



Citation for published version:

Drugowitsch, J & Barry, A 2007, *Generalised mixtures of experts, independent expert training, and learning classifier systems*. Computer Science Technical Reports, no. CSBU-2007-02, University of Bath, Department of Computer Science.

Publication date:
2007

[Link to publication](#)

©The Author April 2007

University of Bath

Alternative formats

If you require this document in an alternative format, please contact:
openaccess@bath.ac.uk

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

**Department of
Computer Science**



UNIVERSITY OF
BATH

Technical Report

Generalised Mixtures of Experts, Independent Expert Training,
and Learning Classifier Systems

Jan Drugowitsch and Alwyn Barry

Copyright ©April 2007 by the authors.

Contact Address:

Department of Computer Science
University of Bath
Bath, BA2 7AY
United Kingdom
URL: <http://www.cs.bath.ac.uk>

ISSN 1740-9497

Generalised Mixtures of Experts, Independent Expert Training, and Learning Classifier Systems

Jan Drugowitsch
J.Drugowitsch@bath.ac.uk

Alwyn M. Barry
A.M.Barry@bath.ac.uk

April 2007

Abstract

We present a generalisation to the Mixtures of Experts model that introduces prior localisation of the experts as part of the model structure, and as such relates them strongly to the evolutionary computation ML method known as Learning Classifier Systems. While the introduced generalisation allows specification of more complex localisation patterns, identifying good models becomes more difficult. We approach this tradeoff by introducing a new training schema that makes fitting a single model computationally less demanding and shifts the importance to searching the space of possible model structures, guided by approximate variational Bayesian inference to fit the model and find the model evidence. We demonstrate model search for simple non-linear curve fitting tasks by sampling from the model posterior, as a proof-of-concept alternative to the genetic algorithm used in Learning Classifier Systems for that purpose.

1 Introduction

While Learning Classifier Systems (LCS) are well established in the field of evolutionary computation, and perform competitively in regression and classification tasks when compared to other ML methods (Bernadó-Mansilla et al., 2002; Butz et al., to appear), their performance lacks explanation due to their purely algorithmic description that does not identify the underlying data model. This makes them hard to access and hinders their further improvement.

In this work we for the first time clearly identify the model underlying LCS by linking them to a generalisation of the Mixtures of Experts (MoE) model (Jacobs et al., 1991; Jordan & Jacobs, 1994). The latter method is well established in the machine learning literature, and trains a set of local models, called *experts*, by localising them in the input space with a gating network. This localisation is achieved by the interdependent training of gating network and experts, resulting in a smoothed linear partitioning of the input space. A variant on MoE replaces the linear partitioning by normalised Gaussian localisation, but with the same underlying principle (Xu et al., 1995).

In contrast, LCS add an additional layer of “forced” localisation (determined by the model structure) to each expert¹, where the shape of localisation is only limited by the choice of representation, allowing the capture of a richer data dependency structure and natural integration of nominal, ordinal and metric data attributes. This freedom comes at the cost of making them harder to train, as one needs to search the potentially complex space of possible expert localisations (that is, the space of possible model structures) in addition to training the model for a fixed set of localisations. While LCS acquire a more efficient model fitting scheme for this task, and search the space of model structures with a genetic algorithm, we describe this fitting scheme and its implications, but leave the adaptation of the genetic algorithm to the presented model as future work. Rather, we perform a less powerful model search by sampling from the model posterior as a proof-of-concept².

To assess the quality of a certain model structure in explaining the given data, we use the model posterior approximated by variational inference on a full Bayesian LCS model, closely related to similar approaches for MoE by Waterhouse et al. (1996), Ueda and Ghahramani (2002) and Bishop and Svensén (2003). Such a principled approach is new to LCS as their lack of model description requires LCS to fall back on a set of heuristics for model assessment.

In the next section we introduce the MoE and our generalisations to it, and discuss in Section 3 the requirement of a model selection mechanism that is realised by Bayesian means with approximate

¹In LCS, the experts are called *classifiers*, despite them usually being regression models.

²For more information on LCS, the interested reader is referred to Butz (2006).

variational inference. Our method of searching the space of possible models and the required modifications to the fitting procedure are described in Section 4, and demonstrated by sampling from the model posterior in Section 5, after which we offer conclusions about our achievements.

2 Generalised Mixtures of Experts

Let us for now assume a fixed model \mathcal{M} with K experts. After Jordan and Jacobs (1994), each expert provides a conditional probability distribution $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_k)$ over an output vector \mathbf{y} given an input vector \mathbf{x} and a set of parameters $\boldsymbol{\theta}_k$ of the k th expert. In MoE these experts are combined by a gating network to form the mixture distribution $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$.

The gating network is best explained from the generative point-of-view. The random gating vector $\mathbf{z} = (z_1, \dots, z_K)^T$ of binary latent gating variables z_k determines which expert generates the observation $\{\mathbf{x}, \mathbf{y}\}$. An observation is always assigned to a single expert, and so \mathbf{z} has a 1-of- K structure, where $g_k(\mathbf{x}) \equiv p(z_k = 1|\mathbf{x}, \mathbf{v}_k)$ denotes the probability of expert k generating $\{\mathbf{x}, \mathbf{y}\}$, and \mathbf{v}_k is the parameter vector for g_k .

2.1 Localisation through Matching

While in the original MoE model there are no prior restrictions on the association of observations to experts, we introduce such a restriction by the concept of matching: an additional binary random variable m_k per expert restricts the set of observations that this expert can have generated to the set that it *matches*, that is, for which $m_k = 1$. This corresponds to each expert in the LCS being restricted to modelling a subset of the input space, and so matching is a property of the model \mathcal{M} that is specified by the matching function $m_k(\mathbf{x}) \equiv p(m_k = 1|\mathbf{x})$ and remains unchanged during the model fitting process. Hence, the model structure is fully specified by the number of experts K and their matching functions $\mathbf{M} = \{m_k\}$, that is, $\mathcal{M} = \{K, \mathbf{M}\}$.

We enforce matching by defining z_k conditioned on m_k to be

$$p(z_k = 1|m_k, \mathbf{x}, \mathbf{v}_k) \propto \begin{cases} \exp(\mathbf{v}_k^T \vartheta(\mathbf{x})) & \text{if } m_k = 1, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

effectively modelling the generative probability for expert k by a generalised linear model of $\vartheta(\mathbf{x})$ if it matches, and locking it to 0 if it does not. The transfer function ϑ over the input vectors is an additional generalisation over the original MoE model, which currently uses $\vartheta(\mathbf{x}) = \mathbf{x}$.

To get the generating probability g_k we marginalise over all m_k and add the normalising term, resulting in

$$\begin{aligned} g_k(\mathbf{x}) &= \sum_{m \in \{0,1\}} p(z_k = 1, m_k = m|\mathbf{x}, \mathbf{v}_k) \\ &= \frac{m_k(\mathbf{x}) \exp(\mathbf{v}_k^T \vartheta(\mathbf{x}))}{\sum_j m_j(\mathbf{x}) \exp(\mathbf{v}_j^T \vartheta(\mathbf{x}))}. \end{aligned} \quad (2)$$

Note that this is the well-known softmax function with augmenting matching functions m_k , which weights the output of the generalised linear model by the degree of matching of the corresponding expert. While we have generalised the MoE model by introducing expert localisation through matching and the additional moderating function ϑ , its original formulation can be recovered by simply setting $m_k(\mathbf{x}) = 1$ and $\vartheta(\mathbf{x}) = \mathbf{x}$ for all \mathbf{x} and k , which causes each expert to match all observations.

2.2 The Data Likelihood

Given an input \mathbf{x} , we stochastically choose among the experts that match this input to generate the observed output \mathbf{y} . This becomes more obvious when using the 1-of- K structure of \mathbf{z} to write

$$p(\mathbf{y}|\mathbf{x}, \mathbf{z}, \boldsymbol{\theta}) = \prod_{k=1}^K p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_k)^{z_k}. \quad (3)$$

Here we do not need to explicitly consider matching, as by Eq. (1) we can only have $z_k = 1$ if $m_k = 1$.

To get the likelihood of a single observation $\{\mathbf{x}, \mathbf{y}\}$ we marginalise over the gating variables \mathbf{z} to obtain

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) &= \sum_{\mathbf{z}} \prod_{k=1}^K p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_k)^{z_k} \prod_{j=1}^K p(z_j = 1|\mathbf{x}, \mathbf{v}_j)^{z_j} \\ &= \sum_{k=1}^K g_k(\mathbf{x}) p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_k). \end{aligned} \quad (4)$$

This shows that the conditional distribution of \mathbf{y} given \mathbf{x} is a mixture distribution of the expert models weighted by the gating network.

For a set of N i.i.d. inputs $\mathbf{X} = \{\mathbf{x}_n\}$ and the corresponding outputs $\mathbf{Y} = \{\mathbf{y}_n\}$, the likelihood is given by

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_n p(\mathbf{y}_n|\mathbf{x}_n, \boldsymbol{\theta}). \quad (5)$$

where every input/output pair has a set of latent gating variables $\mathbf{z}_n = \{z_{nk}\}$ associated with it.

2.3 Model Training

In the original MoE model, the parameters of the experts and the gating network are found by maximising the likelihood by use of the EM-algorithm, where in the E-step one finds the posterior over the gating variables $\{z_n\}$ and in the M-step the complete-data likelihood is maximised with respect to the expert and gating parameters. Independent of the type of experts, the gating parameters are trained by the iteratively re-weighted least-squares (IRLS) algorithm, which can be easily adjusted to handle Eq. (2) of the generalised model.

However, the power of the generalisation lies in the localisation of the experts expressed by the matching function. Thus, in addition to fitting the model to the data, we also want to identify adequate localisation of the experts. As this localisation is fixed for a single model \mathcal{M} , we need to be able to i) compare different models and to ii) efficiently search for better ones.

The problem with the maximum likelihood approach is its tendency to over-fit the data, which causes it to always prefer more complex models over simpler ones, making it inappropriate for model selection. Embedding the model into a Bayesian framework, on the other hand, avoids over-fitting and allows us to derive an expression for the probability of a model given the data, which we can use to compare different models. Even though a fully Bayesian treatment of the generalised MoE is not analytically tractable, we will describe in the next section how to use an approximate variational inference similar to Waterhouse et al. (1996) in order to get a quick and sufficiently accurate estimate of model fit and evidence.

While this allows us to evaluate the quality of a single model structure with respect to the data, we also need to find a good model structure by searching the potentially complex space $\{\mathcal{M}\}$. LCS approach this by evaluating the current model structure and stochastically adding or removing experts or changing their matching function, thus traversing the space of possible models. This requires a frequent re-evaluation of the model evidence for a changed model structure, and how to achieve this efficiently is discussed in Section 4.

3 A Bayesian Approach

The Bayesian approach tries not to estimate the parameters by maximum likelihood, but rather by estimating a posterior distribution by conditioning a prior distribution on the available data. In that way it also mitigates over-fitting of the data because the model parameters are integrated out.

Additionally, if we consider the model structure \mathcal{M} as a random variable, we can find the posterior $p(\mathcal{M}|\mathbf{X}, \mathbf{Y})$ and improve the model structure by maximising this posterior. In order to do so, we can observe that

$$p(\mathcal{M}|\mathbf{X}, \mathbf{Y}) \propto p(\mathbf{Y}|\mathbf{X}, \mathcal{M})p(\mathcal{M}), \quad (6)$$

and hence are only required to find the model evidence $p(\mathbf{Y}|\mathbf{X}, \mathcal{M})$ and specify some prior distribution $p(\mathcal{M})$ over the possible model structures. The latter can be used to deal with symmetries in the model parameterisation, as we will see in Section 5.

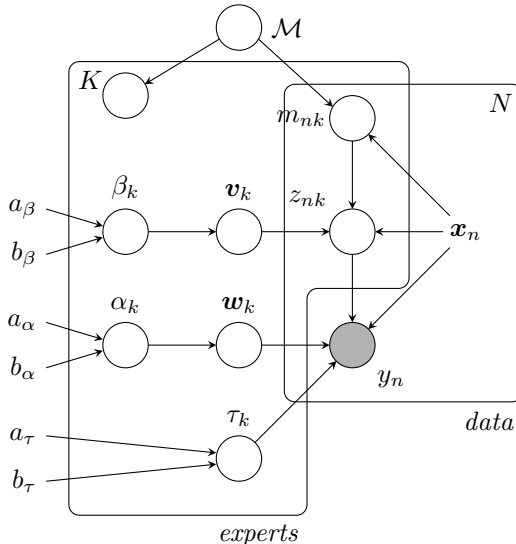


Figure 1: Directed Acyclic Graph showing the Bayesian generalised MoE. The circular nodes are random variables, which are observed when shaded. Labels without nodes are constants. The boxes are “plates”, comprising replicas of the entities inside them. Note that K is a random variable that is conditional on the model structure \mathcal{M} .

3.1 Expert Structure and Priors

For the rest of the paper we will assume univariate linear regression experts, but the method can be easily adapted to multivariate regression or binomial/multinomial classification experts. Additionally, we assume that the input vector \mathbf{x} is augmented by the constant term 1, and therefore automatically introduces the bias term in the regression model. The conditional distribution of \mathbf{y} given \mathbf{x} for expert k is given by

$$p(y|\mathbf{x}, \boldsymbol{\theta}_k) = \mathcal{N}(y|\mathbf{w}_k^T \mathbf{x}, \tau_k^{-1}), \quad (7)$$

where $\boldsymbol{\theta}_k = \{\mathbf{w}_k, \tau_k\}$, \mathbf{w}_k is the weight vector, and τ_k is the noise precision (inverse variance) of the observations.

Adopting priors similar to Bishop and Svensén (2003), the prior on \mathbf{w}_k is conjugate Gaussian, and given by

$$p(\mathbf{w}_k|\alpha_k) = \mathcal{N}(\mathbf{w}_k|0, \alpha_k^{-1} \mathbf{I}) \quad (8)$$

for each expert separately, where the precision hyper-parameter α_k determines the shrinkage on \mathbf{w}_k and \mathbf{I} is the identity matrix. Similarly, for parameters $\{\mathbf{v}_k\}$ of the gating network we define the Gaussian priors

$$p(\mathbf{v}_k|\beta_k) = \mathcal{N}(\mathbf{v}_k|0, \beta_k^{-1} \mathbf{I}), \quad (9)$$

with precision hyper-parameter β_k . The noise precisions $\boldsymbol{\tau} = \{\tau_k\}$ and hyper-parameters $\boldsymbol{\alpha} = \{\alpha_k\}$ and $\boldsymbol{\beta} = \{\beta_k\}$ get assigned conjugate Gamma priors, given by

$$p(\tau_k) = \text{Gam}(\tau_k|a_\tau, b_\tau), \quad (10)$$

$$p(\alpha_k) = \text{Gam}(\alpha_k|a_\alpha, b_\alpha), \quad (11)$$

$$p(\beta_k) = \text{Gam}(\beta_k|a_\beta, b_\beta). \quad (12)$$

We have used $a_\alpha = a_\beta = a_\tau = 10^{-2}$ and $b_\alpha = b_\beta = b_\tau = 10^{-4}$ to get sufficiently broad priors and hyper-priors. As the matching variables $\mathbf{M} = \{m_k\}$ are fixed for each model \mathcal{M} , we do not need to specify priors on them. The directed probabilistic graphical model that describes this Bayesian structure is shown in Figure 1.

While the model is not directly analytically tractable, we will follow the same approach as in Waterhouse et al. (1996) or Bishop and Svensén (2003), using variational Bayesian inference to find an approximate posterior and model evidence.

3.2 Variational Bayesian Inference

Our goal is on one hand to find the variational distribution $q(\mathbf{U})$ that approximates the true posterior $p(\mathbf{U}|\mathbf{Y})$ and on the other hand to get the model evidence $p(\mathbf{Y})$, where $\mathbf{U} = \{\mathbf{W}, \boldsymbol{\tau}, \boldsymbol{\alpha}, \mathbf{Z}, \mathbf{V}, \boldsymbol{\beta}\}$ is the set of hidden variables, and all distributions are implicitly conditional on \mathbf{X} and \mathcal{M} . Variational Bayesian inference is based on the decomposition (Bishop, 2006)

$$\ln p(\mathbf{Y}) = \mathcal{L}(q) + \text{KL}(q\|p), \quad (13)$$

$$\mathcal{L}(q) = \int q(\mathbf{U}) \ln \frac{p(\mathbf{U}, \mathbf{Y})}{q(\mathbf{U})} d\mathbf{U}, \quad (14)$$

$$\text{KL}(q\|p) = - \int q(\mathbf{U}) \ln \frac{p(\mathbf{U}|\mathbf{Y})}{q(\mathbf{U})} d\mathbf{U}, \quad (15)$$

which holds for any choice of q . As the Kullback-Leibler divergence $\text{KL}(q\|p)$ is always non-negative, and 0 if and only if $p(\mathbf{U}|\mathbf{Y}) = q(\mathbf{U})$, $\mathcal{L}(q)$ is a lower bound on $\ln p(\mathbf{Y})$ and only equivalent to the latter if $q(\mathbf{U})$ is the true posterior $p(\mathbf{U}|\mathbf{Y})$. Hence, we can approximate the posterior by maximising the lower bound $\mathcal{L}(q)$, which brings the variational distribution closer to the true posterior and at the same time gives us an approximation of the model evidence.

3.3 Factorised Distributions

To make this approach tractable, we need to choose a family of distributions $q(\mathbf{U})$ that gives an analytical solution. A frequently used approach (for example, (Bishop & Svensén, 2003; Waterhouse et al., 1996)) that is still sufficiently flexible to give a good approximation to the true posterior is to use the set of distributions that factorises with respect to disjoint groups \mathbf{U}_i of variables

$$q(\mathbf{U}) = \prod_i q_i(\mathbf{U}_i), \quad (16)$$

which allows us to maximise $\mathcal{L}(q)$ with respect to each hidden variable separately while keeping the other ones fixed. This results in

$$\ln q_i^*(\mathbf{U}_i) = \mathbb{E}_{i \neq j}(\ln p(\mathbf{U}, \mathbf{Y})) + \text{const.}, \quad (17)$$

when maximising with respect to \mathbf{U}_i , where the expectation is taken with respect to all hidden variables except for \mathbf{U}_i and the constant is the logarithm of the normalisation constant of q_i^* .

3.4 Handling the Softmax Function

If the model has a conjugate-exponential structure, Eq. (17) gives an analytical solution with a distribution form equal to the prior on the corresponding hidden variable. However, in our case the softmax function in Eq. (2) does not conform to the conjugate-exponential structure, and we need to deal with it separately. A possible approach is to replace the softmax function by an exponential lower bound on it, which consequently introduces additional variational variables with respect to which $\mathcal{L}(q)$ also needs to be maximised. This approach was followed in Bishop and Svensén (2003) for the logistic sigmoid function, but currently there is no known exponential lower bound function on the softmax besides a conjectured one in Gibbs (1997)³. Alternatively, we can follow the approach taken in Waterhouse (1997), where $q_V^*(\mathbf{V})$ is approximated by a Laplace approximation $\tilde{q}_V^*(\mathbf{V})$. Despite such approximation invalidating the lower bound nature of $\mathcal{L}(q)$, we have chosen to use it due to the lack of better alternatives.

3.5 Update Equations and Model Posterior

To get the variational update equations, we need to evaluate Eq. (17) for each hidden variable in \mathbf{U} separately, similarly to the derivations in Waterhouse et al. (1996; 1997) and (Ueda & Ghahramani, 2002). All update equations are derived in the appendix.

Similarly, we can find a closed-form expression for $\mathcal{L}(\tilde{q})$ by evaluating Eq. (14), where we can reuse many terms that have already been used for finding the variational update equations. The expression for $\mathcal{L}(\tilde{q})$ is derived similarly to those found in Ueda and Ghahramani (2002) and Bishop (2006) and

³A more general bound was recently developed in Wainwright et al. (2005), but its applicability still needs to be evaluated.

is presented in the appendix. As the variational bound $\mathcal{L}_{\mathcal{M}}(\tilde{q})$ for model \mathcal{M} is an approximate lower bound on the logarithm of the model evidence $\ln p(\mathbf{Y}|\mathbf{X}, \mathcal{M})$, we can find a closed-form approximation to the model posterior Eq. (6) by

$$\ln \tilde{p}(\mathcal{M}|\mathbf{X}, \mathbf{Y}) \geq \mathcal{L}_{\mathcal{M}}(\tilde{q}) + \ln p(\mathcal{M}) + \text{const.}, \quad (18)$$

which we will use to compare different models with the aim of finding better model structures.

4 Searching the Model Structure Space

A model structure \mathcal{M} comprises the number of experts K and their matching functions \mathbf{M} . While these functions could potentially be arbitrarily defined, they are usually parametric to keep the search tractable. Nonetheless, even in the light of such restrictions, model structure search is still a complex task. For example, given 10-dimensional input vectors and matching functions that have 2 scalar parameters (like some specification of an interval) per dimension we have a 21-dimensional, possibly multi-modal, model structure space. Hence, model structure search is potentially computationally expensive, which we need to counter-balance by making the evaluation of a single model structure cheaper.

Once a model structure $\mathcal{M} = (K, \mathbf{M})$ is specified, we need to find its posterior approximation $\tilde{p}(\mathcal{M}|\mathbf{X}, \mathbf{Y})$ by Eq. (18), which involves evaluating the variational bound $\mathcal{L}_{\mathcal{M}}(\tilde{q})$ and, hence, fitting the model to the data. In this process it is important to obtain good solutions to the variational equations by avoiding poor local minima, of which there are many in the MoE model (Bishop & Svensén, 2003). Repeated model fitting with randomised initial conditions is not an option in our case, as i) it is computationally too expensive, and ii) it does not guarantee finding the global optimum. Hence, we will use a different training strategy to bypass the issue of local optima in fitting the model.

4.1 Independent Expert Training

The MoE method achieves localisation of the experts by interdependent training of experts and the gating network. This becomes apparent by inspecting the maximised variational distribution after Eq. (17) for the expert weight vectors \mathbf{W} , which factorises with respect to the experts, and for expert k is given by

$$\begin{aligned} q_{\mathbf{W}}^*(\mathbf{w}_k) &= \mathcal{N}(\mathbf{w}_k | \mathbf{w}_k^*, \boldsymbol{\Sigma}_k^*), \\ \boldsymbol{\Sigma}_k^* &= \left(\mathbb{E}_{\alpha}(\alpha_k) \mathbf{I} + \mathbb{E}_{\tau}(\tau_k) \sum_n r_{nk} \mathbf{x}_n \mathbf{x}_n^T \right)^{-1}, \\ \mathbf{w}_k^* &= \mathbb{E}_{\tau}(\tau_k) \boldsymbol{\Sigma}_k^* \sum_n r_{nk} \mathbf{x}_n y_n, \end{aligned} \quad (19)$$

where r_{nk} is the *responsibility* of expert k for observation n , given by $r_{nk} = \mathbb{E}_Z(z_{nk})$. Hence, the experts are trained by weighted linear regression with a Gaussian shrinkage prior, where the observations are weighted according to the expert's responsibilities. These are found by evaluating $q_Z^*(\mathbf{Z})$, and reveal that the gating network distributes the responsibilities according to the goodness-of-fit of the experts, which depends on the responsibilities of the last iteration. While these interdependencies make localisation of the experts possible, they are also the main cause of the highly multi-modal structure of the variational bound that we want to maximise.

In our generalisation of the MoE model we provide a second layer of localisation that is determined by the model structure rather than applied while fitting it to the data. Hence, we can remove the interdependence of expert and gating network training by replacing the responsibilities r_{nk} in the expert training with the matching functions $m_k(\mathbf{x}_n)$. From the generative point-of-view, this causes a shift from the assumption that each observation was generated by one and only one expert to the idea that the observation is generated by some mixture of independent processes, as a single observation can now be fully attributed to several experts at once (Note that matching does not underlie the 1-of- K assumption). As a result, the role of the gating network is now to combine the experts to best explain the set of observations, rather than to assign each observation to one and only one expert. This brings the whole model interpretation closer to ensemble methods.

For fitting a certain model structure the modification has the important consequence that the experts can now be trained independently of the gating network. Hence, if new experts are added in the process of searching for better model structures, only the new experts need to be fitted in addition

to updating the gating network. When removing experts, the only task is to re-fit the gating network. This reduces the computational complexity of evaluating a model structure immensely.

On the downside, the quality of the fit is reduced if many experts are localised in the same area of the input space. While interdependent expert and gating network training leads to at least a local optimum of the variational bound $\mathcal{L}_{\mathcal{M}}(q)$ of that model structure, removing the interdependency by fixing the responsibilities in expert training causes a worse fit due to the lack of feedback from the gating networks to the experts. Therefore, the modified training scheme shifts the emphasis from expert localisation by the gating network to localisation by the matching functions, and adds additional responsibility to the model search mechanism.

While the change in training policy causes a worse fit of the model when several experts are localised in the same regions of the input space, a locally mostly disjoint set of experts behaves similarly to the original MoE, but with higher degrees of freedom in how to choose the localisation.

4.2 Model Search in LCS

In Learning Classifier Systems, the space of possible models is searched by adding, removing or replacing single experts, based on some quality metric, such as, for example, the noise precision for regression experts. The space is traversed stochastically with a rather complex genetic algorithm that heuristically balances over- and underfitting of the experts and aims at finding an adequate coverage of the input space. The model quality is determined by the interaction between the algorithm and various system parameters, and there is currently no explicit quality metric for a model.

While we provide such a quality metric by the Bayesian model posterior, adjusting the genetic algorithm to use this metric is clearly outside of the scope of this paper, and a topic of further research. Instead, we provide a simple demonstration of how the model space can be searched by means of MCMC.

5 Sampling the Model Posterior

As a proof-of-concept on how to train a model and how to search the space of possible model structures, we have designed a simple Metropolis-Hastings algorithm that samples from the model posterior, similar to the one used for CART model search in Chipman et al. (1998).

5.1 Sampling by Metropolis-Hastings

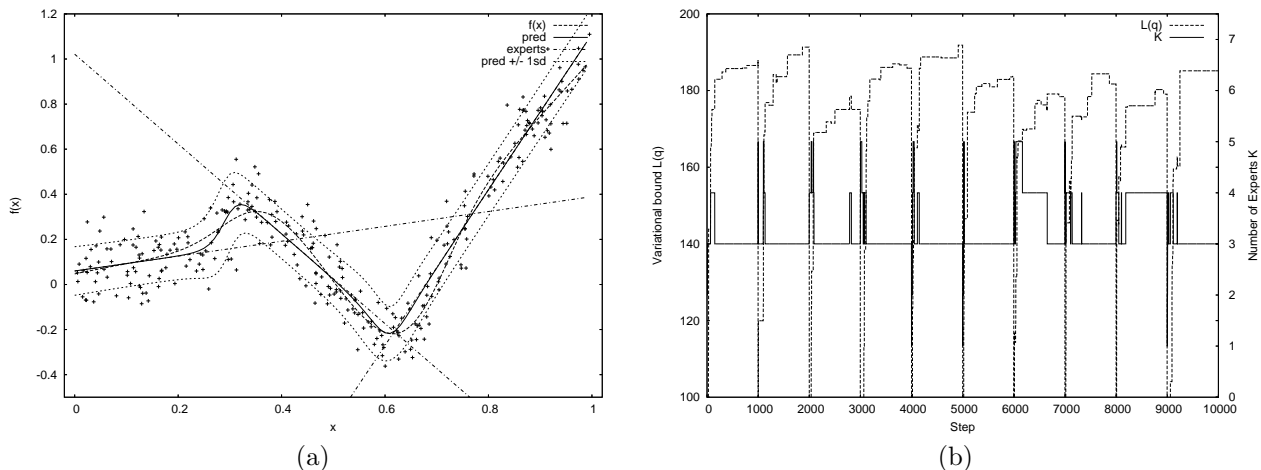


Figure 2: (a) shows the original function $f_1(x)$, generated by the mixture of 3 localised experts, and the data points used for training. The search method correctly identifies all 3 experts, which are shown by the dashed lines, with their localisation and mixed prediction. The error bars show one standard deviation from the model prediction. (b) shows how $\mathcal{L}(\hat{q})$ and K change during the sampling process.

The algorithm starts with an initial model structure \mathcal{M}_0 at $t = 0$, and then proceeds by iterating over the following steps:

1. Sample a candidate model structure \mathcal{M}^* from the Markov chain $p(\mathcal{M}^*|\mathcal{M}_t)$.

2. Accept the candidate model structure by setting $\mathcal{M}_{t+1} = \mathcal{M}^*$ with probability

$$\min \left(\frac{p(\mathcal{M}_t | \mathcal{M}^*) \tilde{p}(\mathcal{M}^* | \mathbf{X}, \mathbf{Y})}{p(\mathcal{M}^* | \mathcal{M}_t) \tilde{p}(\mathcal{M}_t | \mathbf{X}, \mathbf{Y})}, 1 \right), \quad (20)$$

where Eq. (18) is used to compute the model posteriors, or otherwise reject it by setting $\mathcal{M}_{t+1} = \mathcal{M}_t$.

Following this algorithm causes the sequence $\mathcal{M}_0, \mathcal{M}_1, \dots$ to approach the approximated model posterior distribution $\tilde{p}(\mathcal{M} | \mathbf{X}, \mathbf{Y})$, and therefore we will draw more samples from models with a high posterior probability.

We generate the Markov chain $p(\mathcal{M}^* | \mathcal{M}_t)$ by choosing randomly from one of the following actions: **Change:** Pick one expert uniformly at random and randomly re-initialise its matching function parameters; **Add:** Add one expert and randomly initialise its matching function parameters; **Remove:** Pick one expert uniformly at random and remove it. While these actions are rather simplistic and do not exploit the information that is available in the model, they were chosen such that the reverse transition probabilities $p(\mathcal{M}_t | \mathcal{M}^*)$ are easy to evaluate. In fact, for the *Change* action, the reverse transition equals the forward transition probability, and for the *Add* and *Remove* actions there is significant cancellation in the computation of the acceptance probability due to their complementary nature. In all experiments we have chosen *Add* and *Remove* with a probability of 1/4, and *Change* with a probability of 1/2.

The matching function we have used is a Gaussian basis function with a diagonal covariance structure to keep the model search space small. Its matching probability for an input \mathbf{x} (not considering its bias term) is given by

$$m_k(\mathbf{x}) = \exp \left(- \sum_i \frac{1}{2\sigma_i^2} (x_i - \mu_i)^2 \right), \quad (21)$$

where x_i is the i th component of \mathbf{x} , σ_i is the i th diagonal component of the covariance matrix, and μ_i is the i th component of the mean vector $\boldsymbol{\mu}$. Such a matching function introduces symmetries in the model space which we need to account for in the model prior $p(\mathcal{M})$ to avoid a bias towards models with a higher number of experts. Given, for example, a model with K experts, there are $K!$ possible permutations of the same model by rearranging the order of the experts. We have countered this effect by using the model prior $p(\mathcal{M}) \propto 1/K_{\mathcal{M}}!$ when evaluating Eq. (18).

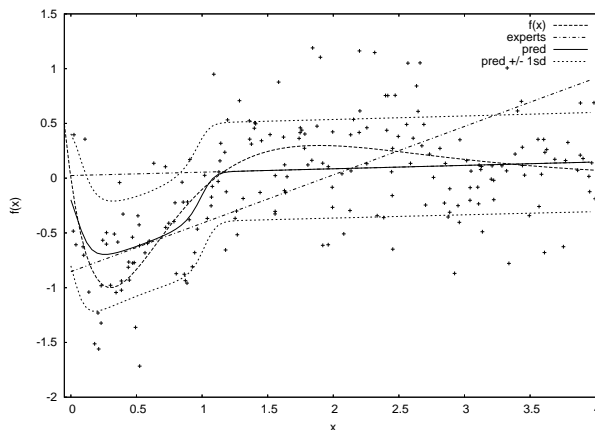


Figure 3: The original function $f_2(x)$, the training data with added noise, the prediction of the separate experts and the mixed prediction. The error bars show a single standard deviation from the model prediction.

5.2 Experiments

To evaluate the method's ability to identify the correct model, we have generated a 1-dimensional function with added Gaussian noise from a set of 3 linear experts, and have used the above procedure to identify the correct number of experts and matching function parameters, with random re-initialisation of the model structure every 1000 steps to escape local peaks in the model posterior. As done in LCS,

we completely rely on the model search for localisation by setting the gating features to $\vartheta(\mathbf{x}) = 1$, such that the gating network only performs a non-localised re-weighting of the experts. As can be seen from the results in Figure 2(a), the MCMC method with Bayesian model selection was able to correctly identify both the number of experts and their localisation in the input space. Figure 2(b) shows the change of the variational bound $\mathcal{L}(q)$ and the current number of experts while searching the model space, and demonstrates that the method quickly finds the correct number of experts after almost every random restart.

In an additional experiment, we have tested the performance on an artificial data set used in Waterhouse et al. (1996) by sampling from $f_2(x) = 4.26(e^{-x} - 4e^{-2x} + 3e^{-3x}) + \epsilon$ over the range $[0, 4]$, where ϵ is Gaussian noise with variance 0.44. With such a high variance no pattern was apparent, and the best model contained only a single expert. Upon reducing the noise variance to 0.2, two experts were identified as shown in Figure 3. Surprisingly, the curve was fitted more compactly with one local expert interleaving the prediction of another global expert, where we would have expected the use of 3 experts that perform a piecewise linear fit.

While the MCMC model search was able to find good models for low-dimensional data, we quickly reached its limits when increasing the dimensionality. The commonly observed pattern was the removal of all except for a single expert, and maintaining that expert. This behaviour is explained by the very low probability of finding good matching functions with random sampling in a high-dimensional model space. Additionally, the underlying generative Markov chain causes two successive models to be highly correlated and therefore does not support large changes in the model: Given that the single expert's matching is well tuned to the data, then the probability of adding another expert that in combination with the previous one performs better is highly unlikely.

Clearly, to improve model search we need to either design a sampling algorithm that uses more of the information on the structure of the data that is available in the model, or redesign the genetic algorithm of current LCS to make use of the model posterior approximate. Either approach is left as future work.

6 Discussion and Conclusion

By introducing a generalisation to the MoE model we have for the first time explicitly identified the model that underlies the LCS method. Compared to the original MoE, LCS allow for more complex localisation pattern of the experts, only limited by the choice of representation for the matching functions. This increased flexibility comes at the cost of having to perform a search in the potentially complex model structure space, which is made computationally cheaper by introducing a training scheme that makes the expert training independent of the gating network.

The link between MoE and LCS immediately enables several improvements to LCS: i) while matching is usually binary in LCS, we can now match to a degree by specifying the probability of matching $m_k(\mathbf{x})$; ii) the experts were so far restricted to linear regression models, but can now be easily extended to any generalised linear model; iii) we do not need to use $\vartheta(\mathbf{x}) = 1$, which allows the gating network to combine equally localised experts better; iv) Bayesian model selection replaces the heuristics in LCS. This clearly opens up a wide range of future work, some of which has been identified throughout the paper.

A Variational Distributions

In order to derive an expression for the variational distributions over the hidden random variables $\mathbf{U} = \{\mathbf{W}, \boldsymbol{\tau}, \boldsymbol{\alpha}, \mathbf{Z}, \mathbf{V}, \boldsymbol{\beta}\}$, we need to evaluate Eq. (17) for each element of \mathbf{U} , and the necessary moments of these. In all derivations we drop the terms that are independent of the hidden variable in question, making use of the factorisation

$$p(\mathbf{U}, \mathbf{Y}) = p(\mathbf{Y}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\tau})p(\mathbf{W}|\boldsymbol{\alpha})p(\boldsymbol{\tau})p(\boldsymbol{\alpha})p(\mathbf{Z}|\mathbf{V})p(\mathbf{V}|\boldsymbol{\beta})p(\boldsymbol{\beta}). \quad (22)$$

All distributions are implicitly conditional on \mathbf{X} and \mathcal{M} .

A.1 Expert weights $q_{\mathbf{W}}(\mathbf{W})$

We need to evaluate

$$\begin{aligned} \ln q_{\mathbf{W}}^*(\mathbf{W}) &= \mathbb{E}_{\boldsymbol{\tau}, \boldsymbol{\alpha}, \mathbf{Z}, \mathbf{V}, \boldsymbol{\beta}}(\ln p(\mathbf{Y}, \mathbf{U})) + \text{const.} \\ &= \mathbb{E}_{\boldsymbol{\tau}, \mathbf{Z}}(\ln p(\mathbf{Y}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\tau})) + \mathbb{E}_{\boldsymbol{\alpha}}(\ln p(\mathbf{W}|\boldsymbol{\alpha})) + \text{const.}, \end{aligned}$$

using Eqs. (3), (7), and (8) to get

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\tau}, \mathbf{Z}}(\ln p(\mathbf{Y}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\tau})) &= \sum_n \sum_k \mathbb{E}_{\mathbf{Z}}(z_{nk}) \mathbb{E}_{\boldsymbol{\tau}}(\ln \mathcal{N}(y_n | \mathbf{w}_k^T \mathbf{x}_n, \tau_k^{-1})) \\ &= \sum_n \sum_k \mathbb{E}_{\mathbf{Z}}(z_{nk}) \mathbb{E}_{\boldsymbol{\tau}}\left(\frac{\tau_k}{2}(y_n - \mathbf{w}_k^T \mathbf{x}_n)^2\right) + \text{const.} \\ &= \sum_k \left(-\frac{\mathbb{E}_{\boldsymbol{\tau}}(\tau_k)}{2} \mathbf{w}_k^T \sum_n (\mathbb{E}_{\mathbf{Z}}(z_{nk}) \mathbf{x}_n \mathbf{x}_n^T) \mathbf{w}_k + \mathbb{E}_{\boldsymbol{\tau}}(\tau_k) \mathbf{w}_k^T \sum_n \mathbb{E}_{\mathbf{Z}}(z_{nk}) \mathbf{x}_n y_n \right) \\ &\quad + \text{const.}, \\ \mathbb{E}_{\boldsymbol{\alpha}}(\ln p(\mathbf{W}|\boldsymbol{\alpha})) &= \sum_k \mathcal{N}(\mathbf{w}_k | 0, \alpha_k^{-1} \mathbf{I}) \\ &= \sum_k \left(-\frac{\mathbb{E}_{\boldsymbol{\alpha}}(\alpha_k)}{2} \mathbf{w}_k^T \mathbf{w}_k \right) + \text{const.} \\ &= \sum_k \left(-\frac{1}{2} \mathbf{w}_k^T (\mathbb{E}_{\boldsymbol{\alpha}}(\alpha_k) \mathbf{I}) \mathbf{w}_k \right) + \text{const.} \end{aligned}$$

Hence, $q_{\mathbf{W}}^*(\mathbf{W})$ factorises w.r.t. k , which allows us to treat each $q_{\mathbf{W}}^*(\mathbf{w}_k)$ separately, resulting in

$$\begin{aligned} \ln q_{\mathbf{W}}^*(\mathbf{w}_k) &= -\frac{1}{2} \mathbf{w}_k^T \left(\mathbb{E}_{\boldsymbol{\alpha}}(\alpha_k) \mathbf{I} + \mathbb{E}_{\boldsymbol{\tau}}(\tau_k) \sum_n \mathbb{E}_{\mathbf{Z}}(z_{nk}) \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w}_k \\ &\quad + \mathbb{E}_{\boldsymbol{\tau}}(\tau_k) \mathbf{w}_k^T \sum_n \mathbb{E}_{\mathbf{Z}}(z_{nk}) \mathbf{x}_n y_n + \text{const.} \end{aligned}$$

By completing the square we can see that hence $q_{\mathbf{W}}^*(\mathbf{w}_k)$ is a Gaussian distribution given by

$$\begin{aligned} q_{\mathbf{W}}^*(\mathbf{w}_k) &= \mathcal{N}(\mathbf{w}_k | \mathbf{w}_k^*, \boldsymbol{\Sigma}_k^*), \\ \boldsymbol{\Sigma}_k^* &= \left(\mathbb{E}_{\boldsymbol{\alpha}}(\alpha_k) \mathbf{I} + \mathbb{E}_{\boldsymbol{\tau}}(\tau_k) \sum_n \mathbb{E}_{\mathbf{Z}}(z_{nk}) \mathbf{x}_n \mathbf{x}_n^T \right)^{-1}, \\ \mathbf{w}_k^* &= \mathbb{E}_{\boldsymbol{\tau}}(\tau_k) \boldsymbol{\Sigma}_k^* \sum_n \mathbb{E}_{\mathbf{Z}}(z_{nk}) \mathbf{x}_n y_n. \end{aligned} \quad (23)$$

A.2 Expert precisions $q_{\boldsymbol{\tau}}(\boldsymbol{\tau})$

We require

$$\begin{aligned} \ln q_{\boldsymbol{\tau}}^*(\boldsymbol{\tau}) &= \mathbb{E}_{\mathbf{W}, \boldsymbol{\alpha}, \mathbf{Z}, \mathbf{V}, \boldsymbol{\beta}}(\ln p(\mathbf{Y}, \mathbf{U})) + \text{const.} \\ &= \mathbb{E}_{\mathbf{W}, \mathbf{Z}}(\ln p(\mathbf{Y}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\tau})) + \ln p(\boldsymbol{\tau}) + \text{const.} \end{aligned}$$

Using Eqs. (3), (7), and (10) we get

$$\begin{aligned}
 \mathbb{E}_{W,Z}(\ln p(\mathbf{Y}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\tau})) &= \sum_n \sum_k \mathbb{E}_Z(z_{nk}) \mathbb{E}_W(\ln \mathcal{N}(y_n | \mathbf{w}_k^T \mathbf{x}_n, \tau_k^{-1})) \\
 &= \sum_n \sum_k \mathbb{E}_Z(z_{nk}) \left(\frac{1}{2} \ln \tau_k - \frac{\tau_k}{2} \mathbb{E}_W((y_n - \mathbf{w}_k^T \mathbf{x}_n)^2) \right) + \text{const.}, \\
 \ln p(\boldsymbol{\tau}) &= \sum_k \ln \text{Gam}(\tau_k | a_\tau, b_\tau) \\
 &= \sum_k ((a_\tau - 1) \ln \tau_k - b_\tau \tau_k) + \text{const.}
 \end{aligned}$$

As for $q_W^*(\mathbf{W})$, this distribution factorises w.r.t. k , and therefore we can evaluate each $q_\tau^*(\tau_k)$ separately, resulting in

$$\begin{aligned}
 \ln q_\tau^*(\tau_k) &= \left(\frac{1}{2} \sum_n \mathbb{E}_Z(z_{nk}) + a_\tau - 1 \right) \ln \tau_k \\
 &\quad + \tau_k \left(b_\tau + \frac{1}{2} \sum_n \mathbb{E}_Z(z_{nk}) \mathbb{E}_W((y_n - \mathbf{w}_k^T \mathbf{x}_n)^2) \right) + \text{const...},
 \end{aligned}$$

which is a Gamma distribution of the form

$$\begin{aligned}
 q_\tau(\tau_k) &= \text{Gam}(\tau_k | a_{\tau_k}^*, b_{\tau_k}^*), \\
 a_{\tau_k}^* &= a_\tau + \frac{1}{2} \sum_n \mathbb{E}_Z(z_{nk}), \\
 b_{\tau_k}^* &= b_\tau + \frac{1}{2} \sum_n \mathbb{E}_Z(z_{nk}) \mathbb{E}_W((y_n + \mathbf{w}_k^T \mathbf{x}_n)^2).
 \end{aligned} \tag{24}$$

A.3 Expert weight priors $q_\alpha(\boldsymbol{\alpha})$

We need to evaluate

$$\begin{aligned}
 \ln q_\alpha^*(\boldsymbol{\alpha}) &= \mathbb{E}_{W,\tau,Z,V,\beta}(\ln p(\mathbf{Y}, \mathbf{U})) + \text{const.} \\
 &= \mathbb{E}_W(\ln p(\mathbf{W}, \boldsymbol{\alpha})) + \ln p(\boldsymbol{\alpha}) + \text{const.}
 \end{aligned}$$

Using Eqs. (8) and (11) we get

$$\begin{aligned}
 \mathbb{E}_W(\ln p(\mathbf{W}, \boldsymbol{\alpha})) &= \sum_k \mathbb{E}_W(\ln \mathcal{N}(\mathbf{w}_k | 0, \alpha_k^{-1} \mathbf{I})) \\
 &= \sum_k \left(\frac{D_w}{2} \ln \alpha_k - \frac{\alpha_k}{2} \mathbb{E}_W(\mathbf{w}_k^T \mathbf{w}_k) \right) + \text{const.}, \\
 \ln p(\boldsymbol{\alpha}) &= \sum_k \ln \text{Gam}(\alpha_k | a_\alpha, b_\alpha) \\
 &= \sum_k ((a_\alpha - 1) \ln \alpha_k - \alpha_k b_\alpha) + \text{const.},
 \end{aligned}$$

where D_w is the size of the weight vector \mathbf{w}_k . Again, this distribution factorises w.r.t. k , and so we get for $q_\alpha^*(\alpha_k)$

$$\ln q_\alpha^*(\alpha_k) = \left(\frac{D_w}{2} + a_\alpha - 1 \right) \ln \alpha_k - \alpha_k \left(b_\alpha + \frac{1}{2} \mathbb{E}_W(\mathbf{w}_k^T \mathbf{w}_k) \right) + \text{const.},$$

which is the Gamma distribution

$$\begin{aligned}
 q_\alpha^*(\alpha_k) &= \text{Gam}(\alpha_k | a_{\alpha_k}^*, b_{\alpha_k}^*), \\
 a_{\alpha_k}^* &= a_\alpha + \frac{D_w}{2}, \\
 b_{\alpha_k}^* &= b_\alpha + \frac{1}{2} \mathbb{E}_W(\mathbf{w}_k^T \mathbf{w}_k).
 \end{aligned} \tag{25}$$

A.4 Gating weights $q_V(\mathbf{V})$

To get $q_V^*(\mathbf{V})$ we need to evaluate

$$\begin{aligned}\ln q_V^*(\mathbf{V}) &= \mathbb{E}_{W,\tau,\alpha,Z,\beta}(\ln p(\mathbf{Y}, \mathbf{U})) + \text{const.} \\ &= \mathbb{E}_Z(\ln p(\mathbf{Z}|\mathbf{V})) + \mathbb{E}_\beta(\ln p(\mathbf{V}|\boldsymbol{\beta})) + \text{const.}\end{aligned}$$

Using Eq. (12) we get

$$\begin{aligned}\mathbb{E}_\beta(\ln p(\mathbf{V}|\boldsymbol{\beta})) &= \sum_k \mathbb{E}_\beta(\ln \mathcal{N}(\mathbf{v}_k|0, \beta_k^{-1}\mathbf{I})) \\ &= \sum_k \left(-\frac{\mathbb{E}_\beta(\beta_k)}{2} \mathbf{v}_k^T \mathbf{v}_k \right) + \text{const.}\end{aligned}$$

To get $\mathbb{E}_Z(\ln p(\mathbf{Z}|\mathbf{V}))$ we can use the 1-of- K structure of \mathbf{z} , resulting in

$$p(\mathbf{Z}|\mathbf{V}) = \prod_n \prod_k p(z_{nk} = 1 | \mathbf{x}_n, \mathbf{v}_k)^{z_{nk}} = \prod_n \prod_k g_k(\mathbf{x}_n)^{z_{nk}}, \quad (26)$$

which gives

$$\mathbb{E}_Z(\ln p(\mathbf{Z}|\mathbf{V})) = \sum_n \sum_k \mathbb{E}_Z(z_{nk}) \ln g_k(\mathbf{x}_n).$$

Therefore, we get for $\ln q_V^*(\mathbf{V})$

$$\ln q_V^*(\mathbf{V}) = \sum_n \sum_k \mathbb{E}_Z(z_{nk}) \ln g_k(\mathbf{x}_n) - \sum_k \frac{\mathbb{E}_\beta(\beta_k)}{2} \mathbf{v}_k^T \mathbf{v}_k + \text{const.},$$

which is not an exponential distribution and therefore spoils our conjugate-exponential structure. Hence, we will proceed as in Waterhouse et al. (1996) and approximate q_V^* by a Laplace approximation, using a Gaussian $\mathcal{N}(\mathbf{V}|\tilde{\mathbf{V}}, \tilde{\boldsymbol{\Lambda}}_V^{-1})$ centred on the mode of q_V^* , and with a similar covariance structure. We get the mode by setting the derivative of $\ln q_V^*(\mathbf{V})$ w.r.t. \mathbf{V} to zero, that is

$$\sum_n (\mathbb{E}_Z(z_{nk}) - g_k(\mathbf{x}_n)) \vartheta(\mathbf{x}_n) - \mathbb{E}_\beta(\beta_k) \mathbf{v}_k = 0, \quad \forall k,$$

for which the solution can be obtained by the IRLS algorithm, as, for example, described in Bishop (2006).

The precision matrix $\tilde{\boldsymbol{\Lambda}}_V$ is block-symmetrical and is obtained in these blocks $(\tilde{\boldsymbol{\Lambda}}_V)_{jk}$ by

$$(\tilde{\boldsymbol{\Lambda}}_V)_{jk} = \frac{\partial \ln q_V^*(\mathbf{V})}{\partial \mathbf{v}_j \partial \mathbf{v}_k} = \left(-\sum_n g_k(\mathbf{x}_n) (\mathbf{I}_{jk} - g_j(\mathbf{x}_n)) \vartheta(\mathbf{x}_n) \vartheta(\mathbf{x}_n)^T \right) - \mathbf{I}_{jk} \mathbb{E}_\beta(\beta_k),$$

where \mathbf{I}_{jk} is the jk th element of the identity matrix, and $\tilde{\mathbf{V}}$ is used to evaluate $g_k(\mathbf{x}_n)$.

This results in the Gaussian approximation to q_V^* , given by

$$\tilde{q}_V^*(\mathbf{V}) = \mathcal{N}(\mathbf{V}|\tilde{\mathbf{V}}, \tilde{\boldsymbol{\Lambda}}_V^{-1}). \quad (27)$$

A.5 Gating weight priors $q_\beta(\boldsymbol{\beta})$

We require

$$\begin{aligned}\ln q_\beta^*(\boldsymbol{\beta}) &= \mathbb{E}_{W,\tau,\alpha,Z,V}(\ln p(\mathbf{Y}, \mathbf{U})) + \text{const.} \\ &= \mathbb{E}_V(\ln p(\mathbf{V}|\boldsymbol{\beta})) + \ln p(\boldsymbol{\beta}) + \text{const.}\end{aligned}$$

Using Eqs. (9) and (12) we get

$$\begin{aligned}\mathbb{E}_V(\ln p(\mathbf{V}|\boldsymbol{\beta})) &= \sum_k \ln \mathcal{N}(\mathbf{v}_k|0, \beta_k^{-1}\mathbf{I}) \\ &= \sum_k \left(\frac{D_v}{2} \ln \beta_k - \frac{\beta_k}{2} \mathbb{E}_V(\mathbf{v}_k^T \mathbf{v}_k) \right) + \text{const.}, \\ \ln p(\boldsymbol{\beta}) &= \sum_k \ln \text{Gam}(\beta_k | a_\beta, b_\beta) \\ &= \sum_k ((a_\beta - 1) \ln \beta_k - \beta_k b_\beta) + \text{const.},\end{aligned}$$

where D_v is the size of \mathbf{v}_k , and which factorise w.r.t. k . Hence, we have

$$\ln q_{\beta}^*(\beta_k) = \left(\frac{D_v}{2} + a_{\beta} - 1 \right) \ln \beta_k - \beta_k \left(b_{\beta} + \frac{1}{2} \mathbb{E}_V(\mathbf{v}_k^T \mathbf{v}_k) \right) + \text{const.},$$

which is the Gamma distribution

$$\begin{aligned} q_{\beta}^*(\beta_k) &= \text{Gam}(\beta_k | a_{\beta_k}^*, b_{\beta_k}^*), \\ a_{\beta_k}^* &= a_{\beta} + \frac{D_v}{2}, \\ b_{\beta_k}^* &= b_{\beta} + \frac{1}{2} \mathbb{E}_V(\mathbf{v}_k^T \mathbf{v}_k). \end{aligned} \tag{28}$$

A.6 Gating $q_Z(\mathbf{Z})$

We need to evaluate

$$\begin{aligned} \ln q_Z^*(\mathbf{Z}) &= \mathbb{E}_{W, \tau, \alpha, V, \beta}(\ln p(\mathbf{Y}, \mathbf{U})) + \text{const.} \\ &= \mathbb{E}_{W, \tau}(\ln p(\mathbf{Y} | \mathbf{Z}, \mathbf{W}, \tau)) + \mathbb{E}_V(\ln p(\mathbf{Z} | \mathbf{V})) + \text{const.} \end{aligned}$$

With the use of Eqs. (3), (7), and (26), we get

$$\begin{aligned} \mathbb{E}_{W, \tau}(\ln p(\mathbf{Y} | \mathbf{Z}, \mathbf{W}, \tau)) &= \sum_n \sum_k z_{nk} \mathbb{E}_{W, \tau}(\ln \mathcal{N}(y_n | \mathbf{w}_k^T \mathbf{x}_n, \tau_k^{-1})) \\ &= \sum_n \sum_k z_{nk} \left(\frac{1}{2} \mathbb{E}_{\tau}(\ln \tau_k) - \frac{\mathbb{E}_{\tau}(\tau_k)}{2} \mathbb{E}_W((y_n - \mathbf{w}_k^T \mathbf{x}_n)^2) \right) + \text{const.}, \\ \mathbb{E}_V(\ln p(\mathbf{Z} | \mathbf{V})) &= \sum_n \sum_k z_{nk} \mathbb{E}_V(\ln g_k(\mathbf{x}_n)). \end{aligned}$$

Hence, we have

$$\ln q_Z^*(\mathbf{Z}) = \sum_n \sum_k z_{nk} \ln \rho_{nk} + \text{const.},$$

with

$$\ln \rho_{nk} = \ln \mathbb{E}_V(\ln g_k(\mathbf{x}_n)) + \frac{1}{2} \mathbb{E}_{\tau}(\ln \tau_k) - \frac{1}{2} \mathbb{E}_{\tau}(\tau_k) \mathbb{E}_W((y_n - \mathbf{w}_k^T \mathbf{x}_n)^2).$$

It follows that

$$q_Z^*(\mathbf{Z}) \propto \prod_n \prod_k \rho_{nk}^{z_{nk}},$$

which, when normalised by $\sum_k z_{nk} = 1$, gives

$$q_Z^*(\mathbf{Z}) = \prod_n \prod_k r_{nk}^{z_{nk}}, \tag{29}$$

$$r_{nk} = \frac{\rho_{nk}}{\sum_j \rho_{nj}}. \tag{30}$$

A.7 Moments

To evaluate the parameters of the variational distributions we require several moments of other distributions that we evaluate below.

The required moments of the expert weights are based on Eq. (23) and are given by

$$\begin{aligned} \mathbb{E}_W(\mathbf{w}_k^T \mathbf{w}_k) &= \mathbf{w}_k^{*T} \mathbf{w}_k^* + \text{Tr}(\mathbf{\Sigma}_k^*), \\ \mathbb{E}_W((y_n - \mathbf{w}_k^T \mathbf{x}_n)^2) &= y_n^2 - 2 \mathbb{E}_W(\mathbf{w}_k)^T \mathbf{x}_n y_n + \mathbf{x}_n^T \mathbb{E}_W(\mathbf{w}_k \mathbf{w}_k^T) \mathbf{x}_n \\ &= y_n^2 - 2 \mathbf{w}_k^{*T} \mathbf{x}_n y_n + \mathbf{x}_n^T (\mathbf{w}_k^* \mathbf{w}_k^{*T} + \mathbf{\Sigma}_k^*) \mathbf{x}_n \\ &= \|y_n - \mathbf{w}_k^T \mathbf{x}_n\|^2 + \mathbf{x}_n^T \mathbf{\Sigma}_k^* \mathbf{x}_n. \end{aligned}$$

The moments of the expert precision Eq. (24) are

$$\begin{aligned} \mathbb{E}_{\tau}(\tau_k) &= \frac{a_{\tau_k}^*}{b_{\tau_k}^*}, \\ \mathbb{E}_{\tau}(\ln \tau_k) &= \psi(a_{\tau_k}^*) - \ln b_{\tau_k}^*, \end{aligned}$$

where $\psi(\cdot)$ is the digamma function. For the expert weight priors Eq. (25) we have the moments

$$\begin{aligned}\mathbb{E}_\alpha(\alpha_k) &= \frac{a_{\alpha_k}^*}{b_{\alpha_k}^*}, \\ \mathbb{E}_\alpha(\ln \alpha_k) &= \psi(a_{\alpha_k}^*) - \ln b_{\alpha_k}^*.\end{aligned}$$

Even though $\mathbb{E}_\alpha(\ln \alpha_k)$ is not required to find the variational distribution, we have evaluated it for use in a later section.

We get the expectation of the gating variables z_{nk} by inspection of Eq. (29) ,

$$\mathbb{E}_Z(z_{nk}) = r_{nk},$$

where r_{nk} is called the *responsibility* of expert k for observation n . The moments of the gating weights are based on the approximation Eq. (27), and are

$$\begin{aligned}\mathbb{E}_V(\ln g_k(\mathbf{x}_n)) &\approx \ln g_k(\mathbf{x}_n)|_{\mathbf{V}=\tilde{\mathbf{V}}}, \\ \mathbb{E}_V(\mathbf{v}_k^T \mathbf{v}_k) &= \tilde{\mathbf{v}}_k^T \tilde{\mathbf{v}}_k + \text{Tr} \left((\tilde{\mathbf{\Lambda}}_V^{-1})_{kk} \right),\end{aligned}$$

where the moment $\mathbb{E}_V(\ln g_k(\mathbf{x}_n))$ cannot be directly evaluated and is therefore crudely approximated by the logarithm of the expectation, as also done in Waterhouse et al. (1996). In $\mathbb{E}_V(\mathbf{v}_k^T \mathbf{v}_k)$, $(\tilde{\mathbf{\Lambda}}_V^{-1})_{kk}$ stands for the k th diagonal block matrix in $\tilde{\mathbf{\Lambda}}_V^{-1}$. The moments of the gating weight prior is after Eq. (28) given by

$$\begin{aligned}\mathbb{E}_\beta(\beta_k) &= \frac{a_{\beta_k}^*}{b_{\beta_k}^*}, \\ \mathbb{E}_\beta(\ln \beta_k) &= \psi(a_{\beta_k}^*) - \ln b_{\beta_k}^*.\end{aligned}$$

Again, $\mathbb{E}_\beta(\ln \beta_k)$ is evaluated for use in a later section.

A.8 Independent Expert Training

When the experts are trained independently of the gating network, then the update equations for $q_W^*(\mathbf{w}_k)$ and $q_\tau^*(\tau_k)$ are changed to

$$\begin{aligned}\Sigma_k^* &= \left(\mathbb{E}_\alpha(\alpha_k) \mathbf{I} + \mathbb{E}_\tau(\tau_k) \sum_n m_k(\mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^T \right)^{-1}, \\ \mathbf{w}_k^* &= \mathbb{E}_\tau(\tau_k) \Sigma_k^* \sum_n m_k(\mathbf{x}_n) \mathbf{x}_n y_n, \\ a_{\tau_k}^* &= a_\tau + \frac{1}{2} \sum_n m_k(\mathbf{x}_n), \\ b_{\tau_k}^* &= b_\tau + \frac{1}{2} \sum_n m_k(\mathbf{x}_n) \mathbb{E}_W((y_n + \mathbf{w}_k^T \mathbf{x}_n)^2),\end{aligned}$$

by replacing $\mathbb{E}_Z(z_{nk})$ with $m_k(\mathbf{x}_n)$. All other update equations stay unchanged.

B Variational Bound $\mathcal{L}(\tilde{q})$

The variational bound $\mathcal{L}(\tilde{q})$ is evaluated after Eq. (14), which is given by

$$\begin{aligned}\mathcal{L}(\tilde{q}) &= \sum_Z \int \cdots \int \tilde{q}(\mathbf{W}, \tau, \alpha, \mathbf{Z}, \mathbf{V}, \beta) \ln \left(\frac{p(\mathbf{Y}, \mathbf{W}, \tau, \alpha, \mathbf{Z}, \mathbf{V}, \beta)}{\tilde{q}(\mathbf{W}, \tau, \alpha, \mathbf{Z}, \mathbf{V}, \beta)} \right) d\mathbf{W} d\tau d\alpha d\mathbf{V} d\beta \\ &= \mathbb{E}_{\mathbf{W}, \tau, \alpha, \mathbf{Z}, \mathbf{V}, \beta}(\ln p(\mathbf{Y}, \mathbf{W}, \tau, \alpha, \mathbf{Z}, \mathbf{V}, \beta)) - \mathbb{E}_{\mathbf{W}, \tau, \alpha, \mathbf{Z}, \mathbf{V}, \beta}(\ln q(\mathbf{W}, \tau, \alpha, \mathbf{Z}, \mathbf{V}, \beta)) \\ &= \mathbb{E}_{\mathbf{W}, \tau, \mathbf{Z}}(\ln p(\mathbf{Y} | \mathbf{Z}, \mathbf{W}, \tau)) + \mathbb{E}_{\mathbf{W}, \alpha}(\ln p(\mathbf{W} | \alpha)) + \mathbb{E}_\tau(\ln p(\tau)) + \mathbb{E}_\alpha(\ln p(\alpha)) \\ &\quad + \mathbb{E}_{\mathbf{Z}, \mathbf{V}}(\ln p(\mathbf{Z} | \mathbf{V})) + \mathbb{E}_{\mathbf{V}, \beta}(\ln p(\mathbf{V} | \beta)) + \mathbb{E}_\beta(\ln p(\beta)) \\ &\quad - \mathbb{E}_{\mathbf{W}}(\ln q(\mathbf{W})) - \mathbb{E}_\tau(\ln q(\tau)) - \mathbb{E}_\alpha(\ln q(\alpha)) \\ &\quad - \mathbb{E}_{\mathbf{Z}}(\ln q(\mathbf{Z})) - \mathbb{E}_{\mathbf{V}}(\ln q(\mathbf{V})) - \mathbb{E}_\beta(\ln q(\beta)).\end{aligned}$$

All distributions are implicitly conditional on \mathbf{X} and \mathcal{M} . The next sections are dedicated to deriving the required moments, after which we will return to deriving the close-form expression for $\mathcal{L}(\tilde{q})$.

B.1 $\mathbb{E}_{W,\tau,Z}(\ln p(\mathbf{Y}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\tau}))$

Using Eqs. (3) and (7), and the moments from Section A.7, we get

$$\begin{aligned}
 \mathbb{E}_{W,\tau,Z}(\ln p(\mathbf{Y}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\tau})) &= \iiint q(\mathbf{W}, \boldsymbol{\tau}) \ln p(\mathbf{Y}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\tau}) d\mathbf{Z} d\mathbf{W} d\boldsymbol{\tau} \\
 &= \sum_n \sum_k \int q(\mathbf{Z}) z_{nk} d\mathbf{Z} \iint q(\mathbf{W}, \boldsymbol{\tau}) \ln \mathcal{N}(y_n | \mathbf{w}_k^T \mathbf{x}_n, \tau_k^{-1}) d\mathbf{W} d\boldsymbol{\tau} \\
 &= \sum_n \sum_k \mathbb{E}_Z(z_{nk}) \left(\frac{1}{2} \mathbb{E}_\tau(\ln \tau_k) - \frac{1}{2} \ln 2\pi - \frac{1}{2} \mathbb{E}_\tau(\tau_k) \mathbb{E}_W((y_n - \mathbf{w}_k^T \mathbf{x}_n)^2) \right) \\
 &= \sum_k \left((\psi(a_{\tau_k}^*) - \ln b_{\tau_k}^* - \ln 2\pi) \frac{1}{2} \sum_k r_{nk} - \frac{1}{2} \frac{a_{\tau_k}^*}{b_{\tau_k}^*} \sum_n r_{nk} \mathbb{E}_W((y_n - \mathbf{w}_k^T \mathbf{x}_n)^2) \right).
 \end{aligned}$$

B.2 $\mathbb{E}_{W,\alpha}(\ln p(\mathbf{W}|\boldsymbol{\alpha}))$ and $\mathbb{E}_W(\ln q(\mathbf{W}))$

Using Eq. (8) and the moments from Section A.7, we get

$$\begin{aligned}
 \mathbb{E}_{W,\alpha}(\ln p(\mathbf{W}|\boldsymbol{\alpha})) &= \iint q(\mathbf{W}, \boldsymbol{\alpha}) \ln p(\mathbf{W}|\boldsymbol{\alpha}) d\mathbf{W} d\boldsymbol{\alpha} \\
 &= \sum_k \iint q(\mathbf{w}_k, \alpha_k) \ln \mathcal{N}(\mathbf{w}_k | 0, \alpha_k^{-1} \mathbf{I}) d\mathbf{w}_k d\alpha_k \\
 &= \sum_k \left(-\frac{D_w}{2} \ln 2\pi + \frac{D_w}{2} \mathbb{E}_\alpha(\ln \alpha_k) - \frac{1}{2} \mathbb{E}_\alpha(\alpha_k) \mathbb{E}_W(\mathbf{w}_k^T \mathbf{w}_k) \right) \\
 &= \sum_k \left(-\frac{D_w}{2} \ln 2\pi + \frac{D_w}{2} (\psi(a_{\alpha_k}^*) - \ln b_{\alpha_k}^*) - \frac{1}{2} \frac{a_{\alpha_k}^*}{b_{\alpha_k}^*} (\mathbf{w}_k^{*T} \mathbf{w}_k^* + \text{Tr}(\boldsymbol{\Sigma}_k^*)) \right).
 \end{aligned}$$

We can evaluate $\mathbb{E}_W(\ln q(\mathbf{W}))$ by using Eq. (23) and observing that

$$-\mathbb{E}_W(\ln q(\mathbf{W})) = -\int q(\mathbf{W}) \ln q(\mathbf{W}) d\mathbf{W} = \sum_k \left(-\int q(\mathbf{w}_k) \ln q(\mathbf{w}_k) d\mathbf{w}_k \right)$$

is the sum of entropies of $q(\mathbf{w}_k)$, and therefore $\mathbb{E}_W(\ln q(\mathbf{W}))$ is given by

$$\mathbb{E}_W(\ln q(\mathbf{W})) = -\sum_k \left(\frac{1}{2} \ln |\boldsymbol{\Sigma}_k^*| + \frac{D_w}{2} (1 + \ln 2\pi) \right).$$

B.3 $\mathbb{E}_\tau(\ln p(\boldsymbol{\tau}))$ and $\mathbb{E}_\tau(\ln q(\boldsymbol{\tau}))$

From Eq. (10) and the moments evaluated in Section A.7 follows

$$\begin{aligned}
 \mathbb{E}_\tau(\ln p(\boldsymbol{\tau})) &= \int q(\boldsymbol{\tau}) \ln p(\boldsymbol{\tau}) d\boldsymbol{\tau} \\
 &= \sum_k \int q(\tau_k) \ln \text{Gam}(\tau_k | a_\tau, b_\tau) d\tau_k \\
 &= \sum_k \left(-\ln \Gamma(a_\tau) + a_\tau \ln b_\tau + (a_\tau - 1) \mathbb{E}_\tau(\ln \tau_k) - b_\tau \mathbb{E}_\tau(\tau_k) \right) \\
 &= \sum_k \left(-\ln \Gamma(a_\tau) + a_\tau \ln b_\tau + (a_\tau - 1) (\psi(a_{\tau_k}^*) - \ln b_{\tau_k}^*) - b_\tau \frac{a_{\tau_k}^*}{b_{\tau_k}^*} \right),
 \end{aligned}$$

where $\Gamma(\cdot)$ is the gamma function. We can again observe that $-\mathbb{E}_\tau(\ln q(\boldsymbol{\tau}))$ is the sum of entropies of $q_\tau^*(\tau_k)$ after Eq. (24), and therefore we get

$$\begin{aligned}
 \mathbb{E}_\tau(\ln q(\boldsymbol{\tau})) &= \int q(\boldsymbol{\tau}) \ln q(\boldsymbol{\tau}) d\boldsymbol{\tau} \\
 &= -\sum_k \left(-\int q(\tau_k) \ln q(\tau_k) d\tau_k \right) \\
 &= -\sum_k \left(\ln \Gamma(a_{\tau_k}^*) - (a_{\tau_k}^* - 1) \psi(a_{\tau_k}^*) - \ln b_{\tau_k}^* + a_{\tau_k}^* \right)
 \end{aligned}$$

B.4 $\mathbb{E}_\alpha(\ln p(\boldsymbol{\alpha}))$ and $\mathbb{E}_\alpha(\ln q(\boldsymbol{\alpha}))$

The derivations are the same as for $\mathbb{E}_\tau(\ln p(\boldsymbol{\tau}))$ and $\mathbb{E}_\tau(\ln q(\boldsymbol{\tau}))$ and result in

$$\begin{aligned}\mathbb{E}_\alpha(\ln p(\boldsymbol{\alpha})) &= \sum_k \left(-\ln \Gamma(a_\alpha) + a_\alpha \ln b_\alpha + (a_\alpha - 1)(\psi(a_{\alpha_k}^*) - \ln b_{\alpha_k}^*) - b_\alpha \frac{a_{\alpha_k}^*}{b_{\alpha_k}^*} \right), \\ \mathbb{E}_\alpha(\ln q(\boldsymbol{\alpha})) &= -\sum_k (\ln \Gamma(a_{\alpha_k}^*) - (a_{\alpha_k}^* - 1)\psi(a_{\alpha_k}^*) - \ln b_{\alpha_k}^* + a_{\alpha_k}^*).\end{aligned}$$

B.5 $\mathbb{E}_{Z,V}(\ln p(\mathbf{Z}|\mathbf{V}))$ and $\mathbb{E}_Z(\ln q(\mathbf{Z}))$

From Eq. (26) and Section A.7 we get

$$\begin{aligned}\mathbb{E}_{Z,V}(\ln p(\mathbf{Z}|\mathbf{V})) &= \iint q(\mathbf{Z}, \mathbf{V}) \ln p(\mathbf{Z}|\mathbf{V}) d\mathbf{Z} d\mathbf{V} \\ &= \sum_n \sum_k \int q(z_{nk}) z_{nk} dz_{nk} \int q(\mathbf{v}_k) \ln g_k(\mathbf{x}_n) d\mathbf{v}_k \\ &= \sum_n \sum_k \mathbb{E}_Z(z_{nk}) \mathbb{E}_V(g_k(\mathbf{x}_n)) \\ &= \sum_n \sum_k r_{nk} \ln g_k(\mathbf{x}_n)|_{\mathbf{v}_k = \tilde{\mathbf{v}}_k}.\end{aligned}$$

Additionally, using Eq. (29),

$$\begin{aligned}\mathbb{E}_Z(\ln q(\mathbf{Z})) &= \int q(\mathbf{Z}) \ln q(\mathbf{Z}) d\mathbf{Z} \\ &= \sum_n \sum_k \mathbb{E}_Z(z_{nk}) \ln r_{nk} \\ &= \sum_n \sum_k r_{nk} \ln r_{nk},\end{aligned}$$

which is the negative entropy of $q(\mathbf{Z})$.

B.6 $\mathbb{E}_{V,\beta}(\ln p(\mathbf{V}|\boldsymbol{\beta}))$ and $\mathbb{E}_V(\ln \tilde{q}(\mathbf{V}))$

Deriving $\mathbb{E}_{V,\beta}(\ln p(\mathbf{V}|\boldsymbol{\beta}))$ is similar to the derivation of $\mathbb{E}_{W,\alpha}(\ln p(\mathbf{W}|\boldsymbol{\alpha}))$ and results in

$$\mathbb{E}_{V,\beta}(\ln p(\mathbf{V}|\boldsymbol{\beta})) = \sum_k \left(-\frac{D_v}{2} \ln 2\pi + \frac{D_v}{2} (\psi(a_{\beta_k}^*) - \ln b_{\beta_k}^*) - \frac{1}{2} \frac{a_{\beta_k}^*}{b_{\beta_k}^*} \left(\tilde{\mathbf{v}}_k^T \tilde{\mathbf{v}}_k + \text{Tr}((\tilde{\boldsymbol{\Lambda}}_V^{-1})_{kk}) \right) \right).$$

To evaluate $\mathbb{E}_V(\ln \tilde{q}(\mathbf{V}))$ we note that $-\mathbb{E}_V(\ln \tilde{q}(\mathbf{V}))$ is the entropy of $\tilde{q}(\mathbf{V})$ after Eq. (27) to get

$$\begin{aligned}\mathbb{E}_V(\ln \tilde{q}(\mathbf{V})) &= -\left(-\int \tilde{q}(\mathbf{V}) \ln \tilde{q}(\mathbf{V}) d\mathbf{V} \right) \\ &= -\left(\frac{1}{2} \ln |\tilde{\boldsymbol{\Lambda}}_V^{-1}| + \frac{KD_v}{2} (1 + \ln 2\pi) \right).\end{aligned}$$

B.7 $\mathbb{E}_\beta(\ln p(\boldsymbol{\beta}))$ and $\mathbb{E}_\beta(\ln q(\boldsymbol{\beta}))$

The derivations are the same as for $\mathbb{E}_\tau(\ln p(\boldsymbol{\tau}))$ and $\mathbb{E}_\tau(\ln q(\boldsymbol{\tau}))$ and result in

$$\begin{aligned}\mathbb{E}_\beta(\ln p(\boldsymbol{\beta})) &= \sum_k \left(-\ln \Gamma(a_\beta) + a_\beta \ln b_\beta + (a_\beta - 1)(\psi(a_{\beta_k}^*) - \ln b_{\beta_k}^*) - b_\beta \frac{a_{\beta_k}^*}{b_{\beta_k}^*} \right), \\ \mathbb{E}_\beta(\ln q(\boldsymbol{\beta})) &= -\sum_k (\ln \Gamma(a_{\beta_k}^*) - (a_{\beta_k}^* - 1)\psi(a_{\beta_k}^*) - \ln b_{\beta_k}^* + a_{\beta_k}^*).\end{aligned}$$

B.8 Closed-Form Equation for $\mathcal{L}(\tilde{q})$

Let $\mathcal{L}_k(q)$ denote the contribution of expert k independent of the gating network to the variational bound $\mathcal{L}(q)$. Using the contribution of expert k to $\mathbb{E}_{W,\tau,Z}(\ln p(\mathbf{Y}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\tau}))$, $\mathbb{E}_{W,\alpha}(\ln p(\mathbf{W}|\boldsymbol{\alpha}))$, $\mathbb{E}_\tau(\ln p(\boldsymbol{\tau}))$, $\mathbb{E}_\alpha(\ln p(\boldsymbol{\alpha}))$, $\mathbb{E}_W(\ln q(\mathbf{W}))$, $\mathbb{E}_\tau(\ln q(\boldsymbol{\tau}))$ and $\mathbb{E}_\alpha(\ln q(\boldsymbol{\alpha}))$, this contribution is given by

$$\begin{aligned} \mathcal{L}_k(q) &= (\psi(a_{\tau_k}^*) - \ln b_{\tau_k}^* - \ln 2\pi) \frac{1}{2} \sum_n r_{nk} - \frac{1}{2} \frac{a_{\tau_k}^*}{b_{\tau_k}^*} \sum_n r_{nk} \mathbb{E}_W((y_n - \mathbf{w}_k^T \mathbf{x}_n)^2) \\ &\quad - \frac{D_w}{2} \ln 2\pi + \frac{D_w}{2} (\psi(a_{\alpha_k}^*) - \ln b_{\alpha_k}^*) - \frac{1}{2} \frac{a_{\alpha_k}^*}{b_{\alpha_k}^*} (\mathbf{w}_k^{*T} \mathbf{w}_k^* + \text{Tr}(\boldsymbol{\Sigma}_k^*)) \\ &\quad - \ln \Gamma(a_\tau) + a_\tau \ln b_\tau + (a_\tau - 1)(\psi(a_{\tau_k}^*) - \ln b_{\tau_k}^*) - b_\tau \frac{a_{\tau_k}^*}{b_{\tau_k}^*} \\ &\quad - \ln \Gamma(a_\alpha) + a_\alpha \ln b_\alpha + (a_\alpha - 1)(\psi(a_{\alpha_k}^*) - \ln b_{\alpha_k}^*) - b_\alpha \frac{a_{\alpha_k}^*}{b_{\alpha_k}^*} \\ &\quad + \frac{1}{2} \ln |\boldsymbol{\Sigma}_k^*| + \frac{D_w}{2} (1 + \ln 2\pi) \\ &\quad + \ln \Gamma(a_{\tau_k}^*) - (a_{\tau_k}^* - 1)\psi(a_{\tau_k}^*) - \ln b_{\tau_k}^* + a_{\tau_k}^* \\ &\quad + \ln \Gamma(a_{\alpha_k}^*) - (a_{\alpha_k}^* - 1)\psi(a_{\alpha_k}^*) - \ln b_{\alpha_k}^* + a_{\alpha_k}^*. \end{aligned}$$

Observing that

$$\begin{aligned} a_{\tau_k}^* - a_\tau &= \frac{1}{2} \sum_n r_{nk}, \\ b_{\tau_k}^* - b_\tau &= \frac{1}{2} \sum_n r_{nk} \mathbb{E}_W((y_n + \mathbf{w}_k^T \mathbf{x}_n)^2), \\ a_{\alpha_k}^* - a_\alpha &= \frac{D_w}{2}, \\ b_{\alpha_k}^* - b_\alpha &= \frac{1}{2} (\mathbf{w}_k^{*T} \mathbf{w}_k^* + \text{Tr}(\boldsymbol{\Sigma}_k^*)), \end{aligned}$$

we can simplify $\mathcal{L}_k(q)$ to

$$\begin{aligned} \mathcal{L}_k(q) &= -\ln \Gamma(a_\tau) + a_\tau \ln b_\tau + \ln \Gamma(a_{\tau_k}^*) - a_{\tau_k}^* \ln b_{\tau_k}^* + (a_\tau - a_{\tau_k}^*) \ln 2\pi \\ &\quad - \ln \Gamma(a_\alpha) + a_\alpha \ln b_\alpha + \ln \Gamma(a_{\alpha_k}^*) - a_{\alpha_k}^* \ln b_{\alpha_k}^* \\ &\quad + \frac{1}{2} \ln |\boldsymbol{\Sigma}_k^*| + \frac{D_w}{2}. \end{aligned}$$

To get the full variational bound, we can use

$$\begin{aligned} a_{\beta_k}^* - a_\beta &= \frac{D_v}{2}, \\ b_{\beta_k}^* - b_\beta &= \frac{1}{2} (\tilde{\mathbf{v}}_k^T \tilde{\mathbf{v}}_k + \text{Tr}((\tilde{\boldsymbol{\Lambda}}_V^{-1})_{kk})), \end{aligned}$$

to get

$$\begin{aligned} \mathcal{L}(\tilde{q}) &= \sum_k \mathcal{L}_k(q) + \sum_k (-\ln \Gamma(a_\beta) + a_\beta \ln b_\beta + \ln \Gamma(a_{\beta_k}^*) - a_{\beta_k}^* \ln b_{\beta_k}^*) \\ &\quad + \frac{1}{2} \ln |\tilde{\boldsymbol{\Lambda}}_V^{-1}| + \frac{KD_v}{2} + \sum_n \sum_k r_{nk} \ln \frac{g_k(\mathbf{x}_n)|_{\mathbf{v}_k=\tilde{\mathbf{v}}_k}}{r_{nk}}. \end{aligned}$$

B.9 $\mathcal{L}(\tilde{q})$ for Independent Expert Training

When the experts are trained independently, their update equations are modified as described in Section A.8. Hence, our simplification based on the expression for $a_{\tau_k}^* - a_\tau$ and $b_{\tau_k}^* - b_\tau$ is not valid

anymore. Therefore, the expression of the contribution of expert k to $\mathcal{L}(\tilde{q})$ becomes

$$\begin{aligned} \mathcal{L}_k(q) &= (\psi(a_{\tau_k}^*) - \ln b_{\tau_k}^* - \ln 2\pi) \frac{1}{2} \sum_n r_{nk} - \frac{1}{2} \frac{a_{\tau_k}^*}{b_{\tau_k}^*} \sum_n r_{nk} (\|y_n - \mathbf{w}_k^T \mathbf{x}_n\|^2 + \mathbf{x}_n^T \boldsymbol{\Sigma}_k^* \mathbf{x}_n) \\ &\quad - \ln \Gamma(a_{\tau_k}) + a_{\tau_k} (\ln b_{\tau_k} - \ln b_{\tau_k}^*) + \ln \Gamma(a_{\tau_k}^*) - b_{\tau_k} \frac{a_{\tau_k}^*}{b_{\tau_k}^*} a_{\tau_k}^* \psi(a_{\tau_k}^*) + a_{\tau_k}^* \\ &\quad + \frac{1}{2} \ln |\boldsymbol{\Sigma}_k^*| + \frac{D_w}{2} \\ &\quad - \ln \Gamma(a_{\alpha_k}) + a_{\alpha_k} \ln b_{\alpha_k} + \ln \Gamma(a_{\alpha_k}^*) - a_{\alpha_k}^* \ln b_{\alpha_k}^*. \end{aligned}$$

The expression for $\mathcal{L}(\tilde{q})$ remains unchanged.

C Predictive Distribution

Given a new input $\hat{\mathbf{x}}$, we assume that $\hat{\mathbf{z}}$ is the associated latent variable, and want to find $p(\hat{y}|\hat{\mathbf{x}}, \mathbf{X}, \mathbf{Y})$ by evaluating

$$p(\hat{y}|\hat{\mathbf{x}}, \mathbf{X}, \mathbf{Y}) = \sum_{\hat{\mathbf{z}}} \iiint p(\hat{y}|\hat{\mathbf{x}}, \hat{\mathbf{z}}, \mathbf{W}, \boldsymbol{\tau}) p(\hat{\mathbf{z}}|\hat{\mathbf{x}}, \mathbf{V}) p(\mathbf{W}, \boldsymbol{\tau}, \mathbf{V}|\mathbf{X}, \mathbf{Y}) d\mathbf{W} d\boldsymbol{\tau} d\mathbf{V}.$$

Using the definition of the various distributions and summing over $\hat{\mathbf{z}}$ we get

$$p(\hat{y}|\hat{\mathbf{x}}, \mathbf{X}, \mathbf{Y}) = \sum_k \iiint g_k(\hat{\mathbf{x}}) \mathcal{N}(\hat{y}|\mathbf{x}_k^T \hat{\mathbf{x}}, \tau_k^{-1}) p(\mathbf{W}, \boldsymbol{\tau}, \mathbf{V}|\mathbf{X}, \mathbf{Y}) d\mathbf{W} d\boldsymbol{\tau} d\mathbf{V},$$

where we will approximate the posterior $p(\mathbf{W}, \boldsymbol{\tau}, \mathbf{V}|\mathbf{X}, \mathbf{Y})$ by the variational distribution $q_W^*(\mathbf{W}) q_{\boldsymbol{\tau}}^*(\boldsymbol{\tau}) q_V^*(\mathbf{V})$. Thus, the predictive distribution becomes

$$p(\hat{y}|\hat{\mathbf{x}}, \mathbf{X}, \mathbf{Y}) = \sum_k \int q_V^*(\mathbf{v}_k) g_k(\hat{\mathbf{x}}) d\mathbf{v}_k \iint q_W^*(\mathbf{w}_k) q_{\boldsymbol{\tau}}^*(\tau_k) \mathcal{N}(\hat{y}|\mathbf{w}_k^T \hat{\mathbf{x}}, \tau_k^{-1}) d\mathbf{w}_k d\tau_k.$$

The first integral is $\mathbb{E}_V(g_k(\hat{\mathbf{x}}))$ which we approximate as in Ueda and Ghahramani (2002) by its MAP estimate

$$\int q_V^*(\mathbf{v}_k) g_k(\hat{\mathbf{x}}) d\mathbf{v}_k \approx g_k(\hat{\mathbf{x}})|_{\mathbf{v}_k=\bar{\mathbf{v}}_k}.$$

The second integral is also solved as in Ueda and Ghahramani (2002) and results in the Student-t distribution

$$\iint q_W^*(\mathbf{w}_k) q_{\boldsymbol{\tau}}^*(\tau_k) \mathcal{N}(\hat{y}|\mathbf{w}_k^T \hat{\mathbf{x}}, \tau_k^{-1}) d\mathbf{w}_k d\tau_k = \text{St} \left(\hat{y} | \mathbf{w}_k^T \hat{\mathbf{x}}, \frac{a_{\tau_k}^*}{b_{\tau_k}^*} (1 + \hat{\mathbf{x}}^T \boldsymbol{\Sigma}_k^* \hat{\mathbf{x}})^{-1}, 2a_{\tau_k}^* \right).$$

Therefore, our predictive distribution is the mixture of Student-t distributions

$$p(\hat{y}|\hat{\mathbf{x}}, \mathbf{X}, \mathbf{Y}) = \sum_k g_k(\hat{\mathbf{x}})|_{\mathbf{v}_k=\bar{\mathbf{v}}_k} \text{St} \left(\hat{y} | \mathbf{w}_k^T \hat{\mathbf{x}}, \frac{a_{\tau_k}^*}{b_{\tau_k}^*} (1 + \hat{\mathbf{x}}^T \boldsymbol{\Sigma}_k^* \hat{\mathbf{x}})^{-1}, 2a_{\tau_k}^* \right),$$

with mean

$$\mathbb{E}(\hat{y}|\hat{\mathbf{x}}, \mathbf{X}, \mathbf{Y}) = \sum_k g_k(\hat{\mathbf{x}})|_{\mathbf{v}_k=\bar{\mathbf{v}}_k} \mathbf{w}_k^T \hat{\mathbf{x}}.$$

References

- Bernadó-Mansilla, E., Llorá, X., & Garrell-Guiu, J. M. (2002). XCS and GALE: A Comparative Study of Two Learning Classifier Systems on Data Mining. *IWLCS '01: Revised Papers from the 4th International Workshop on Advances in Learning Classifier Systems* (pp. 115–132). London, UK: Springer.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer.
- Bishop, C. M., & Svensén, M. (2003). Bayesian Hierarchical Mixtures of Experts. *Proceedings of the 19th Annual Conference on Uncertainty in Artificial Intelligence (UAI-03)* (pp. 57–64). San Francisco, CA: Morgan Kaufmann.
- Butz, M. V. (2006). *Rule-Based Evolutionary Online Learning Systems: A Principled Approach to LCS Analysis and Design*, vol. 191 of *Studies in Fuzziness and Soft Computing*. Springer.
- Butz, M. V., Lanzi, P. L., & Wilson, S. W. (to appear). Function Approximation with XCS: Hyperellipsoidal Conditions, Recursive Least Squares, and Compaction. *IEEE Transactions on Evolutionary Computations*.
- Chipman, H. A., George, E. I., & McCulloch, R. E. (1998). Bayesian cart model search. *Journal of the American Statistical Association*, *93*, 935–948.
- Gibbs, M. N. (1997). *Bayesian gaussian processes for regression and classification*. Doctoral dissertation, University of Cambridge.
- Jacobs, R. A., Jordan, M. I., Nowlan, S., & Hinton, G. E. (1991). Adaptive mixtures of local experts. *Neural Computation*, *3*, 1–12.
- Jordan, M. I., & Jacobs, R. A. (1994). Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, *6*, 181–214.
- Ueda, N., & Ghahramani, Z. (2002). Bayesian model search for mixture models based on optimizing variational bounds. *Neural Networks*, *15*, 1223–1241.
- Wainwright, M., Jaakkola, T., & Willsky, A. (2005). A new class of upper bounds on the log partition function. *IEEE Transactions on Information Theory*, *51*, 2313–2335.
- Waterhouse, S. (1997). *Classification and Regression using Mixtures of Experts*. Doctoral dissertation, Department of Engineering, University of Cambridge.
- Waterhouse, S., MacKay, D., & Robinson, T. (1996). Bayesian Methods for Mixtures of Experts. *Advances in Neural Information Processing Systems 8* (pp. 351–357). Cambridge, MA: MIT Press.
- Xu, L., Jordan, M., & Hinton, G. E. (1995). An Alternative Model for Mixtures of Experts. *Advances in Neural Information Processing Systems 7* (pp. 633–640). Cambridge, MA: MIT Press.