

Again we can form the posterior expectation as a new Gaussian:

$$Q(\mathbf{W}) \propto \prod_{c=1}^C \mathcal{N}_{\mathbf{w}_c}(\tilde{\mathbf{y}}_c \mathbf{K}^{\tilde{\beta}\tilde{\Theta}} \mathbf{V}_c, \mathbf{V}_c) \quad (\text{B.2})$$

where  $\mathbf{V}_c$  is the covariance matrix defined as

$$\mathbf{V}_c = \left( \sum_{i=1}^S \sum_{j=1}^S \tilde{\beta}_i \tilde{\beta}_j \mathbf{K}^{i\tilde{\theta}_i} \mathbf{K}^{j\tilde{\theta}_j} + (\tilde{\mathbf{Z}}_c)^{-1} \right)^{-1} \quad (\text{B.3})$$

and  $\tilde{\mathbf{Z}}_c$  is a diagonal matrix of the expected variances  $\tilde{\zeta}_{c1}, \dots, \tilde{\zeta}_{cN}$  for each class.

### B.3. Q(Z)

$$\begin{aligned} Q(\mathbf{Z}) &\propto \exp\{E_{Q(\mathbf{W})}(\log p(\mathbf{W}|\mathbf{Z}) + \log p(\mathbf{Z}|\tau, v))\} \\ &= \exp\{E_{Q(\mathbf{W})}(\log p(\mathbf{W}|\mathbf{Z}))\} p(\mathbf{Z}|\tau, v) \end{aligned}$$

Analysing the exponential term only:

$$\begin{aligned} &\exp \left\{ E_{Q(\mathbf{W})} \left( \log \prod_{c=1}^C \prod_{n=1}^N \mathcal{N}_{\mathbf{w}_{cn}}(0, \zeta_{cn}) \right) \right\} \\ &= \exp \left\{ E_{Q(\mathbf{W})} \left( \sum_{c=1}^C \sum_{n=1}^N -\frac{1}{2} \log \zeta_{cn} - \frac{1}{2} \mathbf{w}_{cn}^2 \zeta_{cn}^{-1} \right) \right\} \\ &= \prod_{c=1}^C \prod_{n=1}^N \zeta_{cn}^{-1/2} \exp \left( -\frac{1}{2} \frac{\widetilde{\mathbf{w}_{cn}^2}}{\zeta_{cn}} \right) \end{aligned}$$

which combined with the  $p(\mathbf{Z}|\tau, v)$  prior Gamma distribution leads to our expected posterior distribution:

$$\begin{aligned} Q(\mathbf{Z}) &\propto \prod_{c=1}^C \prod_{n=1}^N \zeta_{cn}^{-(\tau+1/2)+1} \frac{v^\tau e^{-\zeta_{cn}^{-1}(v+(1/2)\widetilde{w}_{cn}^2)}}{\Gamma(\tau)} \\ &= \prod_{c=1}^C \prod_{n=1}^N \text{Gamma}\left(\zeta_{cn}^{-1}|\tau + \frac{1}{2}, v + \frac{1}{2}\widetilde{w}_{cn}^2\right) \end{aligned}$$

**B.4.  $Q(\beta), Q(\rho), Q(\Theta), Q(\pi), Q(\chi)$**

Importance sampling techniques [2] are used to approximate these posterior distributions as they are intractable. We present the case of the *mean composite* kernel analytically and leave the cases of the *product* and *binary composite* kernel as a straightforward modifications.

For  $Q(\rho)$  we have  $p(\rho|\beta) \propto p(\beta|\rho)p(\rho|\mu, \lambda)$

The unnormalized posterior is  $Q^*(\rho) = p(\beta|\rho)p(\rho|\mu, \lambda)$  and hence the importance weights are

$$\mathcal{W}(\rho^i) = \frac{Q^*(\rho^i)}{\prod_{s=1}^S \text{Gamma}_{\rho_s^i}(\mu, \lambda)} = \frac{\text{Dir}(\rho^i)}{\sum_{i=1}^I \text{Dir}(\rho^i)}$$

where *Gamma* and *Dir* are the Gamma and Dirichlet distributions, while *I* is the total number of samples of  $\rho$  taken until now from the product gamma distributions and *i*' denotes the current (last) sample. So now we can estimate any function *f* of  $\rho$  based on

$$\tilde{f}(\rho) = \sum_{i=1}^I f(\rho) \mathcal{W}(\rho^i)$$

In the same manner as above but now for  $Q(\beta)$  and  $Q(\Theta)$  we can use the unnormalized posteriors  $Q^*(\beta)$  and  $Q^*(\Theta)$ , where  $p(\beta|\rho, \mathbf{Y}, \mathbf{W}) \propto p(\mathbf{Y}|\mathbf{W}, \beta)p(\beta|\rho)$  and  $p(\Theta|\omega, \phi, \mathbf{Y}, \mathbf{W}) \propto p(\mathbf{Y}|\mathbf{W}, \Theta)p(\Theta|\omega, \phi)$  with importance weights defined as

$$\mathcal{W}(\beta^i) = \frac{Q^*(\beta^i)}{\sum_{i=1}^I \frac{Q^*(\beta^i)}{\text{Dir}(\beta^i)}} = \frac{\prod_{n=1}^N \mathcal{N}_{\tilde{y}_n}(\tilde{\mathbf{W}}\mathbf{k}_n^{\beta^i}, \mathbf{I})}{\sum_{i=1}^I \prod_{n=1}^N \mathcal{N}_{\tilde{y}_n}(\tilde{\mathbf{W}}\mathbf{k}_n^{\beta^i}, \mathbf{I})}$$

and

$$\mathcal{W}(\Theta^i) = \frac{Q^*(\Theta^i)}{\sum_{i=1}^I \frac{Q^*(\Theta^i)}{\text{Dir}(\Theta^i)}} = \frac{\prod_{n=1}^N \mathcal{N}_{\tilde{y}_n}(\tilde{\mathbf{W}}\mathbf{k}_n^{\Theta^i}, \mathbf{I})}{\sum_{i=1}^I \prod_{n=1}^N \mathcal{N}_{\tilde{y}_n}(\tilde{\mathbf{W}}\mathbf{k}_n^{\Theta^i}, \mathbf{I})}$$

and again we can estimate any function *g* of  $\beta$  and *h* of  $\Theta$  as

$$\tilde{g}(\beta) = \sum_{i=1}^I g(\beta) \mathcal{W}(\beta^i) \quad \text{and} \quad \tilde{h}(\Theta) = \sum_{i=1}^I g(\Theta) \mathcal{W}(\Theta^i)$$

For the *product composite* kernel case we proceed in the same manner by importance sampling for  $Q(\pi)$  and  $Q(\chi)$ .

**Appendix C. Posterior expectations of each  $\mathbf{y}_n$**

As we have shown  $Q(\mathbf{Y}) \propto \prod_{n=1}^N \delta(y_{i,n} > y_{k,n} \forall k \neq i) \delta(t_n = i) \mathcal{N}_{\tilde{\mathbf{y}}_n}(\tilde{\mathbf{W}}\mathbf{k}_n^{\tilde{\beta}}, \mathbf{I})$ . Hence  $Q(\mathbf{y}_n)$  is a truncated multivariate Gaussian

distribution and we need to calculate the correction to the normalizing term  $\mathcal{Z}_n$  caused by the truncation. Thus, the posterior expectation can be expressed as

$$Q(\mathbf{y}_n) = \mathcal{Z}_n^{-1} \prod_{c=1}^C \mathcal{N}_{y_{cn}}^{t_n}(\tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}}, 1)$$

where the superscript  $t_n$  indicates the truncation needed so that the appropriate dimension *i* (since  $t_n = i \iff y_{in} > y_{jn} \forall j \neq i$ ) is the largest.

Now,  $\mathcal{Z}_n = P(\mathbf{y}_n \in \mathcal{C})$  where  $\mathcal{C} = \{\mathbf{y}_n : y_{in} > y_{jn}\}$  hence

$$\begin{aligned} \mathcal{Z}_n &= \int_{-\infty}^{+\infty} \mathcal{N}_{y_{in}}(\tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}}, 1) \prod_{j \neq i} \int_{-\infty}^{y_{in}} \mathcal{N}_{y_{jn}}(\tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}}, 1) dy_{jn} dy_{in} \\ &= \mathcal{E}_{p(u)} \left\{ \prod_{j \neq i} \Phi(u + \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}} - \tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}}) \right\} \end{aligned}$$

with  $p(u) = \mathcal{N}_u(0, 1)$ . The posterior expectation of  $y_{cn}$  for all  $c \neq i$  (the auxiliary variables associated with the rest of the classes except the one that object *n* belongs to) is given by

$$\begin{aligned} \tilde{y}_{cn} &= \mathcal{Z}_n^{-1} \int_{-\infty}^{+\infty} y_{cn} \prod_{j=1}^C \mathcal{N}_{y_{jn}}(\tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}}, 1) dy_{jn} \\ &= \mathcal{Z}_n^{-1} \int_{-\infty}^{+\infty} \int_{-\infty}^{y_{in}} y_{cn} \mathcal{N}_{y_{cn}}(\tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}}, 1) \\ &\quad \prod_{j \neq i, c} \mathcal{N}_{y_{jn}}(\tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}}, 1) \Phi(y_{in} - \tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}}) dy_{cn} dy_{in} \\ &= \tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}} - \mathcal{Z}_n^{-1} \mathcal{E}_{p(u)} \left\{ \mathcal{N}_u(\tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}} - \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}}, 1) \right. \\ &\quad \left. \prod_{j \neq i, c} \Phi(u + \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}} - \tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}}) \right\} \end{aligned}$$

For the *i*th class the posterior expectation  $y_{in}$  (the auxiliary variable associated with the known class of the *n*th object) is given by

$$\begin{aligned} \tilde{y}_{in} &= \mathcal{Z}_n^{-1} \int_{-\infty}^{+\infty} y_{in} \mathcal{N}_{y_{in}}(\tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}}, 1) \prod_{j \neq i} \Phi(y_{in} - \tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}}) dy_{in} \\ &= \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}} + \mathcal{Z}_n^{-1} \mathcal{E}_{p(u)} \left\{ u \prod_{j \neq i} \Phi(u + \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}} - \tilde{\mathbf{w}}_j \mathbf{k}_n^{\tilde{\beta}}) \right\} \\ &= \tilde{\mathbf{w}}_i \mathbf{k}_n^{\tilde{\beta}} + \sum_{c \neq i} (\tilde{\mathbf{w}}_c \mathbf{k}_n^{\tilde{\beta}} - \tilde{y}_{cn}) \end{aligned}$$

where we have made use of the fact that for a variable  $u \sim \mathcal{N}(0, 1)$  and any differentiable function  $g(u)$ ,  $\mathcal{E}\{ug(u)\} = \mathcal{E}\{g'(u)\}$ .

**Appendix D. Predictive distribution**

In order to make a prediction  $t^*$  for a new point  $\mathbf{x}^*$  we need to know:

$$\begin{aligned} p(t^* = c | \mathbf{x}^*, \mathbf{X}, \mathbf{t}) &= \int p(t^* = c | \mathbf{y}^*) p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{X}, \mathbf{t}) d\mathbf{y}^* \\ &= \int \delta_c^* p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{X}, \mathbf{t}) d\mathbf{y}^* \end{aligned}$$

Hence we need to evaluate

$$\begin{aligned} p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{X}, \mathbf{t}) &= \int p(\mathbf{y}^* | \mathbf{W}, \mathbf{x}^*) p(\mathbf{W} | \mathbf{X}, \mathbf{t}) d\mathbf{W} \\ &= \prod_{c=1}^C \int \mathcal{N}_{\mathbf{w}_c \mathbf{K}^c(\mathbf{y}^*, \mathbf{I})} \mathcal{N}_{\mathbf{w}_c}(\tilde{\mathbf{Y}}_c \mathbf{K} \mathbf{V}_c, \mathbf{V}_c) d\mathbf{w}_c \end{aligned}$$

We proceed by analysing the integral, gathering all the terms depending on  $\mathbf{w}_c$ , completing the square twice and reforming to

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{t}) = \prod_{c=1}^C \int \mathcal{N}_{\mathbf{y}_c^*}(\tilde{\mathbf{y}}_c \mathbf{K} \Lambda_c \mathbf{K}^* \tilde{\mathcal{V}}_c^*, \tilde{\mathcal{V}}_c^*) \times \mathcal{N}_{\mathbf{w}_c}((\mathbf{y}_c^* \mathbf{K}^{*T} + \tilde{\mathbf{y}}_c \mathbf{K}) \Lambda_c, \Lambda_c) \quad (\text{D.1})$$

with

$$\tilde{\mathcal{V}}_c^* = (\mathbf{I} - \mathbf{K}^{*T} \Theta_c \mathbf{K}^*)^{-1} \quad (N_{\text{test}} \times N_{\text{test}})$$

and

$$\Lambda_c = (\mathbf{K}^* \mathbf{K}^{*T} + \mathbf{V}_c^{-1})^{-1} \quad (N \times N) \quad (\text{D.2})$$

Finally we can simplify  $\tilde{\mathcal{V}}_c^*$  by applying the Woodbury identity and reduce its form to

$$\tilde{\mathcal{V}}_c^* = (\mathbf{I} + \mathbf{K}^{*T} \mathbf{V}_c \mathbf{K}^*)^{-1} \quad (N_{\text{test}} \times N_{\text{test}})$$

Now the Gaussian distribution wrt  $\mathbf{w}_c$  integrates to one and we are left with

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{t}) = \prod_{c=1}^C \mathcal{N}_{\mathbf{y}_c^*}(\tilde{\mathbf{m}}_c^*, \tilde{\mathcal{V}}_c^*) \quad (\text{D.3})$$

where  $\tilde{\mathbf{m}}_c^* = \tilde{\mathbf{y}}_c \mathbf{K} \Lambda_c \mathbf{K}^* \tilde{\mathcal{V}}_c^* \quad (1 \times N_{\text{test}})$

Hence we can go back to the predictive distribution and consider the case of a single test point (i.e.  $N_{\text{test}} = 1$ ) with associated scalars  $\tilde{m}_c^*$  and  $\tilde{v}_c^*$

$$\begin{aligned} p(t^* = c | \mathbf{x}^*, \mathbf{X}, \mathbf{t}) &= \int \delta_c^* \prod_{c=1}^C \mathcal{N}_{\mathbf{y}_c^*}(\tilde{m}_c^*, \tilde{v}_c^*) dy_c^* \\ &= \int_{-\infty}^{+\infty} \mathcal{N}_{\mathbf{y}_c^*}(\tilde{m}_c^*, \tilde{v}_c^*) \prod_{j \neq c} \int_{-\infty}^{y_c^*} \mathcal{N}_{\mathbf{y}_j^*}(\tilde{m}_j^*, \tilde{v}_j^*) dy_j^* dy_c^* \\ &= \int_{-\infty}^{+\infty} \mathcal{N}_{\mathbf{y}_c^* - \tilde{m}_c^*}(0, \tilde{v}_c^*) \prod_{j \neq c} \int_{-\infty}^{y_c^* - \tilde{m}_j^*} \mathcal{N}_{\mathbf{y}_j^* - \tilde{m}_j^*}(0, \tilde{v}_j^*) dy_j^* dy_c^* \end{aligned}$$

Setting  $u = (y_c^* - \tilde{m}_c^*) \tilde{v}_c^{*-1}$  and  $x = (y_j^* - \tilde{m}_j^*) \tilde{v}_j^{*-1}$  we have

$$\begin{aligned} p(t^* = c | \mathbf{x}^*, \mathbf{X}, \mathbf{t}) &= \int_{-\infty}^{+\infty} \mathcal{N}_u(0, 1) \prod_{j \neq c} \int_{-\infty}^{(u \tilde{v}_c^* + \tilde{m}_c^* - \tilde{m}_j^*) \tilde{v}_j^{*-1}} \mathcal{N}_x(0, 1) dx du \\ &= E_{p(u)} \left\{ \prod_{j \neq c} \Phi \left[ \frac{1}{\tilde{v}_j^*} (u \tilde{v}_c^* + \tilde{m}_c^* - \tilde{m}_j^*) \right] \right\} \end{aligned}$$

## Appendix E. Lower bound

From Jensen's inequality in Eq. (5) and by conditioning on current values of  $\beta, \Theta, \mathbf{Z}, \rho$  we can derive the variational lower bound using the relevant components

$$\mathcal{E}_{Q(\mathbf{Y})Q(\mathbf{W})} \{ \log p(\mathbf{Y} | \mathbf{W}, \beta, \mathbf{X}) \} \quad (\text{E.1})$$

$$+ \mathcal{E}_{Q(\mathbf{Y})Q(\mathbf{W})} \{ \log p(\mathbf{W} | \mathbf{Z}, \mathbf{X}) \} \quad (\text{E.2})$$

$$- \mathcal{E}_{Q(\mathbf{Y})} \{ \log Q(\mathbf{Y}) \} \quad (\text{E.3})$$

$$- \mathcal{E}_{Q(\mathbf{W})} \{ \log Q(\mathbf{W}) \} \quad (\text{E.4})$$

which, by noting that the expectation of a quadratic form under a Gaussian is another quadratic form plus a constant, leads to the

following expression for the lower bound:

$$\begin{aligned} & - \frac{NC}{2} \log 2\pi - \frac{1}{2} \sum_{c=1}^C \sum_{n=1}^N \{ \tilde{y}_{cn}^2 + \mathbf{k}_n^T \tilde{\mathbf{w}}_c \tilde{\mathbf{w}}_c \mathbf{k}_n - 2 \tilde{y}_{cn} \tilde{\mathbf{w}}_c \mathbf{k}_n \} \\ & - \frac{NC}{2} \log 2\pi - \frac{1}{2} \sum_{c=1}^C \log |\mathbf{Z}_c| - \frac{1}{2} \sum_{c=1}^C \tilde{\mathbf{w}}_c \mathbf{Z}_c^{-1} \tilde{\mathbf{w}}_c^T \\ & - \frac{1}{2} \sum_{c=1}^C \text{Tr}[\mathbf{Z}_c^{-1} \mathbf{V}_c] + \sum_{n=1}^N \log \mathcal{L}_n + \frac{NC}{2} \log 2\pi \\ & + \frac{1}{2} \sum_{n=1}^N \sum_{c=1}^C \{ \tilde{y}_{cn}^2 - 2 \tilde{y}_{cn} \tilde{\mathbf{w}}_c \mathbf{k}_n + \mathbf{k}_n^T \tilde{\mathbf{w}}_c \tilde{\mathbf{w}}_c \mathbf{k}_n \} \\ & + \frac{NC}{2} \log 2\pi + \frac{1}{2} \sum_{c=1}^C \log |\mathbf{V}_c| + \frac{NC}{2} \end{aligned}$$

which simplifies to our final expression

$$\text{Lower Bound} = \frac{NC}{2} + \frac{1}{2} \sum_{c=1}^C \log |\mathbf{V}_c| + \sum_{n=1}^N \log \mathcal{L}_n \quad (\text{E.5})$$

$$- \frac{1}{2} \sum_{c=1}^C \text{Tr}[\mathbf{Z}_c^{-1} \mathbf{V}_c] - \frac{1}{2} \sum_{c=1}^C \tilde{\mathbf{w}}_c \mathbf{Z}_c^{-1} \tilde{\mathbf{w}}_c^T \quad (\text{E.6})$$

$$- \frac{1}{2} \sum_{c=1}^C \log |\mathbf{Z}_c| - \frac{1}{2} \sum_{c=1}^C \sum_{n=1}^N \mathbf{k}_n^T \mathbf{V}_c \mathbf{k}_n \quad (\text{E.7})$$

## References

- [1] J. Albert, S. Chib, Bayesian analysis of binary and polychotomous response data, *Journal of the American Statistical Association* 88 (1993) 669–679.
- [2] C. Andrieu, An introduction to MCMC for machine learning, *Machine Learning* 50 (2003) 5–43.
- [3] L. Bai, L. Shen, Combining wavelets with HMM for face recognition, in: *Proceedings of the 23rd Artificial Intelligence Conference*, 2003.
- [4] M.J. Beal, Variational algorithms for approximate Bayesian inference, Ph.D. Thesis, The Gatsby Computational Neuroscience Unit, University College London, 2003.
- [5] J.O. Berger, *Statistical Decision Theory and Bayesian Analysis*, Springer Series in Statistics, Springer, Berlin, 1985.
- [6] C.M. Bishop, *Pattern Recognition and Machine Learning*, Springer, New York, USA, 2006.
- [7] T. Damoulas, M.A. Girolami, Probabilistic multi-class multi-kernel learning: on protein fold recognition and remote homology detection, *Bioinformatics* 24 (10) (2008) 1264–1270.
- [8] T. Damoulas, M.A. Girolami, Pattern recognition with a Bayesian kernel combination machine, *Pattern Recognition Letters* 30 (1) (2009) 46–54.
- [9] T. Damoulas, Y. Ying, M.A. Girolami, C. Campbel, Inferring sparse kernel combinations and relevant vectors: an application to subcellular localization of proteins, in: *IEEE, International Conference on Machine Learning and Applications (ICMLA '08)*, 2008.
- [10] N. de Freitas, P. Højten-Sørensen, M. Jordan, S. Russell, Variational MCMC, in: *Proceedings of the 17th conference in Uncertainty in Artificial Intelligence*, 2001.
- [11] D.G.T. Denison, C.C. Holmes, B.K. Mallick, A.F.M. Smith, *Bayesian Methods for Nonlinear Classification and Regression*, Wiley Series in Probability and Statistics, West Sussex, UK, 2002.
- [12] T.G. Dietterich, Ensemble methods in machine learning, in: *Proceedings of the 1st International Workshop on Multiple Classifier Systems*, 2000, pp. 1–15.
- [13] C. Ding, I. Dubchak, Multi-class protein fold recognition using support vector machines and neural networks, *Bioinformatics* 17 (4) (2001) 349–358.
- [14] M. Girolami, S. Rogers, Hierarchic Bayesian models for kernel learning, in: *Proceedings of the 22nd International Conference on Machine Learning*, 2005, pp. 241–248.
- [15] M. Girolami, S. Rogers, Variational Bayesian multinomial probit regression with Gaussian process priors, *Neural Computation* 18 (8) (2006) 1790–1817.
- [16] M. Girolami, M. Zhong, Data integration for classification problems employing Gaussian process priors, in: B. Schölkopf, J. Platt, T. Hoffman (Eds.), *Advances in Neural Information Processing Systems*, vol. 19, MIT Press, Cambridge, MA, 2007, pp. 465–472.
- [17] W.K. Hastings, Monte Carlo sampling methods using Markov chains and their applications, *Biometrika* 57 (1970) 97–109.
- [18] J. Kittler, M. Hatef, R.P.W. Duin, J. Matas, On combining classifiers, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20 (3) (1998) 226–239.
- [19] L.I. Kuncheva, *Combining Pattern Classifiers, Methods and Algorithms*, Wiley, New York, 2004.

- [20] G.R.G. Lanckriet, N. Cristianini, P. Bartlett, L.E. Ghaoui, M.I. Jordan, Learning the kernel matrix with semidefinite programming, *Journal of Machine Learning Research* 5 (2004) 27–72.
- [21] W.-J. Lee, S. Verzakov, R.P. Duin, Kernel combination versus classifier combination, in: 7th International Workshop on Multiple Classifier Systems, 2007.
- [22] D.P. Lewis, T. Jebara, W.S. Noble, Nonstationary kernel combination, in: 23rd International Conference on Machine Learning, 2006.
- [23] D.P. Lewis, T. Jebara, W.S. Noble, Support vector machine learning from heterogeneous data: an empirical analysis using protein sequence and structure, *Bioinformatics* 22 (22) (2006) 2753–2760.
- [24] D. MacKay, *Information Theory, Inference and Learning Algorithms*, Cambridge University Press, Cambridge, 2003.
- [25] D.J.C. MacKay, The evidence framework applied to classification networks, *Neural Computation* 4 (5) (1992) 698–714.
- [26] S. Manocha, M.A. Girolami, An empirical analysis of the probabilistic k-nearest neighbour classifier, *Pattern Recognition Letters* 28 (13) (2007) 1818–1824.
- [27] R.M. Neal, Regression and classification using Gaussian process priors, *Bayesian Statistics* 6 (1998) 475–501.
- [28] C.S. Ong, A.J. Smola, R.C. Williamson, Learning the kernel with hyperkernels, *Journal of Machine Learning Research* 6 (2005) 1043–1071.
- [29] C.E. Rasmussen, C.K.I. Williams, *Gaussian Processes for Machine Learning*, MIT Press, Cambridge, MA, USA, 2006.
- [30] B. Schölkopf, A. Smola, *Learning with Kernels*, The MIT Press, Cambridge, MA, USA, 2002.
- [31] J. Shawe-Taylor, N. Cristianini, *Kernel Methods for Pattern Analysis*, Cambridge University Press, Cambridge, England, UK, 2004.
- [32] H.-B. Shen, K.-C. Chou, Ensemble classifier for protein fold pattern recognition, *Bioinformatics* 22 (14) (2006) 1717–1722.
- [33] S. Sonnenburg, G. Rätsch, C. Schäfer, A general and efficient multiple kernel learning algorithm, in: Y. Weiss, B. Schölkopf, J. Platt (Eds.), *Advances in Neural Information Processing Systems*, vol. 18, MIT Press, Cambridge, MA, 2006, pp. 1273–1280.
- [34] D.M.J. Tax, M. van Breukelen, R.P.W. Duin, J. Kittler, Combining multiple classifiers by averaging or by multiplying?, *Pattern Recognition* 33 (2000) 1475–1485.
- [35] M.E. Tipping, Sparse Bayesian learning and the relevance vector machine, *Journal of Machine Learning Research* 1 (2001) 211–244.
- [36] D.H. Wolpert, Stacked generalization, *Neural Networks* 5 (2) (1992) 241–259.

**About the Author**—THEODOROS DAMOULAS was awarded the NCR Ltd. Ph.D. Scholarship in 2006 in the Department of Computing Science at the University of Glasgow. He is a member of the Inference Research Group (<http://www.dcs.gla.ac.uk/inference>) and his research interests are in statistical machine learning and pattern recognition. He holds a 1st class M.Eng. degree in Mechanical Engineering from the University of Manchester and an M.Sc. in Informatics (Distinction) from the University of Edinburgh.

**About the Author**—MARK A. GIROLAMI is Professor of Computing and Inferential Science and holds a joint appointment in the Department of Computing Science and the Department of Statistics at the University of Glasgow. His various research interests lie at the boundary between computing, statistical, and biological science. He currently holds a prestigious five year long (2007 until 2012) Advanced Research Fellowship from the Engineering & Physical Sciences Research Council (EPSRC). In 2009 the International Society of Photo-Optical Instrumentation Engineers (SPIE) presented him with their Pioneer Award for his contributions to Biomedical Signal Processing and the development of analytical tools for EEG and fMRI which revolutionized the automated processing of such signals. In 2005 he was awarded an MRC Discipline Hopping Award in the Department of Biochemistry and during 2000 he was the TEKES visiting professor at the Laboratory of Computing and Information Science in Helsinki University of Technology. In 1998 and 1999 Professor Girolami was a research fellow at the Laboratory for Advanced Brain Signal Processing in the Brain Science Institute, RIKEN, Wako-Shi, Japan. He has been a visiting researcher at the Computational Neurobiology Laboratory (CNL) of the Salk Institute. The main focus of the seven post-doctoral researchers and eight Ph.D. students in his group is in the development of statistical inferential methodology to support reasoning about biological mechanisms within a systems biology context.