

Econometric Applications of Hierarchical Mixture of Experts

Lucas C. Dowiak

January 14, 2021

PhD Program in Economics, City University of New York, Graduate Center,
New York, NY, 10016, *Email: ldowiak@gradcenter.cuny.edu*

Abstract

In this article, a novel mixture model is studied. Named the hierarchical mixture of experts (HME) in the machine learning literature, the mixture model utilizes a set of covariates and a tree-based architecture to efficiently allocate each observation to the appropriate local regression. The nature of the conditional weighting scheme provides the researcher a natural interpretation of how the local (and latent) sub-populations are formed. The model is demonstrated by estimating a Mincer earning function using census data. Marginal effects, robust standard errors, a tree-growing algorithm, and a modest extension are also discussed.

Keywords: Hierarchical mixture of experts, expectation maximization

JEL Classification:

1 Introduction

The concepts of mixture models and mixture distributions are old hat in the economics field. Hamilton 1989 and Goldfeld and Quandt 1973 are a few of the pioneering works for time series and cross sectional regression, respectively. We are

also deep into the age of machine learning, and it's reigning champion, the artificial neural network, has been successfully adapted and studied in the context of applied econometrics. This article adds to the small body of literature that employs a novel neural network architecture to model the weights of a mixture model. In doing so, we leverage the highly flexible nature of a neural network but maintain interpretability and the means to quantify marginal effects. The model under investigation is called the Hierarchical Mixture of Experts (HME), a class of mixture models whose defining feature is its conditional weighting scheme. The model's origin story traces back to R. A. Jacobs et al. 1991. The authors use a single multinomial classifier to assign, in a probabilistic sense, input patterns to local *experts*. These experts are almost always some flavor of regression or classification model. The multinomial structure that assigns inputs to experts is referred to as the *gating network*. The authors employ this mixture of experts (ME) framework to model vowel discrimination in a speech recognition context. Shortly after, M. I. Jordan and Robert A. Jacobs 1992 generalize this single-layer multinomial gating network to one with an arbitrary number of layers. M. Jordan and R. Jacobs 1993 then demonstrate an Expectation-Maximization approach to model estimation that is capable handling the additional complexity the generalization requires during optimization. The result of this extension is a gating network that takes on a tree-like structure, stemming from an initial multinomial split and filtering down through additional multinomial partitions of the input space. HME models nest ME models as special case. Pushing a little further, one additional case is studied as well. As the depth of an HME grows, so too must the number of experts. If we have a symmetric HME network, this growth is geometric with respect to the network's depth. With this in mind, a further extension can be considered where each expert is not unique, but a member of a fixed set of experts. We refer to this additional model as a Hierarchical Mixture of Repeated Experts (HMRE). Figure (1) provides an example of each of the variations of this class of model.

This article investigates the adoption of ME and HME models to an applied econometric framework, with particular attention focused on interpretation of the gating network and robust inference of parameter estimates. The outline for the rest of this manuscript is as follows: the remainder of this section provides a brief literature review and section 2 describes the model in formal detail. Section 3 discusses the expectation-maximization approach to estimation while section 4 concerns itself with robust inference of the estimated parameters. Section 5 provides detail on how to derive the marginal effects of the model's covariates. In section 6, we provide a very simple demonstration of the HME in action and then move on to a more economically

relevant example of applying the HME model to a Mincer wage equation in section 7. Section 8 concludes.

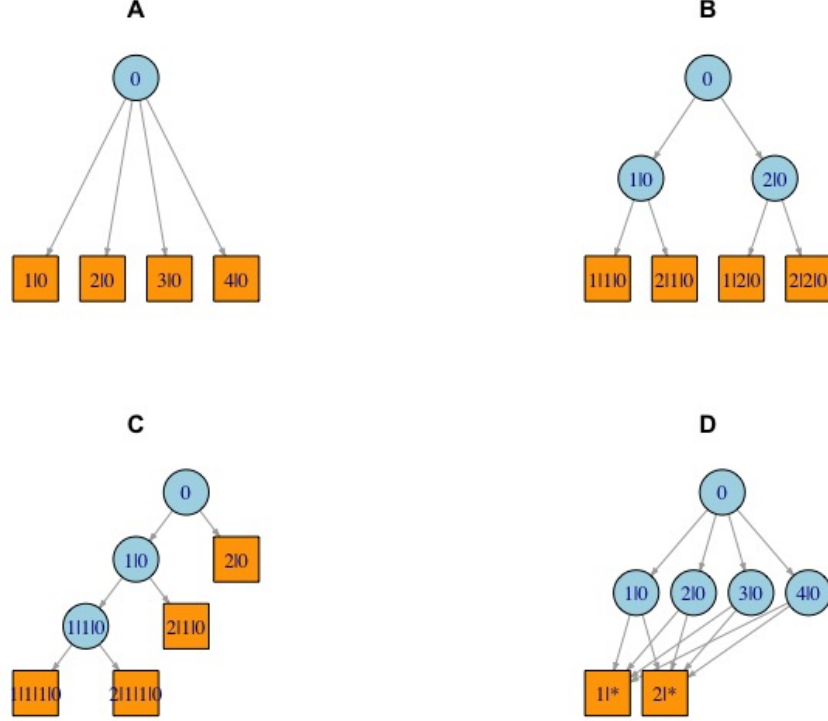


Figure 1: Networks **A** - **D** depict various network architectures that are discussed in this article. For all four networks, gating nodes are represented as blue circles and experts as orange rectangles. Network **A** illustrates the original Mixture of Experts (ME) architecture with a single multinomial split leading to a set of experts one layer down. Networks **B** and **C** both represent different flavors of a Hierarchical Mixture of Experts (HME). Network **B** is a symmetric network of depth 2 with successive binary splits. Network **C** is an asymmetric network of depth 3 with successive binary splits. Network **D** is an example of the Hierarchical Mixture of Repeated Experts (HMRE) architecture. Notice that multiple paths exist from the root node 0 to each expert. Compare this to networks **A** - **C**, where there is only one unique path from the root node to each expert.

1.1 Relevant Literature

ME and HME frameworks have been utilized for both time series and cross-sectional analysis. Within the cross-sectional literature, S. Waterhouse and A. Robinson 1995 puts forth a method to grow an HME from a single split from the root node. The authors are influenced by the popular technique used for classification and regression trees Brieman et al. 1984 and apply it to an HME structure. Once the gating structure to an HME tree has been grown, the authors suggest an additional trimming algorithm to prevent overfitting. Fritsch, Finke, and Waibel 1997 consider S. Waterhouse and A. Robinson 1995 and alter their growing algorithm with a mind to scaling the model to handle thousands of experts. M. Jordan and Xu 1995 provide an extended discussion on the convergence of the model used by M. Jordan and R. Jacobs 1993. The authors also suggest algorithmic improvements to help with estimation. Continuing the theoretical discussing, Jiang and M. A. Tanner 1999 cover convergence rates of an HME model where experts are from the exponential family with generalized linear mean functions. Jiang and M. A. Tanner 2000 provide regularity conditions on the HME structure for for a mixture of general linear models estimated by maximum likelihood to produce consistent and asymptotically normal estimates of the mean response. The conditions are validated for poisson, gamma, gaussian, and binomial experts.

Alternatively, Weigend, Mangeas, and Srivastava 1995 provide a detailed discussion examining ME applied in a time series context and provide valuable insights to avoid overfitting the model to the data, a common problem in neural network applications. Huerta, Jiang, and M. A. Tanner 2003 extend Weigend, Mangeas, and Srivastava 1995 to an HME framework. Using five and a half decades of monthly US industrial production data, the authors allow the series to choose between two models, one modeled as a random walk and the other as trend stationary. In addition, they present a Bayesian approach to estimation. Carvalho and M. Tanner 2003 lay out the necessary regularity conditions to perform hypothesis tests on stationary ME time series of generalized linear models (ME-GLM) using Wald tests. The dual cases of a well-specified and a misspecified model are considered. The authors restrict their analysis to ME-GLM models involving lagged dependent and lagged external covariate variables only. Generalization to include lagged conditional mean values is left for another time. Carvalho and M. Tanner 2005 take a similar approach to Carvalho and M. Tanner 2003 but apply their analysis to a purely auto-regressive context restricted to Gaussian models. The authors extend arguments in Carvalho and M. Tanner 2003 to non-stationary series and provide simulated evidence that using the BIC is helpful in selecting the appropriate number of experts to include.

Carvalho and M. Tanner 2006 re-focus the discussion on ME of time series regressions restricted to exponential family distributions. Distilling the available literature at the time, the authors cover the important topics of estimation and asymptotic properties in the maximum likelihood framework, selection of the number of experts, model validation and fitting. Carvalho and Skoulakis 2010 applies ME of a single time series. Using stock returns the authors structure the gating network using lagged dependent variables and an 'external' covariate capturing a measure of the trade volume at that time.

In this article estimation and inference is from a maximum likelihood perspective and will remain the primary focus. Estimation of ME and HME models from a Bayesian has received considerable amount of attention as well. S. R. Waterhouse, MacKay, and A. J. Robinson 1995 provided an initial approach to estimating a ME by combining gaussian priors on the gating and expert parameters with gamma hyperparameter priors in an approximating ensemble to the true joint density of the model. Optimization of the parameter vector for the approximating density occurs a block of parameters at a time. Ueda and Ghahramani 2002 improve on S. R. Waterhouse, MacKay, and A. J. Robinson 1995 by optimizing for the appropriate number of experts in addition to model parameters. Bishop and Svenson 2003 find previous bayesian approaches to estimating an HME lacking. Using variational inference, the authors provide a complete bayesian estimation approach to the log marginal likelihood. With an eye to prediction, the author's advocate that their approach makes the HME model easier to estimate without overfitting.

1.2 Additional Articles to Include

Neal and Pfeiffer 2001 cross section

Blei, Kucukelbir, and McAuliffe 2016 A review of variational inference applied to generalized linear models and basic examples.

Carvalho and Skoulakis 2005

2 Model

We start by presenting the HME as a standard mixture model. For a given input and output pair (\mathbf{x}_t, y_t) , each expert provides a probabilistic model relating input \mathbf{x}_t to output y_t :

$$P_t^m \equiv P^m(y_t|\mathbf{x}_t; \boldsymbol{\beta}^m), \quad m = 1, 2, \dots, M \quad (1)$$

where m is one of the M component experts in the mixture. The experts are combined with associated weights into a mixture distribution

$$P(y_t|\mathbf{x}_t; \boldsymbol{\beta}) = \sum_{m=1}^M \mathbb{I}(m|t) P^m(y_t|\mathbf{x}_t; \boldsymbol{\beta}^m) \quad (2)$$

Here, $\mathbb{I}_t(m)$ is the probability that the input unit t belongs to expert m and has the usual restrictions: $0 \leq \mathbb{I}(m|t) \leq 1$ for each m and $\sum_m \mathbb{I}(m|t) = 1$. The gating network of the model applies a particular functional form to model $\mathbb{I}(m|t)$, which includes a second set of covariates \mathbf{z}_t and parameter vector $\boldsymbol{\omega}$:

$$P(y_t|\mathbf{x}_t, \mathbf{z}_t; \boldsymbol{\beta}, \boldsymbol{\Omega}) = \sum_{m=1}^M \mathbb{I}(m|\mathbf{z}_t; \boldsymbol{\Omega}) P^m(y_t|\mathbf{x}_t; \boldsymbol{\beta}^m) \quad (3)$$

2.1 Gating Network and $\mathbb{I}(m|\mathbf{Z}, \boldsymbol{\Omega})$

The gating network model is structured as a collection of nodes in a tree structure that branches out in successive layers. The location of these nodes will be referred to by their address a . The root node resides at the apex of the tree and has the address 0. The root node then splits into J different nodes, one level down the tree. The addresses for these J new nodes are $1|0, 2|0, \dots, J|0$. This type of naming convention continues as the rest of network is traversed. At its most general, each gating node can yield an arbitrary number of splits. While a fully generalized gating network is conceptually attractive, it presents practical challenges for implementation. In this paper we address several architectures for the gating network, each with its own set of structural restrictions on the shape of the network and the number of splits each gating node can take. For arbitrary gating node at address a , we use a multinomial logistic regression to model the split in direction i to be:

$$g_t^{a,i} \equiv g_t^{a,i}(\mathbf{z}_t, \boldsymbol{\omega}^a) = \frac{\exp(\mathbf{z}_t^\top \boldsymbol{\omega}^{a,i})}{\sum_{j=1}^J \exp(\mathbf{z}_t^\top \boldsymbol{\omega}^{a,j})} \quad (4)$$

The parameters in equation (4) are subject to the usual identification restrictions. For the remainder of the article, we choose to set $\boldsymbol{\omega}^{a,J} = \mathbf{0}$ for every gating node. It is important to keep track of the product path an input vector travels from one node to another. If the observation index is suppressed, the product path from one node (say the root node 0) to another (say $k|\dots|j|i$) can be defined as

$$\pi_{g^0 \longleftrightarrow g^k | \dots | j | i | 0} = \begin{cases} g^{0,i} g^{i|0,j} \dots g^{\dots | j | i | 0, k} & \text{if path is feasible} \\ 1 & \text{otherwise} \end{cases} \quad (5)$$

If one of the nodes is an expert, then we can define the mixture weight of expert m for input pattern i to be the product of the path taken from the root node to expert m :

$$\mathbb{P}(m | \mathbf{Z}, \boldsymbol{\Omega}) = \pi_{g^0 \longleftrightarrow P^m} \quad (6)$$

For network architectures with multiple paths from the root node to the same expert (see bottom right of figure (1)), we can index these multiples paths by l so that

$$\mathbb{P}(m | \mathbf{Z}, \boldsymbol{\Omega}) = \sum_l \pi_{g^0 \xleftrightarrow{l} P^m} \quad (7)$$

By collecting and summing all possible paths from the root node to each expert, the conditional probability given in equation (3) can be expanded and expressed as:

$$\begin{aligned} P(y_t | \mathbf{x}_t, \mathbf{z}_t; \boldsymbol{\Omega}, \boldsymbol{\beta}) &= \sum_m \mathbb{P}(m | \mathbf{z}_t, \boldsymbol{\Omega}) P^m(y_t | \mathbf{x}_t; \boldsymbol{\beta}^m) \\ &= \sum_m P^m(y_t | \mathbf{x}_t; \boldsymbol{\beta}^m) \sum_l \pi_{g^0 \xleftrightarrow{l} P^m} \end{aligned} \quad (8)$$

If we concatenate the parameters of the gating network with the parameters of the experts as $\boldsymbol{\theta} = [\boldsymbol{\Omega} \ \boldsymbol{\beta}]$, then the product of these individual probabilities across the full sample size T yields the likelihood function.

$$\mathcal{L}(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}, \mathbf{Z}) = \prod_t \sum_m P^m(y_t | \mathbf{x}_t; \boldsymbol{\beta}^m) \sum_l \pi_{g^0 \xleftrightarrow{l} P^m} \quad (9)$$

And taking its log yields the log likelihood

$$\mathbf{l}(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}, \mathbf{Z}) = \sum_t \log \sum_m P^m(y_t | \mathbf{x}_t; \boldsymbol{\beta}^m) \sum_l \pi_{g^0 \xleftrightarrow{l} P^m} \quad (10)$$

The functional form of the log likelihood (10) does not lend itself easily to direct optimization, but a well established technique using expectation maximization (Dempster, Laird, and Rubin 1977) to estimate mixture models is available. This was the primary insight of M. Jordan and R. Jacobs 1993's original paper.

3 The EM Set-Up

The EM approach to estimating an HME model starts by suggesting that if a researcher had perfect information, each input vector \mathbf{x}_t could be matched to the expert P^m that generated it with certainty. If a set of indicator variables is introduced that captures this certainty, an *augmented* version of the likelihood in equation (9) can be put forward. Define the indicator set as:

$$I_t(m) = \begin{cases} 1 & \text{if observation } t \text{ is generated by expert } m \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

We can then reformulate the likelihood equation

$$\mathcal{L}_c(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}, \mathbf{Z}) = \prod_t \prod_m \left[P^m(y_t|\mathbf{x}_t; \boldsymbol{\beta}^m) \sum_l \pi_{g^0 \xleftrightarrow{l} P^m} \right]^{I_t(m)} \quad (12)$$

leading to the complete-data log-likelihood

$$l_c(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}, \mathbf{Z}) = \sum_t \sum_m I_t(m) \left[\log P^m(y_t|\mathbf{x}_t; \boldsymbol{\beta}^m) + \log \sum_l \pi_{g^0 \xleftrightarrow{l} P^m} \right] \quad (13)$$

As mentioned previously, summing over multiple paths l in equation (13) is only necessary in the HMRE case. For the ME and HME cases, $l = 1$, simplifying the second log in (13) to $\log(\pi_{g^0 \xleftrightarrow{1} P^m})$. Going forward, we will focus our analysis on the ME and HME specifications with work on the HMRE case left for another time.

3.1 E-Step

The E-step of the algorithm performs an expectation over the complete log-likelihood equation (13), where the expectation includes the additional information contained in the expert regressions. One of the results of this expectation is the creation of second set of weights h^a that parallel the weights from the gating network g^a discussed in section (2.1). For an HME model:

$$\begin{aligned}
Q(\boldsymbol{\theta}) &= \mathbb{E} [\mathbf{l}_c(\boldsymbol{\theta}|\mathbf{y}, \mathbf{X}, \mathbf{Z})] = \sum_t \sum_m \mathbb{E} [I_t(m)] [\log P^m(y_t|\mathbf{x}_t; \boldsymbol{\beta}^m) + \log \pi_{g_t^0 \longleftrightarrow P_t^m}] \\
&= \sum_t \sum_m \pi_{h_t^0 \longleftrightarrow P_t^m} [\log P^m(y_t|\mathbf{x}_t; \boldsymbol{\beta}^m) + \log \pi_{g_t^0 \longleftrightarrow P_t^m}] \\
&= \sum_t Q_t(\boldsymbol{\theta})
\end{aligned} \tag{14}$$

Here $\pi_{h^0 \longleftrightarrow h^k, \dots | j|i|0}$ is analogous to equation (5)

$$\pi_{h^0 \longleftrightarrow h^k | \dots | j|i|0} = \begin{cases} h^{0,i} h^{i|0,j} \dots h^{\dots | j|i|0,k} & \text{if path is feasible} \\ 1 & \text{otherwise} \end{cases} \tag{15}$$

and the $h^{a,i}$ are arrived at using Bayes' theorem.

$$h_t^{a,i} = \frac{g^{a,i} \sum_k P_t^k \pi_{g_t^{i|a} \longleftrightarrow P^k}}{\sum_j g^{a,j} \sum_m P_t^m \pi_{g_t^{j|a} \longleftrightarrow P^m}} \tag{16}$$

So, now we have two different forms of weights, g 's and h 's. The way the g 's are formed in equation (4), they are only functions of the nodes in the gating network, separate from the expert regressions and the information they contain. For this reason, M. Jordan and R. Jacobs 1993 refer to g 's as *priors*. The h 's draw from both the gating network and the expert regressions and are referred to as *posterior* weights.

3.2 M-Step

One of the more attractive features of using EM to optimize a HME is how the log-likelihood function compartmentalizes into a set of independent functions which can be individually optimized. After taking the expectation of the log-likelihood function (14), the parameters governing each expert and each gating network can be grouped together and optimized on their own. For the experts we have:

$$\arg \max_{\boldsymbol{\beta}^m} \sum_t \pi_{h_t^0 \longleftrightarrow P_t^m} \log P^m(y_t|\mathbf{x}_t; \boldsymbol{\beta}^m) \tag{17}$$

And for the gating nodes:

$$\arg \max_{\boldsymbol{\omega}^a} \sum_t \pi_{h_t^0 \longleftrightarrow h_t^a} \log g(\mathbf{z}_t, \boldsymbol{\omega}^a) \tag{18}$$

4 Inference

When considering inference, it's worth thinking about what would motivate a researcher to turn to an HME model in the first place. At times, a researcher may suspect that a latent structure exists within the data and that a single regression $y_t = \mathbf{x}_t^\top \boldsymbol{\beta}$ may mask a critical change in relationship depending on membership to some unknown sub-group j of the data $y_{tj} = \mathbf{x}_t^\top \boldsymbol{\beta}_j$. A wide variety of time series, especially those with longer histories, experience changes in behaviour over time. They can be subjected to sharp one-off changes in value or more gradual changes of behavior over time. Regardless of the context, any latent structural change in the data generating process may also introduce some hidden form of heterogeneity to the error terms. Rather than taking a firm stance on any concealed structure, an HME setup ideally limits the work the researcher needs to do to specifying a set of well-chosen conditioning variables \mathbf{Z} to feed through the gating network. This limited workload may come at a cost, though. By allowing the gating network to find it's own mixture allocations, the odds of arriving at a misspecified model becomes a concern. To guard against this, we use a sandwich estimator for the variance-covariance matrix:

$$\mathbf{V}(\boldsymbol{\theta}) = \mathbf{H}^{-1}(\boldsymbol{\theta}) \mathbf{G}(\boldsymbol{\theta}) \mathbf{H}^{-1}(\boldsymbol{\theta}) \quad (19)$$

where $\mathbf{G}(\boldsymbol{\theta})$ is the sum of the outer products of the score vectors

$$\mathbf{G}(\boldsymbol{\theta}) = \sum_t \mathbf{S}_t(\boldsymbol{\theta}) \mathbf{S}_t(\boldsymbol{\theta})^\top \quad (20)$$

and $\mathbf{H}(\boldsymbol{\theta})$ is the empirical Hessian:

$$\mathbf{H}(\boldsymbol{\theta}) = \frac{1}{T} \sum_t \mathbf{H}_t(\boldsymbol{\theta}) \quad (21)$$

We discuss the score vectors and the Hessians in more detail in the following two sections.

4.1 The Score

For the score vector, we concatenate the scores of each gating node and those of each local expert regressions.

$$\mathbf{S}_t(\boldsymbol{\theta}) = [\mathbf{S}_t(\boldsymbol{\Omega}) \ \mathbf{S}_t(\boldsymbol{\beta})] \quad (22)$$

Starting with parameters of the gating network, the full vector can be partitioned in some logical order into the sub-vectors of each node’s individual score.

$$\mathbf{S}_t(\boldsymbol{\Omega}) = [\mathbf{S}_t(\boldsymbol{\omega}^0) \ \mathbf{S}_t(\boldsymbol{\omega}^{1|0}) \ \mathbf{S}_t(\boldsymbol{\omega}^{2|0}) \ \dots] \quad (23)$$

$$\mathbf{S}_t(\boldsymbol{\omega}^a) = [\mathbf{S}_t(\boldsymbol{\omega}^{a,1}) \ \dots \ \mathbf{S}_t(\boldsymbol{\omega}^{a,J-1})] \quad (24)$$

For a generic gating node a we can define the individual score for sample t as:

$$\mathbf{S}_t(\boldsymbol{\omega}^{a,i}) = \frac{\partial Q_t}{\partial \boldsymbol{\omega}^{a,i}} = \pi_{h_t^0 \longleftrightarrow h_t^a} (1 - g_t^{a,i}) \mathbf{z}_t \quad (25)$$

Turning our attention to the expert regressions, the exact functional form of the score vector depends on the type of regression we wish to run. In most cases, all experts in an HME model are from the same family (Huerta, Jiang, and M. A. Tanner 2003 is a notable exception). When all experts share the same functional form, it’s standard to accept the restriction that no experts in the HME model produce the same parameter vector $\boldsymbol{\beta}^j \neq \boldsymbol{\beta}^k$. Such an HME is defined by Jiang and M. A. Tanner 2000 as being *irreducible*. The irreducibility of an HME plays a critical role in guaranteeing the convergence of the model.

In this article, each HME discussed will employ a set of experts running a standard linear regression model with Gaussian errors. To aide with model optimization, the specification of the parameter vector for each regression $\boldsymbol{\beta}^m = [\beta_0^m \ \dots \ \beta_k^m \ \phi^m]$ takes on a unique form where we model the log variance explicitly $\phi = \log \sigma^2$.

$$P^m(y_t | \mathbf{x}_t; \boldsymbol{\beta}^m, \phi^m) = (2\pi \exp(\phi^m))^{-\frac{1}{2}} \exp \left(-\frac{(y_t - \mathbf{x}_t^\top \boldsymbol{\beta}^m)^2}{2 \exp(\phi^m)} \right) \quad (26)$$

In this case the score vector for any particular regression expert is:

$$\mathbf{S}_t(\boldsymbol{\beta}^m) = \left(\frac{\partial Q_t}{\partial \boldsymbol{\beta}^m}, \frac{\partial Q_t}{\partial \phi^m} \right)^\top \quad (27)$$

$$\mathbf{S}_t(\boldsymbol{\beta}) = [\mathbf{S}_t(\boldsymbol{\beta}^1) \ \dots \ \mathbf{S}_t(\boldsymbol{\beta}^M)] \quad (28)$$

with:

$$\frac{\partial Q_t}{\partial \boldsymbol{\beta}^m} = \pi_{h_t^0 \longleftrightarrow f_t^m} \frac{(y_t - \mathbf{x}_t^\top \boldsymbol{\beta}^m)}{\exp(\phi^m)} \mathbf{x}_t \quad (29)$$

and

$$\frac{\partial Q_t}{\partial \phi^m} = \frac{\pi_{h_t^0 \longleftrightarrow f_t^m}}{2} \left(\frac{(y_t - \mathbf{x}_t^\top \boldsymbol{\beta}^m)^2}{\exp(\phi^m)} - 1 \right) \quad (30)$$

4.2 The Hessian

The hessian is equally straight-forward. Starting with equation (25), the hessian for each gating node is:

$$\mathbf{H}_t(\boldsymbol{\omega}^a) \equiv \frac{\partial^2 Q}{\partial \boldsymbol{\omega}^{a,i} \partial \boldsymbol{\omega}^{a,j}} = \pi_{h_t^0 \longleftrightarrow h_t^a} \boldsymbol{\Gamma}_t^a \otimes \mathbf{z}_t \mathbf{z}_t^\top \quad (31)$$

where \otimes is the kronecker product and:

$$\boldsymbol{\Gamma}_t^a = \begin{bmatrix} -g_t^{a,1}(1 - g_t^{a,1}) & g_t^{a,1}g_t^{a,2} & \dots & g_t^{a,1}g_t^{a,J-1} \\ g_t^{a,1}g_t^{a,2} & -g_t^{a,2}(1 - g_t^{a,2}) & \dots & g_t^{a,2}g_t^{a,J-1} \\ \vdots & \vdots & \ddots & \vdots \\ g_t^{a,1}g_t^{a,J-1} & g_t^{a,2}g_t^{a,J-1} & \dots & -g_t^{a,J-1}(1 - g_t^{a,J-1}) \end{bmatrix} \quad (32)$$

For each expert regression:

$$\mathbf{H}_t(\boldsymbol{\beta}^m) = \frac{\pi_{h_t^0 \longleftrightarrow f_t^m}}{\exp(\phi^m)} \begin{bmatrix} \mathbf{x}_t \mathbf{x}_t^\top & \mathbf{x}_t \epsilon_t^m \\ \mathbf{x}_t^\top \epsilon_t^m & \frac{1}{2}(\epsilon_t^m)^2 \end{bmatrix} \quad (33)$$

where we have set $\epsilon_t^m = y_t - \mathbf{x}_t^\top \boldsymbol{\beta}^m$ to ease the notational burden. Staying consistent with the score vector, we sum the hessian matrices across observations:

$$\mathbf{H}(\boldsymbol{\omega}^a) = \sum_t^T \mathbf{H}_t(\boldsymbol{\omega}^a) \quad (34)$$

$$\mathbf{H}(\boldsymbol{\beta}^m) = \sum_t^T \mathbf{H}_t(\boldsymbol{\beta}^m) \quad (35)$$

5 Marginal Effects

Due to the complexity of the model's structure and the ability to place covariates in either the gating network, the expert regressions, or both, viewing the relationship between the covariates and the dependent variable through their marginal effects may provide a simplifying lens of the model's governing principles. Just as for logistic and multinomial regression, the marginal effects of an HME model have a closed form solution. Starting with equation (3) we replace the expert distributions P_t^m with the regression's output y_t^m and use the relationship in equation (6) to arrive at:

$$y_t = f(\mathbf{x}_t, \mathbf{z}_t; \boldsymbol{\beta}, \boldsymbol{\Omega}) = \sum_{m=1}^M \pi_{g_t^0 \longleftrightarrow y_t^m} y_t^m \quad (36)$$

The functional form of the marginal effect depends on where the variables appear in the model. Our existing notation labels the covariates in gating network as \mathbf{Z} and the covariates in the expert regressions as \mathbf{X} . As seen later, the variables belonging to \mathbf{Z} and \mathbf{X} do not need to be mutually exclusive. There is also no requirement that they differ at all. In light of this, a few more notational definitions are needed to cover all the cases:

- $\mathbf{T} = \mathbf{Z} \cup \mathbf{X}$
- $\mathbf{W} = \mathbf{Z} \cap \mathbf{X}$
- $\mathbf{U}_Z = \mathbf{Z} \setminus \mathbf{X}$
- $\mathbf{U}_X = \mathbf{X} \setminus \mathbf{Z}$

The full list of variables considered in the model is labeled \mathbf{T} . Covariates that appear in both the gating network and the expert regressions are collect in \mathbf{W} . \mathbf{U}_Z and \mathbf{U}_X are used to label variables that appear only in the gating network or only in the expert regressions, respectively. With this notation, we can express the full marginal effects of the HME by where the explanatory variables appear in the model.

$$\frac{\partial y_t}{\partial \mathbf{T}} \equiv \boldsymbol{\Delta} = \sum_{m=1}^M \boldsymbol{\Delta}^m = \sum_{m=1}^M \left[\frac{\partial y_t^m}{\partial \mathbf{U}_Z} \quad \frac{\partial y_t^m}{\partial \mathbf{W}} \quad \frac{\partial y_t^m}{\partial \mathbf{U}_X} \right] \quad (37)$$

with the functional form of the each covariate group in (37) defined as:

$$\frac{\partial y_t^m}{\partial \mathbf{U}_Z} = \frac{\partial \pi_{g_t^0 \longleftrightarrow y_t^m}}{\partial \mathbf{U}_Z} y_t^m \quad (38)$$

$$\frac{\partial y_t^m}{\partial \mathbf{U}_X} = \pi_{g_t^0 \longleftrightarrow y_t^m} \frac{\partial y_t^m}{\partial \mathbf{U}_X} \quad (39)$$

$$\frac{\partial y_t^m}{\partial \mathbf{W}} = \frac{\partial \pi_{g_t^0 \longleftrightarrow y_t^m}}{\partial \mathbf{W}} y_t^m + \pi_{g_t^0 \longleftrightarrow y_t^m} \frac{\partial y_t^m}{\partial \mathbf{W}} \quad (40)$$

Not matter how complex the model becomes, the researcher can always interpret the estimated HME through a single vector of marginal effects of \mathbf{T} .

Of the four components in equations (38) - (40), three have already been established: y_t^m is the output from local expert m , $\pi_{g_t^0 \longleftrightarrow y_t^m}$ is the prior weight for input t for local expert m , and $\frac{\partial y_t^m}{\partial \mathbf{X}}$ is the marginal effect of the local expert m . What is left is the partial derivative of the gating network with respect to a variables in that network $\frac{\partial \pi_{g_t^0 \longleftrightarrow y_t^m}}{\partial \mathbf{Z}}$. Starting with equation (5), we take the partial with respect to gating matrix \mathbf{Z} :

$$\delta^m \equiv \frac{\partial \pi_{g_t^0 \longleftrightarrow y_t^m}}{\partial \mathbf{Z}} = \frac{\partial g^{0,i} g^{i|0,j} \dots g^{k|\dots|j|i|0,m}}{\partial \mathbf{Z}} \quad (41)$$

and applying the product rule gives us:

$$\begin{aligned} \delta^m &= \frac{\partial g^{0,i}}{\partial \mathbf{Z}} g^{i|0,j} \dots g^{k|\dots|j|i|0,m} \\ &+ g^{0,i} \frac{\partial g^{i|0,j}}{\partial \mathbf{Z}} \dots g^{k|\dots|j|i|0,m} \\ &+ \dots \\ &+ g^{0,i} g^{i|0,j} \dots \frac{\partial g^{k|\dots|j|i|0,m}}{\partial \mathbf{Z}} \end{aligned} \quad (42)$$

and since:

$$\frac{\partial g^{a,i}}{\partial \mathbf{Z}} = g^{a,i} \left(\omega^{a,i} - \sum_j g^{a,j} \omega^{a,j} \right) = g^{a,i} (\omega^{a,i} - \bar{\omega}^a) \quad (43)$$

we can substitute equation (43) into (42) to arrive at:

$$\delta^m = \pi_{g_t^0 \longleftrightarrow y_t^m} (\omega^{0,i} + \omega^{i|0,j} + \dots + \omega^{k|\dots|j|i|0,m} - (\bar{\omega}^0 + \bar{\omega}^{i|0} + \dots + \bar{\omega}^{k|\dots|j|i|0})) \quad (44)$$

5.1 Delta method

Using the delta method, we can approximate standard errors for the marginal effects of the HME model. Starting with equation (37) from the previous section, we break down the gradient of the marginal effects with respect to the parameters by those in the gating network, $\boldsymbol{\Omega}$, and the parameters in the expert regression, $\boldsymbol{\beta}$. These results are collected in table 1.

	$\underline{U_Z}$	\underline{W}	$\underline{U_X}$
$\frac{\partial \Delta_t^m}{\partial \omega^a}$	$\frac{\partial \delta_t^m}{\partial \omega^a} y_t^m$	$\frac{\partial \delta_t^m}{\partial \omega^a} y_t^m + \frac{\partial \pi_{g_t^0 \leftrightarrow y_t^m}}{\partial \omega^a} \frac{\partial y_t^m}{\partial \mathbf{W}}$	$\mathbf{0}$
$\frac{\partial \Delta_t^m}{\partial \beta^m}$	$\mathbf{0}$	$\delta_t^m \frac{\partial y_t^m}{\partial \beta^m} + \pi_{g_t^0 \leftrightarrow y_t^m} \frac{\partial^2 y_t^m}{\partial \mathbf{W} \partial \beta^m}$	$\pi_{g_t^0 \leftrightarrow y_t^m} \frac{\partial^2 y_t^m}{\partial U_X \partial \beta^m}$

Table 1: Delta Method Gradient Cases

Again, many of the expressions in table 1 have already been defined in previous sections. The two expressions new to this section are $\frac{\partial^2 y_t^m}{\partial \mathbf{X} \partial \beta^m}$ and $\frac{\partial \delta_t^m}{\partial \omega^{a,i}}$. For the standard OLS regressions that are considered in this paper, $\frac{\partial^2 y_t^m}{\partial \mathbf{X} \partial \beta^m} = \mathbf{1}$. Conceptually, $\frac{\partial \delta_t^m}{\partial \omega^{a,i}}$ describes how the marginal effects of the gating network change in response to small changes in the parameters of $\mathbf{\Omega}$. The value of $\frac{\partial \delta_t^m}{\partial \omega^{a,i}}$ depends on what role $\omega^{a,i}$ plays in navigating an input pattern from the root node to the expert m . For instance, say that we're at the root node, and it's our mission is to traverse the gating network down to expert m . When we arrive at node a , if the direction we need to take to reach expert m is along path i , then we'll call $\omega^{a,i}$ an *explicit* parameter set with respect to expert m . If taking path i leads to a different expert, then $\omega^{a,i}$ will be referred to as an *implicit* parameter set.

For an explicit path

$$\frac{\partial \delta_t^m}{\partial \omega_p^{a,i}} = \pi_{g^0 \leftrightarrow f^m} [(1 - g^{a,i}) + [W_p^m(1 - g^{a,i}) - G_p^{a,i}] Z_p] \quad (45a)$$

$$\frac{\partial \pi_{g^0 \leftrightarrow f^m}}{\partial \omega_p^{a,i}} = \pi_{g^0 \leftrightarrow f^m} (1 - g^{a,i}) Z_p \quad (45b)$$

and for an implicit path

$$\frac{\partial \delta_t^m}{\partial \omega_p^{a,j}} = \pi_{g^0 \leftrightarrow f^m} [-g^{a,j} + [-W_p^m(1 - g^{a,j}) - G_p^{a,j}] Z_p] \quad (46a)$$

$$\frac{\partial \pi_{g^0 \leftrightarrow f^m}}{\partial \omega_p^{a,j}} = -\pi_{g^0 \leftrightarrow f^m} g^{a,j} Z_p \quad (46b)$$

where

$$W_p^m = [\omega_p^{0,i} + \dots + \omega_p^{k|\dots|j|i|0,m} - (\bar{\omega}_p^0 + \dots + \bar{\omega}_p^{k|\dots|j|i|0})] \quad (47)$$

$$G_p^{a,i} = \left\{ g^{a,i}(1 - g^{a,i})\omega_p^{a,i} - \sum_{j \neq i} g^{a,i}g^{a,j}\omega_p^{a,j} \right\} \quad (48)$$

The intermediate step for equation 45a

$$\begin{aligned} \frac{\partial \delta_t^m}{\partial \omega_p^{a,i}} = & (1 - g^{a,i})\pi_{g^0 \longleftrightarrow f^m} [\omega^{0,i} + \dots + \omega^{k|\dots|j|i|0,m} - (\bar{\omega}^0 + \dots + \bar{\omega}^{k|\dots|j|i|0})] Z_p + \\ & \pi_{g^0 \longleftrightarrow f^m} \left[(1 - g^{a,i}) - \left\{ g^{a,i}(1 - g^{a,i})\omega_p^{a,i} - \sum_{j \neq i} g^{a,i}g^{a,j}\omega_p^{a,j} \right\} Z_p \right] \end{aligned}$$

Standard errors for the marginal effects for the HME models can then be constructed with the robust variance-covariance matrix from equation (19) and the collection of equations from (37) to (46).

$$Asy.Var [\hat{\Delta}] = \sum_{n=1}^M \left(\frac{1}{T} \sum_{t=1}^T \frac{\partial \Delta_t}{\partial \theta_n} \right) \mathbf{V}(\hat{\theta}) \left(\frac{1}{T} \sum_{t=1}^T \frac{\partial \Delta_t}{\partial \theta_n} \right)^\top \quad (49)$$

Note that equation (49) estimates the precision of the marginal effects of \mathbf{T}_k on the full model, which relies on equations (38), (39), and (40). If desired, we can isolate the marginal effects of the variables in the gating network (\mathbf{Z}_k) on the outcome, which would take a slight modification of equation (38).

$$\Delta_{Z_t} \equiv \sum^m \frac{\partial y_t^m}{\partial \mathbf{Z}} = \sum^m \frac{\partial \pi_{g_t^0 \longleftrightarrow y_t^m}}{\partial \mathbf{Z}} y_t^m \quad (50)$$

Similarly, we can isolate the marginal effects of the variables in the expert regressions on the outcome, which would take a slight modification of equation (39).

$$\Delta_{X_t} \equiv \sum^m \frac{\partial y_t^m}{\partial \mathbf{X}} = \sum^m \pi_{g_t^0 \longleftrightarrow y_t^m} \frac{\partial y_t^m}{\partial \mathbf{X}} \quad (51)$$

Substituting Δ_{Z_t} or Δ_{X_t} into (49) will yield analagous estimates of the precision of (50) and (51).

6 A simple example

In order to provide a concrete example of the concepts discussed previously, the ME and HME models are demonstrated on a small and well known dataset collected by Edgar Anderson (Anderson 1936) and popularized in the statistics literature by Ronald Fisher (Fisher 1936). Anderson collected 50 measurements each from three different species of iris flowers; the width and length of both the petal and the sepal. Figure 2 provides a basic view of the species specific clustering inherent in the data.

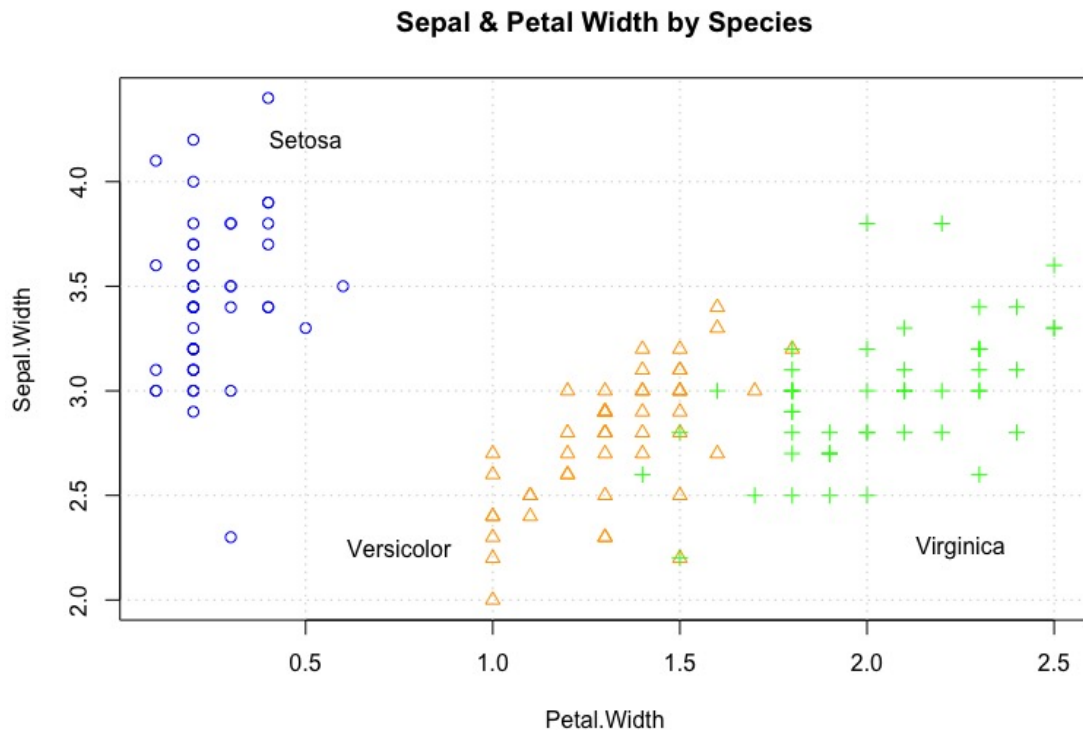


Figure 2: Three different iris species: Setosa (blue circles), Versicolor (orange triangles), Virginia (green crosses). Sepal width is on the vertical axis and petal width on the horizontal axis.

The work below uses the ME and HME architectures to estimate a flower’s sepal width using only its petal width as a predictor. The petal width will be used as the sole covariate in the local linear expert regressions (\mathbf{X}) as well as in the gating network (\mathbf{Z}).

$$sepal.width_i = \beta_0 + \beta_1 * petal.width_i + \varepsilon_i \mid \omega_0 + \omega_1 * petal.width_i \quad (52)$$

The goal is to have the gating network of the models identify the inherent species-specific clustering without explicit knowledge of each observation’s species classification, and then fit an appropriate local regression to the self-identified clusters. As a benchmark, an OLS model is run where a flower’s petal width is interacted with it’s species, resulting in a species-specific estimation of sepal width.

$$sepal.width_{is} = \beta_{0,s} + \beta_{1,s} * petal.width_{is} + \varepsilon_{is} \quad (53)$$

Two sets of regressions are run. Since the Versicolor and Virginica species can be viewed as one larger cluster, a two-expert ME model is run and compared to a benchmark OLS where Versicolor and Virginica are labelled as the same species. A second set of regressions are run with three mixture experts. When moving to the three expert model, there is now a choice on what kind of gating architecture to employ. We can go deep by adding a gating network with depth two (HME), or we can go wide by keeping the depth of the gating network at one (ME). Again, for comparative purposes, a benchmark OLS regression is estimated for each species separately. Results are collected in table 2. Coefficients for local experts in the two expert ME regression match closely with the OLS benchmark. The strong separation between the Setosa and Versicolor/Virginica clusters makes it easy for the ME gating network to discriminate between the two using just the Petal Width dimension. This task becomes a little more complicated when considering all three species at the same time since there exists some overlap between the Versicolor and Virginica clusters. When comparing the coefficients of the local regressions (see table 2), the HME architecture clearly outperforms the ME architecture. While the ME model does obtain a larger log-likelihood value than the OLS estimate, it fails to identify the three separate species that are known to exist. The HME model, on the other hand, naturally picks up on the three underlying clusters while also providing a superior likelihood value. This speaks to one of the major caveats of using this class of model. The likelihood value of an ME or HME can always be improved by adding more and more experts, but this improvement should not be confused with the model gaining a finer understanding of the underlying data generating process. It simply starts to over-fit to the data at hand.

Table 2: Iris Dataset - OLS vs ME vs HME

	2 Expert Mixture				3 Expert Mixture					
	OLS		ME		OLS		HME		ME	
	Coef.	SE	Coef.	SE	Coef.	SE	Coef.	SE	Coef.	SE
<hr/> Setosa <hr/>										
Const.	3.22	0.11**	3.22	0.13**	3.22	0.11**	3.22	0.13**	3.45	0.13**
Petal.Width	0.84	0.42*	0.95	0.49**	0.84	0.41*	0.94	0.49	0.39	0.46
<hr/> Virginica <hr/>										
Const.	—	—	—	—	1.70	0.32**	1.96	0.12**	3.02	0.05**
Petal.Width	—	—	—	—	0.63	0.16**	0.50	0.06**	0.21	0.31
<hr/> Versicolor <hr/>										
Const.	—	—	—	—	1.37	0.29**	1.15	0.12**	2.13	0.09**
Petal.Width	—	—	—	—	1.05	0.22**	1.29	0.09**	0.44	0.06**
<hr/> Virg + Versi <hr/>										
Const.	2.13	0.13**	2.13	0.09**	—	—	—	—	—	—
Petal.Width	0.44	0.07**	0.44	0.06**	—	—	—	—	—	—
<hr/> AME <hr/>										
Petal.Width	0.57	—	0.49	—	0.84	—	0.57	—	0.62	—
Log-Like	-35.5	—	-31.9	—	-29.3	—	-21.8	—	-27.8	—
N	150	—	150	—	150	—	150	—	150	—

** $p < 0.01$, * $p < 0.05$

OLS regressions are modeled using equation (53)

ME regressions are modeled using equation (52) and architecture **A** from figure 1

HME regressions are modeled using equation (52) and architecture **C** from figure 1

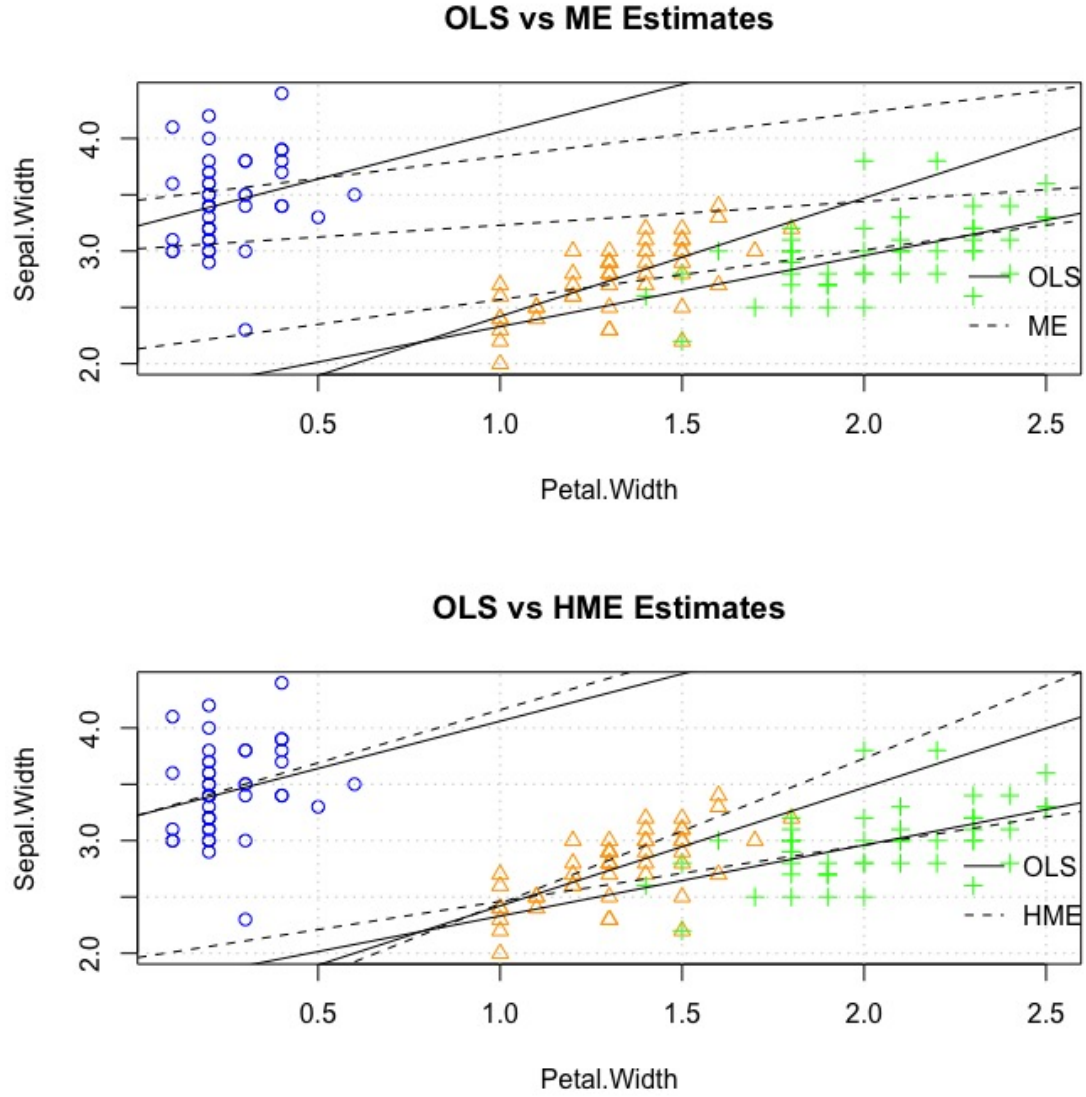


Figure 3: Comparison of the fitted experts between the ME and HME architectures applied to the Iris dataset. OLS regression estimates are drawn in solid lines. Although the HME and ME both achieve superior log-likelihood values compared to OLS, only the HME is able to identify the three iris species clusters.

7 A Mincer Wage Equation

For a more economically relevant example, we turn our attention to a common topic in labor economics: the income return on an additional year of education. At times called the "Mincer wage equation", our version of it will be:

$$\log(wage) = \beta_0 + \beta_1 * Age + \beta_2 * Age^2 + \beta_3 * YrsEdu + \beta_4 \mathbf{X} + \varepsilon \quad (54)$$

with \mathbf{X} containing a set of individual-specific variables as well as a set of occupation-specific attributes. Our data will come from two sources. First, from the 2000 Census, we devise a measure of the hourly (log) wage. In addition to income, we also collect information on age, years of education (YrsEdu), job occupations codes, and a set of demographic identifiers indicating the race of the individuals contained in the census sample. For the occupational codes, we use the Standard Occupation Classification (SOC) codes from the Occupation Information Network (ONet). Each occupation is associated with a set of knowledge and skill-based attributes describing what qualities are necessary to perform each job suitably. A federally sponsored source, ONet details, "the knowledge, skills, and abilities required as well as how the work is performed in terms of tasks, work activities, and other descriptors" (*Occupational Information Network (O*NET)* 2019).

To link the occupational codes in the census data to the SOC codes used by ONet, we use the cross walk provided by Sarah Porter (Porter 2019). This mapping is not one-to-one. When more than one SOC code points to a single census code, we take the average of the SOC codes. After a quick but careful scan of the job attributes available on ONet, the following four were selected. The footnootes provide the full classification hierarchy listed on the website.

1. Social Perceptiveness ¹
2. Design ²
3. Data Analytics ³
4. Creative Thinking ⁴

The guiding principle for attribute selection was to choose a small but diverse set of attributes that contrast well, with each attribute embodying a human skill valued

¹Skills - Social Skills - Social Perceptiveness

²Work Activities - Mental Processes - Analyzing Data or Information

³Knowledge - Design

⁴Work Activities - Mental Processes - Thinking Creatively

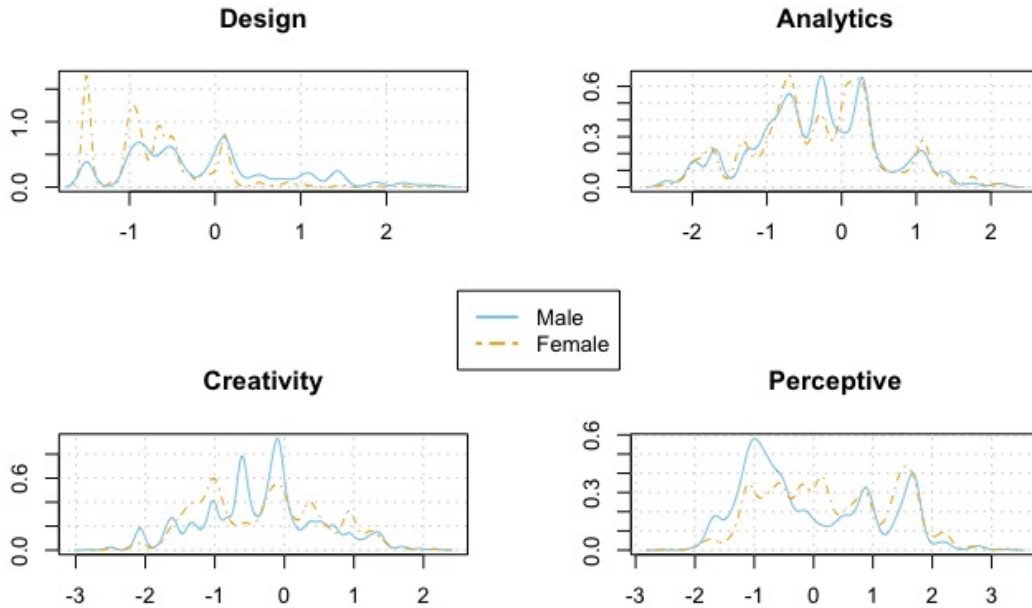


Figure 4: Density estimates of ONet job characteristics broken down by sex. The job characteristics have been mean centered and scaled to have unit variance.

across industry, culture, and society. For these selected attributes, ONet grades their relevance on a 100 point scale. Each attribute contains two scales, an "importance" scale and a "level" scale. The importance scale denotes how critical the attribute is to the occupation while the level indicates how much the skill is required or needed to perform the occupation. To unify the two measures, we follow the paper Prof Wijverberg gave me and take a cobb-douglass style average with a $2/3$'s weight for importance and a $1/3$ weight for the level scale.

The total number of individuals in the Census data numbers 105,796. After applying our crosswalk, only 75,957 cases remain with complete information across both datasets. Of those 75,957, roughly ten percent (7,315) are randomly held-out and used as a test set to gauge out-of-sample forecast performance across model specifications. This leaves 68,642 individuals left as a training set. A statistical summary of the covariates is provided in table 3.

A natural question to consider as a researcher is where to put the variable(s) of interest while performing an HME estimation. Jiang and M. A. Tanner 2000

Table 3: Summary Statistics

	25%	Mean	50%	75%
Wage (hr)	9.20	15.82	13.32	19.44
Yrs Edu	12.00	13.78	14.00	16.00
Age	30.00	39.15	39.00	48.00
Age16	14.00	23.15	23.00	32.00
Female	—	40.47	—	—
Af Amer	—	8.62	—	—
Indian	—	1.05	—	—
White	—	77.00	—	—
Hispanic	—	10.00	—	—
Asian	—	3.36	—	—
Creative	46.65	53.30	53.82	58.94
Design	15.00	30.58	26.33	38.97
Analytic	44.01	52.63	52.68	62.18
Perceptive	41.15	50.49	46.03	59.84

N = 68,642

provide their proof of model consistency for HME of GLM’s for the case where all covariates appear in the gating network as well as the experts. We will call this the *full* specifications:

$$\log(wage) = Age + YrsEdu + Sex + Race + Occ \mid Age + YrsEdu + Sex + Race + Occ \quad (55)$$

We will compare this *full* specification to two others. A *mid* specification where the local experts contain age and years of education while removing demographic indicators:

$$\log(wage) = Age + YrsEdu \mid Age + YrsEdu + Sex + Race + Occ \quad (56)$$

And finally a *minimal* specification where our core variable of interest, years of education, appears solely in the gating network.

$$\log(wage) = Age \mid Age + YrsEdu + Sex + Race + Occ \quad (57)$$

For comparative purposes, we estimate several different regressions across three different dimensions: model architectures (ME vs HME), the number of experts,

and the regression specification (equations (55) - (57)). Table 4 presents a view of these results across those dimensions. After looking at the results, two themes emerge. First, there is a clear advantage to using the HME structure if the aim is to maximize the likelihood value. The HME structure shows consistent improvement across specifications as the number of experts increase, while the ME struggles to improve the likelihood value if there is only one gating split. This increase in efficiency is most likely due to the HME’s more refined gating architecture, whose recursive partitioning is more effective at finding the next improvement in the parameter vector than the single multinomial split in the ME. As for the second theme, it’s best to give the expert regressions as much information as possible. The Full specification clearly outperforms the Mid specification, which outperforms the Min specification. Referencing table 4, if one holds the architecture and the number of experts constant, the performance metrics show clear improvement as the regression specification adds more explanatory variables.

Turning attention to the main variable of focus, table 5 provides a comparison of the average marginal effect for *YrsEdu* across the same dimensions explored for the performance metrics. There is a noticeable change across model specifications. Compared to the OLS coefficient of 0.76, the Min specification, which includes *YrsEdu* only in the gating network, underestimates the returns to education. The Mid specification, which includes *Age* and *YrsEdu* in the expert regressions as well as the gating network, overestimates the returns to education in all the models except the HME with four and five experts. The Full specification, which has the entire suite of variables in both places, matches most closely to the OLS estimate across the estimated models.

For our census sample, estimating up to five experts is pretty extreme. It’s rather unlikely that there exists more than one distinct cluster, let alone two⁵. Because of this, a deeper analysis of the regression results are only explored for the three models that have the least complexity/experts. We first estimate equation (54) for a two expert model. At this specification, there is no distinction between the HME and ME. A three expert model is then estimated for these two respective architectures to asses if different conclusions to the estimated Mincer equations arise. Results for these regressions are collected in tables 6, 8, and 10 and complimented by graphs 7, 9, and 11, which provide mean and median values for the subset of individuals in the census sample that are attributed to each expert based on the value of their posterior weights⁶.

⁵Testing if a (H)ME model is even necessary would be a valuable addition to this paper

⁶For example, observation i is assigned to expert j if the posterior vector’s largest value is the j -th index: $\arg \max \mathbf{h}_i = h_{ij}$

Broadly speaking, all three models explored share the same macro view of the data. On the right side of tables 6, 8, and 10 are a group of columns titled '(H)ME Marginal Effects'. Here the marginal effects of the model can be broken down and attributed to the gating network or the expert regressions. "Both", "Experts", and "Gates" refers to marginal effects referenced by equations (40), (39), and (38), respectively. The values are fairly consistent across variables and model architectures with the coefficients for *Age* and it's square a modest exception, ranging from 0.028 (HME) to 0.042 (2-Expert ME) for *Age*. Notice also that the marginal effects from the expert regressions are the lion's share of total marginal effect, ranging from one to two orders of magnitude larger than marginal effects for the gating network.

When left to segment the data set on it's own, the two expert ME model estimates two different wage equations, one for the majority of the population that tends to be older, whiter, and more educated (see table 7), and a second smaller popular that is more diverse, significantly young, with less education on average. The difference between the average age of the two populations is noticeable and might play a role behind the marginal effects for *Age* moving around as much as it does. This younger cohort breaking off from the bulk of the sample repeats for both 3-expert models as well. Interestingly, the 3-expert models then share a further partition of the sample around age, with a third older cohort separating itself from the smaller sample. When taken together with the regression diagnostics, the model suggests that returns to education evolve over an individual's lifetime. When young, the returns to education are at their smallest and then expand during main earning years of middle age. Returns then dip slightly as individuals come closer to retirement age.

When looking at the occupational attributes there is similar agreement between the estimated models. The marginal effects for all three are in close proximity between the ME and HME models. Those individuals who specialize in performing analytics enjoy the greatest hourly rate (0.126 - 0.128). Design (0.074 to 0.081) and Perceptive (0.053 to 0.057) attributes get a smaller bump to the their hourly wage while Creative types (-0.044 to -0.043) clearly have alternative motivation than monetary gain.

8 Conclusion

In this article, a novel mixture model is explored that borrows equally from the economic and deep learning fields. A flexible (and optionally deep) gating network is used to learn the latent structure of a dataset and then apply local regressions to that latent structure. Robust standard errors and closed form expressions for marginal effects were developed and demonstrated on two different datasets.

Table 4: Comparing Complexity, Architecture, and Regression Specification

Specification	Architecture	Experts	Performance Metrics			
			Log-Lik	AIC	BIC	MSE
Full	ME	2	-0.541	1.082	1.088	0.182
	ME	3	-0.526	1.053	1.062	0.182
	ME	4	-0.537	1.078	1.091	0.181
	ME	5	-0.535	1.073	1.089	0.182
	HME	3	-0.525	1.052	1.061	0.182
	HME	4	-0.515	1.034	1.047	0.181
	HME	5	<u>-0.505</u>	<u>1.015</u>	<u>1.031</u>	<u>0.178</u>
Mid	ME	2	-0.560	1.120	1.123	0.185
	ME	3	-0.558	1.117	1.123	0.186
	ME	4	-0.581	1.163	1.171	0.192
	ME	5	-0.590	1.182	1.192	0.199
	HME	3	-0.541	1.083	1.088	0.184
	HME	4	-0.528	1.057	1.065	0.183
	HME	5	<i>-0.519</i>	<i>1.039</i>	<i>1.050</i>	<i>0.182</i>
Min	ME	2	-0.596	1.192	1.195	0.192
	ME	3	-0.587	1.176	1.181	0.192
	ME	4	-0.629	1.260	1.268	0.211
	ME	5	-0.564	1.131	1.140	0.189
	HME	3	-0.581	1.163	1.168	0.190
	HME	4	-0.546	1.094	1.101	0.182
	HME	5	<i>-0.524</i>	<i>1.049</i>	<i>1.059</i>	<i>0.182</i>

Note: Log-Likelihood, AIC, and BIC are divided by the sample size: 68,642. Italicized entries are the winning values within specification while underlined entries are the best values across all three specifications.

Note: The MSE is calculated from a hold-out test set with sample size: 7,315

Note: After looking at the results, two themes emerge. **One**, there is a clear advantage to using the HME structure if the aim is to maximize the likelihood value. The HME structure shows consistent improvement across specifications as the number of experts increases, while the ME struggles to match this consistency. **Two**, give the expert regressions as much information as possible. The Full specification clearly outperforms the Mid and Min specifications across the board.

Table 5: Returns to Years of Education

Depth	Experts	Avg. Marginal Effect		
		Min	Mid	Full
ME	2	0.051	0.082	0.076
ME	3	0.051	0.081	0.074
ME	4	0.039	0.085	0.075
ME	5	0.063	0.095	0.076
HME	3	0.063	0.080	0.073
HME	4	0.063	0.078	0.073
HME	5	0.068	0.075	0.069

Note: OLS coef: 0.076

Note: There is a noticeable change across in the marginal return to an extra year of education. Compared to the OLS coefficient of 0.76, the Min specification, which includes *YrsEdu* only in the gating network, underestimates the returns to education. The Mid specification, which includes *Age* and *YrsEdu* in the expert regressions as well as the gating network, overestimates the returns to education in all the models except the HME with four and five experts. The Full specification, which has the entire suite of variables in both places, matches most closely to the OLS estimate across the estimated models.

Table 6: Regression Results: Two-Expert, Full Parameter Specification

	ME Regressions ¹		OLS ²		ME Marginal Effects ³		
	Coef.	Coef.	Coef.		Both	Experts	Gates
Intercept	1.231 *	1.494 *	1.241 *		1.225 *	1.260 *	-0.040
Age16	0.032 *	0.068 *	0.035 *		0.042	0.038 *	0.004
Age16sq	-0.000 *	-0.002 *	-0.001 *		-0.001	-0.001 *	-0.000
YrsEduc	0.082 *	0.036 *	0.076 *		0.076 *	0.075 *	0.000
Female	-0.244 *	-0.032 *	-0.215 *		-0.209 *	-0.207 *	-0.002
Af Amer	-0.076 *	-0.045 *	-0.076 *		-0.076 *	-0.071 *	-0.005
Indian	-0.081 *	1.390 *	-0.091 *		-0.085 +	-0.079 *	-0.005
Asian	-0.045 *	0.036 *	-0.032 *		-0.024	-0.028 *	0.003
Hisp	-0.121 *	-0.082 *	-0.106 *		-0.112 *	-0.112 *	-0.000
Creativity	-0.054 *	-0.008 *	-0.046 *		-0.044 *	-0.045 *	0.002
Design	0.080 *	0.078 *	0.082 *		0.081 *	0.080 *	0.001
Analytics	0.133 *	0.112 *	0.131 *		0.126 *	0.129 *	-0.003
Perceptive	0.063 *	-0.013 *	0.058 *		0.053 *	0.049 *	0.004
Log-Variance	-1.651 *	-2.682 *	—		—	—	—
Share ⁴ :	0.826	0.174	1.000		—	—	—

Signif. Codes: 0.01 '**', 0.05 '+', 0.1 '-'

Log-Likelihood: ME -0.541, OLS -0.558

¹ Fitted coefficients from the two-expert model with the full parameter specification from equation (55)

² Fitted coefficients from an OLS regression. These coefficient values can be compared to the HME coefficients to their left as well as to the marginal values to their right

³ Marginal effects for the HME model. Both, Experts, and Gates refers to marginal effects referenced by equations (49), (51), and (50), respectively.

⁴ The share is calculated by summing the posterior weights across observations for each expert.

Table 7: Sample Mean Comparison: Two-Expert ME

Share: ¹	(0.826)		(0.174)	
	Mean	Median	Mean	Median
Wage (hr)	2.679	2.681	2.175	2.197
Age	25.814	25.000	6.965	7.000
Age16	759.812	625.000	62.478	49.000
Female	0.408	0.000	0.386	0.000
Af Amer	0.084	0.000	0.101	0.000
Indian	0.009	0.000	0.018	0.000
White	0.778	1.000	0.698	1.000
Hispanic	0.037	0.000	0.028	0.000
Asian	0.091	0.000	0.155	0.000
YrsEduc	13.916	14.000	12.974	12.000
Creative	-0.191	-0.137	-0.464	-0.542
Design	-0.344	-0.535	-0.442	-0.635
Analytic	-0.196	-0.247	-0.499	-0.550
Perceptive	0.230	0.127	-0.233	-0.532
N	—	58,939	—	9,703

¹ The share is calculated by summing the posterior weights across observations for each expert.

Note: Mean and median values are applied to individuals in the census sample that are classified based on the value of their posterior weights. For example, observation i is assigned to expert j if the posterior vector's largest value is the j -th index: $\arg \max \mathbf{h}_i = h_{ij}$

Table 8: Regression Results: Wide Three-Expert, Full Parameter Specification

	ME Regressions ¹						OLS ²		ME Marginal Effects ³				
	Coef.		Coef.		Coef.		Coef.		Both		Experts		Gates
Intercept	1.379	*	1.574	*	0.562	*	1.241	*	1.367	1.340	*	0.032	
Age16	0.021	*	0.045	*	0.060	*	0.035	*	0.029	0.027	*	0.002	
Age16sq	-0.000	*	-0.001	*	-0.001	*	-0.001	*	-0.000	-0.000	*	0.000	
YrsEduc	0.082	*	0.032	*	0.080	*	0.076	*	0.074	0.077	*	-0.002	
Female	-0.251	*	-0.022	*	-0.149	*	-0.215	*	-0.206	-0.218	*	0.012	
Af Amer	-0.084	*	-0.056	*	-0.054	-	-0.076	*	-0.076	-0.078	*	0.002	
Indian	-0.105	*	-0.046	*	0.010		-0.091	*	-0.091	-0.090	*	-0.002	
Asian	-0.030	*	0.057	*	-0.091	*	-0.032	*	-0.024	-0.025	*	0.001	
Hisp	-0.136	*	-0.061	*	0.071	+	-0.106	*	-0.107	-0.111	*	0.004	
Creativity	-0.038	*	-0.022	*	-0.177	*	-0.046	*	-0.044	-0.047	*	0.003	
Design	0.080	*	0.080	*	-0.037	*	0.082	*	0.075	0.071	*	0.004	
Analytics	0.123	*	0.110	*	0.196	*	0.131	*	0.128	0.128	*	0.000	
Perceptive	0.060	*	-0.008	*	0.168	*	0.058	*	0.057	0.061	*	-0.004	
Log-Variance	-1.893	*	-2.891	*	-0.627	*	—						
Share ⁴ :	0.809		0.111		0.080		1.000		—	—		—	

Signif. Codes: 0.01 '**', 0.05 '+', 0.1 '-'

Log-Likelihood: ME -0.526, OLS -0.558

¹ Fitted coefficients from the three-expert model with the full parameter specification from equation (55)

² Fitted coefficients from an OLS regression. These coefficient values can be compared to the HME coefficients to their left as well as to the marginal values to their right

³ Marginal effects for the HME model. Both, Experts, and Gates refers to marginal effects referenced by equations (49), (51), and (50), respectively.

⁴ The share is calculated by summing the posterior weights across observations for each expert.

Table 9: Sample Mean Comparison: Wide Three-Expert HME

Share: ¹	(0.809)		(0.111)		(0.080)	
	Mean	Median	Mean	Median	Mean	Median
Wage (hr)	2.664	2.667	2.106	2.096	2.549	2.221
Age	24.916	24.000	5.830	6.000	27.827	28.000
Age16	722.355	576.000	41.789	36.000	904.103	784.000
Female	0.420	0.000	0.301	0.000	0.250	0.000
Af Amer	0.090	0.000	0.060	0.000	0.064	0.000
Indian	0.010	0.000	0.012	0.000	0.010	0.000
Hispanic	0.036	0.000	0.020	0.000	0.102	0.000
Asian	0.100	0.000	0.114	0.000	0.045	0.000
YrsEduc	13.802	14.000	13.101	12.000	15.837	16.000
Creative	-0.209	-0.141	-0.422	-0.456	-0.195	-0.282
Design	-0.344	-0.535	-0.387	-0.535	-0.765	-0.860
Analytic	-0.218	-0.264	-0.472	-0.412	-0.072	0.049
Perceptive	0.177	0.127	-0.122	-0.455	0.851	0.877
N	–	60,396	–	6,603	–	1,643

¹ The share is calculated by summing the posterior weights across observations for each expert.

Note: Mean and median values are applied to individuals in the census sample that are classified based on the value of their posterior weights. For example, observation i is assigned to expert j if the posterior vector's largest value is the j -th index: $\arg \max \mathbf{h}_i = h_{ij}$

Table 10: Regression Results: Deep Three-Expert, Full Parameter Specification

	HME Regressions ¹						OLS ²		HME Marginal Effects ³					
	Coef.		Coef.		Coef.		Coef.		Both		Experts		Gates	
Intercept	1.404	*	1.559	*	0.898	*	1.241	*	1.393	1.382	*	0.011		
Age16	0.020	*	0.050	*	0.044	*	0.035	*	0.028	0.026	*	0.003		
Age16sq	-0.000	*	-0.001	*	-0.001	*	-0.001	*	-0.000	-0.000	*	0.000		
YrsEduc	0.082	*	0.034	*	0.074	*	0.076	*	0.073	0.075	*	-0.001		
Female	-0.257	*	-0.034	*	-0.131	*	-0.215	*	-0.209	-0.217	*	0.008		
Af Amer	-0.086	*	-0.048	*	-0.041		-0.076	*	-0.076	-0.077	*	0.001		
Indian	-0.113	*	-0.057	*	0.043		-0.091	*	-0.100	-0.093	*	-0.007		
Asian	-0.033	*	0.058	*	-0.062	+	-0.032	*	-0.025	-0.023	*	-0.001		
Hisp	-0.143	*	-0.066	*	0.077	*	-0.106	*	-0.111	-0.114	*	0.003		
Creativity	-0.042	*	-0.021	*	-0.136	*	-0.046	*	-0.043	-0.047	*	0.004		
Design	0.080	*	0.068	*	-0.048	*	0.082	*	0.074	0.068	*	0.006		
Analytics	0.124	*	0.112	*	0.183	*	0.131	*	0.128	0.127	*	0.000		
Perceptive	0.063	*	-0.003		0.135	*	0.058	*	0.056	0.061	*	-0.004		
Log-Variance	-1.895	*	-2.791	*	-0.622	*	—							
Share ⁴ :	0.783		0.133		0.084		1.000		—	—		—		

Signif. Codes: 0.01 '**', 0.05 '+', 0.1 '.'

Log-Likelihood: HME -0.525, OLS -0.558

¹ Fitted coefficients from the three-expert model with the full parameter specification from equation (55)

² Fitted coefficients from an OLS regression. These coefficient values can be compared to the HME coefficients to their left as well as to the marginal values to their right

³ Marginal effects for the HME model. Both, Experts, and Gates refers to marginal effects referenced by equations (49), (51), and (50), respectively.

⁴ The share is calculated by summing the posterior weights across observations for each expert.

References

- Anderson, Edgar (1936). “The species problem in iris”. In: *Annals of the Missouri Botanical Gardens* 23.3, pp. 457–509.
- Bishop, Christopher and Markus Svenson (2003). “Bayesian Hierarchical Mixtures of Experts”. In: *Proceedings Nineteenth Conference on Uncertainty in Artificial Intelligence*. Morgan Kaufmann, pp. 57–64.
- Blei, David M., Alp Kucukelbir, and McAuliffe (2016). “Variational Inference: A Review for Statisticians”. In: *ArXiv e-prints*. eprint: 1601.00670.
- Brieman, L. et al. (1984). *Classification and Regression Trees*. Wadsworth and Brooks/Cole.
- Carvalho, Alexandre and Georgios Skoulakis (2005). “Ergodicity and existence of moments for local mixtures of linear autoregressions”. In: *Statistics and Probability Letters* 71.3, pp. 313–322.
- (2010). “Time Series Mixtures of Generalized t Experts: ML Estimation and an Application to stock return density forecasting”. In: *Econometric Reviews* 29.5-6, pp. 642–687. DOI: 10.1080/07474938.2010.481987.
- Carvalho, Alexandre and Martin Tanner (2003). “Hypothesis testing in mixture-of-experts of generalized linear time series”. In: *2003 IEEE International Conference on Computational Intelligence for Financial Engineering, 2003. Proceedings*. Pp. 285–292.
- (2005). “Mixtures-of-experts of autoregressive time series: asymptotic normality and model specification”. In: *IEEE Transactions on Neural Networks* 16.1, pp. 39–56. ISSN: 1045-9227.
- (2006). “Modeling nonlinearities with mixtures-of-experts of time series models”. In: *International Journal of Mathematics and Mathematical Sciences* 2006.
- Dempster, A. P., N. M. Laird, and D. B. Rubin (1977). “Maximum likelihood from incomplete data via the em algorithm”. In: *Journal of the Royal Statistical Society. Series B*. 39.1, pp. 1–38.
- Fisher, R.A. (1936). “The use of multiple measurements in taxonomic problems”. In: *Annals of Eugenics* 7.2, pp. 179–188.
- Fritsch, Jürgen, Michael Finke, and Alex Waibel (1997). “Adaptively Growing Hierarchical Mixtures of Experts”. In: *Advances in Neural Information Processing Systems 9*. Ed. by M. C. Mozer, M. I. Jordan, and T. Petsche. MIT Press, pp. 459–465. URL: <http://papers.nips.cc/paper/1279-adaptively-growing-hierarchical-mixtures-of-experts.pdf>.
- Goldfeld, Stephan M. and Richard E. Quandt (1973). “A Markov Model for Regime Switching”. In: *Journal of Econometrics* 1 (1), pp. 3–16.

Table 11: Sample Mean Comparison: Deep Three-Expert HME

Share: ¹	(0.783)		(0.133)		(0.084)	
	Mean	Median	Mean	Median	Mean	Median
Wage (hr)	2.683	2.676	2.137	2.140	2.432	2.075
Age	25.523	25.000	6.494	7.000	26.913	27.000
Age16	746.250	625.000	50.858	49.000	873.924	729.000
Female	0.414	0.000	0.358	0.000	0.313	0.000
Af Amer	0.088	0.000	0.073	0.000	0.075	0.000
Indian	0.010	0.000	0.016	0.000	0.018	0.000
White	0.770	1.000	0.753	1.000	0.749	1.000
Hispanic	0.036	0.000	0.027	0.000	0.101	0.000
Asian	0.096	0.000	0.131	0.000	0.057	0.000
YrsEduc	13.846	14.000	13.077	12.000	15.378	16.000
Creative	-0.198	-0.137	-0.444	-0.508	-0.201	-0.282
Design	-0.330	-0.530	-0.477	-0.635	-0.757	-0.859
Analytic	-0.206	-0.253	-0.471	-0.412	-0.161	-0.007
Perceptive	0.185	0.127	-0.082	-0.308	0.756	0.877
N	–	58,429	–	8,674	–	1,539

¹ The share is calculated by summing the posterior weights across observations for each expert.

Note: Mean and median values are applied to individuals in the census sample that are classified based on the value of their posterior weights. For example, observation i is assigned to expert j if the posterior vector's largest value is the j -th index: $\arg \max \mathbf{h}_i = h_{ij}$

- Hamilton, J.D. (1989). “A new approach to the economic analysis of nonstationary time series and the business cycle”. In: *Econometrica* 57, pp. 357–384.
- Huerta, Gabriel, Wenxin Jiang, and Martin A. Tanner (2003). “Time series modeling via hierarchical mixtures”. In: *Statistica Sinica* 13.
- Jacobs, R. A. et al. (1991). “Adaptive mixture of local experts”. In: *Neural Computation* 3, pp. 79–82.
- Jiang, Wenxin and Martin A. Tanner (1999). “Hierarchical Mixture-of-Experts for Exponential Family Regression Models: Approximation and Maximum Likelihood Estimation”. In: *The Annals of Statistics* 27.3, pp. 987–1011.
- (2000). “On the asymptotic normality of hierarchical mixtures-of-experts for generalized linear models”. In: 46.3, pp. 1005–1013. ISSN: 0018-9448.
- Jordan, M. and R. Jacobs (1993). “Hierarchical mixtures of experts and the EM algorithm”. In: *Proceedings of 1993 International Joint Conference on Neural Networks*.
- Jordan, M. and L. Xu (1995). “Convergence results for the em approach to mixtures-of-experts architectures”. In: *Neural Networks* 8 (9), pp. 1409–1431.
- Jordan, Michael I. and Robert A. Jacobs (1992). “Hierarchies of adaptive experts”. In: *Advances in Neural Information Processing Systems 4*. Ed. by J. E. Moody, S. J. Hanson, and R. P. Lippmann. Morgan-Kaufmann, pp. 985–992. URL: <http://papers.nips.cc/paper/514-hierarchies-of-adaptive-experts.pdf>.
- Neal, Jeffries and Ruth Pfeiffer (2001). “A mixture model for the probability distribution of rain rate”. In: *Environmetrics* 12.1, pp. 1–10.
- Occupational Information Network (O*NET) (2019). URL: <https://www.doleta.gov/programs/onet/> (visited on 01/28/2019).
- Porter, Sarah (2019). *Census to ONet Mapping*. URL: http://econterms.net/pbmeyer/research/occs/wiki/index.php?title=Crosswalk_by_Sarah_Porter_to_map_1980_codes_forward_in_SAS (visited on 01/28/2019).
- Ueda, N. and Z. Ghahramani (2002). “Bayesian model search for mixture models based on optimizing variational bounds”. In: *Neural Networks* 15.10, pp. 1223–1241.
- Waterhouse, S.R. and A.J. Robinson (1995). “Constructive Algorithms for Hierarchical Mixture of Experts”. In: *Advances in Neural Information Processing Systems* 8.
- Waterhouse, Steve R., David MacKay, and Anthony J. Robinson (1995). “Bayesian Methods for Mixtures of Experts”. In: *NIPS*.
- Weigend, A., M. Mangeas, and A. Srivastava (1995). “Nonlinear gated experts for time series: discovering regimes and avoiding overfitting”. In: *International Journal of Neural Systems* 6, pp. 373–399.