

# Assigning Features using Additive Clustering

**Wheeler Ruml**

Division of Engineering and Applied Sciences  
Harvard University  
33 Oxford Street  
Cambridge, MA 02138  
ruml@eecs.harvard.edu

April 25, 2001\*

## Abstract

If the promise of computational modeling is to be fully realized in higher-level cognitive domains such as language processing, principled methods must be developed to construct the semantic representations that serve as these models' input and/or output. In this paper, we propose the use of an established formalism from mathematical psychology, *additive clustering*, as a means of automatically assigning discrete features to objects using only pairwise similarity data. Similar approaches have not been widely adopted in the past, as existing methods for the unsupervised learning of such models do not scale well to large problems. We propose a new algorithm for additive clustering, based on heuristic combinatorial optimization. Through extensive empirical tests on both human and synthetic data, we find that the new algorithm is more effective than previous methods and that it also scales well to larger problems. By making additive clustering practical, we take a significant step toward scaling connectionist models beyond hand-coded examples.

## 1 Introduction

Many cognitive models posit mental representations based on discrete substructures. Even connectionist models whose processing involves manipulation of real-valued activations typically represent objects as patterns of 0s

---

\*This report summarizes work done from the summer of 1998 through January, 2001.

and 1s across a set of units (Noelle, Cottrell, and Wilms, 1997). Often, individual units are taken to represent specific features of the objects and two representations will share features to the degree to which the two objects are similar. While this arrangement is intuitively appealing, it can be difficult to construct the features to be used in such a model. Using random feature assignments clouds the relationship between the model and the objects it is intended to represent, diminishing the model’s value. But as Clouse and Cottrell (1996) point out, hand-crafted representations are tedious to construct and it can be difficult to precisely articulate or justify the principles that guided their design. These difficulties effectively limit the number of objects that can be encoded.

In this paper, we investigate methods for automatically synthesizing features directly from the pairwise object similarities that the model is intended to respect. (These similarities need not be obtained from psychological testing: one might use word co-occurrence statistics as a proxy for lexical semantic similarity, for example (Resnik, 1995).) This automatic approach eliminates the manual burden of selecting and assigning features while providing an explicit design criterion that objectively connects the representations to empirical data. After formalizing the problem, we will review existing algorithms that have been proposed for solving it. We will then investigate a new approach, based on combinatorial optimization. When using a novel heuristic search technique, we find that the new approach, despite its simplicity, performs better than previous algorithms and that, perhaps more important, it maintains its effectiveness on larger problems.

## 1.1 Additive Clustering

There are many ways to formalize the problem of constructing discrete features from similarity information. We will use what may be the oldest and most established, the *additive clustering* model of Shepard and Arabie (1979). In this framework, clusters represent arbitrarily overlapping discrete features. Each of the  $k$  features has a non-negative real-valued weight  $w_k$ , and the similarity between two objects  $i$  and  $j$  is just the sum of the weights of the features they share. If  $f_{ik}$  is 1 if object  $i$  has feature  $k$  and 0 otherwise, and  $c$  is a real-valued constant, then the similarity of  $i$  and  $j$  is modeled as

$$\hat{s}_{ij} = \sum_k w_k f_{ik} f_{jk} + c .$$

This class of models is very expressive, encompassing non-hierarchical as well as hierarchical arrangements of clusters. (An example model, derived

Table 1: An 8-feature model derived from consonant confusability data. With  $c = 0.024$ , the VAF is 91.8%.

Wt.	Objects with feature	Interpretation
.350	$f\theta$	front unvoiced fricatives
.243	$dg$	back voiced stops
.197	$p\ k$	unvoiced stops (without t)
.182	$b\ v\partial$	front voiced
.162	$ptk$	unvoiced stops
.127	$mn$	nasals
.075	$dgv\partial z\check{z}$	voiced (without b)
.049	$ptkf\theta s\check{s}$	unvoiced

using the `ewindclus-klb` algorithm described below, is shown in Table 1.) Additive clustering is asymmetric in the sense that only the shared features of two objects contribute to their similarity, not the ones they both lack. (This is the more general formulation, as an additional feature containing the set complement of the original feature could always be used to produce such an effect.) Additive clustering was extended by Carroll and Arabie (1983) to the case in which one has similarity data gathered from multiple sources. The assumption is that the underlying cluster memberships remain the same across the conditions, but that separate weights and constants can account for source variation. Carroll and Arabie term this model INDCLUS, for individual differences clustering.

With a model formalism in hand, we can then phrase the problem of constructing feature assignments as simply finding the INDCLUS model that best matches the given similarity data using the desired number of features. The fit of a model (comprising  $F, W, c$ ) to a matrix,  $S$ , can be quantified by the variance accounted for (VAF), which compares the model’s accuracy to just guessing the mean similarity:

$$\text{VAF} = 1 - \frac{\sum_{i,j}(s_{ij} - \hat{s}_{ij})^2}{\sum_{i,j}(s_{ij} - \bar{s})^2}$$

A VAF of 0 can always be achieved by setting all  $w_k$  to 0 and  $c$  to  $\bar{s}$ .

Similar approaches to feature construction for cognitive modeling have been proposed previously. Clouse and Cottrell (1996) present a framework, based on an analogy to multi-dimensional scaling, that encompasses many models. The results they present are based on a model that assumes that feature absence also implies similarity, but we will consider an INDCLUS

variation of their algorithm below. Also, they use an intermediate vector representation of each object, derived using principal components analysis, rather than working directly from the similarities. Lee (1998) also investigates the case in which feature absence implies similarity, although he briefly discusses extending his algorithm to INDCLUS models. We will consider a variation of a more recent INDCLUS algorithm of his below (Lee, submitted).

## 2 Previous Algorithms

Additive clustering is a difficult 0-1 quadratic programming problem and only heuristic approaches, which do not guarantee an optimal model, have been proposed. Many different approaches have been taken:

**Subsets:** Shepard and Arabie (1979) proposed an early algorithm based on subset analysis that was clearly superseded by Arabie’s later work below. Hojo (1983) also proposed an algorithm along these lines. We will not consider these algorithms further.

**Non-discrete Approximation:** Arabie and Carroll (1980) and Carroll and Arabie (1983) propose the two-stage `indclus` algorithm. In the first stage, cluster memberships are treated as real values and optimized for each cluster in turn by gradient descent. At the same time, a penalty term for non-0-1 values is gradually increased. Afterwards, a combinatorial clean-up stage tries all possible changes to 1 or 2 cluster memberships. Experiments reported below use the original code, modified slightly to handle large instances. Runs that crashed or produced VAFs  $\leq 0$  were ignored.

**Asymmetric Approximation:** In the `sindclus` algorithm, Chaturvedi and Carroll (1994) optimize an asymmetric model with two sets of cluster memberships, having the form  $\hat{s}_{ij} = \sum_k w_k f_{ik} g_{jk} + c$ . By considering each cluster in turn, this formulation allows a fast method for determining each of  $F$ ,  $G$ , and  $w$  given the other two. In practice,  $F$  and  $G$  often become identical, yielding an INDCLUS model. Experiments reported below use both a version of the original implementation that has been modified to handle large instances and a reimplemented version that differs in its behavior at boundary cases (handling 0 weights, empty clusters, ties). Runs of the original code that crashed or seemed to hang were ignored. Models from runs in which  $F$  and  $G$  did not converge were each converted into several INDCLUS models by taking only  $F$ , only  $G$ , their intersection, or their

union. The weights and constants of each model were optimized using constrained least-squares linear regression (Stark and Parker, 1995), ensuring non-negative cluster weights, and the one with the highest VAF was used.

**Alternating Clusters:** Kiers (1997) proposes an element-wise simplified `sindclus` algorithm, which we abbreviate as `ewindclus`. Like `sindclus`, it considers each cluster in turn, alternating between the weights and the cluster memberships, although only one set of clusters is maintained. Weights are set by a simple regression and memberships are determined by a gradient function that assumes object independence and fixed weights. The experiments reported below use a new implementation, similar to the reimplementation of `sindclus`.

**Expectation Maximization:** Tenenbaum (1996) reformulates INDCLUS fitting in probabilistic terms as a problem with multiple hidden factorial causes, and proposes a combination of the EM algorithm, Gibbs sampling, and simulated annealing to solve it. The experiments below use a modified version of the original implementation which we will notate as `em-indclus`. It terminates early if 10 iterations of EM pass without a change in the solution quality. (A comparison with the original code showed this modification to give equivalent results using less running time.)

Unfortunately, it is not clear which of these approaches is the best. Most published comparisons of additive clustering algorithms use only a small number of test problems (or only artificial data) and report only the best solution found within an unspecified amount of time. Because each algorithm often returns solutions of widely varying quality, this leaves it unclear which algorithm gives the best results on a typical run. Furthermore, different algorithms require very different running times, and multiple runs of a fast algorithm with high variance in solution quality may produce a better result in the same time as a single run of a more predictable algorithm. The next section reports on a new empirical comparison that addresses these concerns.

## 2.1 Evaluation of Previous Algorithms

We compared `indclus`, both implementations of `sindclus`, `ewindclus`, and `em-indclus` on 3 sets of problems. (For compatibility with the `em-indclus` code, only problems with a single source matrix were used.) The first set is

Table 2: Benchmark problems.

Name	$n$	$k$	First use
animals-small	10	3	Mechelen and Storms (1995)
numbers	10	8	Shepard and Arabie (1979)
workers	14	7	Shepard and Arabie (1979)
consonants	16	8	Shepard and Arabie (1979)
animals	30	5	Chaturvedi and Carroll (1994)
letters	26	12	Shepard and Arabie (1979)

a collection of 6 typical data sets from psychological experiments that have been used in previous additive clustering work. These instances are summarized in Table 2, which lists the number of objects ( $n$ ), the number of features used to fit the data ( $k$ ), and the paper in which additive clustering was first applied to that set. The second set of problems contains noiseless synthetic data derived from INDCLUS models with 8, 16, 32, 64, and 128 objects. In a rough approximation of the human data, the number of clusters was set to  $2\log_2(n)$ , and as in previous INDCLUS work, each object was inserted in each cluster with probability 0.5. A single similarity matrix was generated from each model using weights and constants uniformly distributed between 1 and 6. The third set of problems was derived from the second by adding gaussian noise with a variance of 10% of the variance of the similarity data and enforcing symmetry. Each algorithm was run at least 50 times on each data set. To avoid biasing our conclusions in favor of methods requiring more computation time, those results were then used to derive the distribution of results that would be expected if all algorithms were run simultaneously and those that finished early were re-run repeatedly until the slowest algorithm finished its first run, with any re-runs in progress at that point discarded.<sup>1</sup>

The time-equated distributions of solutions produced by each algorithm on each of the human data sets are shown in Figure 1. (`em-indclus` took much longer than the other algorithms and for clarity its performance is shown separately in Figure 2.) Each box represents the middle 50% of a distribution, a line marks the median, vertical whiskers from each box cover all values within 4 quartiles of the median, and a small grey box represents the 95% confidence interval around the mean. Small circles mark outliers.

<sup>1</sup>Depending as it does on running time, this comparison remains imprecise due to variations in the degree of code tuning and the quality of the compilers used, and the need to normalize timings between the multiple machines used in the tests.

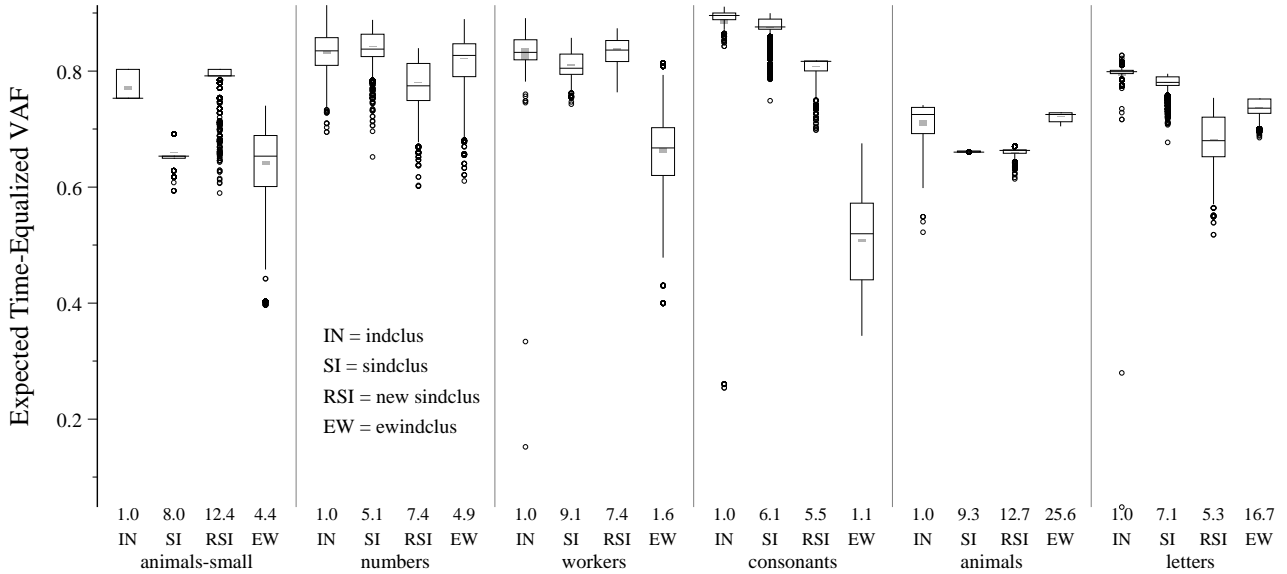


Figure 1: The performance of several previously proposed algorithms on data sets from psychological experiments.

The small numbers below the bars indicate how many runs were necessary on average to achieve time parity with the slowest algorithm on that data set. On most instances, there is remarkable variance in the VAF achieved by each algorithm.<sup>2</sup> Overall, despite the variety of approaches that have been brought to bear over the years, the original `indclus` algorithm appears to be the best. `Animals-small` is the only data set on which its median performance is not the best, and its overall distribution of results is consistently competitive. It is revealing to note the differences in performance between the original and reimplemented versions of `sindclus`. Small changes in the handling of boundary cases make a large difference in the performance of the algorithm.

Surprisingly, on the synthetic data sets (not shown), the relative performance of the algorithms was quite different, and almost the same on the noisy data as on the noise-free data. (This suggests that the randomly generated data sets that are commonly used to evaluate INDCLUS algorithms do

<sup>2</sup>Figure 2 shows one anomaly: no `em-indclus` run on `animals` resulted in a  $\text{VAF} \geq 0$ . This also occurred on all synthetic problems with 32 or more objects (although very good solutions were found on the smaller problems). Tenenbaum (personal communication) suggests that the default annealing schedule in the `em-indclus` code may need to be modified for these problems.

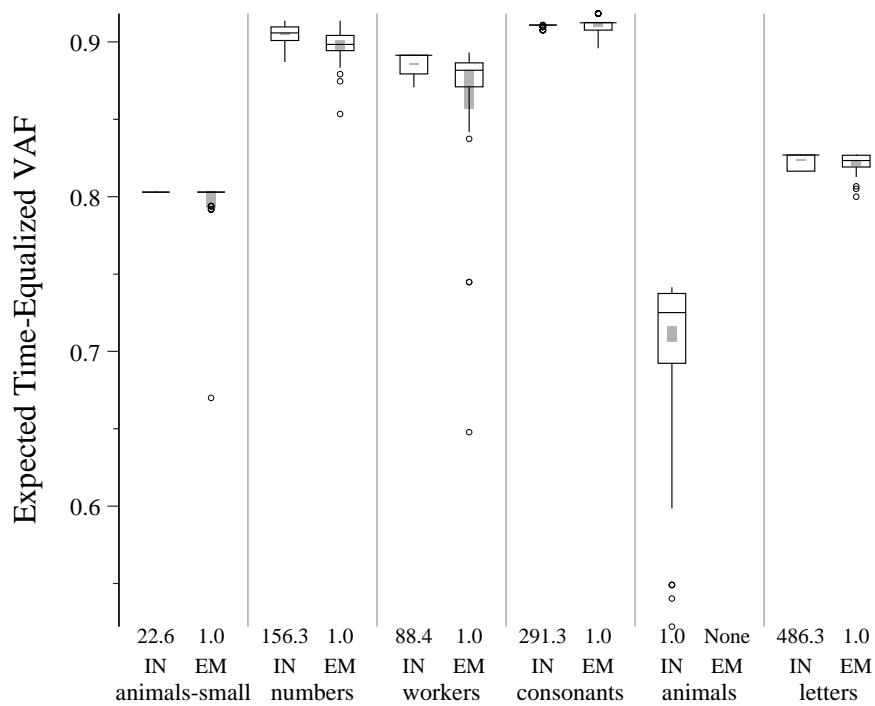


Figure 2: The performance of `indclus` (IN) and `em-indclus` (EM) on the human data sets.



not accurately reflect the problems of interest to practitioners.) `ewindclus` performed best here, although it was only occasionally able to recover the original models from the noise-free data.

Overall, it appears that current methods of additive clustering are quite sensitive to the type of problem they are run on and that there is no assurance that they can reliably scale to problems larger than those examined in the first papers on the topic. In an attempt to address these problems, we turn now to a new approach.

### 3 A Purely Combinatorial Approach

One common theme in `indclus`, `sindclus`, and `ewindclus` is their computation of each cluster and its weight in turn, at each step fitting only the residual similarity not accounted for by the other clusters. Assuming that clusters are independent allows simple equations for updating the cluster weights, but forces memberships to be considered in a predetermined order and allows weights to become obsolete. Inspired in part by recent work of Lee (submitted), we propose an orthogonal decomposition of the problem.<sup>3</sup> Instead of computing the elements and weight of each cluster in succession, we first consider all the memberships and then derive all the weights using constrained regression. And where previous algorithms recompute all the memberships of one cluster simultaneously (and therefore independently), we will change memberships one by one in a dynamically determined order using simple heuristic search techniques, recomputing the weights after each step. From this perspective, `INDCLUS` becomes a purely combinatorial optimization problem.

We will evaluate three different algorithms based on this approach, all of which attempt to improve a random initial model. The first, `indclus-hc`, is a simple hill-climbing strategy which attempts to toggle memberships in an arbitrary order and the first one resulting in an improved model is accepted. The algorithm terminates when no membership can be changed to give an improvement. This strategy was suggested by Clouse and Cottrell (1996), although here we are using the `INDCLUS` model of similarity. In the second algorithm, `indclus-pbil`, the `PBIL` algorithm of Baluja (1997) is used to search for appropriate memberships. This is an efficient simplification of the strategy suggested by Lee (submitted), whose proposal includes elements concerned with automatically controlling model complexity. We use the

---

<sup>3</sup>Lee's recent proposals take a similar combinatorial approach, although his algorithms use it within a more complex two-stage or incremental design.

parameter settings he suggests but only allow the algorithm to generate 10,000 solutions.

While the two previous approaches do not use any problem-specific information beyond solution quality, the third algorithm uses the gradient function from the `ewindclus` algorithm to guide the search. The move strategy is a novel combination of gradient ascent and the method of Kernighan and Lin (1970) which we call ‘KL break-out’. It proceeds by gradient ascent, changing the entry in  $F$  whose `ewindclus` gradient points most strongly to the opposite of its current value. When the ascent no longer results in an improvement, a local maximum has been reached. Motivated by results suggesting that good maxima tend to cluster (Boese, Kahng, and Muddu, 1994; Ruml et al., 1996), the algorithm tries to break out of the current basin of attraction and find a nearby maximum rather than start from another random model. It selects the least damaging variable to change, using the heuristic as in the ascent, but now it locks each variable after changing it. The pool of unlocked variables shrinks, thus forcing the algorithm out of the local maximum and into another part of the space. To determine if it has escaped, a new gradient ascent is attempted after each locking step. If the ascent surpasses the previous maximum, the current break-out attempt is abandoned and the ascent is pursued. If the break-out procedure changes all variables without any ascent finding a better maximum, the algorithm terminates. This strategy, which we will call `ewindclus-klb`, surpassed the original KL method in time-equated tests. It is also conceptually simple and has no parameters that need to be tuned.

### 3.1 Evaluation of the Combinatorial Algorithms

The time-equated performance of the combinatorial algorithms on the human data sets is shown in Figure 3, with `indclus`, the best of the previous algorithms, shown for comparison. As one might expect, adding heuristic guidance to the search helps it enormously: `ewindclus-klb` surpasses the other combinatorial algorithms on every problem. It performs better than `indclus` on three of the human data sets (top panel), equals its performance on two, and performs worse on one data set, letters. The variance of `indclus` on letters is very small, and the results suggest that `ewindclus-klb` is the better choice on this data set if one can afford the time to take the best of 20 runs. (Experiments using 7 additional human data sets found that letters represented the weakest performance of `ewindclus-klb`.)

Performance of a plain KL strategy (not shown) surpassed or equaled `indclus` on all but two problems (consonants and letters), indicating that

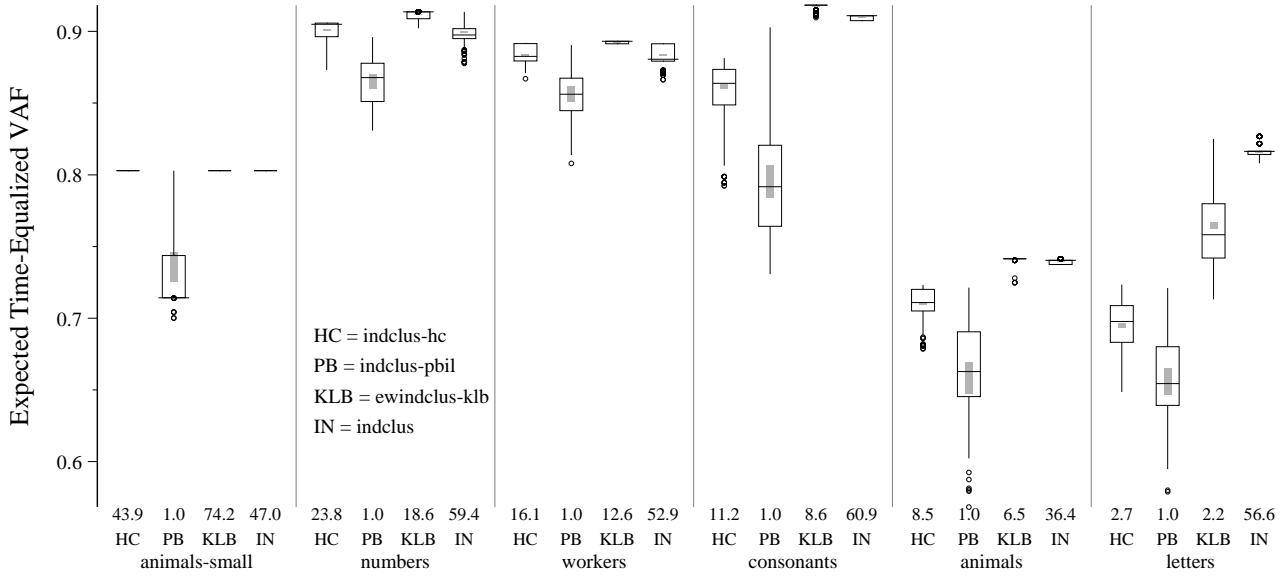


Figure 3: The performance of the combinatorial algorithms on human data sets.

the combinatorial approach, in tandem with heuristic guidance, is powerful even without the new ‘KL break-out’ strategy. It is also interesting to note that `indclus` itself is unique among previous algorithms in incorporating a purely combinatorial search phase. Our results suggest that this portion of the algorithm may account for a substantial part of its effectiveness.

While we have already seen that synthetic data does not predict the relative performance of algorithms on human data very well, it does provide a test of how well they scale to larger problems. On noise-free synthetic data, `ewindclus-klb` reliably recovered the original model on all data sets. It was also the best performer on the noisy synthetic data (a comparison with `indclus` is presented in Figure 4.<sup>4</sup>) These results show that, in addition to performing best on the human data, the combinatorial approach retains its effectiveness on larger problems.

<sup>4</sup>The `indclus` code crashed on 94% of runs with 128 objects. Because these results reflect only runs that completed, they are sensitive to any correlation that might exist between final solution quality and code malfunction.

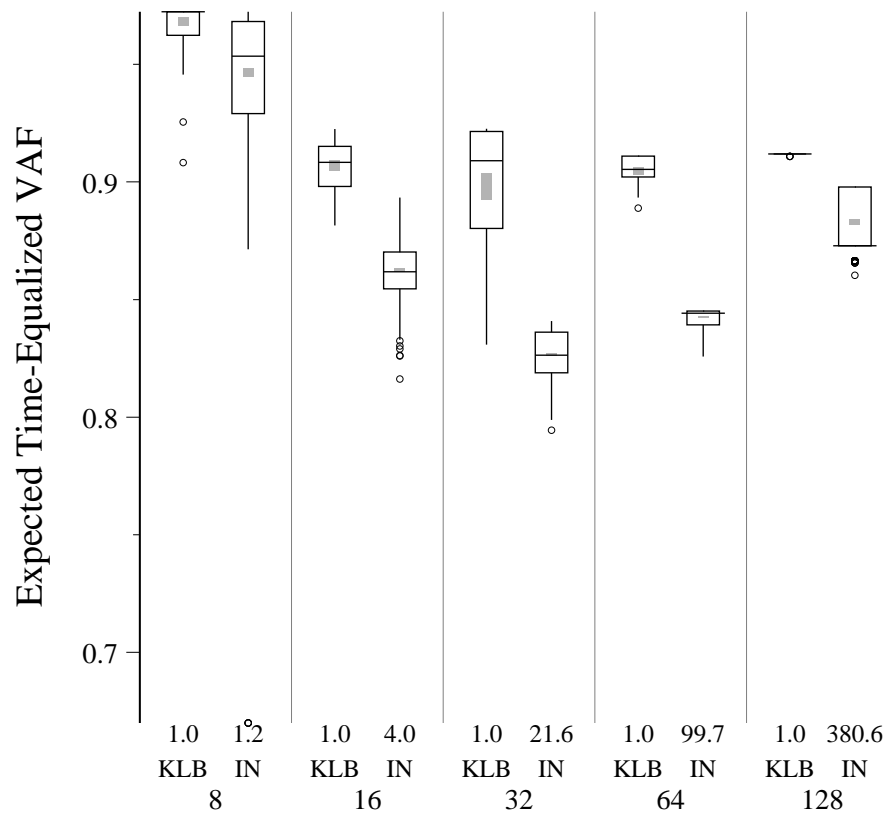


Figure 4: ewindclus-klb (KLB) and indclus (IN) on noisy synthetic data sets of increasing size.

## 4 Conclusions

We developed a connection between constructing feature-based representations for cognitive modeling and the unsupervised learning of INDCLUS models from similarity data. In an empirical comparison sensitive to variance in solution quality and computation time, we found that several recently proposed methods for recovering such models perform worse than the original `indclus` algorithm of Arabie and Carroll (1980). Noting that `indclus` does not scale well to large problems, we suggested a purely combinatorial approach that is simpler than previous proposals. When provided with gradient information, this approach was comparable to `indclus`, and when using a novel move strategy, it surpassed all other algorithms. Perhaps even more important, the new method seems to perform well on large problems.

For simplicity, we considered algorithms in isolation, although it would be interesting to explore the potential for using the output of one algorithm as the starting point for another. We also did not consider the tolerance of the algorithms for missing data, or the question of determining how many features are justifiable for a particular data set.

While this work has extended the reach of the additive clustering paradigm to large problems, it is directly applicable only to cognitive models whose representations encode similarity as shared features. (The cluster weights can be represented by duplicating strong features or by varying connection weights.) However, the simplicity of the combinatorial approach should make it straightforward to extend to models in which the absence of features can enhance similarity.

## 5 Acknowledgments

Thanks to Josh Tenenbaum, Michael Lee, Stuart Shieber, Avi Pfeffer, Emil Gilliam, Alfonso Caramazza, and the Harvard AI Group for stimulating discussions; to Josh, Anil Chaturvedi, Henk Kiers, J. Douglas Carroll, and Phipps Arabie for providing source code for their algorithms; Josh, Michael, and Phipps for providing data sets; and Michael for sharing unpublished work. This work was supported in part by the NSF under grants CDA-94-01024 and IRI-9618848.

## References

- Arabie, Phipps and J. Douglas Carroll. 1980. MAPCLUS: A mathematical programming approach to fitting the adclus model. *Psychometrika*, 45(2):211–235, June.
- Baluja, Shumeet. 1997. Genetic algorithms and explicit search statistics. In Michael C. Mozer, Michael I. Jordan, and Thomas Petsche, editors, *Advances in Neural Information Processing Systems 9*.
- Boese, Kenneth D., Andrew B. Kahng, and Sudhakar Muddu. 1994. A new adaptive multi-start technique for combinatorial global optimizations. *Operations Research Letters*, 16:101–113.
- Carroll, J. Douglas and Phipps Arabie. 1983. INDCLUS: An individual differences generalization of the ADCLUS model and the MAPCLUS algorithm. *Psychometrika*, 48(2):157–169, June.
- Chaturvedi, Anil and J. Douglas Carroll. 1994. An alternating combinatorial optimization approach to fitting the INDCLUS and generalized INDCLUS models. *Journal of Classification*, 11:155–170.
- Clouse, Daniel S. and Garrison W. Cottrell. 1996. Discrete multi-dimensional scaling. In *Proceedings of the Eighteenth Annual Conference of the Cognitive Science Society*, pages 290–294.
- Hojo, Hiroshi. 1983. A maximum likelihood method for additive clustering and its applications. *Japanese Psychological Research*, 25(4):191–201.
- Kernighan, B. and S. Lin. 1970. An efficient heuristic procedure for partitioning graphs. *The Bell System Technical Journal*, 49(2):291–307, February.
- Kiers, Henk A. L. 1997. A modification of the SINDCLUS algorithm for fitting the ADCLUS and INDCLUS models. *Journal of Classification*, 14(2):297–310.
- Lee, Michael D. 1998. Neural feature abstraction from judgements of similarity. *Neural Computation*, 10:1815–1830.
- Lee, Michael D. submitted. A simple method for generating additive clustering models with limited complexity. Manuscript available on the web.

- Mechelen, I. Van and G. Storms. 1995. Analysis of similarity data and Tversky's contrast model. *Psychologica Belgica*, 35(2-3):85-102.
- Noelle, David C., Garrison W. Cottrell, and Fred R. Wilms. 1997. Extreme attraction: On the discrete representation preference of attractor networks. In M. G. Shafto and P. Langley, editors, *Proceedings of the Nineteenth Annual Conference of the Cognitive Science Society*, page 1000. Lawrence Erlbaum.
- Resnik, Philip. 1995. Using information content to evaluate semantic similarity in a taxonomy. In *Proceedings of IJCAI-95*, pages 448-453.
- Ruml, Wheeler, J. Thomas Ngo, Joe Marks, and Stuart Shieber. 1996. Easily searched encodings for number partitioning. *Journal of Optimization Theory and Applications*, 89(2).
- Shepard, Roger N. and Phipps Arabie. 1979. Additive clustering: Representation of similarities as combinations of discrete overlapping properties. *Psychological Review*, 86(2):87-123, March.
- Stark, Philip B. and Robert L. Parker. 1995. Bounded-variable least-squares: An algorithm and applications. *Computational Statistics*, 10:129-141.
- Tenenbaum, Joshua B. 1996. Learning the structure of similarity. In D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, editors, *Advances in Neural Information Processing Systems 8*.