

## Adaptive Mixtures of Local Experts

Robert A. Jacobs

Michael I. Jordan

Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology,  
Cambridge, MA 02139 USA

Steven J. Nowlan

Geoffrey E. Hinton

Department of Computer Science, University of Toronto,  
Toronto, Canada M5S 1A4

We present a new supervised learning procedure for systems composed of many separate networks, each of which learns to handle a subset of the complete set of training cases. The new procedure can be viewed either as a modular version of a multilayer supervised network, or as an associative version of competitive learning. It therefore provides a new link between these two apparently different approaches. We demonstrate that the learning procedure divides up a vowel discrimination task into appropriate subtasks, each of which can be solved by a very simple expert network.

separate networks, each of which learns to handle a subset of the complete set of training cases

### 1 Making Associative Learning Competitive

---

If backpropagation is used to train a single, multilayer network to perform different subtasks on different occasions, there will generally be strong interference effects that lead to slow learning and poor generalization. If we know in advance that a set of training cases may be naturally divided into subsets that correspond to distinct subtasks, interference can be reduced by using a system composed of several different "expert" networks plus a gating network that decides which of the experts should be used for each training case.<sup>1</sup> Hampshire and Waibel (1989) have described a system of this kind that can be used when the division into subtasks is known prior to training, and Jacobs *et al.* (1990) have described a related system that learns how to allocate cases to experts. The idea behind such a system is that the gating network allocates a new case to one or a few experts, and, if the output is incorrect, the weight changes are localized to these experts (and the gating network).

This idea is different from the 'world model' which people are chasing as of now. Not really breakdown and separation for now. However, with the MoE & model merging, we seems to be heading to a new direction here, resurrecting this early idea from 90s.

---

<sup>1</sup>This idea was first presented by Jacobs and Hinton at the Connectionist Summer School in Pittsburgh in 1988.

So there is no interference with the weights of other experts that specialize in quite different cases. The experts are therefore local in the sense that the weights in one expert are decoupled from the weights in other experts. In addition they will often be local in the sense that each expert will be allocated to only a small local region of the space of possible input vectors.

each expert is allocated to only a small local region of the space of possible input vectors, representation vector & k-means clustering

Unfortunately, both Hampshire and Waibel and Jacobs *et al.* use an error function that does not encourage localization. They assume that the final output of the whole system is a linear combination of the outputs of the local experts, with the gating network determining the proportion of each local output in the linear combination. So the final error on case  $c$  is

$$E^c = \|\mathbf{d}^c - \sum_i p_i^c \mathbf{o}_i^c\|^2 \quad (1.1)$$

where  $\mathbf{o}_i^c$  is the output vector of expert  $i$  on case  $c$ ,  $p_i^c$  is the proportional contribution of expert  $i$  to the combined output vector, and  $\mathbf{d}^c$  is the desired output vector in case  $c$ .

so the fundamental issue here, is that firstly during backprop the optimization of the gating function will simply be buried by the optimization signal of the bigger experts parameters. This means the gating will never be really optimized, if it's optimized together with the entire network.

This error measure compares the desired output with a blend of the outputs of the local experts, so, to minimize the error, each local expert must make its output cancel the residual error that is left by the combined effects of all the other experts. When the weights in one expert change, the residual error changes, and so the error derivatives for all the other local experts change.<sup>2</sup> This strong coupling between the experts causes them to cooperate nicely, but tends to lead to solutions in which many experts are used for each case. It is possible to encourage competition by adding penalty terms to the objective function to encourage solutions in which only one expert is active (Jacobs *et al.* 1990), but a simpler remedy is to redefine the error function so that the local experts are encouraged to compete rather than cooperate.

That is the real problem of learnable gating functions: they are too small compared to the expert network and therefore its learned parameters are too sub-optimal compared to the sub experts, and yet the gating function is so important such that without a proper module, the sub-expert can NOT even hope to become a real domain specified expert, since a horrible gating weight will only makes them work more towards predicting the error made by all the other experts, and essentially waste all its energy to learn to cancel the residual error that is left by the combined effects of all the other experts.

Instead of linearly combining the outputs of the separate experts, we imagine that the gating network makes a stochastic decision about which single expert to use on each occasion (see Fig. 1). The error is then the expected value of the squared difference between the desired and actual output vectors

$$E^c = \langle \|\mathbf{d}^c - \mathbf{o}_i^c\|^2 \rangle = \sum_i p_i^c \|\mathbf{d}^c - \mathbf{o}_i^c\|^2 \quad (1.2)$$

Notice that in this new error function, each expert is required to produce the whole of the output vector rather than a residual. As a result, the goal of a local expert on a given training case is not directly affected by the weights within other local experts. There is still some indirect

then the learning of a local expert on a given training case is not directly affected by the weights within other local experts

<sup>2</sup>For Hampshire and Waibel, this problem does not arise because they do not learn the task decomposition. They train each expert separately on its own preassigned subtask.

Weighted sum of all experts create a strong coupling between the experts, and with a tough-to-back-prop gating module, expert learns to cancel the residual error left by the combined effects of all the other experts, instead of focusing on solving the real problem...

Competitive learning encourage competition by adding penalty terms to encourage solutions in which only one expert is active

Weighted sum by nature encourage cooperate while by definition of 'expert', we want them to compete (!!)

When the weights in one expert change, the residual error changes, and so the error derivative for all the other local experts change.

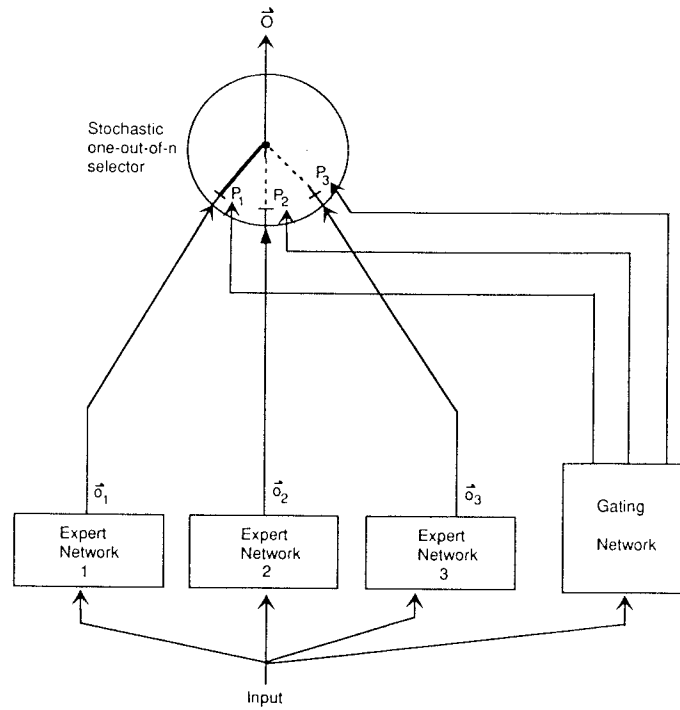


Figure 1: A system of expert and gating networks. Each expert is a feed-forward network and all experts receive the same input and have the same number of outputs. The gating network is also feedforward, and typically receives the same input as the expert networks. It has normalized outputs  $p_j = \exp(x_j) / \sum_i \exp(x_i)$ , where  $x_j$  is the total weighted input received by output unit  $j$  of the gating network. The selector acts like a multiple input, single output stochastic switch; the probability that the switch will select the output from expert  $j$  is  $p_j$ .

coupling because if some other expert changes its weights, it may cause the gating network to alter the responsibilities that get assigned to the experts, but at least these responsibility changes cannot alter the sign of the error that a local expert senses on a given training case. If both the gating network and the local experts are trained by gradient descent in this new error function, the system tends to devote a single expert to each training case. Whenever an expert gives less error than the weighted average of the errors of all the experts (using the outputs of the gating network to decide how to weight each expert's error) its responsibility for that case

*but in a naive weighted average gating, the sign of the error that a local expert senses on a given training case could easily alter, when other expert changes their weights, since a local expert is supposed to predict the residual from the collective decision of all the other experts, which require strong coupling.*

*In this stochastic route-1 scenario, at least the sign of the error will not be altered when other expert changes their weight.*

will be increased, and whenever it does worse than the weighted average its responsibility will be decreased.

The error function in equation 1.2 works in practice but in the simulations reported below we used a different error function which gives better performance:

$$E^c = -\log \sum_i p_i^c e^{-\frac{1}{2} \|d^c - o_i^c\|^2} \quad (1.3)$$

so we have a 'mixture of gaussian' distribution prediction from the mixture of expert model, and loss function defined to be the negative log probability is more effective. and we could do the stochastic routing during inference time, to use only one expert here, or rather the observation in this paper is that stochastic routing works best with mixture of gaussian model formulation, which sounds very intuitive and correct

The error defined in equation 1.3 is simply the negative log probability of generating the desired output vector under the mixture of gaussians model described at the end of the next section. To see why this error function works better, it is helpful to compare the derivatives of the two error functions with respect to the output of an expert. From equation 1.2

So the right view of this is a mixture of gaussian model, and we aim to minimize the negative log probability of generating the desired output vector ... Interesting

we get

$$\frac{\partial E^c}{\partial o_i^c} = -2p_i^c (d^c - o_i^c) \quad (1.4)$$

this looks very similar to the REINFORCE gradient calculation... unification right here... instead of just calculating the gradient according to the predicted dumb gating prob, we should reinforce ourselves with the actual error signal, we should actually check which one actually work, relative to the other ones, according to our little intuition that is, so we effectively reinforce our judgement, according to some sort of competitive schemes.

while from equation 1.3 we get

$$\frac{\partial E^c}{\partial o_i^c} = - \left[ \frac{p_i^c e^{-\frac{1}{2} \|d^c - o_i^c\|^2}}{\sum_j p_j^c e^{-\frac{1}{2} \|d^c - o_j^c\|^2}} \right] (d^c - o_i^c) \quad (1.5)$$

In equation 1.4 the term  $p_i^c$  is used to weight the derivative for expert  $i$ . In equation 1.5 we use a weighting term that takes into account how well expert  $i$  does relative to other experts. This is a more useful measure of the relevance of expert  $i$  to training case  $c$ , especially early in the training. Suppose, for example, that the gating network initially gives equal weights to all experts and  $\|d^c - o_i^c\| > 1$  for all the experts. Equation 1.4 will adapt the best-fitting expert the slowest, whereas equation 1.5 will adapt it the fastest.

whenever you put exponential, you put contrast, relative big difference are now exponentially bigger

## 2 Making Competitive Learning Associative

It is natural to think that the "data" vectors on which a competitive network is trained play a role similar to the *input* vectors of an associative network that maps input vectors to output vectors. This correspondence is assumed in models that use competitive learning as a preprocessing stage within an associative network (Moody and Darken 1989). A quite different view is that the data vectors used in competitive learning correspond to the *output* vectors of an associative network. The competitive network can then be viewed as an inputless stochastic generator of output vectors and competitive learning can be viewed as a procedure for making the network generate output vectors with a distribution that matches the distribution of the "data" vectors. The weight vector of each competitive hidden unit represents the mean of a multidimensional gaussian

distribution, and output vectors are generated by first picking a hidden unit and then picking an output vector from the gaussian distribution determined by the weight vector of the chosen hidden unit. The log probability of generating any particular output vector  $\mathbf{o}'$  is then

$$\log P' = \log \sum_i p_i k e^{-\frac{1}{2} \|\boldsymbol{\mu}_i - \mathbf{o}'\|^2} \quad (2.1)$$

*so essentially the suggesting, obvious suggestion is to treat all your prediction as parameter for a mixture-of-gaussian model here, the sampling should be done in the same formula, as the training process*

where  $i$  is an index over the hidden units,  $\boldsymbol{\mu}_i$  is the “weight” vector of the hidden unit,  $k$  is a normalizing constant, and  $p_i$  is the probability of picking hidden unit  $i$ , so the  $p_i$  are constrained to sum to 1. In the statistics literature (McLachlan and Basford 1988), the  $p_i$  are called “mixing proportions.”

“Soft” competitive learning modifies the weights (and also the variances and the mixing proportions) so as to increase the product of the probabilities (i.e., the likelihood) of generating the output vectors in the training set (Nowlan 1990a). “Hard” competitive learning is a simple approximation to soft competitive learning in which we ignore the possibility that a data vector could be generated by several different hidden units. Instead, we assume that it must be generated by the hidden unit with the closest weight vector, so only this weight vector needs to be modified to increase the probability of generating the data vector.

If we view a competitive network as generating output vectors, it is not immediately obvious what role input vectors could play. However, competitive learning can be generalized in much the same way as Barto (1985) generalized learning automata by adding an input vector and making the actions of the automaton be conditional on the input vector. We replace each hidden unit in a competitive network by an entire expert network whose output vector specifies the mean of a multidimensional gaussian distribution. So the means are now a function of the current input vector and are represented by activity levels rather than weights. In addition, we use a gating network which allows the mixing proportions of the experts to be determined by the input vector. This gives us a system of competing local experts with the error function defined in equation 1.3. We could also introduce a mechanism to allow the input vector to dynamically determine the covariance matrix for the distribution defined by each expert network, but we have not yet experimented with this possibility.

### 3 Application to Multispeaker Vowel Recognition

The mixture of experts model was evaluated on a speaker independent, four-class, vowel discrimination problem (Nowlan 1990b). The data consisted of the first and second formants of the vowels [i], [I], [a], and [A] (usually denoted [Λ]) from 75 speakers (males, females, and children) uttered in a **hVd** context (Peterson and Barney 1952). The data forms two

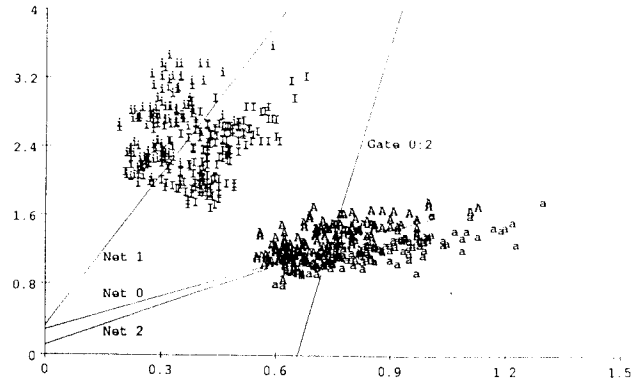


Figure 2: Data for vowel discrimination problem, and expert and gating network decision lines. The horizontal axis is the first formant value, and the vertical axis is the second formant value (the formant values have been linearly scaled by dividing by a factor of 1000). Each example is labeled with its corresponding vowel symbol. Vowels [i] and [I] form one overlapping pair of classes, vowels [a] and [A] form the other pair. The lines labeled Net 0, 1, and 2 represent the decision lines for 3 expert networks. On one side of these lines the output of the corresponding expert is less than 0.5, on the other side the output is greater than 0.5. Although the mixture in this case contained 4 experts, one of these experts made no significant contribution to the final mixture since its mixing proportion  $p_i$  was effectively 0 for all cases. The line labeled Gate 0:2 indicates the decision between expert 0 and expert 2 made by the gating network. To the left of this line  $p_2 > p_0$ , to the right of this line  $p_0 > p_2$ . The boundary between classes [a] and [A] is formed by the combination of the left part of Net 2's decision line and the right part of Net 0's decision line. Although the system tends to use as few experts as it can to solve a problem, it is also sensitive to specific problem features such as the slightly curved boundary between classes [a] and [A].

pairs of overlapping classes, and different experts learn to concentrate on one pair of classes or the other (Fig. 2).

We compared standard backpropagation networks containing a single hidden layer of 6 or 12 units with mixtures of 4 or 8 very simple experts. The architecture of each expert was restricted so it could form only a linear decision surface, which is defined as the set of input vectors for which the expert gives an output of exactly 0.5. All models were trained with data from the first 50 speakers and tested with data from the remaining 25 speakers. The small number of parameters for each expert allows excellent generalization performance (Table 1), and permits

System	Train % correct	Test % correct	Average number of epochs	SD
4 Experts	88	90	1124	23
8 Experts	88	90	1083	12
BP 6 Hid	88	90	2209	83
BP 12 Hid	88	90	2435	124

Table 1: Summary of Performance on Vowel Discrimination Task. Results are based on 25 simulations for each of the alternative models. The first column of the table indicates the system simulated. The second column gives the percent of training cases classified correctly by the final set of weights, while the third column indicates the percent of testing cases classified correctly. The last two columns contain the average number of epochs required to reach the error criterion, and the standard deviation of the distribution of convergence times. Although the squared error was used to decide when to stop training, the criterion for correct performance is based on a weighted average of the outputs of all the experts. Each expert assigns a probability distribution over the classes and these distributions are combined using proportions given by the gating network. The most probable class is then taken to be the response of the system. The identical performance of all the systems is due to the fact that, with this data set, the set of misclassified examples is not sensitive to small changes in the decision surfaces. Also, the test set is easier than the training set.

a graphic representation of the process of task decomposition (Figure 3). The number of hidden units in the backpropagation networks was chosen to give roughly equal numbers of parameters for the backpropagation networks and mixture models. All simulations were performed using a simple gradient descent algorithm with fixed step size  $\epsilon$ . To simplify the comparisons, no momentum or other acceleration techniques were used. The value of  $\epsilon$  for each system was chosen by performing a limited exploration of the convergence from the same initial conditions for a range of  $\epsilon$ . Batch training was used with one weight update for each pass through the training set (*epoch*). Each system was trained until an average squared error of 0.08 over the training set was obtained.

The mixtures of experts reach the error criterion significantly faster than the backpropagation networks ( $p \gg 0.999$ ), requiring only about half as many epochs on average (Table 1). The learning time for the mixture model also scales well as the number of experts is increased: The mixture of 8 experts has a small, but statistically significant ( $p > 0.95$ ), advantage in the average number of epochs required to reach the error criterion. In contrast, the 12 hidden unit backpropagation network requires more epochs ( $p > 0.95$ ) to reach the error criterion than the network with 6



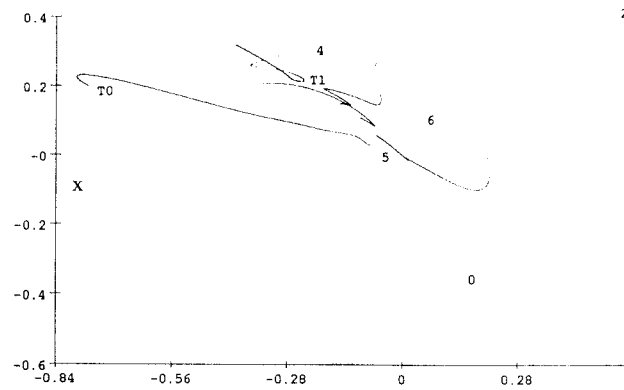


Figure 3: The trajectories of the decision lines of some experts during one simulation. The horizontal axis is the first formant value, and the vertical axis is the second formant value. Each trajectory is represented by a sequence of dots, one per epoch, each dot marking the intersection of the expert's decision line and the normal to that line passing through the origin. For clarity, only 5 of the 8 experts are shown and the number of the expert is shown at the start of the trajectory. The point labeled T0 indicates the optimal decision line for a single expert trained to discriminate [i] from [I]. Similarly, T1 represents the optimal decision line to discriminate [a] from [A]. The point labeled X is the decision line learned by a single expert trained with data from all 4 classes, and represents a type of *average* solution.

hidden units (Table 1). All statistical comparisons are based on a *t* test with 48 degrees of freedom and a pooled variance estimator.

Figure 3 shows how the decision lines of different experts move around as the system learns to allocate pieces of the task to different experts. The system begins in an unbiased state, with the gating network assigning equal mixing proportions to all experts in all cases. As a result, each expert tends to get errors from roughly equal numbers of cases in all 4 classes, and all experts head towards the point X, which represents the optimal decision line for an expert that must deal with all the cases. Once one or more experts begin to receive more error from cases in one class pair than the other, this symmetry is broken and the trajectories begin to diverge as different experts concentrate on one class pair or the other. In this simulation, expert 5 learns to concentrate on discriminating classes [i] and [I] so its decision line approaches the optimal line for this discrimination (T0). Experts 4 and 6 both concentrate on discriminating classes [a] and [A], so their trajectories approach the

*With uniform gating weight initialization (which is the natural thing to do without any prior knowledge), all the experts receive roughly the same error signal from roughly the same set of data (i.i.d. sampled data & error), they therefore learn to update their parameter similarly, so the real question that needs to be addressed with these gating layer, is actually how do you ACCELERATE the asymmetry in their learning trajectory, or how do you speed up the divergence between each module's learning process, which all boils down to trying to create an exponential divergence signal, reacting to the smallest difference between these expert modules and make sure we augment / amplify this small difference and eventually make them into domain experts.*

*At the start, expert work similarity to each other (but actually I feel like their initialized weight would be different so they should have difference in their performance to begin with in fact). But to become a domain expert, we would need to reinforce their comparative divergence, and not their absolute difference — as that could be very small at the beginning, and we will get stuck if we only care about absolute difference therein. This is why I respect the true intelligence of this paper, which just sparkle out the insight and pick the exponential ratio between their relative performance, and use that to perform the gradient update similar to the REINFORCE scheme, which is adopted in Quiet-STAR. And observe that in that way, the entire mixture of expert system is actually just a mixture of Gaussian system ....*

*I still think the issue of gradient mute exists no matter what gating & loss you would ends up picking, somehow it's a natural feeling when you have 1 B expert param, and 1 k gating param, and some offline, heuristic / discriminative gating can be adopted,*



optimal single line (T1) and then split to form a piecewise linear approximation to the slightly curved optimal decision surface (see Fig. 2). Only experts 4, 5, and 6 are active in the final mixture. This solution is typical — in all simulations with mixtures of 4 or 8 experts all but 2 or 3 experts had mixing proportions that were effectively 0 for all cases.

### Acknowledgments

---

Jordan and Jacobs were funded by grants from Siemens and the McDonnell-Pew program in Cognitive Neuroscience. Hinton and Nowlan were funded by grants from the Ontario Information Technology Research Center and the Canadian Natural Science and Engineering Research Council. Hinton is a fellow of the Canadian Institute for Advanced Research.

### References

---

- Barto, A. G. 1985. Learning by statistical cooperation of self-interested neuron-like computing elements. *Human Neurobiol.* **4**, 229–256.
- Hampshire, J., and Waibel, A. 1989. *The meta-pi network: Building distributed knowledge representations for robust pattern recognition*. Tech. Rep. CMU-CS-89-166, Carnegie Mellon University, Pittsburgh, PA.
- Jacobs, R. A., and Jordan, M. I. 1991. Learning piecewise control strategies in a modular connectionist architecture, in preparation.
- Jacobs, R. A., Jordan, M. I., and Barto, A. G. 1991. Task decomposition through competition in a modular connectionist architecture: The what and where vision tasks. *Cog. Sci.*, in press.
- McLachlan, G. J., and Basford, K. E. 1988. *Mixture Models: Inference and Applications to Clustering*. Marcel Dekker, New York.
- Moody, J., and Darken, C. 1989. Fast learning in networks of locally-tuned processing units. *Neural Comp.* **1**(2), 281–294.
- Nowlan, S. J. 1990a. Maximum likelihood competitive learning. In *Advances in Neural Information Processing Systems 2*, D. S. Touretzky, ed., pp. 574–582. Morgan Kaufmann, San Mateo, CA.
- Nowlan, S. J. 1990b. *Competing experts: An experimental investigation of associative mixture models*. Tech. Rep. CRG-TR-90-5, University of Toronto, Toronto, Canada.
- Peterson, G. E., and Barney, H. L. 1952. Control methods used in a study of the vowels. *J. Acoust. Soc. Am.* **24**, 175–184.

---

Received 27 July 1990; accepted 1 November 90.