

# SPARSE BAYESIAN HIERARCHICAL MIXTURE OF EXPERTS

*Iman Mossavat, Oliver Amft*

Signal Processing Systems Group, Department of Electrical Engineering  
Eindhoven University of Technology, The Netherlands

## ABSTRACT

Hierarchical mixture of experts (HME) is a widely adopted probabilistic divide-and-conquer regression model. We extend the variational inference algorithm for HME by using automatic relevance determination (ARD) priors. Unlike Gaussian priors, ARD allows for a few model parameters to take on large values, while forcing others to zero. Thus, using ARD priors encourages sparse models. Sparsity is known to be advantageous to the generalization capability as well as interpretability of the models. We present the variational inference algorithm for sparse HME in detail. Subsequently, we evaluate the sparse HME approach in building objective speech quality assessment algorithms, that are required to determine the quality of service in telecommunication networks.

**Index Terms**— Mixture, sparse, variational, inference, Bayesian, speech quality.

## 1. INTRODUCTION

Hierarchical mixture of experts (HME) is a probabilistic divide-and-conquer regression method [1]. HME divides the feature space into regions, and assigns a linear regressor to each region. HME is a conditional mixture of linear regressors. Mixture components are called *experts*. Experts are combined in the mixture using weights, called mixing coefficients. Mixing coefficients depend on the input feature vector, and are determined by the *gating network*: a tree structure with binary classifiers, or *gates*, at its internal nodes. The gating network task is to divide the feature space into regions by adapting the mixing coefficients based on the feature vector.

Since HME is a probabilistic method, decision boundaries among the regions are ‘soft’. Thus a given data point may lie simultaneously in several regions. The following advantages make HME more favorable than other divide-and-conquer regression methods with ‘hard’ decision boundaries such as CART [2]: the capability to model multi-modal distributions as explained in [3], less abrupt fluctuations near decision boundaries [4], and smaller variance as described in [1].

CART models allow easier interpretation, which is an advantage over HME. CART divides the space by specifying

thresholds on certain feature values, whereas HME uses probabilistic binary classifiers in its gating network. Existing inference algorithms for HME models, e.g. [1, 3] do not yield sparse results, making HME models hard to interpret.

In this work, we extend the variational Bayesian HME of Bishop and Svensén [3], in which they approximate the posterior using the *mean-field* method [4]. We specify new priors for HME based on the automatic relevance determination (ARD) mechanism [5, 6]. ARD allows for experts and gates to be sparse. Sparsity often promotes interpretability as well as better generalization. We present the details of the approximate inference algorithm for the sparse HME.

The new sparse HME is subsequently evaluated by developing a non-intrusive objective speech quality assessment (QA) algorithm, which is an estimator of the perceived speech quality. The state of the art algorithm in non-intrusive speech QA is represented by the ITU-T recommendation P.563 [7], which exploits a divide-and-conquer strategy. P.563 classifies speech signals transmitted over a telecommunication network based on the talker sex as well as how the network distorted the speech, e.g. high-level background noise, speech robotization, etc. P.563 uses hard binary classifiers, thus its output changes abruptly near distortion class boundaries. P.563 is limited to assigning only one class to each signal, thus it cannot reliably predict the quality when both male and female talkers are present in the speech signal [7]. In [8], we use the variational Bayesian HME of [3] to model the speech quality data in the ITU-T coded-speech data-set, Supplement 23 [9]. HME soft decision boundaries effectively addresses the non-smoothness limitation of P.563, and achieves competitive prediction performance with P.563. However, even with two experts, the HME model is hard to interpret.

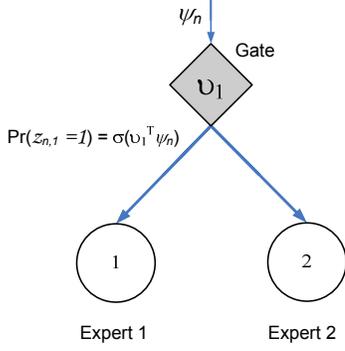
## 2. HME MODEL

HME defines a mixture distribution over  $y$ , the target variable, conditioned on the  $D$ -dimensional feature vector  $\psi$ , which is

$$\Pr(y|\psi, \theta) = \sum_{k=1}^C g_k(\psi) \mathcal{N}(y|\omega_k^T \psi, \tau_k^{-1}), \quad (1)$$

where  $\theta$  denote the model parameters,  $C$  denotes the number of experts, and  $g_k(\psi)$  denotes the mixing-coefficient for the  $k$ -th

Iman Mossavat is funded by STW project 07605: Personalization of Hearing Aids Through Bayesian Preference Elicitation (HearClip).



**Fig. 1.** The gating network of a HME model with one gate and two experts. The gate is a binary classifier. For  $n$ -th input  $\psi_n$ , the latent binary variable  $z_{n,1}$  indicates the gate output, with  $z_{n,1} = 1$  corresponding to selecting the expert on the left. The gate selects the output branch on its left by probability  $\sigma(\mathbf{v}_1^T \psi_n)$ , where  $\sigma(\cdot)$  is the logistic function. The mixing-coefficient for each expert is the product of the link probabilities on the path that connects the root node to the expert.

expert. Experts correspond to Gaussian distributions, denoted by  $\mathcal{N}$ . The mean of the  $k$ -th expert is given by  $\omega_k^T \psi$  and its precision (inverse of variance) by  $\tau_k$ . The gating network determines  $g_k(\psi)$  as described in the following section.

Suppose the HME model has  $C$  experts and  $M$  gates. From this point forward, we use  $k = 1, \dots, C$ ,  $l = 1, \dots, M$ , and  $n = 1, \dots, N$  to denote respectively the expert, the gate, and the data point indices. In addition, the  $i$ -th element of vector  $\mathbf{x}$  is denoted by  $x(i)$ .

## 2.1. Gating Network

We define the latent variable  $z_{n,l} \in \{0, 1\}$ , corresponding to the  $l$ -th gate and  $n$ -th data point  $\psi_n$ , such that  $z_{n,l} = 1$  indicates the  $l$ -th gate left-side branch is chosen for  $\psi_n$ . The mixing-coefficients in Equation (1) are defined as

$$g_k(\psi_n) = \Pr(\zeta_{n,k} = 1), \quad (2)$$

where

$$\zeta_{n,k} = \prod_{l=1}^M z_{n,l}^{\mathbf{S}^L(k,l)} (1 - z_{n,l})^{\mathbf{S}^R(k,l)}, \quad (3)$$

where the gating network topology is specified by binary matrices  $\mathbf{S}^L$  and  $\mathbf{S}^R$ , where  $\mathbf{S}^L(k, l) = 1$  if  $k$ -th expert is on the ‘left’ sub-tree of  $l$ -th gate, and zero otherwise. Similarly,  $\mathbf{S}^R(k, l) = 1$  if  $k$ -th expert is on the ‘right’ sub-tree of  $l$ -th gate, and zero otherwise.

It is straightforward to verify that  $\zeta_{n,k} = 1$ , if and only if, for the point  $\psi_n$  the gates connect the root to the  $k$ -th expert. The probability distribution of  $z_{n,l}$  is

$$\Pr(z_{n,l} | \psi_n, \mathbf{v}_l) = \sigma(\mathbf{v}_l^T \psi_n)^{z_{n,l}} (1 - \sigma(\mathbf{v}_l^T \psi_n))^{(1-z_{n,l})}, \quad (4)$$

where  $\sigma(x) = 1 / (1 + \exp(-x))$  is the sigmoid function, and  $\mathbf{v}_l$  is the weight vector of  $l$ -th gate. We illustrate the gating network of a HME model with two experts in Figure 1.

## 3. ARD PRIOR

To infer sparse experts and gates parameters from the data, we incorporate ARD in the priors for  $\omega_k$  and  $\mathbf{v}_l$ . The hierarchical prior for the  $l$ -th gate is defined as

$$\Pr(\mathbf{v}_l | \beta_l) = \prod_{i=1}^D \mathcal{N}(\mathbf{v}_l(i) | 0, \beta_l^{-1}(i)) \quad (5)$$

$$\Pr(\beta_l) = \prod_{i=1}^D \text{Gamma}(\beta_l(i) | a_0, b_0), \quad (6)$$

where  $D$  is the dimensionality of the feature vector  $\psi$  and Gamma denotes the Gamma distribution. By  $\beta_l(i)$  we denote the precision of the  $i$ -th feature weight, i.e.,  $\mathbf{v}_l(i)$ . A large value for  $\beta_l(i)$  results in the  $\mathbf{v}_l(i)$  prior to be concentrated around zero, thus forcing  $\mathbf{v}_l(i)$  to small values. Learning  $\beta_l(i)$  along with  $\mathbf{v}_l(i)$  results in determination of feature relevance in the  $l$ -th gate and encourages sparsity ([4], Chapter 7). Similarly, the  $k$ -th expert prior is

$$\Pr(\omega_k | \tau_k, \alpha_k) = \prod_{i=1}^D \mathcal{N}(\omega_k(i) | 0, \tau_k^{-1} \alpha_k^{-1}(i)) \quad (7)$$

$$\Pr(\tau_k | a_0, b_0) = \text{Gamma}(\tau_k | a_0, b_0) \quad (8)$$

$$\Pr(\alpha_k | a_0, b_0) = \prod_{i=1}^D \text{Gamma}(\alpha_k(i) | a_0, b_0), \quad (9)$$

where  $\alpha_k(i)$  is the shrinkage parameter corresponding to  $\omega_k(i)$ . Jointly inferring  $\alpha_k(i)$  and  $\omega_k(i)$  results in determination of feature relevance for the  $k$ -th expert.

## 4. VARIATIONAL INFERENCE

For non-sparse HME, an efficient Bayesian inference algorithm, based on mean-field approximation method, is given in [3]. We first review the mean-field approximation briefly [4], and then proceed to present our inference algorithm that is working with priors specified in the previous section. As in [3], we approximate the sigmoid function  $\sigma(x)$  such that the conjugacy of the model is restored.

By  $\theta$ ,  $\mathcal{D}$ , and  $\mathcal{M}$  we denote the model parameters, the data set and the model, respectively. In HME, the model  $\mathcal{M}$  specifies the gating network topology and the hyper-priors. To approximate the posterior  $\Pr(\theta | \mathcal{D}, \mathcal{M})$ , we start by writing the log model evidence  $\ln \Pr(\mathcal{D} | \mathcal{M})$  as

$$\ln \Pr(\mathcal{D} | \mathcal{M}) = \mathcal{L}(q) + \text{KL}(q || p) \quad (10)$$

$$\mathcal{L}(q) = \int q(\theta) \ln \frac{\Pr(\theta, \mathcal{D} | \mathcal{M})}{\Pr(\theta | \mathcal{M})} d\theta, \quad (11)$$

where  $q(\boldsymbol{\theta})$  denotes an arbitrary distribution. By  $\text{KL}(q||p)$  we denote the Kullback-Leibler (KL) divergence between distributions  $q(\boldsymbol{\theta})$  and the posterior  $\text{Pr}(\boldsymbol{\theta}|\mathcal{D}, \mathcal{M})$ . KL divergence is always non-negative, thus  $\mathcal{L}(q) \leq \ln \text{Pr}(\mathcal{D}|\mathcal{M})$ , and  $\mathcal{L}(q)$  is a variational *lower-bound* to log-model-evidence. The goal here is to develop a computationally tractable algorithm, which accurately approximates the posterior by  $q(\boldsymbol{\theta})$ . For sparse HME an efficient iterative algorithm for finding the optimal  $q(\boldsymbol{\theta})$  exists, once we ‘assume’ that  $q(\boldsymbol{\theta})$  factorizes as

$$q(\boldsymbol{\theta}) = \prod_{k=1}^C q(\boldsymbol{\omega}_k, \tau_k) q(\boldsymbol{\alpha}_k) \prod_{l=1}^M q(\mathbf{v}_l) q(\boldsymbol{\beta}_l) \prod_{n=1}^N q(z_{n,l}). \quad (12)$$

Equation (12) specifies a set of distributions. It is possible to efficiently find the member of this set,  $q^*(\boldsymbol{\theta})$ , which maximizes the lower-bound  $\mathcal{L}(q^*)$ . The assumed factorization in Equation (12) makes efficient optimization of  $\mathcal{L}(q)$  possible. In each iteration of the  $\mathcal{L}(q)$  optimization, one of the factors is updated while other factors are constant. Given a conjugate-exponential structure, update equations have closed form formulae. We present the update equations, but the formula for  $\mathcal{L}(q)$  and the derivations are omitted due to space limitation.

#### 4.1. Update Equations

Computing the update equations for the gates are complicated by the fact that the conjugate-exponential structure of the HME model is thwarted by the sigmoid function in Equation (4). We use the parametric approximation of the sigmoid function in [3]. For each gate  $l$  and data point  $\boldsymbol{\psi}_n$ , we define a variational parameter  $\gamma_{n,l}$ . In each iteration we update  $\gamma_{n,l}$  along with the other factors in Equation (12). In the following we use  $i = 1, \dots, D$  to denote the feature index. The terms of optimal  $q^*(\boldsymbol{\theta})$  in Equation (12) are

$$q^*(z_{n,l}) = \sigma(h_{n,l})^{z_{n,l}} (1 - \sigma(h_{n,l}))^{1-z_{n,l}} \quad (13)$$

$$q^*(\boldsymbol{\omega}_k|\tau_k) = \mathcal{N}(\boldsymbol{\omega}_k|\bar{\boldsymbol{\omega}}_k, \tau_k^{-1}\mathbf{V}_k) \quad (14)$$

$$q^*(\tau_k) = \text{Gamma}(\tau_k|a_{\tau_k}, b_{\tau_k}) \quad (15)$$

$$q^*(\mathbf{v}_l) = \mathcal{N}(\mathbf{v}_l|\bar{\mathbf{v}}_l, \boldsymbol{\Lambda}_l) \quad (16)$$

$$q^*(\boldsymbol{\alpha}_k) = \prod_{i=1}^D \text{Gamma}(\boldsymbol{\alpha}_k(i)|\mathbf{a}_{\boldsymbol{\alpha}_k}(i), \mathbf{b}_{\boldsymbol{\alpha}_k}(i)) \quad (17)$$

$$q^*(\boldsymbol{\beta}_l) = \prod_{i=1}^D \text{Gamma}(\boldsymbol{\beta}_l(i)|\mathbf{a}_{\boldsymbol{\beta}_l}(i), \mathbf{b}_{\boldsymbol{\beta}_l}(i)), \quad (18)$$

where the update equations are

$$\gamma_{n,l} = \boldsymbol{\psi}_n^T \boldsymbol{\Lambda}_l \boldsymbol{\psi}_n + (\bar{\mathbf{v}}_l^T \boldsymbol{\psi}_n)^2 \quad (19)$$

$$h_{n,l} = \sum_{k \in \mathcal{E}_l^R} \zeta_{n,k}^l \mathbb{E}[\log \mathcal{N}(y_n|\boldsymbol{\omega}_k^T \boldsymbol{\psi}_n, \tau_k^{-1})] - \dots \\ \sum_{k \in \mathcal{E}_l^L} \zeta_{n,k}^l \mathbb{E}[\log \mathcal{N}(y_n|\boldsymbol{\omega}_k^T \boldsymbol{\psi}_n, \tau_k^{-1})] \quad (20)$$

$$\mathbf{V}_k^{-1} = \mathbb{E}[\mathbf{A}_k] + \sum_{n=1}^N \mathbb{E}[\zeta_{n,k}] \boldsymbol{\psi}_n \boldsymbol{\psi}_n^T \quad (21)$$

$$\bar{\boldsymbol{\omega}}_k = \mathbf{V}_k \sum_{n=1}^N \mathbb{E}[\zeta_{n,k}] y_n \boldsymbol{\psi}_n \quad (22)$$

$$a_{\tau_k} = a_0 + 0.5 \sum_{n=1}^N \mathbb{E}[\zeta_{n,k}] \quad (23)$$

$$b_{\tau_k} = b_0 + 0.5 \left( \sum_{n=1}^N \mathbb{E}[\zeta_{n,k}] (\bar{\boldsymbol{\omega}}_k^T \boldsymbol{\psi}_n - y_n)^2 + \dots \right. \\ \left. \bar{\boldsymbol{\omega}}_k^T \mathbb{E}[\mathbf{A}_k] \bar{\boldsymbol{\omega}}_k \right) \quad (24)$$

$$\boldsymbol{\Lambda}_l^{-1} = \mathbb{E}[\mathbf{B}_l] + 2 \sum_{n=1}^N \lambda(\gamma_{n,l}) \boldsymbol{\psi}_n \boldsymbol{\psi}_n^T \quad (25)$$

$$\bar{\mathbf{v}}_l = \boldsymbol{\Lambda}_l \sum_{n=1}^N (\mathbb{E}[z_{n,l}] - 0.5) \boldsymbol{\psi}_n \quad (26)$$

$$\mathbf{a}_{\boldsymbol{\alpha}_k}(i) = a_0 + 0.5 \quad (27)$$

$$\mathbf{b}_{\boldsymbol{\alpha}_k}(i) = b_0 + 0.5 \mathbb{E}[\tau_k] \bar{\boldsymbol{\omega}}_k(i)^2 + \mathbf{V}_k(i, i) \quad (28)$$

$$\mathbf{a}_{\boldsymbol{\beta}_l}(i) = a_0 + 0.5 \quad (29)$$

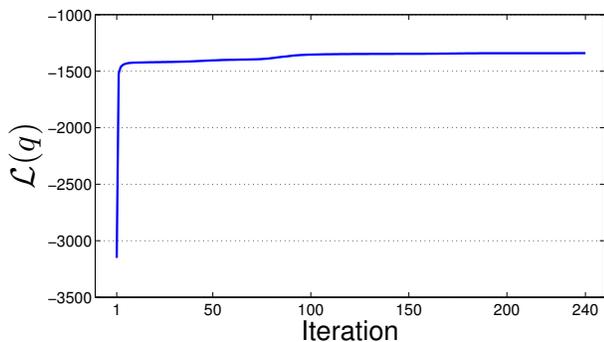
$$\mathbf{b}_{\boldsymbol{\beta}_l}(i) = b_0 + 0.5 \bar{\mathbf{v}}_l(i)^2, \quad (30)$$

where  $\mathbb{E}$  denotes the expectation with respect to  $q^*(\boldsymbol{\theta})$ ,  $\mathcal{E}_l^R$  and  $\mathcal{E}_l^L$  denote the set of experts on the right-hand-side and left-hand-side of the  $l$ -th gate respectively, and  $\zeta_{n,k}^l$  is computed as  $\zeta_{n,k}$  in Equation (3) with terms corresponding to  $l$ -th gate being omitted. By  $\mathbf{B}_l$  and  $\mathbf{A}_k$  we denote diagonal matrices with  $\boldsymbol{\beta}_l$  and  $\boldsymbol{\alpha}_k$  on their diagonal, respectively. The  $i$ -th element on the diagonal of  $\mathbf{V}_k$  is denoted by  $\mathbf{V}_k(i, i)$ . Finally, the expectation of  $\mathbb{E}[\zeta_{n,k}]$  is computed as

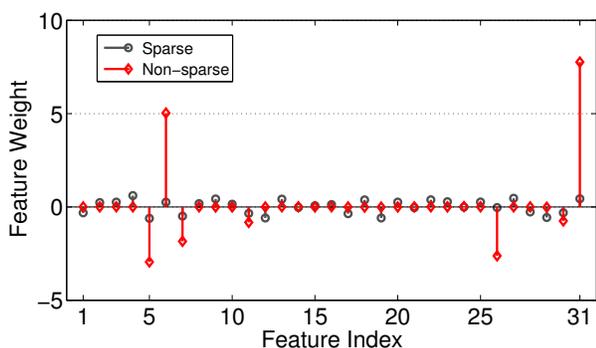
$$\mathbb{E}[\zeta_{n,k}] = \prod_{l=1}^M \sigma(h_{n,l})^{\mathbf{S}^L(k,l)} (1 - \sigma(h_{n,l}))^{\mathbf{S}^R(k,l)} \quad (31)$$

Considering that  $q^*(\tau_k)$ ,  $q^*(\boldsymbol{\alpha}_k)$ , and  $q^*(\boldsymbol{\beta}_l)$  are Gamma distributed, we can readily compute the expectation involving variables  $\tau_k$ ,  $\boldsymbol{\alpha}_k$ , and  $\boldsymbol{\beta}_l$ .

ARD is reflected in Equations (28) and (30), where  $\mathbf{b}_{\boldsymbol{\alpha}_k}(i)$  and  $\mathbf{b}_{\boldsymbol{\beta}_l}(i)$  specify the prior over the  $i$ -th feature in the  $k$ -th expert, and the  $l$ -th gate, respectively.



**Fig. 2.** Convergence of lower-bound to log-model-evidence,  $\mathcal{L}(q)$ , for random initialization.



**Fig. 3.** Feature weights of the only gate in the network for non-sparse and sparse HME models. The features are normalized to zero-mean and unit variance.

## 5. EXPERIMENTAL EVALUATION

To evaluate the effectiveness of sparse HME, we compared its performance against P.563 as well as the variational Bayesian HME of [3]. We used the ITU-T coded-speech data-set, Supplement 23 [9]. More details about our data-set and experimental method is given in [8]. A measure called the condition-averaged correlation-coefficient (CACC) is often used to quantify the accuracy of quality assessment algorithms in cross-validation experiments. For brevity we omit details of the experiments, which are identical to [8]. The CACC performance of the new sparse HME matched that of the non-sparse HME, both yielding 0.88. This result confirms an improvement compared to the P.563, which yielded 0.87. We used Gamma hyper-parameters  $a_0 = 10^{-2}$  and  $b_0 = 10^{-4}$  to obtain broad support for the priors.

Figure 2 shows the convergence of the lower bound to log model evidence  $\mathcal{L}(q)$  as terms in Equation (12) are updated by our algorithm. The monotonic convergence of  $\mathcal{L}(q)$  confirms the correctness of the iterative algorithm and its implementation.

We trained a HME model with a gating network comprising one gate and two experts. The feature weights of the

gate are shown in Figure 3 for the sparse HME as well as the non-sparse HME. Sparse HME pushed most features to zero, while allowing a few features to take on large values. In contrast, non-sparse HME allowed several small non-zero feature weights. This result confirmed that indeed sparsity in the gate was obtained.

## 6. CONCLUSION

We presented the prior and a variational inference algorithm for the HME model that results in sparse gates and experts. Sparse HME was successfully deployed to build non-intrusive speech quality assessment algorithm. While sparse HME matched the performance of non-sparse HME, the gate mechanism used fewer number of features.

## 7. REFERENCES

- [1] M.I. Jordan and R.A. Jacobs, “Hierarchical mixtures of experts and the EM algorithm,” *Neural Computation*, vol. 6, no. 2, pp. 181–214, 1994.
- [2] L. Breiman, *Classification and regression trees*, Chapman & Hall/CRC, 1984.
- [3] C.M. Bishop and M. Svensén, “Bayesian hierarchical mixtures of experts,” in *Proc. Nineteenth Conference on Uncertainty in Artificial intell.*, 2003, pp. 57–64.
- [4] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [5] D.J.C. MacKay, “A practical Bayesian framework for backpropagation networks,” *Neural Computation*, vol. 4, no. 3, pp. 448–472, 1992.
- [6] R.M. Neal, *Bayesian Learning for Neural Networks*, Springer Verlag, 1996.
- [7] “Single-ended method for objective speech quality assessment in narrow-band telephony applications,” 2004, ITU, Geneva, Switzerland, 2004, ITU-T Rec. P.563.
- [8] S. I. Mossavat, O. Amft, B. de Vries, P. N. Petkov, and W. B. Kleijn, “A Bayesian hierarchical mixture of experts approach to estimate speech quality,” in *QoMEX*, 2010.
- [9] “ITU-T coded-speech database,” ITU, Geneva, Switzerland, 1998, ITU-T Rec. P.Supp1. 23.